

History of main

June 6, 2011

Supervised Topic Modeling in Clinical Text

Perotte A Li Y Pivovarov R Weiskopf N Wood F
Columbia University, New York, NY 10027, USA
{ajp9009, yil7003, rip7002, ngw7001}@dbm.s.columbia.edu

Abstract

Current medical record keeping technology relies heavily upon human capacity to effectively summarize and infer information from free-text physician notes. We propose a novel method to suggest diagnostic code assignment for patient visits, based upon narrative medical notes. We applied a supervised latent Dirichlet allocation model to a corpus of free-text medical notes from New York - Presbyterian Hospital to infer a set of specific ICD-9 codes for each patient note. Evaluation of the predictions were conducted by comparison to a gold-standard set of ICD-9s assigned to a set of patient notes.

1 Introduction

Despite the growing emphasis on meaningful use of technology in medicine, many aspects of medical record-keeping remain a manual process. Diagnostic coding for billing and insurance purposes is often handled by professional medical coders who must explore a patient's extensive clinical record before assigning the proper codes. So while electronic health records (EHRs) should be adopted by most medical institutions within the next several years, largely due to the provisions of HITECH under the American Recovery and Reinvestment Act [4], there has been little forward movement in automating medical coding.

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9 CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [15].

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [14, 8, 13, 5], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [12] but the most recent and promising work on the subject was inspired by the 2007 Medical NLP Challenge International Challenge: Classifying Clinical Free Text Using Natural Language Processing (NLP) techniques, included word sense disambiguation and named-entity recognition [6, 7]. The data set in the challenge was annotated with only a few documents that were not part of all of the documents in radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [11]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

An automated process would ideally produce a more complete and accurate diagnosis lists. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Methods

2.1 Data

Our data set was gathered from the clinical data warehouse of New York - Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patients chief complaint, diagnostic findings, therapy administered, patient response to the therapy, and the plans made for the patient until discharge. These notes are used uniformly to structure the discharge summaries as part of a controlled terminology which is the international standard diagnostic classification for epidemiological, health management, and clinical purposes (<http://www.who.int/classifications/icd/en/>). The codes are classified in a root-to-tree structure, with each edge representing an is-a relationship between parent and child; such that the parent diagnosis subsumes the child diagnosis. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary. For the purposes of sLDA, ICD-9 codes will be used as labels for discharge summaries.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within

(13)

(14)

(15)

1

2

3

4

5

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) = \frac{\exp\left\{\sum_{i=1}^I (\eta_i^T z_i) y_{m,i} - A(\eta_i^T z_i)\right\} \left(k_{m,n} + \alpha_k\right)}{\sum_{k=1}^K \exp\left\{\sum_{i=1}^I (\eta_i^T z_i) y_{m,i} - A(\eta_i^T z_i)\right\} \left(k_{m,\bullet} + \alpha_k\right)}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - \eta_i^T z_i)] \frac{k_{m,n}}{\sum_{k=1}^K k_{m,k}}$$

Supervised Topic Modeling in Clinical Text

Perotte A Li Y Pivovarov R Weiskopf N Wood F
Columbia University, New York, NY 10027, USA
{ajp9009, yil7003, rip7002, ngw7001}@dbms1.columbia.edu

Abstract

Current medical record keeping technology relies heavily upon human capacity to effectively summarize and infer information from free-text physician notes. We propose a novel method to suggest diagnostic code assignment for patient visits, based upon narrative medical notes. We applied a supervised latent Dirichlet allocation model to a corpus of free-text medical notes from New York - Presbyterian Hospital to infer a set of specific ICD-9 codes for each patient note. Evaluation of the predictions were conducted by comparison to a gold-standard set of ICD-9s assigned to a set of patient notes.

1 Introduction

Despite the growing emphasis on meaningful use of technology in medicine, many aspects of medical record-keeping remain a manual process. Diagnostic coding for billing and insurance purposes is often handled by professional medical coders who must explore a patient's extensive clinical record before assigning the proper codes. So while electronic health records (EHRs) should be adopted by most medical institutions within the next several years, largely due to the provisions of HITECH under the American Recovery and Reinvestment Act [4], there has been little forward movement in automating medical coding.

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9 CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [15].

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [14, 8, 13, 5], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [12] but the most recent and promising work was inspired by the 2007 Medical NLP Challenge International Challenge: Classifying Clinical Free Text Using Natural Language Processing (NLPChallenge website). The main challenges included word sense disambiguation and named entity recognition. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [2].

An automated process would ideally produce a more complete and accurate diagnosis lists. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Methods

2.1 Data

Our data set was gathered from the clinical data warehouse of New York - Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patients chief complaint, diagnostic findings, therapy and admission response to the chief complaint, and the plans made during the course of the visit. These notes are used to map the structure of discharge summaries to a tree hierarchy which is the international standard diagnostic classification for epidemiological, health management, and clinical purposes (<http://www.who.int/classifications/icd/en/>). The codes are classified in a root-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary. For the purposes of sLDA, ICD-9 codes will be used as labels for discharge summaries.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within

(12)

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. Before beginning data processing, we generated a PHI-free dataset (see Data Pre-processing below).

2.2 Pre-Processing

Patient discharge summaries and their associated ICD-9 diagnoses are stored in two different places in the NewYork - Presbyterian data warehouse, and so had to be linked together before being fed into the sLDA algorithm. Each discharge note and set of diagnoses were assigned a patient unique identifier (PUID) and a visit unique identifier (VUID), allowing the two types of data to be linked.

Natural Language Processing (NLP) techniques were used to process the free-text discharge summaries. First, the Natural Language Toolkit (<http://www.nltk.org/>) was used to tokenize the text. Next, feature selection was performed using a term frequency - inverse document frequency (TF-IDF) algorithm on the entire document set and sorting the words by their TF-IDF values. The top 10,000 words were manually evaluated to eliminate all potentially identifying information. Finally, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a node was present, it could be assumed to be present as well (e.g. if a parent had malignant hypertension, it could be assumed that they also had normal hypertension). Second, if a diagnosis was observed to absent, it could be assumed that all of its descendants were also absent (e.g. if a parent did not have essential hypertension, it could be assumed that they did not have malignant hypertension). Unfortunately, ICD-9 code observations never include observations of disease absence. ICD-9 codes are only documented when the condition is observed to be present. Additionally, ICD-9 codes are known to have relatively low sensitivity; conditions that are present are often not documented in a set of ICD-9 codes (Surjan, 1999). Given these facts, we made the following assumptions regarding each visit: recorded diagnoses and their ancestors were labeled as true; diagnoses that were observed as false for a patient but not at the current visit were labeled as unobserved; and diagnoses that had never been listed for a patient were labeled as false for all of that patient's visits. This last assumption captures the belief that parts of the ICD-9 hierarchy that are never observed for a particular patient are almost certain to be false. Additionally, for computational purposes, we decided not to include an ICD-9 code at all if neither it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

2.3 Supervised Topic Models

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [3].

This model, supervised latent dirichlet allocation (sLDA), builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free-text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [2].

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an is-a hierarchy [1]. An is-a hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors (i.e. there are no loops). In this particular case, the ICD-9 code hierarchy is also partially a prefix tree where the labels for certain nodes are prefixes for child nodes. Given that this rule does not apply to all nodes in the hierarchy, we did not use this feature to determine the structure of the hierarchy. Instead we acquired a dataset that explicitly defined the relationships between the nodes of the hierarchy [1]. In documentation of ICD-9 codes for billing purposes, only a subset of the nodes can be used, however the nodes higher in the hierarchy contain semantic information about the categories of codes that are their descendants. For this reason, we included these nodes in our model.

2.4 Generative Model

Given the number of topics, K, the global prior over topic proportions, α' , and the prior over topics, γ , the generative process for documents and responses is as follows:

1. For each topic:

- (a) Draw a distribution over words $\beta_k \sim Dir(u, \gamma)$
- 2. For each ICD9 Code:
 - (a) Draw regression coefficient $\eta_i \mid \mu, \sigma \sim N(\mu, \sigma)$
 - 3. Draw a prior over topic proportions $m \mid \alpha' \sim Dir(u, \alpha')$
 - 4. For each document:
 - (a) Draw topic proportions $\theta_d \mid \alpha \sim Dir(m, \alpha)$
 - (b) For each word:
 - i. Draw topic assignment $z_{a,d} \mid \theta_d \sim Mult(\theta_d)$
 - ii. Draw word $w_{n,a} \mid z_{a,d}, \beta_1 \sim Mult(\beta_{a,d})$
 - (c) For each ICD-9 code from the root of the hierarchy and recursively descending the tree:

2.7.2 $p(\eta_i \mid \mathbf{z}, \mathbf{Y}, \mu)$ or $p(\eta_i \mid \mathbf{z}, \mathbf{a}, \mu)$ in the augmented probit regression model

Given that η_i and $a_{m,i}$ are distributed normally, this posterior distribution is also normal. In the case for general exponential family distributions, η_i can remain a parameter without a prior, fit with maximum likelihood in the usual fashion.

$$p(\eta_i \mid \mathbf{z}, \mathbf{Y}, \mu) = \mathcal{N}(\eta_i \mid \hat{\mu}_i, \hat{\mathbf{S}}_i)$$

$$\hat{\mu}_i = \mathbf{S}_i \mathbf{Z}^T \mathbf{a}_{\bullet,i}$$

$$\hat{\mathbf{S}}_i^{-1} = \mathbf{I} + \mathbf{Z}^T \mathbf{Z}$$

2.7.3 $p(a_{m,i} \mid \mathbf{z}, \mathbf{Y}, \eta)$ in the augmented probit regression model

In the augmented probit regression model, the posterior distribution of a_i is distributed according to a truncated normal distribution.

$$p(a_{m,i} \mid \mathbf{z}, \mathbf{Y}, \eta) = truncN(a_{m,i} \mid \eta_i^T \mathbf{z}, 1, y_{m,i})$$

2.7.4 $p(y_{m,i} \mid \eta, \mathbf{a}, \xi)$

In our model, response variables are not always observed and are treated as latent and sampled where appropriate. There are two factors influencing predictions of the response variable, $y_{m,i}$. There is an undirected model enforcing the aforementioned constraints and providing a prior and there is the probit regression.

$$p(y_{m,i} \mid \eta, \mathbf{a}, \xi) \propto \psi(\mathbf{Y}) \delta(sign(a_{m,i}) - y_{m,i}) \mathcal{N}(a_{m,i} \mid \eta_i^T \mathbf{z}, 1)$$

$$p(y_{m,i} \mid \eta, \mathbf{a}, \xi) \propto \psi(\mathbf{Y}) truncN(a_{m,i} \mid \eta_i^T \mathbf{z}, 1, y_{m,i})$$

Again, this conditional distribution can be evaluated through enumerations and normalization.

$$p(y_{m,i} \mid \eta, \mathbf{a}, \xi) = \frac{\psi(\mathbf{Y}) truncN(a_{m,i} \mid \eta_i^T \mathbf{z}, 1, y_{m,i})}{\sum_{y_{m,i}} \psi(\mathbf{Y}) truncN(a_{m,i} \mid \eta_i^T \mathbf{z}, 1, y_{m,i})}$$

3 Results

4 Conclusion

References

- [1] International classification of disease. <http://bioportal.bioontology.org/ontologies/35686>, May 2008.
- [2] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [3] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [4] D. Blumenthal. Stimulating the adoption of health information technology. *The New England Journal of Medicine*, 360(15):1477–1479, 2009.
- [5] P. Brown, DG Cochrane, and JR Allegri. The ngram cc classifier: A novel method of automatically creating cc classifiers based on icd9 groupings. *Advances in Disease Surveillance*, 1:30, 2006.
- [6] K. Crammer, M. Dredze, K. Ganchev, PP Talukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [7] R. Farkas and G. Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [8] HR FreitasJúnior, B. RibeiroNeto, RF Vale, AHF Laender, and LRS Lima. Categorization-driven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.
- [9] I. Goldstein, A. Arzumanyan, and O. Uzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [10] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [11] LV Li, S Yu, S Niculescu, and J Bi. Large scale diagnostic code classification for medical patient records. 2008.
- [12] RB Rao, S Sandhya, RS Niculescu, C Gerstenblatt, and H Rao. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 416–425, 2003.
- [13] B RibeiroNeto, AHF Laender, and LRS Lima. An experimental study in automatically categorizing medical documents. *Journal of the American Society for Information Science and Technology*, 52(5):391–401, 2001.
- [14] P. Ruch, J. Gobell, I. Tshiriri, and A. Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [15] G. Surjan. Questions on validity of international classification of diseases-coded diagnoses. *International journal of medical informatics*, 54(2):77–95, 1999.

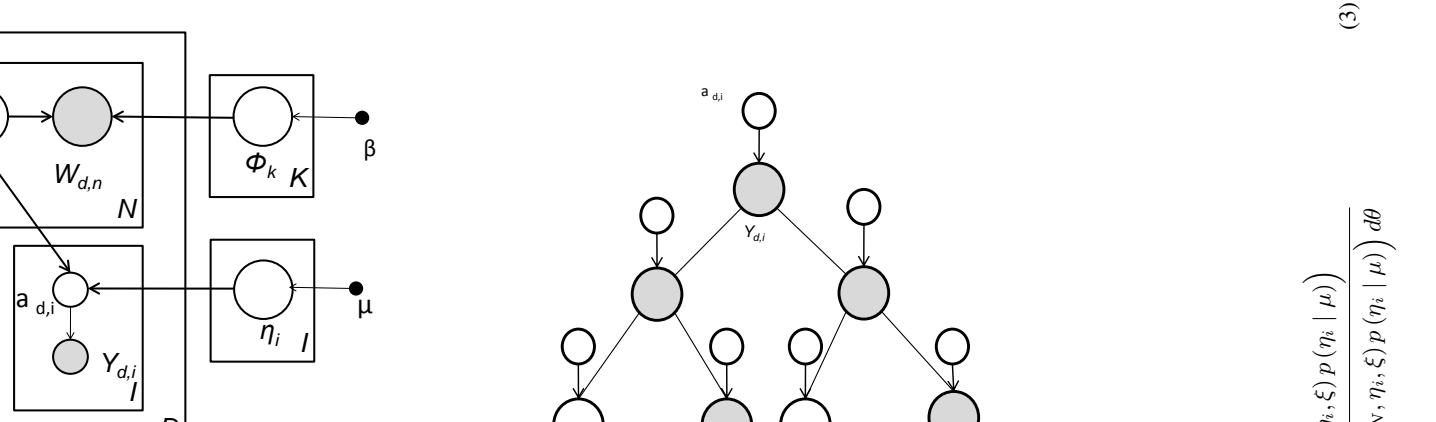


Figure 1: adapted sLDA model

Supervised Topic Modeling in Clinical Text

Perotte A Li Y Pivovarov R Weiskopf N Wood F
Columbia University, New York, NY 10027, USA
{ajp9009, yil7003, rip7002, ngw7001}@dbms1.columbia.edu

Abstract

Current medical record keeping technology relies heavily upon human capacity to effectively summarize and infer information from free-text physician notes. We propose a novel method to suggest diagnostic code assignment for patient visits, based upon narrative medical notes. We applied a supervised latent Dirichlet allocation model for a corpus of free-text medical notes from New York - Presbyterian Hospital to infer a set of specific ICD-9 codes for each patient note. Evaluation of the predictions were conducted by comparison to a gold-standard set of ICD-9s assigned to a set of patient notes.

1 Introduction

Despite the growing emphasis on meaningful use of technology in medicine, many aspects of medical record-keeping remain a manual process. Diagnostic coding for billing and insurance purposes is often handled by professional medical coders who must explore a patient's extensive clinical record before assigning the proper codes. So while electronic health records (EHRs) should be adopted by most medical institutions within the next several years, largely due to the provisions of HITECH under the American Recovery and Reinvestment Act [?], there has been little forward movement in automating medical coding.

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9 CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [?].

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [? ? ? ?], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [?], but the most recent and promising work on the subject was inspired by the 2007 Medical NLP Challenge International Challenge: Classifying Clinical Free Text Using Natural Language Processing (website). The most recent and promising approaches included word vector and bag-of-words models [? ? ?]. The data set used in this challenge was annotated only with documents that were in ICD-9 format and all of the documents had pathology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [?]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

An automated process would ideally produce a more complete and accurate diagnosis lists. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Methods

2.1 Data

Our data set was gathered from the clinical data warehouse of New York - Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patients chief complaint, diagnostic findings, therapy and patient response to the changes, and the plans made during the course of treatment unless otherwise. This data is mapped to a structure of discharge summaries, which are part of a controlled terminology which is the international standard diagnostic classification for epidemiological, health management, and clinical purposes (<http://www.who.int/classifications/icd/en/>). The codes are classified in a rootless tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary. For the purposes of sLDA, ICD-9 codes will be used as labels for discharge summaries.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within

a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. Before beginning data processing, we generated a PHI-free dataset (see Data Pre-processing below).

2.2 Pre-Processing

Patient discharge summaries and their associated ICD-9 diagnoses are stored in two different places in the NewYork - Presbyterian data warehouse, and so had to be linked together before being fed into the sLDA algorithm. Each discharge note and set of diagnoses were assigned a patient unique identifier (PUD) and a visit unique identifier (VUID), allowing the two types of data to be linked.

Natural Language Processing (NLP) techniques were used to process the free-text discharge summaries. First, the Natural Language Toolkit (<http://www.nltk.org/>) was used to tokenize the text. Next, feature selection was performed using a term frequency - inverse document frequency (TF-IDF) algorithm on the entire document set and sorting the words by their TF-IDF values. The top 10,000 words were manually evaluated to eliminate all potentially identifying information. Finally, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a node is an ancestor of its descendants it could be assumed to be present as well (e.g. if a parent had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was observed to be absent, it could be assumed that all of its descendants were also absent (e.g. if a parent did not have essential hypertension, it could be assumed that they did not have malignant hypertension). Unfortunately, ICD-9 code observations never include observations of disease absence. ICD-9 codes are only documented when the condition is observed to be present. Additionally, ICD-9 codes are known to have relatively low sensitivity; conditions that are present are often not documented in a set of ICD-9 codes (Surjan, 1999). Given these facts, we made the following assumptions regarding each visit: recorded diagnoses and their ancestors were labeled as true; diagnoses that were observed at some time for a patient but not at the current visit were labeled as unobserved; and diagnoses that had never been listed for a patient were labeled as false for all of that patients visits. This last assumption captures the belief that parts of the ICD-9 hierarchy that are never observed for a particular patient are almost certain to be false. Additionally, for computational purposes, we decided not to include an ICD-9 code at all if neither it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

2.3 Supervised Topic Models

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [?].

This model, supervised latent dirichlet allocation (sLDA), builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [?].

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an is-a hierarchy [?]. An is-a hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors (ie. there are no loops). In this particular case, the ICD-9 code hierarchy is also partially a prefix tree where the labels for certain nodes are prefixes for child nodes. Given that this rule does not apply to all nodes in the hierarchy, we did not use this feature to determine the structure of the hierarchy. Instead we acquired a dataset that explicitly defined the relationships between the nodes of the hierarchy [?]. In documentation of ICD-9 codes for billing purposes, only a subset of the nodes can be used, however the nodes higher in the hierarchy contain semantic information about the categories of codes that are their descendants. For this reason, we included these nodes in our model.

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 4. The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler for the supervised topic model.

2.4 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

1. For each topic:
 - (a) Draw a distribution over words $\beta_k \sim Dir(u, \gamma)$
2. For each ICD9 Code:
 - (a) Draw regression coefficient $\eta_i \mid \mu, \sigma \sim N(\mu, \sigma)$
 3. Draw a prior over topic proportions $m \mid \alpha' \sim Dir(u, \alpha')$
 4. For each document:
 - (a) Draw topic proportions $\theta_d \mid \alpha \sim Dir(m, \alpha)$
 - (b) For each word:
 - i. Draw topic assignment $z_{a,d} \mid \theta_d \sim Mult(\theta_d)$
 - ii. Draw word $w_{n,a} \mid z_{a,d}, \beta_{1:K} \sim Mult(\beta_{a,d})$
 - (c) For each ICD-9 code from the root of the hierarchy and recursively descending the tree:

2.7.2 $p(\eta_i \mid z, \mathbf{Y}, \mu)$ or $p(\eta_i \mid z, \mathbf{a}, \mu)$ in the augmented probit regression model

Given that η_i and $a_{m,i}$ are distributed normally, this posterior distribution is also normal. In the case for general exponential family distributions, η_i can remain a parameter without a prior, fit with maximum likelihood in the usual fashion.

$$p(\eta_i \mid z, \mathbf{a}, \mu) = N(\eta_i \mid \hat{\mu}_i, \hat{S}_i) \quad (15)$$

$$\hat{\mu}_i = \hat{\mathbf{S}}_i^T \mathbf{a}_{\bullet,i} \quad (16)$$

$$\hat{S}_i^{-1} = \mathbf{I} + \mathbf{Z}^T \mathbf{Z} \quad (17)$$

2.7.3 $p(a_{m,i} \mid z, \mathbf{Y}, \eta)$ in the augmented probit regression model

In the augmented probit regression model, the posterior distribution of a_i is distributed according to a truncated normal distribution.

$$p(a_{m,i} \mid z, \mathbf{Y}, \eta) = truncN(a_{m,i} \mid \eta_i^T z, 1, y_{m,i}) \quad (18)$$

In our model, response variables are not always observed and are treated as latent and sampled where appropriate. There are two factors influencing predictions of the response variable, $y_{m,i}$. There is an undirected model enforcing the aforementioned constraints and providing a prior and there is the probit regression.

$$p(y_{m,i} \mid \eta, \mathbf{a}, \xi) \propto \psi(\mathbf{Y}) \delta(sign(a_{m,i}) - y_{m,i}) N(a_{m,i} \mid \eta_i^T z, 1) \quad (19)$$

$$p(y_{m,i} \mid \eta, \mathbf{a}, \xi) \propto \psi(\mathbf{Y}) truncN(a_{m,i} \mid \eta_i^T z, 1, y_{m,i}) \quad (20)$$

Again, this conditional distribution can be evaluated through enumerations and normalization.

$$p(y_{m,i} \mid \eta, \mathbf{a}, \xi) = \frac{\psi(\mathbf{Y}) truncN(a_{m,i} \mid \eta_i^T z, 1, y_{m,i})}{\sum_{y_{m,i}} \psi(\mathbf{Y}) truncN(a_{m,i} \mid \eta_i^T z, 1, y_{m,i})} \quad (21)$$

3 Results

4 Conclusion

$$p(z_{m,n} \mid z_{-(m,n)}, \mathbf{w}, \eta, \mathbf{a}, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - 1) \mathcal{N}(y_{m,i} \mid \eta_i^T z, 1)] \quad (12)$$

$$p(z_{m,n} \mid z_{-(m,n)}, \mathbf{w}, \eta, \mathbf{a}, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - 1) \mathcal{N}(y_{m,i} \mid \eta_i^T z, 1)] \quad (13)$$

$$p(z_{m,n} \mid z_{-(m,n)}, \mathbf{w}, \eta, \mathbf{a}, \alpha, \beta, \mu) \propto \prod_{i=1}^I [\delta(\text{sign}(a_{m,i}) - 1) \mathcal{N}(y_{m,i} \mid \eta_i^T z, 1)] \quad (14)$$

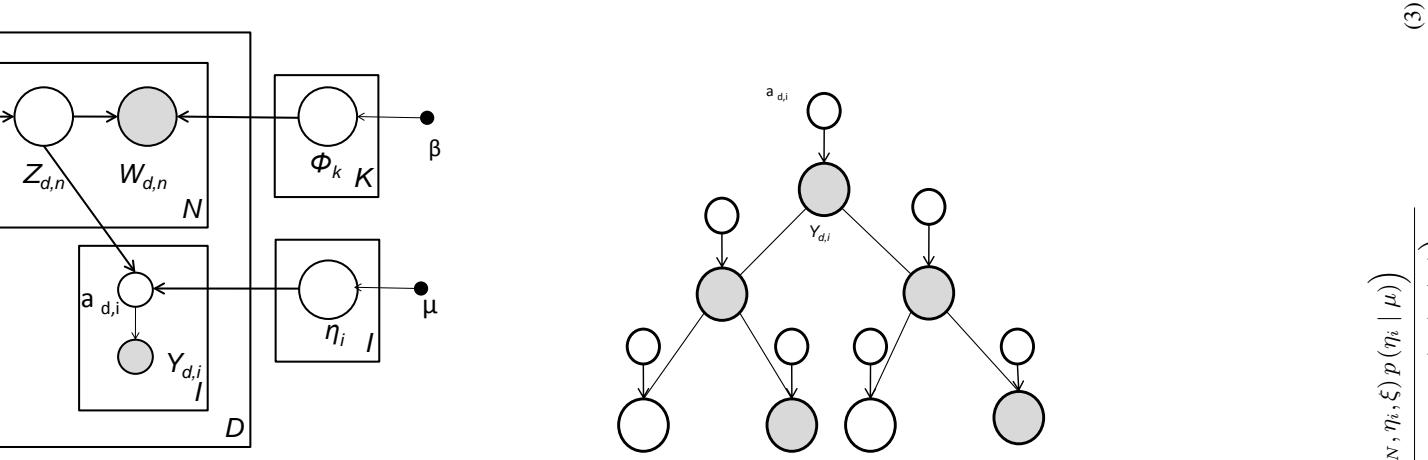


Figure 1: adapted sLDA model

i. Draw a response variable $y_i \mid z, \eta_i, y_{parent} \sim \Phi(\eta_i^T z, y_{parent})$ where $z = N^{-1} \sum_{n=1}^N z_n$ and Φ refers to a conditional probit model.

We will employ a data augmentation scheme with auxiliary variables a_i in the probit model where:

$$y_i \sim \begin{cases} 1, & a_i > 0 \text{ and } y_{parent} = 1 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

$$a_i \sim \mathcal{N}(\eta_i^T z, 1) \quad (2)$$

2.5 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 4.

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler for the supervised topic model.

2.6 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 4. The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler for the supervised topic model.

2.7 Gibbs Sampling

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

2.7 Gibbs Sampling

To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

2.7.1 $p(z_{m,n} \mid z_{-(m,n)}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, m . The conditional probability with respect to this latent variable is proportional to the joint distribution up to a constant.

$$p(z_{m,n} \mid z_{-(m,n)}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu, \xi) \quad (8)$$

Due to the factorization of this model we can rewrite the joint distribution as the following:

$$\propto \prod_{i=1}^I [p(\eta_i \mid p(y_{m,i} \mid z_{m,i}, \mathbf{w}, \eta, \alpha, \beta, \mu))] \quad (9)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [?].

$$\propto \prod_{i=1}^I [p(y_{m,i} \mid z_{m,i}, \eta_i, \xi)] \left(\frac{n_{k,-(m,n)} + \alpha_k}{n_{k,-(m,n)} + \beta_{k,w_{m,n}}} \right) \quad (10)$$

Here, $n_{k,-(m,n)}$ represents the count of word v in document m assigned to topic k omitting the $(m, n)^{th}$ word count. For exponential family distributions, the normalization constant, $h(y_{m,i})$, does not depend on $z_{m,n}$.

$$\propto \exp \left\{ \sum_{i=1}^I (\eta_i^T z) y_{m,i} - A(\eta_i^T z) \right\} \left(\frac{n_{k,-(m,n)} + \alpha_k}{n_{k,-(m,n)} + \beta_{k,w_{m,n}}} \right) \quad (11)$$

Given this expression, $p(z_{m,n} \mid z_{-(m,n)}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu)$ can be sampled through enumeration as seen in Equation 13. In the case of probit regression, the expression for ?? evaluates to Equation 14. Equivalently we can parameterize the model with an auxiliary variable a_i , resulting in Equation 15.

3 Results

3.1 Model

3.2 Evaluation

3.3 Discussion

3.4 Conclusion

3.5 Future Work

Supervised Topic Modeling in Clinical Text

Perotte A Li Y Pivovarov R Weiskopf N Wood F
Columbia University, New York, NY 10027, USA
{ajp9009, yil7003, rip7002, ngw7001}@dbms1.columbia.edu

Abstract

Current medical record keeping technology relies heavily upon human capacity to effectively summarize and infer information from free-text physician notes. We propose a novel method to suggest diagnostic code assignment for patient visits, based upon narrative medical notes. We applied a supervised latent Dirichlet allocation model to a corpus of free-text medical notes from New York - Presbyterian Hospital to infer a set of specific ICD-9 codes for each patient note. Evaluation of the predictions were conducted by comparison to a gold-standard set of ICD-9s assigned to a set of patient notes.

1 Introduction

Despite the growing emphasis on meaningful use of technology in medicine, many aspects of medical record-keeping remain a manual process. Diagnostic coding for billing and insurance purposes is often handled by professional medical coders who must explore a patient's extensive clinical record before assigning the proper codes. So while electronic health records (EHRs) should be adopted by most medical institutions within the next several years, largely due to the provisions of HITECH under the American Recovery and Reinvestment Act [?], there has been little forward movement in automating medical coding.

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9 CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [?].

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [? ? ? ?], fully automatic assignment of ICD-9 codes to medical texts became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [?], but most recent and promising work on the subject was inspired by the 2007 Medical NLP Challenge International Challenge: Classifying Clinical Free Text Using Natural Language Processing (website). The top performing systems included word and bag-based algorithms, [? ? ? ?]. The data set for this challenge was annotated with only the documents that were in the training set of radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [?]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

An automated process would ideally produce a more complete and accurate diagnosis lists. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Methods

2.1 Data

Our data set was gathered from the clinical data warehouse of New York - Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patients chief complaint, diagnostic findings, therapy and patient response to the therapy, and the plans for further treatment unless under discharge. This data set is used to build a tree structure to represent the hierarchy of diseases and are part of a controlled terminology which is the international standard diagnostic classification for epidemiological, health management, and clinical purposes (<http://www.who.int/classifications/icd/en/>). The codes are classified in a root-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary. For the purposes of sLDA, ICD-9 codes will be used as labels for discharge summaries.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within

(12)

1

(13)

(14)

2

2.7.2 $p(\eta_i | \mathbf{z}, \mathbf{Y}, \mu)$ or $p(\eta_i | \mathbf{z}, \mathbf{a}, \mu)$ in the augmented probit regression model

Given that η_i and $a_{m,i}$ are distributed normally, this posterior distribution is also normal. In the case for general exponential family distributions, η_i can remain a parameter without a prior, fit with maximum likelihood in the usual fashion.

$$p(\eta_i | \mathbf{z}, \mathbf{Y}, \mu) = \mathcal{N}(\eta_i | \hat{\mu}_i, \hat{\mathbf{S}}_i)$$

$$\hat{\mathbf{S}}_i^{-1} = \mathbf{I} + \mathbf{Z}^T \mathbf{Z}$$

2.7.3 $p(a_{m,i} | \mathbf{z}, \mathbf{Y}, \eta)$ in the augmented probit regression model

In the augmented probit regression model, the posterior distribution of a_i is distributed according to a truncated normal distribution.

$$p(a_{m,i} | \mathbf{z}, \mathbf{Y}, \eta) = \text{truncN}(a_{m,i} | \eta_i^T \mathbf{z}, 1, y_{m,i})$$

2.7.4 $p(y_{m,i} | \eta, \mathbf{a}, \xi)$

In our model, response variables are not always observed and are treated as latent and sampled where appropriate. There are two factors influencing predictions of the response variable, $y_{m,i}$. There is an undirected model enforcing the aforementioned constraints and providing a prior and there is the probit regression.

$$p(y_{m,i} | \eta, \mathbf{a}, \xi) \propto \psi(\mathbf{Y}) \delta(\text{sign}(a_{m,i}) = y_{m,i}) \mathcal{N}(a_{m,i} | \eta_i^T \mathbf{z}, 1)$$

$$p(y_{m,i} | \eta, \mathbf{a}, \xi) \propto \psi(\mathbf{Y}) \text{truncN}(a_{m,i} | \eta_i^T \mathbf{z}, 1, y_{m,i})$$

Again, this conditional distribution can be evaluated through enumerations and normalization.

$$p(y_{m,i} | \eta, \mathbf{a}, \xi) = \frac{\psi(\mathbf{Y}) \text{truncN}(a_{m,i} | \eta_i^T \mathbf{z}, 1, y_{m,i})}{\sum_{y_{m,i}} \psi(\mathbf{Y}) \text{truncN}(a_{m,i} | \eta_i^T \mathbf{z}, 1, y_{m,i})} \quad (18)$$

3 Results

4 Conclusion

References

- [1] International classification of disease. <http://bioportal.bioontology.org/ontologies/35686>, May 2008.
- [2] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [3] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [4] D. Blumenthal. Stimulating the adoption of health information technology. *The New England Journal of Medicine*, 360(15):1477–1479, 2009.
- [5] P Brown, DG Cochane, and JR Allegri. The ngram cc classifier: A novel method of automatically creating cc classifiers based on icd9 groupings. *Advances in Disease Surveillance*, 1:30, 2006.
- [6] K Crammer, M Dredze, K Ganchev, PP Talukdar, and S Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [7] R Farkas and G Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [8] HR FreitasJunior, B RibeiroNeto, RF Vale, AHF Laender, and LRS Lima. Categorization-driven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.
- [9] I Goldstein, A Arzumanyan, and O Uzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [10] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [11] LY Lita, S Yu, S Niculescu, and J Bi. Large scale diagnostic code classification for medical patient records. 2008.
- [12] RB Rao, S Sandilya, RS Niculescu, C Gerstenblatt, and H Rao. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 416–425, 2003.
- [13] B RibeiroNeto, AHF Laender, and LRS Lima. An experimental study in automatically categorizing medical documents. *Journal of the American Society for Information Science and Technology*, 52(5):391–401, 2001.
- [14] P Ruch, J Gobell, I Thashirii, and A Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [15] G Surjan. Questions on validity of international classification of diseases-coded diagnoses. *International journal of medical informatics*, 54(2):77–95, 1999.

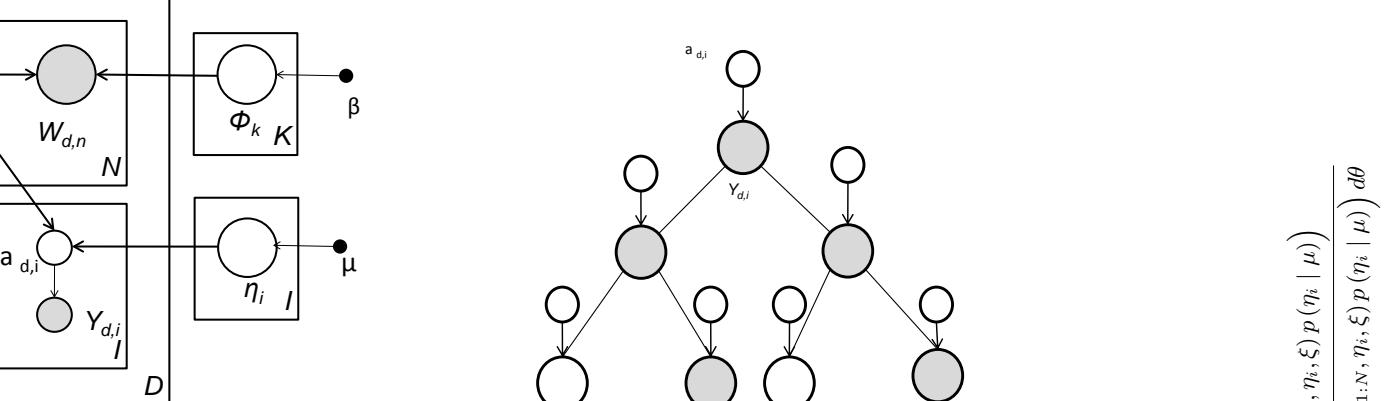


Figure 1: adapted sLDA model

i. Draw a response variable $y_i | z, \eta_i, y_{\text{parent}} \sim \Phi(\eta_i^T z, y_{\text{parent}})$ where $z = N^{-1} \sum_{n=1}^N z_n$ and Φ refers to a conditional probit model.

We will employ a data augmentation scheme with auxiliary variables a_i in the probit model where:

$$y_i \sim \begin{cases} 1, & a_i > 0 \text{ and } y_{\text{parent}} = 1 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

$$a_i \sim \mathcal{N}(\eta_i^T z, 1) \quad (2)$$

2.5 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 4.

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler for the supervised topic model.

2.6 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5–8.

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 4.

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler for the supervised topic model.

2.7 Gibbs Sampling

To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \quad (3)$$

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu, \xi) \quad (4)$$

$$p(w_{n,d} | z_{m,n}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \quad (5)$$

$$p(a_{d,n} | z_{m,n}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \quad (6)$$

$$p(\eta_i | z_{m,n}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu, \xi) \quad (7)$$

$$p(\theta_d | z_{m,n}, \mathbf{Y}, \mathbf{w}, \eta, \alpha, \beta, \mu) \quad (8)$$

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

<p

Supervised Topic Modeling in Clinical Text

Perotte A Li Y Pivovarov R Weiskopf N Wood F
Columbia University, New York, NY 10027, USA
{ajp9009, yil7003, rip7002, ngw7001}@columbia.edu

Abstract

Current medical record keeping technology relies heavily upon human capacity to effectively summarize and infer information from free-text physician notes. We propose a novel method to suggest diagnostic code assignment for patient visits, based upon narrative medical notes. We applied a supervised latent Dirichlet allocation model to a corpus of free-text medical notes from New York - Presbyterian Hospital to infer a set of specific ICD-9 codes for each patient note. Evaluation of the predictions were conducted by comparison to a gold-standard set of ICD-9s assigned to a set of patient notes.

1 Introduction

Despite the growing emphasis on meaningful use of technology in medicine, many aspects of medical record-keeping remain a manual process. Diagnostic coding for billing and insurance purposes is often handled by professional medical coders who must explore a patient's extensive clinical record before assigning the proper codes. So while electronic health records (EHRs) should be adopted by most medical institutions within the next several years, largely due to the provisions of HITECH under the American Recovery and Reinvestment Act [4], there has been little forward movement in automating medical coding.

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the the principal diagnoses [15].

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [14, 8, 13, 5], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [16].

A subset of earlier work was inspired by the 2007 Medical NLP Challenge: "International Challenge: Classifying Clinical Free Text Using Natural Language Processing" [17]. Most of the classification strategies included word and rule-based algorithms [9, 6, 7]. The best system used a rule-based classifier composed of a decision tree and a support vector machine [12]. It used all of the documents in the pathology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [11]. This paper proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

An automated process would ideally produce a more complete and accurate diagnosis lists. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Background

2.1 Data

Our data set was gathered from the clinical data warehouse of New York - Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient's chief complaint, diagnostic findings, therapy and admission plan, a response to the emergency, and relevant plans for follow-up care and discharge. The ICD-9 codes are assigned to structure the discharge summaries; they are part of a controlled terminology which is the international standard for diagnostic classification for epidemiological, health management, and clinical purposes (<http://www.who.int/classifications/icd/en/>). The codes are classified in a root-tree structure, with each edge representing an is-a relationship between parent and child; such that the parent diagnosis subsumes the child diagnosis. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary. For the purposes of sLDA, ICD-9 codes will be used as labels for discharge summaries.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within Equations 5-8.

5.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. Before beginning data processing, we generated a PHI-free dataset (see Data Pre-processing below).

2.2 Pre-Processing

Patient discharge summaries and their associated ICD-9 diagnoses are stored in two different places in the New York - Presbyterian data warehouse, and so had to be linked together before being fed into the sLDA algorithm. Each discharge note and set of diagnoses were assigned a patient unique identifier (PUIID) and a visit unique identifier (VUID), allowing the two types of data to be linked.

Natural Language Processing (NLP) techniques were used to process the free-text discharge summaries. First, the Natural Language Toolkit (<http://www.nltk.org>) was used to tokenize the text. Next, feature selection was performed using a term frequency - inverse document frequency (TF-IDF) algorithm on the entire document set and sorting the words by their TF-IDF values. The top 10,000 words were manually evaluated to eliminate all potentially identifying information. Finally, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had benign hypertension). Second, if a diagnosis was observed to be absent, it could be assumed that all of its descendants were also absent (e.g., if a patient did not have malignant hypertension, it could be assumed that they did not have benign hypertension). Unfortunately, ICD-9 code observations never include observations of disease absence. ICD-9 codes are only documented when the condition is observed to be present. Additionally, ICD-9 codes are known to have relatively low sensitivity; conditions that are present are often not documented in a set of ICD-9 codes (Surjan, 1999). Given these facts, we made the following assumptions regarding each visit: recorded diagnoses and their ancestors were labeled as true; diagnoses that were observed at some time for a patient but not at the current visit were labeled as unobserved; and diagnoses that had never been listed for a patient were labeled as false for all of the patient's visits. This last assumption captures the belief that parts of the ICD-9 hierarchy that are never observed for a particular patient are almost certain to be false. Additionally, for computational purposes, we decided not to include an ICD-9 code at all if either it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

2.3 ICD-9 Code Hierarchy

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an is-a hierarchy [1]. An is-a hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors (i.e., there are no loops). In this particular case, the ICD-9 code hierarchy is also partially a prefix tree where the labels for certain nodes are prefixes for child nodes. Given that this rule does not apply to all nodes in the hierarchy, we did not use this feature to determine the structure of the hierarchy. Instead we acquired a dataset that explicitly defined the relationships between the nodes of the hierarchy [1]. In documentation of ICD-9 codes for billing purposes, only a subset of the nodes can be used, however the nodes higher in the hierarchy contain semantic information about the categories of codes that are their descendants.

For this reason, we included these nodes in our model.

3 Methods

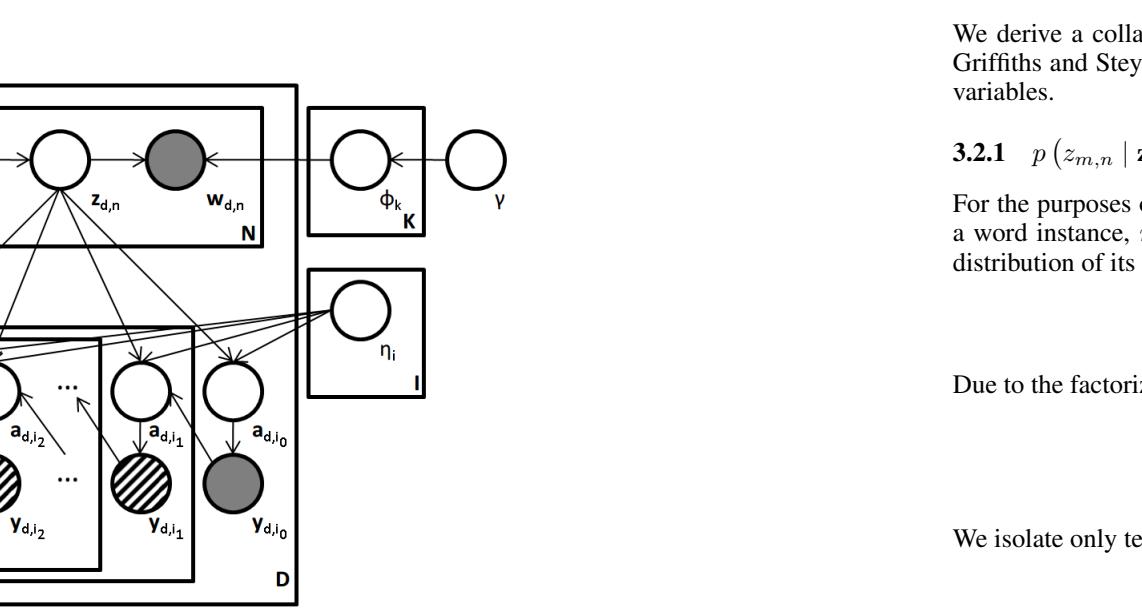


Figure 1: adapted sLDA model

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, k:

$$\begin{aligned} p(\theta | z_{d,n}) &= \frac{p(\theta; \alpha)}{\int_0^{\infty} p(\theta; \alpha) \sum_{k=1}^K \left(\prod_{n=1}^N p(y_{d,n} | z_{d,n}, \theta, \phi_{k,z}) \right) \left(\prod_{n=1}^N p(w_{d,n} | z_{d,n}, \theta, \phi_{w,z}) \right) \left(\prod_{n=1}^N p(z_{d,n} | \theta) \right) \left(\prod_{n=1}^N p(\phi_{k,z} | \theta) \right) d\theta} \\ (13) \quad &= \frac{p(\theta; \alpha)}{\int_0^{\infty} p(\theta; \alpha) \sum_{k=1}^K \left(\prod_{n=1}^N p(y_{d,n} | z_{d,n}, \theta, \phi_{k,z}) \right) \partial \theta d\theta} \\ (14) \quad &= \prod_{k=1}^K \left[\int_{\phi_{k,z}}^{\infty} p(\phi_{k,z}; \beta) \prod_{n=1}^N p(y_{d,n} | z_{d,n}, \theta, \phi_{k,z}) \right] \prod_{k=1}^K \left[\int_{\phi_{w,z}}^{\infty} p(\phi_{w,z}; \beta) \prod_{n=1}^N p(w_{d,n} | z_{d,n}, \theta, \phi_{w,z}) \right] \prod_{k=1}^K \left[\int_{\theta_1}^{\theta_2} p(\theta; \alpha) d\theta \right] \\ (15) \quad &= \prod_{k=1}^K \left[\int_{\phi_{k,z}}^{\infty} p(\phi_{k,z}; \beta) \prod_{n=1}^N p(y_{d,n} | z_{d,n}, \theta, \phi_{k,z}) \right] \prod_{k=1}^K \left[\int_{\phi_{w,z}}^{\infty} p(\phi_{w,z}; \beta) \prod_{n=1}^N p(w_{d,n} | z_{d,n}, \theta, \phi_{w,z}) \right] \prod_{k=1}^K \left[\int_{\theta_1}^{\theta_2} p(\theta; \alpha) d\theta \right] \\ (16) \quad &= \prod_{k=1}^K \left[\int_{\phi_{k,z}}^{\infty} p(\phi_{k,z}; \beta) \prod_{n=1}^N p(y_{d,n} | z_{d,n}, \theta, \phi_{k,z}) \right] \prod_{k=1}^K \left[\int_{\phi_{w,z}}^{\infty} p(\phi_{w,z}; \beta) \prod_{n=1}^N p(w_{d,n} | z_{d,n}, \theta, \phi_{w,z}) \right] \prod_{k=1}^K \left[\int_{\theta_1}^{\theta_2} p(\theta; \alpha) d\theta \right] \\ (17) \quad &= \prod_{k=1}^K \left[\int_{\phi_{k,z}}^{\infty} p(\phi_{k,z}; \beta) \prod_{n=1}^N h(y_n) \exp \left(\left(\frac{z_{d,n}}{\theta} \right)^T y_n - A \left(\frac{z_{d,n}}{\theta} \right) \right) \right] \prod_{k=1}^K \left[\int_{\phi_{w,z}}^{\infty} p(\phi_{w,z}; \beta) \prod_{n=1}^N h(w_n) \exp \left(\left(\frac{z_{d,n}}{\theta} \right)^T w_n - A \left(\frac{z_{d,n}}{\theta} \right) \right) \right] \prod_{k=1}^K \left[\int_{\theta_1}^{\theta_2} p(\theta; \alpha) d\theta \right] \end{aligned}$$

a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. Before beginning data processing, we generated a PHI-free dataset (see Data Pre-processing below).

2. For each ICD9 code, c, at all levels in the tree, l:

(a) Draw regression coefficient $\eta_{i_l,c} | \sigma \sim N_K(-1, \sigma)$

3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$

4. For each document, d:

(a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$

(b) For each word, n:

i. Draw topic assignment $z_{d,n} | \theta_d \sim Mult_K(\theta_d)$

ii. Draw word $w_{d,n} | z_{d,n}, \beta_{i_l,c} \sim Mult_V(\beta_{i_l,c})$

(c) For each level of the ICD-9 code tree, l:

i. For each ICD9 code at this level, c:

A. Draw a latent variable

$$a_{d,i_l,c} \sim \left\{ \begin{array}{ll} \mathcal{N} \left(\frac{z^T \eta_{i_l,c} + 1}{\text{trunc} \mathcal{N}^+ (z^T \eta_{i_l,c})} \right), & y_{parent_{i_l,c}} = 1 \\ \text{trunc} \mathcal{N}^- (z^T \eta_{i_l,c}), & y_{parent_{i_l,c}} = -1 \end{array} \right.$$

where $z = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_c\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,c}\}$

B. Draw a response variable $y_{d,i_l,c} | a_{d,i_l,c} \sim \left\{ \begin{array}{ll} 1, & a_{d,i_l,c} > 0 \\ -1, & \text{otherwise} \end{array} \right.$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{1:N}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \theta, \alpha, \beta, \alpha', \gamma) \propto p(\theta, z_{1:N}, \phi_{1:K}, \eta_{1:N}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \theta, \alpha, \beta, \alpha', \gamma) \quad (1)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{1:N}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \theta, \alpha, \beta, \alpha', \gamma)}{\int_{\theta, \phi, \eta, \alpha, \beta, \alpha', \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{1:N}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \theta, \alpha, \beta, \alpha', \gamma)} \quad (2)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [10]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

3.2.1 $p(z_{m,n} | z_{-(m,n)}, \alpha, \beta, \eta, \alpha', \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z, for a word instance, n, in a document instance, d. The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant:

$$p(z_{d,n} | z_{-(d,n)}, \alpha, \beta, \eta, \alpha', \gamma) \propto p(z_{d,n}, z_{-(d,n)}, \alpha, \beta, \eta, \alpha', \gamma) \quad (3)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_l \in Z} p(a_{i_l,c} | z_{d,n}, \eta_{i_l,c}) p(z_{d,n}, z_{-(d,n)}, \alpha, \beta, \eta, \alpha', \gamma) \quad (4)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [10]:

$$\propto \prod_{i_l \in Z} \exp \left\{ -\frac{(z^T \eta_{i_l,c} - a_{i_l,c})^2}{2} \right\} p(a_{i_l,c} | \eta_{i_l,c}) \left(n_{k-(d,n)}^{k-(d,n)} + \alpha \beta_k \right) \frac{n_{k-(d,n)}^{k-(d,n)} + \gamma}{\sum_{v=1}^V \left(n_{k-(d,n)}^{k-(d,n)} + \alpha \beta_v \right)} \quad (5)$$

Here, $n_{d,v}^{k-(d,n)}$ represents the count of word v in document d assigned to topic k omitting the $(d,n)^{th}$ word count.

Given Equation ??, $p(z_{d,n} | z_{-(d,n)}, \alpha, \beta, \eta, \alpha', \gamma)$ can be sampled through enumeration.

All hyperparameters were given broad gamma priors ($\lambda = \{shape = 2, scale = 1000\}$) and sampled via the Metropolis-Hastings algorithm.

3.3 Prediction

3.4 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [3].

This model, supervised latent dirichlet allocation (sLDA), builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free-text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [2].

Other models - predicting document links, other supervised latent variable models

4 Results

5 Conclusion

References

- [1] International classification of disease. <http://bioportal.bioontology.org/ontologies/35686>, May 2008.
- [2] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing Systems*, 20:121–128, 2008.
- [

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

- Benefits of combining human categorization information into “topic models”

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the primary diagnoses [15].

An automated process would ideally produce a more complete and accurate diagnosis lists. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinical notes. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Background

3 Methods

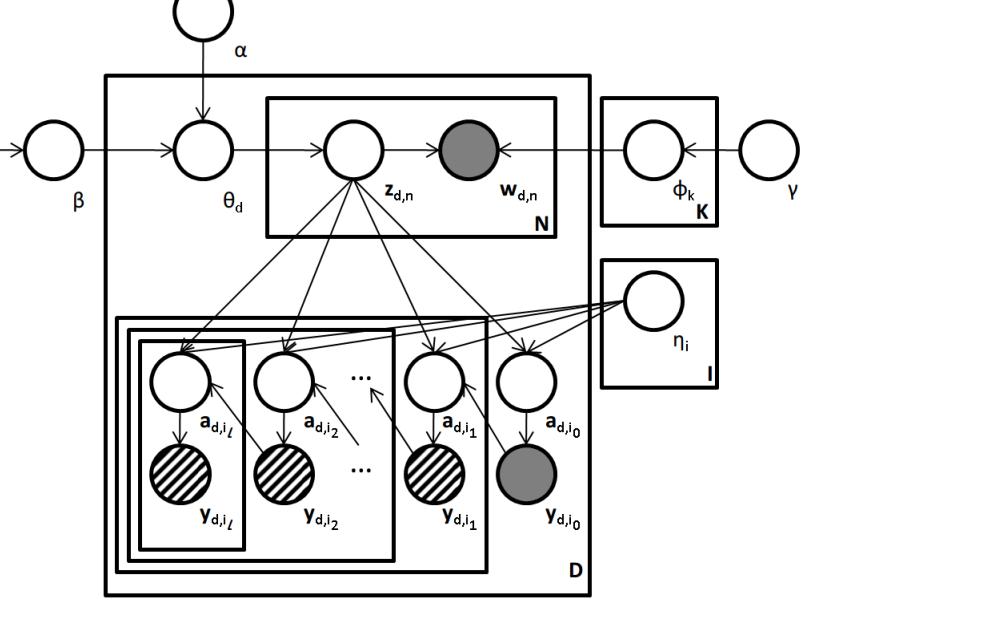


Figure 1: adapted sLDA model

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1

6.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

- For each topic k :
 (a) Draw a distribution over words $\phi_{k,n} \sim Dir(\eta_k, \xi)$
- For each ICD9 code, c , at all levels in the tree:
 (a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
- Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
- For each document, d :
 (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$

(b) For each word, n :

- Draw topic assignment $z_{n,d} \sim Mult_K(\theta_d)$
- Draw word $w_{n,d} | z_{n,d}, \beta_{i,K} \sim Mult_V(\beta_{i,n})$

(c) For each level of the ICD-9 code tree, l :

i. For each ICD-9 code at this level, c :

A. Draw a latent variable

$$a_{d,i_{l,c}} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{i,c}, 1), & y_{parent_{l,c}} = 1 \\ truncN^-(\bar{z}^T \eta_{i,c}, 1), & y_{parent_{l,c}} = -1 \end{cases}$$

where $\bar{z} = N^{-1} \sum_{i=1}^N z_i$ and $Z = \{i_0, i_1, \dots, i_C\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$

$$B. Draw a response variable $y_{d,i_{l,c}} | a_{d,i_{l,c}} \sim \begin{cases} 1, & a_{d,i_{l,c}} > 0 \\ -1, & otherwise \end{cases}$$$

C. Draw a latent variable

$$a_{d,i_{l,c}} \sim \mathcal{N}(\bar{z}^T \eta_{i,c}, 1)$$

where $\bar{z} = N^{-1} \sum_{i=1}^N z_i$ and $Z = \{i_0, i_1, \dots, i_C\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$

D. Draw a response variable $y_{d,i_{l,c}} | a_{d,i_{l,c}} \sim \begin{cases} 1, & a_{d,i_{l,c}} > 0 \\ -1, & otherwise \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

(1)

$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathbb{Z}, a_{i,c} \in \mathbb{Z}, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{i,c} \in \mathbb{Z}, \sigma)$

$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathbb{Z}, a_{i,c} \in \mathbb{Z}, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{i,c} \in \mathbb{Z}, \sigma)}{\int_{\theta, \phi, \alpha, \eta, \alpha', \beta, \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathbb{Z}, a_{i,c} \in \mathbb{Z}, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{i,c} \in \mathbb{Z}, \sigma)}$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalizing factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is simple to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [10]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

3.2.1 $p(z_{m,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

For the purpose of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

(2)

$p(z_{m,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

(3)

$\propto \prod_{i_l \in \mathbb{Z}} p(a_{i_l,c}) p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma)$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [10].

(4)

$\propto \prod_{i_l \in \mathbb{Z}} \exp\left(-\frac{(z_{m,n} - a_{i_l,c})^2}{2}\right) p(a_{i_l,c} | \eta_{i_l,c}) \left(\frac{n_{i_l,c} - (a_{i_l,c} + \alpha_k)}{\sum_{i_l=1}^V n_{i_l,c} + \alpha_k}\right)$

(13)

(14)

(15)

(16)

(17)

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

- Benefits of combining human categorization information into “topic models”

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the primary diagnoses [14].

An automated process would ideally produce a more complete and accurate diagnosis lists. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinical notes. A Latent Dirichlet Allocation (LDA) model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Background

3 Methods

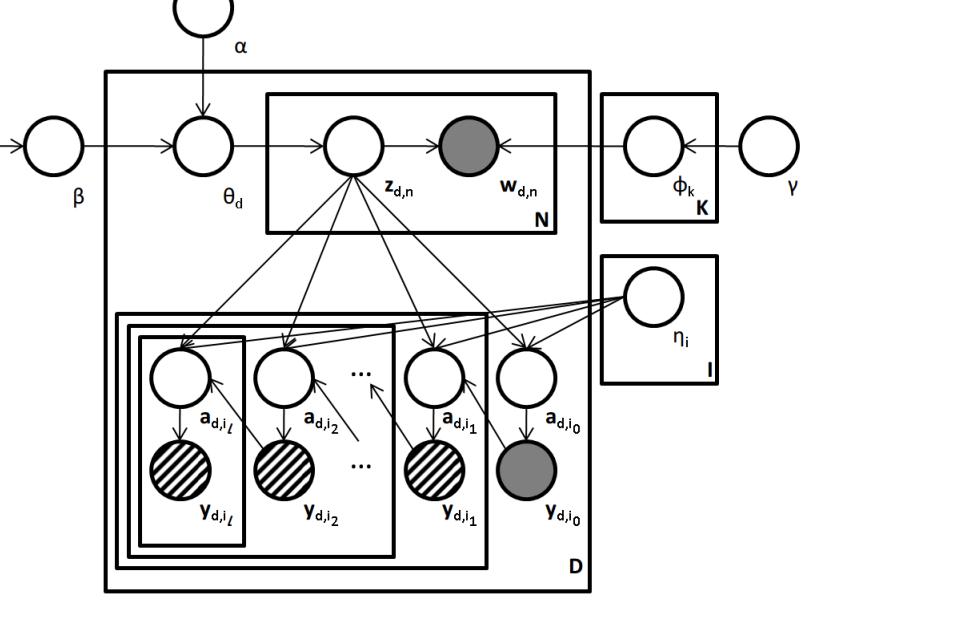


Figure 1: adapted sLDA model

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

6.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

- For each topic k :
 (a) Draw a distribution over words $\phi_{k,n} \sim Dir(\gamma, \beta)$
- For each ICD9 code, c , at all levels in the tree:
 (a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
- Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
- For each document, d :
 (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$

(b) For each word, n :

- Draw topic assignment $z_{n,d} | \theta_d \sim Mult_K(\theta_d)$
- Draw $w_{n,d} | z_{n,d}, \beta_{1,K} \sim Mult_V(\beta_{1,K})$

(c) For each level of the ICD-9 code tree, l :

i. For each ICD-9 code at this level, c :

A. Draw a latent variable

$$a_{d,i_{l,c}} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{i,c}, 1), & y_{parent_{l,c}} = 1 \\ \text{truncN}^-(\bar{z}^T \eta_{i,c}, 1), & y_{parent_{l,c}} = -1 \end{cases}$$

where $\bar{z} = N^{-1} \sum_{i=1}^N z_i$ and $Z = \{i_0, i_1, \dots, i_C\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$

- Draw a response variable $y_{d,i_{l,c}} | a_{d,i_{l,c}} \sim \begin{cases} 1, & a_{d,i_{l,c}} > 0 \\ -1, & \text{otherwise} \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$\begin{aligned} p(z_{1:N}, \phi_{1:K}, \eta_{i,c} | \theta_d, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{i,c} | \sigma, \lambda) &= \frac{p(z_{1:N}, \phi_{1:K}, \eta_{i,c} | \theta_d, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{i,c} | \sigma, \lambda)}{\int_{\theta_d, \phi_{1:K}, \eta_{i,c}} p(z_{1:N}, \phi_{1:K}, \eta_{i,c} | \theta_d, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{i,c} | \sigma, \lambda)} \quad (1) \end{aligned}$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalizing factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is simple to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

$$3.2.1 \quad p(z_{m,n} | z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$$

For the purpose of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

$$p(z_{m,n} | z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_{l,c} \in Z} p(a_{d,i_{l,c}} | \mathbf{z}, \eta_{i,c}) p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [9].

$$\propto \prod_{i_{l,c} \in Z} \exp \left\{ -\frac{(\bar{z}^T \eta_{i,c} - a_{d,i_{l,c}})^2}{2} \right\} \binom{n_{i_{l,c}}^{k_{i_{l,c}}(d,n)} + \alpha \beta_{i_{l,c}}}{\sum_{v=1}^V \binom{n_{i_{l,c}}^{k_{i_{l,c}}(d,n)} + \gamma}{n_{i_{l,c}}^{k_{i_{l,c}}(d,n)} + \alpha_v} \quad (4)$$

$$(13) \quad \int_{\theta_d} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K}$$

$$(14) \quad = \int_{\theta_d} \int_{\phi_{1:K}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K}$$

$$(15) \quad = \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c}$$

$$(16) \quad = \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}}$$

$$(17) \quad = \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}}$$

$$2 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}}$$

$$3 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}}$$

$$4 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)}$$

$$5 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)}$$

$$6 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)}$$

$$7 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)}$$

$$8 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)}$$

$$9 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)}$$

$$10 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)}$$

$$11 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)}$$

$$12 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)}$$

$$13 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)} dy_{(d,n)}$$

$$14 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} p(\theta_d | z_{1:N}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma, \xi) d\theta_d d\phi_{1:K} d\eta_{i,c} da_{d,i_{l,c}} dy_{d,i_{l,c}} dy_{parent_{l,c}} dy_{(d,n)} dy_{(d,n)}$$

$$15 \quad \int_{\theta_d} \int_{\phi_{1:K}} \int_{\eta_{i,c}} \int_{a_{d,i_{l,c}}} \int_{y_{d,i_{l,c}}} \int_{y_{parent_{l,c}}} \int_{y_{(d,n)}} \int_{y_{(d,n)}} \int_{y_{$$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

- Benefits of combining human categorization information into “topic models”
- LDA solved free text
- supervised LDA improves LDA (extra info) and allows new inference (predict links, etc.)
- amazon, freshdirect, netflix, dmoz, pandora (music genome)

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the primary diagnoses [14].

An automated process would ideally produce a more complete and accurate diagnosis lists. Also, this model will reveal information about the medical records themselves. For example, we may gain understanding of why a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of episodic discharges. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Background

3 Methods

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, k :
 - (a) Draw a distribution over words $\phi_k \sim Dir(\gamma, 1)$
 2. For each ICD9 code, c , at all levels in the tree, t :
 - (a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
 3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
 4. For each document, d :
 - (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$
 - (b) For each word, n :
 - i. Draw topic assignment $z_{n,d} | \theta_d \sim Mult_K(\theta_d)$
 - ii. Draw word $w_{n,d} | z_{n,d}, \beta_{i,K} \sim Mult_V(\beta_{i,K})$
 - (c) For each level of the ICD9 code tree, t :
 - i. For each ICD9 code at this level, c :
 - A. Draw a latent variable
 $a_{d,i,c} \sim \begin{cases} N(z^T \eta_{i,c}, 1), & y_{d,i,c} = 1 \\ truncN(z^T \eta_{i,c}, 1), & y_{d,i,c} = -1 \end{cases}$
 - where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_C\}$ and $i_t = \{i_{t,0}, i_{t,1}, \dots, i_{t,C_t}\}$
 - B. Draw a response variable $y_{d,i,c} | a_{d,i,c} \sim \begin{cases} 1, & a_{d,i,c} > 0 \\ -1, & otherwise \end{cases}$

6.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

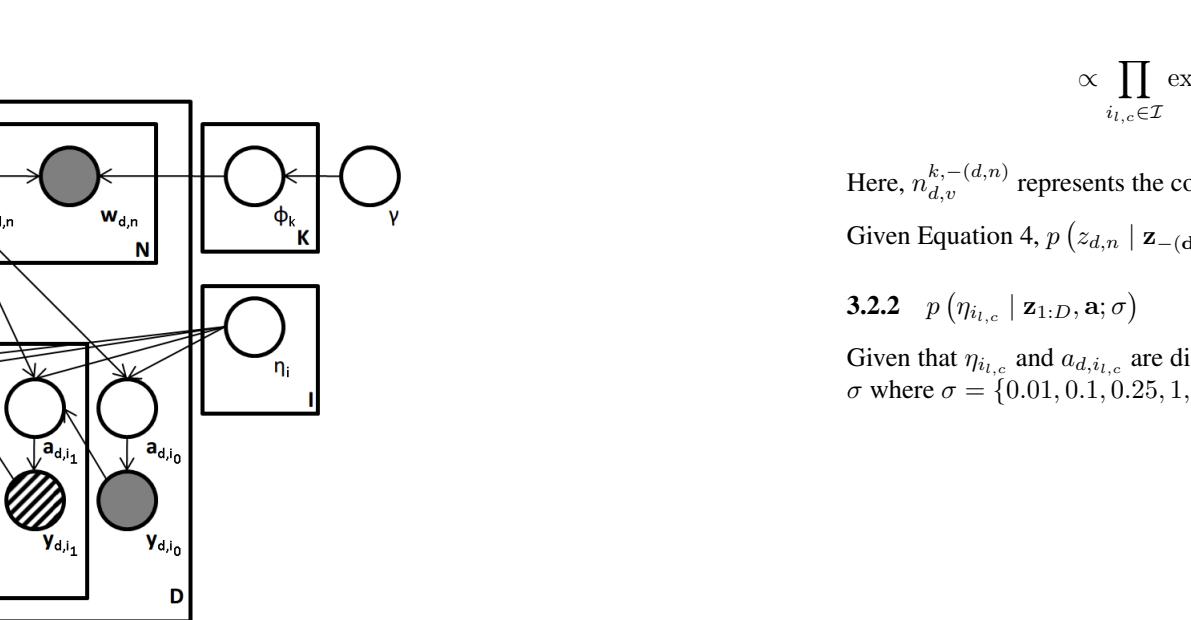


Figure 1: adapted sLDA model

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathcal{Z}, a_{i,c} \in \mathcal{A}, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{i,c} \in \mathcal{Y}; \sigma, \lambda) \quad (1)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathcal{Z}, a_{i,c} \in \mathcal{A}, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{i,c} \in \mathcal{Y}; \sigma, \lambda)}{\int_{\theta, \phi, \eta, \alpha, \alpha', \beta, \gamma} \int_{z_{1:N}} \int_{a_{i,c} \in \mathcal{A}} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathcal{Z}, a_{i,c} \in \mathcal{A}, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{i,c} \in \mathcal{Y}; \sigma, \lambda)}$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

$$3.2.1 p(z_{m,n} | \mathbf{z}_{-(m,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

$$p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i,j \in \mathcal{I}} p(a_{i,j} | \mathbf{z}, \eta_{i,j}) p(z_{d,n} | \mathbf{z}, \eta_{d,n}) \quad (3)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [9].

3.3.2 p(a_{d,i,c} | z, Y, \eta, \alpha, \beta, \gamma)

Given Equation 4, $p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.3.2.2 p(\eta_{i,c} | z_{1:D}, \mathbf{a}, \sigma)

Given that $\eta_{i,c}$ and $a_{d,i,c}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$\propto \prod_{i,c \in \mathcal{I}} \exp \left\{ - \frac{(z^T \eta_{i,c} - a_{d,i,c})^2}{2} \right\} p(a_{d,i,c} | \eta_{i,c}) \left(n_{d,(c)}^{k_{-(d,n)}} + \alpha \beta_k \right) \frac{n_{d,(c)}^{k_{-(d,n)}} + \gamma}{\sum_{v=1}^V (n_{d,(c)}^{k_{-(d,n)}} + \alpha_v)} \quad (4)$$

Here, $n_{d,(c)}^{k_{-(d,n)}}$ represents the count of word v in document d assigned to topic k omitting the $(d, n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

$$3.2.2 p(\eta_{i,c} | z_{1:D}, \mathbf{a}, \sigma)$$

Given that $\eta_{i,c}$ and $a_{d,i,c}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$3.2.2.1 p(\eta_{i,c} | z_{1:D}, \mathbf{a}, \sigma) = \mathcal{N}(\eta_{i,c} | \hat{\mu}_i, \hat{\Sigma}_i) \quad (5)$$

$$\hat{\mu}_i = \hat{\Sigma}_i (-1\sigma^{-1} + \bar{Z}^T \mathbf{a}_{(i),i,c})$$

$$\hat{\Sigma}_i^{-1} = \mathbf{I}\sigma^{-1} + \bar{Z}^T \bar{Z}$$

$$3.2.3 p(a_{d,i,c} | \mathbf{z}, \mathbf{Y}, \eta) \text{ and } p(y_{m,i} | \mathbf{a})$$

In the augmented probit regression model, the posterior distribution of $a_{d,i,c}$ is distributed according to a truncated normal distribution where the response variable is observed.

$$p(a_{d,i,c} | \mathbf{z}, \mathbf{Y}, \eta) = \begin{cases} truncN^+(a_{d,i,c}, \eta^T \bar{z}, 1, y_{d,i,c}) & if \ y_{d,i,c} = 1 \\ truncN^-(a_{d,i,c}, \eta^T \bar{z}, 1, y_{d,i,c}) & if \ y_{d,i,c} = -1 \end{cases} \quad (6)$$

However, if $y_{d,i,c}$ is unobserved then $a_{d,i,c}$ must be sampled jointly with $y_{d,i,c}$ to ensure that the Markov chain is ergodic. Suppose that $a_{d,i,c}$ is sampled to have a negative value and $y_{d,i,c}$ is appropriately sampled at -1. Although there exist states where $a_{d,i,c} > 0$ and $y_{d,i,c} = 1$, they will never be reached by such a Markov chain since $p(a_{d,i,c} < 0 | y_{d,i,c} = -1) = 1$ and $p(y_{d,i,c} = -1 | a_{d,i,c} < 0) = 1$. Therefore, to ensure ergodicity, $a_{d,i,c}$ and $y_{d,i,c}$ must be sampled from the joint distribution as shown in Equation ??.

$$p(a_{d,i,c}, y_{d,i,c} | \mathbf{z}, \mathbf{Y}, \eta) \propto p(y_{d,i,c} | \mathbf{a}, \mathbf{y}_{-(l,c)}) p(a_{d,i,c} | \mathbf{z}, \mathbf{Y}, \eta) \quad (7)$$

$$p(a_{d,i,c}, y_{d,i,c} | \mathbf{z}, \mathbf{Y}, \eta) = \delta(sign(a_{d,i,c}) = y_{d,i,c}) p(y_{d,i,c} | y_{parents_{l,c}}) \prod_{i_t \in children_{l,c}} p(y_{i_t} | y_{d,i,c}) \quad (8)$$

$$p(y_{d,i,c} = -1 | y_{parents_{l,c}}) = \begin{cases} 1, & y_{parents_{l,c}} = -1 \\ 0, & y_{parents_{l,c}} = 1 \end{cases} \quad (9)$$

$$p(a_{d,i,c}, y_{d,i,c} | \mathbf{z}, \mathbf{Y}, \eta) = \begin{cases} N(a_{d,i,c} | \bar{z}^T \eta_{d,i,c}, 1) p(y_{d,i,c} | a_{d,i,c}), & y_{parents_{l,c}} = 1, \forall y_{i_t} \in y_{children_{l,c}}, y_{i_t} = -1 \\ truncN^-(a_{d,i,c} | \bar{z}^T \eta_{d,i,c}, 1), \delta(y_{d,i,c} = -1), & y_{parents_{l,c}} = -1 \\ truncN^+(a_{d,i,c} | \bar{z}^T \eta_{d,i,c}, 1), \delta(y_{d,i,c} = 1), & \exists y_{i_t} \in y_{children_{l,c}} \setminus y_{d,i,c} = 1 \\ 0, & otherwise \end{cases} \quad (10)$$

where $\mathbf{Y}_{-(d,i,c)}$ denotes all of the response variables excluding the response variable being sampled.

3.2.4 p(\beta | z, \alpha', \alpha)

In our model, we place a hierarchical Dirichlet prior over topic assignments. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion. This flexible distribution allows for an asymmetric prior over document level distributions over topics [Wallach et al. [16]].

Posterior inference was performed using the “direct assignment” method of Teh et al. [15].

$$\beta \sim Dir(m_{(.),1}, m_{(.),2}, \dots, m_{(.),K}) \quad (11)$$

$$p(m_{d,k} | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m_{(.),k})^{\alpha \beta_k} \quad (12)$$

represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [3].

This model, supervised latent dirichlet allocation (sLDA), builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [2].

Other models - predicting document links, other supervised latent variable models

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [13, 7, 12, 4], fully automatic assignment of ICD-9 codes became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [11] but the most recent and promising work on the subject was inspired by the 2007 Medical NLP Challenge: “International Challenge: Classifying Clinical Free Text Using Natural Language Processing” (website). Most of the classification strategies included word matching and rule-based algorithms. [8, 5, 6]. The data set given to the participants consisted only of documents that were 1-2 lines each and all of the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [10]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

5 Results

6 Conclusion

References

- [1] International classification of disease. <http://bioportal.bioontology.org/ontologies/35686>, May 2008.
- [2] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing Systems*, 20:121–128, 2008.
- [3] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [4] P. Brown, DG Cocke, and JR Allegre. The ngram cc classifier: A novel method of automatically creating cc classifiers based on icd9 groupings. *Advances in Disease Surveillance*, 1:30, 2006.
- [5] K. Crammer, M. Dredze, K. Ganchev, PF Tkalukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [6] R. Farkas and G. Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMCS bioinformatics*, 9(Suppl 3):S10, 2008.
- [7] HR Freitas, B. RibeiroNeto, AHF Laender, and LRS De Lima. Categorizationdriven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.
- [8] I. Goldstein, A. Arzumanyan, and O. Uzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [9] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl):15228–15235, 2004.
- [10] LV Lita, S Yu, S Niculescu, and JBL. Large scale diagnostic code classification for medical patient records. 2008.
- [11] RB Rao, S Sandilya, SR Niculescu, C Germond, and H Rao. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 416–425, 2003.
- [12] B RibeiroNeto, AHF Laender, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and technology*, 52(5):391–401, 2001.
- [13] P. Ruch, J. Golbini, I. Thabirin, and A. Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [14] G. Surjan. Questions on validation of international classification of diseases-coded diagnoses. *International journal of medical informatics*, 54(2):77–95, 1999.
- [15] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [16] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking lda: Why priors matter. In Y. Bengio, D. Schuurmans, J. Laff

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

There exist surprisingly many sources of unstructured text data that have been partially or completely categorized by human editors. Examples include hierarchical directories of webpages [7], large hierarchically annotated product catalogs (e.g. [7]) as available from [7]), manually annotated patient medical records, and more. In this work we show how to combine these two sources of information in a single model that allows us to, amongst other things, automatically annotate and/or categorize new text documents, effectively inserting them into the category tree.

- Benefits of combining human categorization information into “topic models”
- LDA solved free text
- supervised LDA improves LDA (extra info) and allows new inference (predict links, etc.)
- amazon, freshdirect, netflix, dmoz, pandora (music genome)

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (s-LDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [14].

An automated process would ideally produce a more complete and accurate diagnostic list. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Background

3 Methods

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, k :
 (a) Draw a distribution over words $\theta_k \sim Dir(\gamma, 1 - \gamma)$
2. For each ICD9 code, c , at all levels in the tree, l :
 (a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
4. For each document, d :
 (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir(\gamma, 1 - \gamma)$
 (b) For each word, n :
 i. Draw topic assignment $z_{d,n} | \theta_d \sim Mult_K(\theta_d)$
 ii. Draw word $w_{d,n} | z_{d,n}, \beta_i, \alpha \sim Mult_V(\beta_{z_n})$
 (c) For each level of the ICD-9 code tree, l :
 i. For each ICD-9 code at this level, c :

1

6.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

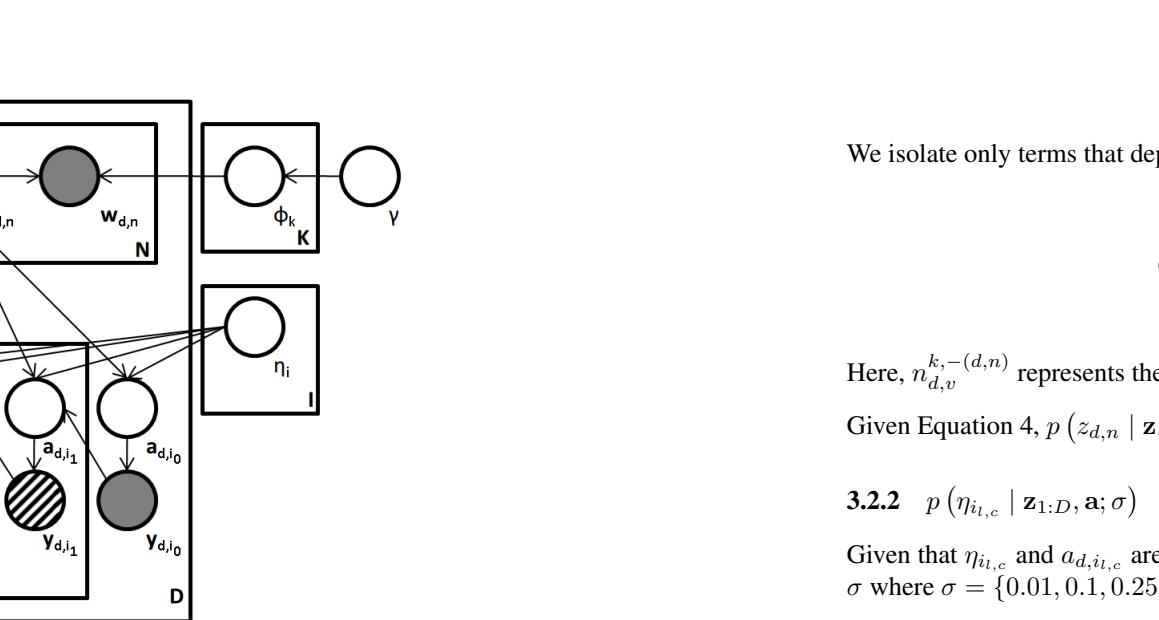


Figure 1: adapted sLDA model

A. Draw a latent variable

$$a_{d,i_c} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{i_c}, 1), & y_{i_c} = 1 \\ \text{truncN}(\bar{z}^T \eta_{i_c}, 1), & y_{i_c} = -1 \end{cases}$$

where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $\mathcal{I} = \{i_1, i_2, \dots, i_{l,1}, \dots, i_{l,1}\}$

$$\text{B. Draw a response variable } y_{d,i_c} \sim \begin{cases} 1, & a_{d,i_c} > 0 \\ -1, & \text{otherwise} \end{cases}$$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 1.

$$\begin{aligned} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, a_{d,i_c}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{i_c} \in \mathcal{I}; \sigma, \lambda) \\ = \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, a_{d,i_c}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{i_c} \in \mathcal{I}; \sigma, \lambda)}{\sum_{\theta, \phi, \eta, a, \alpha, \beta, \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, a_{d,i_c}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{i_c} \in \mathcal{I}; \sigma, \lambda)} \end{aligned} \quad (1)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

$$3.2.1 \quad p(z_{m,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z_m for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant:

$$p(z_{m,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$(13) \quad \int_{\theta_{1:D}} p(\theta_{1:D} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\phi_{1:K}} p(\phi_{1:K} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\eta_{i_c}} p(\eta_{i_c} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{a_{d,i_c}} p(a_{d,i_c} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\theta_d} p(\theta_d | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{z_{d,n}} p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{w_{d,n}} p(w_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{y_{d,i_c}} p(y_{d,i_c} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}} p(\mathbf{a} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}} p(\mathbf{w} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}} p(\mathbf{z} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{a}_{-(d,n)}} p(\mathbf{a}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{w}_{-(d,n)}} p(\mathbf{w}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}_{-(d,n)}} p(\mathbf{y}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{z}_{-(d,n)}} p(\mathbf{z}_{-(d,n)} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{y}} p(\mathbf{y} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \int_{\mathbf{$$

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

There exist surprisingly many sources of unstructured text data that have been partially or completely categorized by human editors. Examples include hierarchical directories of webpages [7], large hierarchically annotated product catalogs (e.g. [7]) as available from [7]), manually annotated patient medical records, and much more. In this work we show how to combine these two sources of information in a single model that allows us to, amongst other things, automatically annotate and/or categorize new text documents, effectively inserting them into the category tree.

- Benefits of combining human categorization information into “topic models”
- LDA solved free text
- supervised LDA improves LDA (extra info) and allows new inference (predict links, etc.)
- amazon, freshdirect, netflix, dmoz, pandora (music genome)

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [14].

An automated process would ideally produce a more complete and accurate diagnostic list. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Background

3 Methods

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, k :
 (a) Draw a distribution over words $\theta_k \sim Dir(\gamma, 1 - \gamma)$
2. For each ICD9 code, c , at all levels in the tree, l :
 (a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
4. For each document, d :
 (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir(\gamma, 1 - \gamma)$
 (b) For each word, n :
 i. Draw topic assignment $z_{d,n} | \theta_d \sim Mult_K(\theta_d)$
 ii. Draw word $w_{d,n} | z_{d,n}, \beta_i, \alpha \sim Mult_V(\beta_{z_n})$
 (c) For each level of the ICD-9 code tree, l :
 i. For each ICD-9 code at this level, c :

1

6.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

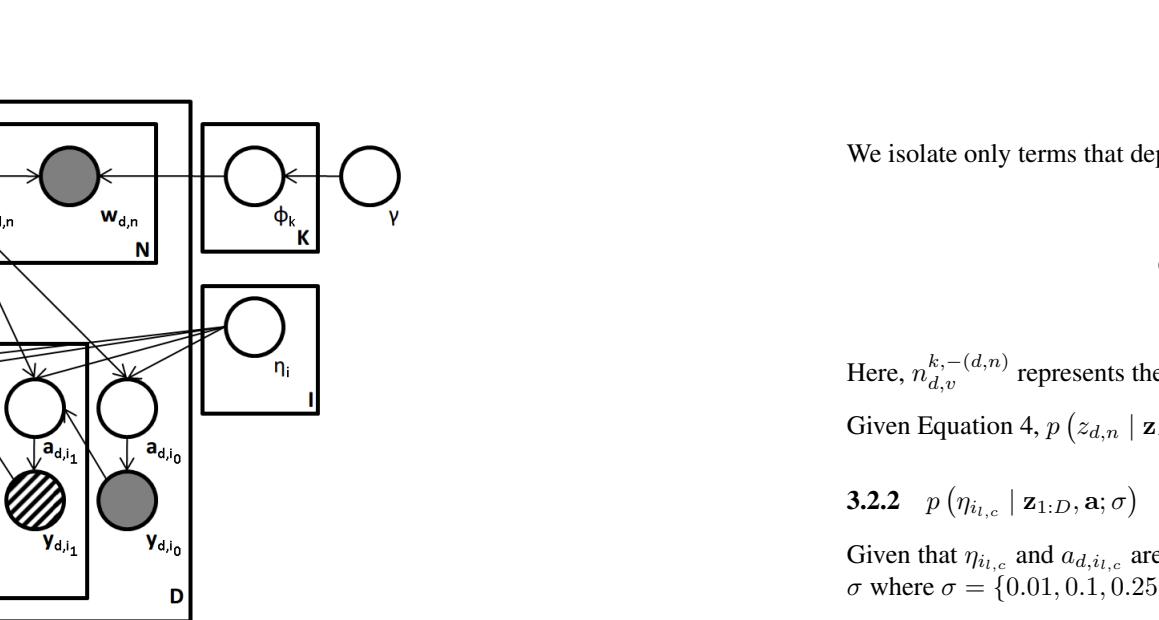


Figure 1: adapted sLDA model

A. Draw a latent variable

$$a_{d,i_c} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{i_c}, 1) & y_{parent,c} = 1 \\ \text{truncN}(\bar{z}^T \eta_{i_c}, 1) & y_{parent,c} = -1 \end{cases}$$

where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $\mathcal{I} = \{i_1, i_2, \dots, i_{C_l}\}$

$$y_{d,i_c} \sim \begin{cases} 1, & a_{d,i_c} > 0 \\ -1, & \text{otherwise} \end{cases}$$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 1.

$$\begin{aligned} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, a_{d,i_c}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{i_c}, \bar{z}; \sigma, \lambda) \\ = \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, a_{d,i_c}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{i_c}, \bar{z}; \sigma, \lambda)}{\sum_{\theta, \phi, \eta, a, \alpha, \beta, \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, a_{d,i_c}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{i_c}, \bar{z}; \sigma, \lambda)} \end{aligned} \quad (1)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

3.2.1 $p(z_{m,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z_m for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

$$p(z_{m,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\begin{aligned} (13) & p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, a_{d,i_c}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{i_c}, \bar{z}; \sigma, \lambda) \\ (14) & = \int_0^\infty p(\theta | \alpha) \sum_{k=1}^K \left(\prod_{n=1}^N p(z_{n,k} | \theta_n, \phi_{n,k}, \xi) \right) \left(\prod_{n=1}^N p(w_{n,k} | z_{n,k}, \eta_{n,k}, \xi) \right) \left(\prod_{n=1}^N p(a_{n,k} | \theta_n, \phi_{n,k}, \xi) \right) d\theta \end{aligned}$$

$$(15) \quad p(\mathbf{Y}, \mathbf{w}, \mathbf{z}, \mathbf{a}, \phi_{1:K}, \beta, \mu, \xi) = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,1:k} | \phi_{1:k-1}, \beta) \prod_{n=1}^N p(z_{n,k} | \theta_n, \phi_{n,k}, \xi) d\phi_{1:K}$$

$$(16) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] d\phi_{1:K}$$

$$(17) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(18) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(19) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(20) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(21) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(22) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(23) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(24) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(25) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(26) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(27) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(28) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(29) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(30) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(31) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\phi_{n,k})} \right]$$

$$(32) \quad = \prod_{k=1}^K \left[\frac{\Gamma(\sum_{n=1}^N \Gamma(\theta_n))}{\prod_{n=1}^N \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\theta_n))}{\prod_{k=1}^K \Gamma(\theta_n)} \right] \prod_{n=1}^N \left[\frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{n,k}))}{\prod_{k=1}^K \Gamma(\$$

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

There exist surprisingly many sources of unstructured text data that have been partially or completely categorized by human editors. Examples include hierarchical directories of webpages [7], large hierarchically annotated product catalogs (e.g. [7]) as available from [7]), manually annotated patient medical records, and more. In this work we show how to combine these two sources of information in a single model that allows us to, amongst other things, automatically annotate and/or categorize new text documents, effectively inserting them into the category tree.

- Benefits of combining human categorization information into “topic models”
- LDA solved free text
- supervised LDA improves LDA (extra info) and allows new inference (predict links, etc.)
- amazon, freshdirect, netflix, dmoz, pandora (music genome)

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (s-LDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [14].

An automated process would ideally produce a more complete and accurate diagnostic list. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the sLDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Background

3 Methods

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, k :
 (a) Draw a distribution over words $\theta_k \sim Dir(\gamma, 1 - \gamma)$
2. For each ICD9 code, c , at all levels in the tree, l :
 (a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
4. For each document, d :
 (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir(\gamma, 1 - \gamma)$
 (b) For each word, n :
 i. Draw topic assignment $z_{d,n} | \theta_d \sim Mult_K(\theta_d)$
 ii. Draw word $w_{d,n} | z_{d,n}, \beta_i, \alpha \sim Mult_V(\beta_{z_n})$
5. For each level of the ICD-9 code tree, l :
 i. For each ICD-9 code at this level, c :

1

6.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

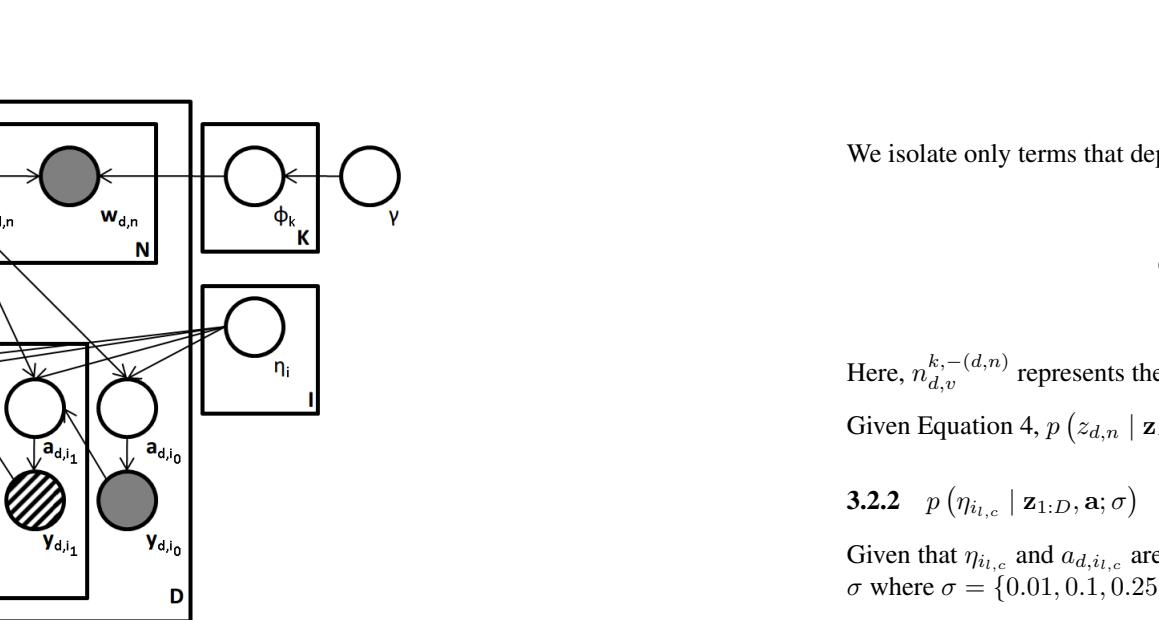


Figure 1: adapted sLDA model

A. Draw a latent variable

$$a_{d,i_c} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{i_c}, 1) & y_{parent,c} = 1 \\ \text{truncN}(\bar{z}^T \eta_{i_c}, 1) & y_{parent,c} = -1 \end{cases}$$

where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $\mathcal{I} = \{i_1, i_2, \dots, i_{C_l}\}$

$$a_{d,i_c} \sim \begin{cases} 1, & a_{d,i_c} > 0 \\ -1, & \text{otherwise} \end{cases}$$

B. Draw a response variable

$$y_{d,i_c} | a_{d,i_c} \sim \begin{cases} 1, & a_{d,i_c} > 0 \\ -1, & \text{otherwise} \end{cases}$$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 1.

$$\begin{aligned} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \bar{z}; \sigma, \lambda) \\ = \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{1:N}, \bar{z}; \sigma, \lambda)}{\sum_{\theta, \phi, \eta, \alpha, \beta, \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{1:N}, \bar{z}; \sigma, \lambda)} \end{aligned} \quad (1)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

3.2.1 $p(z_{m,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z_m for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant:

$$p(z_{m,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\begin{aligned} (13) & p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \bar{z}; \sigma, \lambda) \\ & = \frac{p(\theta | \alpha) (\prod_{k=1}^K p(\phi_{k|K} | \theta) p(\phi_{k|K})) (\prod_{k=1}^K p(\eta_{i_c} | \theta, \eta_{i_c}, \phi_{k|K})) (\prod_{k=1}^K p(w_{d,n} | \theta, z_{d,n}, \phi_{k|K}))}{\int_{\theta} p(\theta) \sum_{k=1}^K (\prod_{n=1}^N p(z_{n,k} | \theta, \eta_{i_c}, \phi_{k|K})) (\prod_{n=1}^N p(w_{d,n} | \theta, z_{d,n}, \phi_{k|K})) d\theta} \end{aligned}$$

$$(14) \quad \int_{\theta} \prod_{k=1}^K p(\phi_{k|K} | \theta) \prod_{n=1}^N p(z_{n,k} | \theta, \eta_{i_c}, \phi_{k|K}) d\theta = \prod_{k=1}^K \int_{\phi_{k|K}} p(\phi_{k|K} | \theta) d\phi_{k|K} \prod_{n=1}^N \int_{\theta} p(z_{n,k} | \theta, \eta_{i_c}, \phi_{k|K}) d\theta$$

$$(15) \quad \prod_{k=1}^K \int_{\phi_{k|K}} p(\phi_{k|K} | \theta) d\phi_{k|K} = \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(16) \quad \prod_{n=1}^N \int_{\theta} p(z_{n,k} | \theta, \eta_{i_c}, \phi_{k|K}) d\theta = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(17) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(18) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(19) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(20) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(21) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(22) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(23) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(24) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(25) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(26) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))} = \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}{\prod_{n=1}^N \Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))}$$

$$(27) \quad \prod_{n=1}^N \frac{\Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \beta_k))}{\prod_{k=1}^K \Gamma(\sum_{k=1}^K \Gamma(\phi_{k|K}^n + \alpha_k))} \prod_{k=1}^K \frac{\Gamma(\sum_{n=1}^N \Gamma(\phi_{k|K}^n + \alpha_k))$$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured clinical data that has been at least in part manually categorized. Examples include patient records that are not linked to diagnoses and curated hierarchical directories of the same [2]. Product descriptions and catalogs (e.g. [2, 1]) as available from [2] and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records and the International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned to them[1]). In this work we show how to combine these two sources of information using a single model that allows one, among other things, to automatically categorize new text documents, suggest labels that might be inaccurate, and compute improved similarities between documents for information retrieval purposes. The models and techniques that we develop in this paper are applicable in other domains as well, for instance, unstructured representations of data that have been hierarchically categorized.

In this work we extend supervised latent Dirichlet allocation (sLDA) to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) augmented with per document “supervision” often taking the form of a single numerical or categorical “label.” More generally this “supervision” can be seen as extra data generated about a document; for instance its quality or relevance (e.g. online reviews), marks given to written work (e.g. graded exams), or the number of times a web document is linked. These labels are usually modeled as having been generated conditioned on the topic of topics found in a document. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. As one concrete example consider web retailers. They often have both a browsable hierarchy and free-text descriptions of all products they sell. The situation of each product in the product hierarchy (often multiply classified) can be seen as a form of multiple labeling and as a similar products based on the similarity of their textual description is an *u*. An equivalent challenge, particularly for larger retailers, is to situate the merchandise in as many categories as possible.

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [14].

• Benefits of combining human categorization information into “topic models”
 • LDA solved free text
 • supervised LDA improves LDA (extra info) and allows new inference (predict links, etc.)
 • amazon, freshdirect, netflix, dmoz, pandora (music genome)

2 Background

3 Methods
 Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, *k*:
 (a) Draw a distribution over words $\phi_k \sim Dir(\alpha, \beta)$
2. For each ICD9 code, *c*, at all levels in the tree, *l*:
 (a) Draw regression coefficient $a_{l,i,c} | \sigma \sim \mathcal{N}_K(-1, \sigma)$
3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$

1

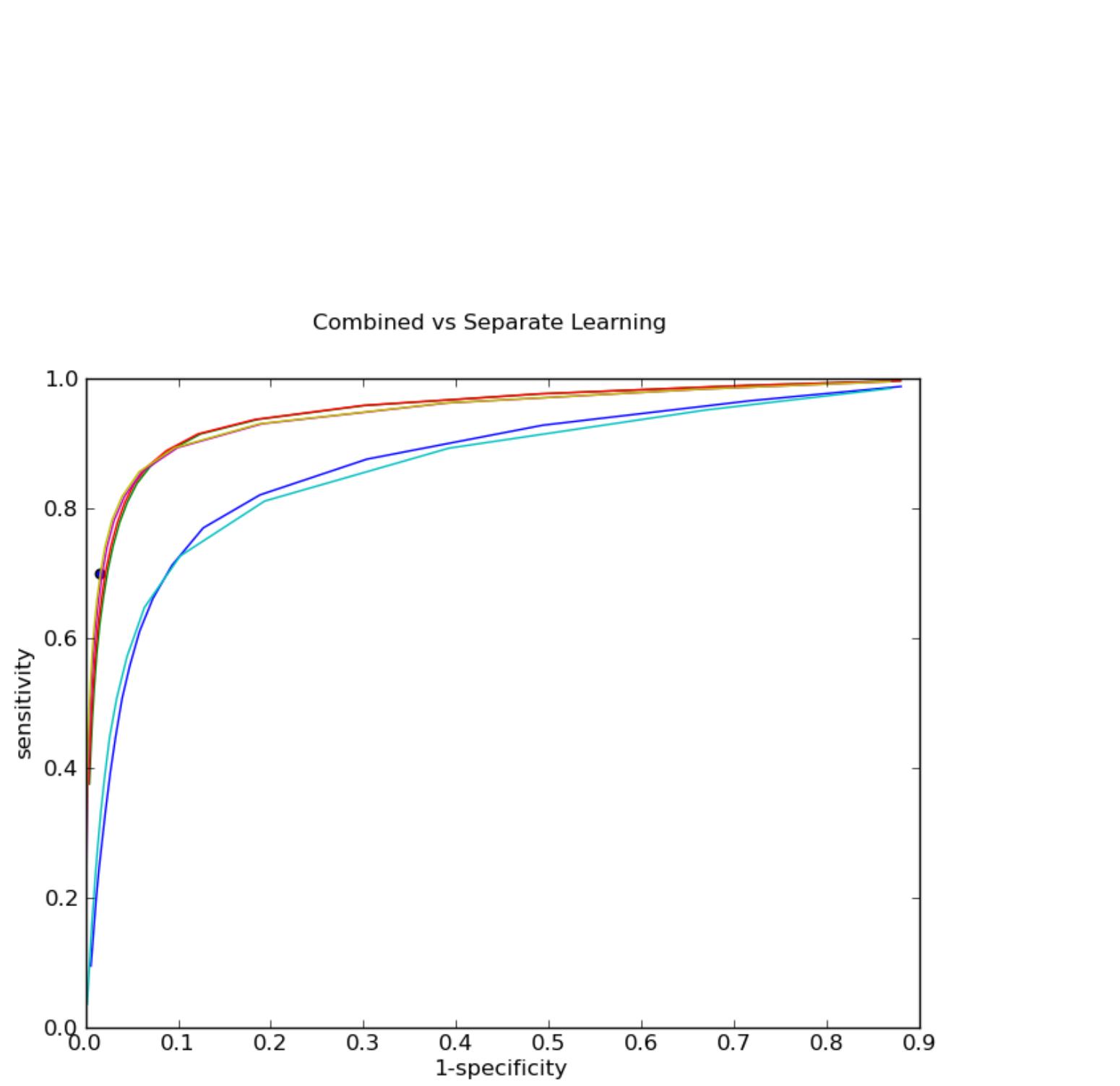


Figure 2: default

2

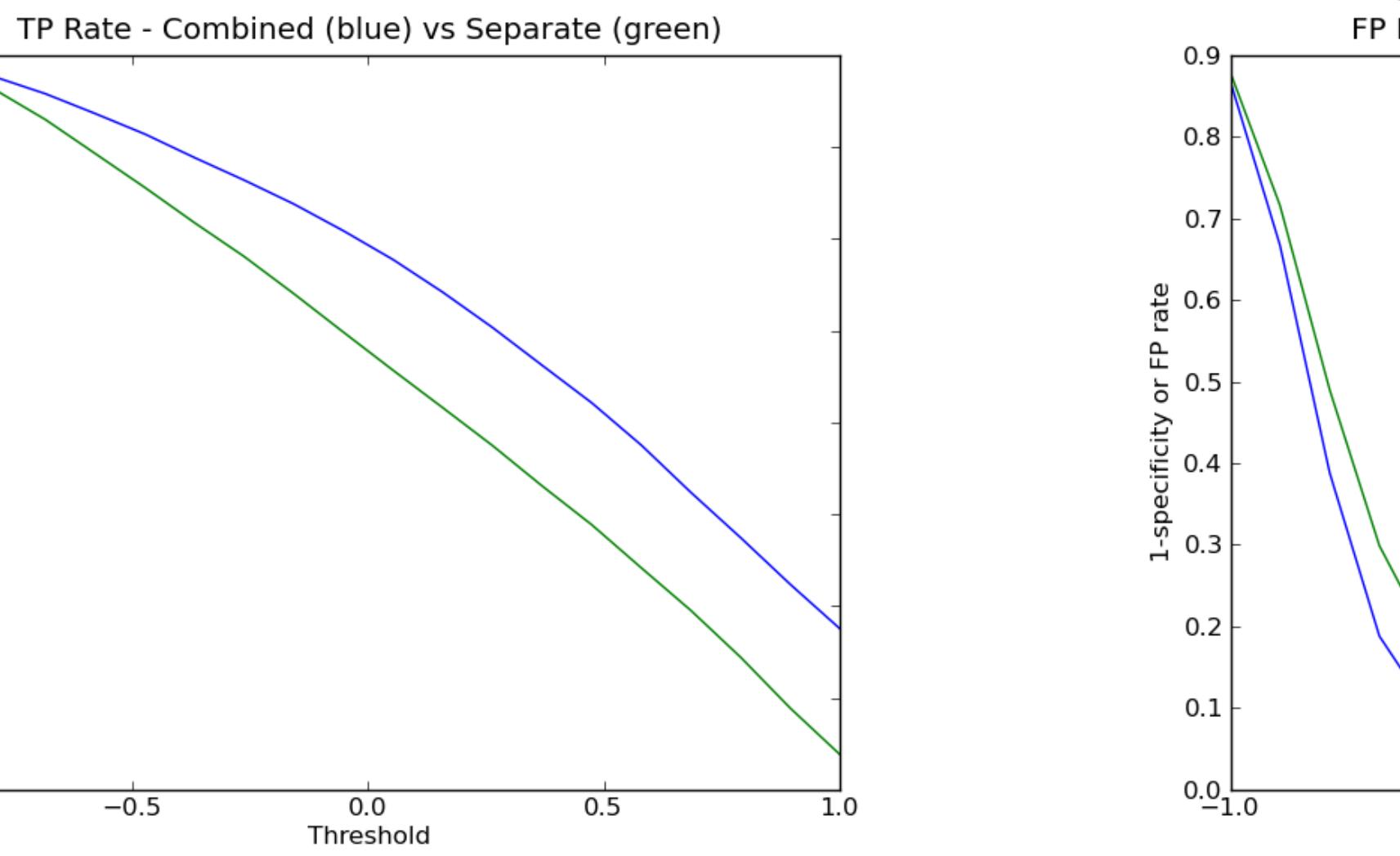


Figure 3: default

3

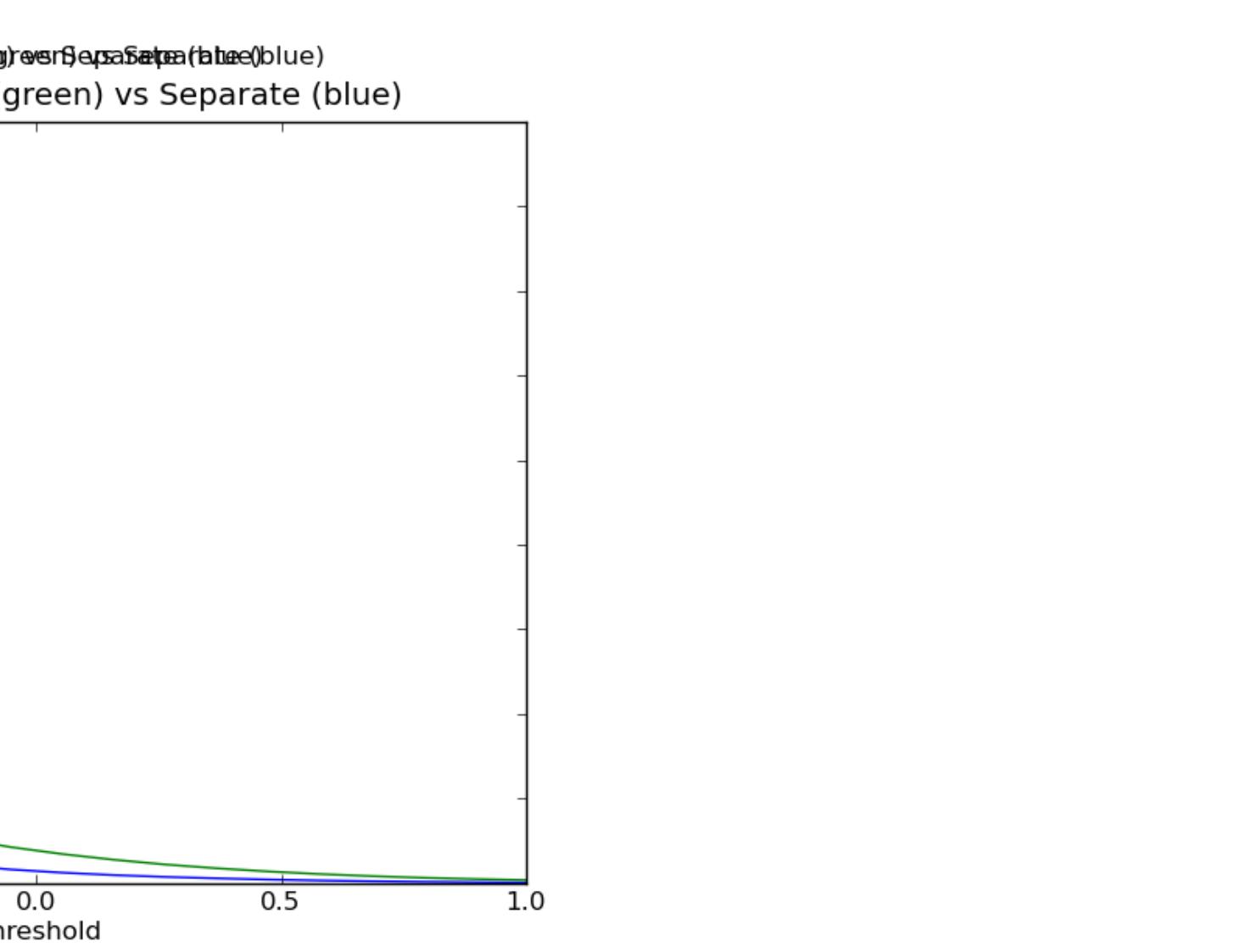


Figure 4: default

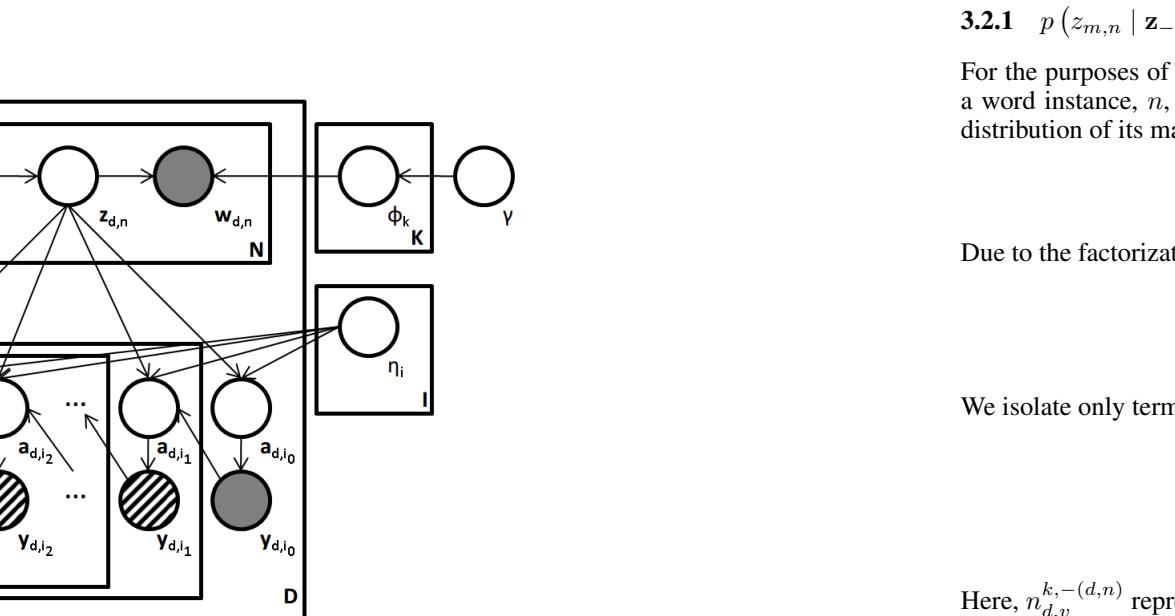


Figure 1: adapted sLDA model

4. For each document, *d*:

- (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$
- (b) For each word, *n*:
 - i. Draw topic assignment $z_{d,n} | \theta_d \sim Mult_K(\theta_d)$
 - ii. Draw word $w_{d,n} | z_{d,n}, \beta_{1,K} \sim Mult(\beta_{z_{d,n}})$
- (c) For each level of the ICD-9 code tree, *l*:
 - i. For each ICD-9 code at this level, *c*:
 - A. Draw a latent variable $a_{l,i,c} \sim \begin{cases} \mathcal{N}(\bar{z}_{l,i,c}, 1), & y_{l,i,c} = 1 \\ \text{truncN}^+(\bar{z}_{l,i,c}, 1), & y_{l,i,c} = -1 \end{cases}$

where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_C\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$

B. Draw a response variable $y_{l,i,c} | a_{l,i,c} \sim \begin{cases} 1, & a_{l,i,c} > 0 \\ -1, & \text{otherwise} \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,i,c} | \mathbf{z}, \mathbf{Y}, \mathbf{a}, \mathbf{w}, \mathbf{y}_{l,i,c}; \sigma, \lambda) \quad (1)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,i,c} | \mathbf{z}, \mathbf{Y}, \mathbf{a}, \mathbf{w}, \mathbf{y}_{l,i,c}; \sigma, \lambda)}{\int_{\theta, z_{1:N}, \phi_{1:K}, \eta_{l,i,c}} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,i,c} | \mathbf{z}, \mathbf{Y}, \mathbf{a}, \mathbf{w}, \mathbf{y}_{l,i,c}; \sigma, \lambda)} \quad (2)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

3.2.1 $p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, *z*, for a word instance, *n*, in a document instance, *d*. The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant:

$$p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_{l,c} \in \mathcal{Z}} p(a_{l,i_{l,c}} | \mathbf{z}, \eta_{l,i_{l,c}}) p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{d,n}$ and absorb all other constant terms into the normalization constant [9].

$$\propto \prod_{i_{l,c} \in \mathcal{Z}} \exp \left\{ -\frac{(\bar{z}^T \eta_{l,i_{l,c}} - a_{l,i_{l,c}})^2}{2} \right\} \left(n_{d,-(d,n)}^{k_{d,-(d,n)}} + \alpha \right) \sum_{i_{l,c} \in \mathcal{Z}} \left(\frac{n_{d,i_{l,c}}^{k_{d,i_{l,c}}} + \gamma}{n_{d,-(d,n)} + \gamma} \right) \quad (4)$$

Here, $n_{d,v}^{k_{d,v}}$ represents the count of word *v* in document *d* assigned to topic *k* omitting the $(d, n)^{\text{th}}$ word count.

Given Equation 4, $p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.2.2 $p(\eta_{l,i_{l,c}} | \mathbf{z}_{1:D}, \mathbf{a}, \sigma)$

Given that $\eta_{l,i_{l,c}}$ and $a_{l,i_{l,c}}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$p(\eta_{l,i_{l,c}} | \mathbf{z}_{1:D}, \mathbf{a}, \sigma) = \mathcal{N}(\eta_{l,i_{l,c}} | \bar{\mu}_l, \bar{\Sigma}_l) \quad (5)$$

$$\bar{\mu}_l = \bar{\Sigma}_l (-1\sigma^{-1} + \bar{Z}^T \mathbf{a}_{(l,i_{l,c})})$$

$$\bar{\Sigma}_l^{-1} = \mathbf{I}\sigma^{-1} + \bar{Z}^T \bar{Z}$$

3.2.3 $p(a_{l,i_{l,c}} | \mathbf{z}, \mathbf{Y}, \eta)$ and $p(y_{l,i_{l,c}} | \mathbf{a})$

In the augmented probit regression model, the posterior distribution of $a_{l,i_{l,c}}$ is distributed according to a truncated normal distribution where the response variable is observed.

$$p(a_{l,i_{l,c}} | \mathbf{z}, \mathbf{Y}, \eta) = \begin{cases} \text{truncN}^+(a_{l,i_{l,c}} | \eta_l^T \bar{z}_{l,i_{l,c}}, 1) & \text{if } y_{l,i_{l,c}} = 1 \\ \text{truncN}^-(a_{l,i_{l,c}} | \eta_l^T \bar{z}_{l,i_{l,c}}, 1) & \text{if } y_{l,i_{l,c}} = -1 \end{cases} \quad (6)$$

However, if $y_{l,i_{l,c}}$ is unobserved then $a_{l,i_{l,c}}$ must be sampled jointly with $y_{l,i_{l,c}}$ to ensure that the Markov chain is ergodic. Suppose that $a_{l,i_{l,c}}$ is sampled to have a negative value and $y_{l,i_{l,c}}$ is appropriately sampled at -1. Although there exist valid states where $a_{l,i_{l,c}} > 0$ and $y_{l,i_{l,c}} = 1$, they will never be reached by such a Markov chain since $p(a_{l,i_{l,c}} < 0 | y_{l,i_{l,c}} = -1) = 1$ and $p(y_{l,i_{l,c}} = 1 | a_{l,i_{l,c}} < 0) = 1$. Therefore, to ensure ergodicity, $a_{l,i_{l,c}}$ and $y_{l,i_{l,c}}$ must be sampled from the joint distribution as shown in Equation ??.

$$p(a_{l,i_{l,c}}, y_{l,i_{l,c}} | \mathbf{z}, \mathbf{Y}, \eta) \propto p(a_{l,i_{l,c}} | \mathbf{z}, \mathbf{Y}, \eta) p(y_{l,i_{l,c}} | a_{l,i_{l,c}}, \mathbf{z}, \mathbf{Y}, \eta) \quad (7)$$

$$p(y_{l,i_{l,c}} | \mathbf{a}, \mathbf{y}_{-(l,c)}) = \delta(\text{sign}(a_{l,i_{l,c}}) = y_{l,i_{l,c}}) p(y_{l,i_{l,c}} | y_{l,i_{l,c}}) \prod_{i_{l,c} \in \mathcal{Z} \setminus l} p(y_{l,i_{l,c}} | y_{l,i_{l,c}}) \quad (8)$$

$$p(y_{l,i_{l,c}} = -1 | y_{l,i_{l,c}}) = \begin{cases} 1, & y_{l,i_{l,c}} = -1 \\ 0.5, & y_{l,i_{l,c}} = 1 \end{cases} \quad (9)$$

$$p(a_{l,i_{l,c}}, y_{l,i_{l,c}} | \mathbf{z}, \mathbf{Y}, \eta) \quad (10)$$

$$= \begin{cases} \mathcal{N}(a_{l,i_{l,c}} | \bar{z}^T \eta_{l,i_{l,c}}, 1) p(y_{l,i_{l,c}} | a_{l,i_{l,c}}, \mathbf{z}, \mathbf{Y}, \eta), & y_{l,i_{l,c}} = 1, \forall y_{l,i_{l,c}} \in \mathcal{Z} \setminus l, y_{l,i_{l,c}} \neq -1 \\ \text{truncN}^-(a_{l,i_{l,c}} | \bar{z}^T \eta_{l,i_{l,c}}, 1) \delta(y_{l,i_{l,c}} = -1), & y_{l,i_{l,c}} = -1 \\ \text{truncN}^+(a_{l,i_{l,c}} | \bar{z}^T \eta_{l,i_{l,c}}, 1) \delta(y_{l,i_{l,c}} = 1), & y_{l,i_{l,c}} \in \mathcal{Z} \setminus l, y_{l,i_{l,c}} \neq 1 \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

where $\mathbf{Y}_{-(l,c)}$ denotes all of the response variables excluding the response variable being sampled.

3

4.2 Pre-Processing

Patient discharge summaries and their associated ICD-9 diagnoses are stored in two different places in the New York - Presbyterian data warehouse, and so had to be linked together before being fed into the sLDA algorithm. Each discharge note and set of diagnoses were assigned a patient unique identifier (PUIID), allowing the two types of data to be linked.

Natural Language Processing (NLP) techniques were used to process the free-text discharge summaries. First, the Natural Language Toolkit (<http://www.nltk.org/>) was used to tokenize the text. Next, feature selection was performed using a term frequency - inverse document frequency (TF-IDF) algorithm on the entire document set and sorting the words by their TF-IDF values. The top 10,000 words were manually evaluated to eliminate all potentially identifying information. Finally, each discharge summary was converted to a bag-of-words,

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was observed to be absent, all of its descendants could be assumed to be absent as well (e.g., if a patient had a history of hypertension, it could be assumed that they did not have malignant hypertension). Unfortunately, ICD-9 codes observations never include observations of disease absence, therefore to ensure ergodicity, $a_{l,i_{l,c}}$ and $y_{l,i_{l,c}}$ must be sampled from the joint distribution as shown in Equation ??.

ICD-9 code at all if neither it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

4.3 ICD-9 Code Hierarchy

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an is-a hierarchy [1]. An is-a hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors

4

6.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

5

5 Results

6 Conclusion

References

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

There exist surprisingly many sources of unstructured text data that have been partially or completely categorized by human editors. Examples include hierarchical directories of webpages [7], large hierarchically annotated product catalogs (e.g. [7]) as available from [7]), manually annotated patient medical records, and more. In this work we show how to combine these two sources of information in a single model that allows us to, amongst other things, automatically annotate and/or categorize new text documents, effectively inserting them into the category tree.

- Benefits of combining human categorization information into “topic models”
- LDA solved free text
- supervised LDA improves LDA (extra info) and allows new inference (predict links, etc.)
- amazon, freshdirect, netflix, dmoz, pandora (music genome)

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (s-LDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [14].

An automated process would ideally produce a more complete and accurate diagnostic lists. Also, this model will reveal information about the medical records themselves. For example, we may gain an understanding of what a specific code actually means in terms of clinical narratives. Similarly, viewing the distribution of topics over discharge summaries may reveal information about the latent structure of clinician documentation. Lastly, the s-LDA model would provide a novel approach to dealing with the problem of high dimensionality when representing narrative text in a vector space specifically by reducing dimensions from an entire vocabulary of potentially tens of thousands of words to a set of several dozen topics.

2 Background

3 Methods

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, k :
 (a) Draw a distribution over words $\theta_k | \beta, \alpha \sim Dir(\beta, \alpha)$
2. For each ICD9 code, c , at all levels in the tree, I :
 (a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N(-1, \sigma)$
3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
4. For each document, d :
 (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir(\beta, \alpha)$
 (b) For each word, n :
 i. Draw topic assignment $z_{d,n} | \theta_d \sim Mult_K(\theta_d)$
 ii. Draw word $w_{d,n} | z_{d,n}, \beta_{i,K} \sim Mult_V(\beta_{i,K})$
- (c) For each level of the ICD-9 code tree, I :
 i. For each ICD-9 code at this level, c :

Figure 1: adapted sLDA model

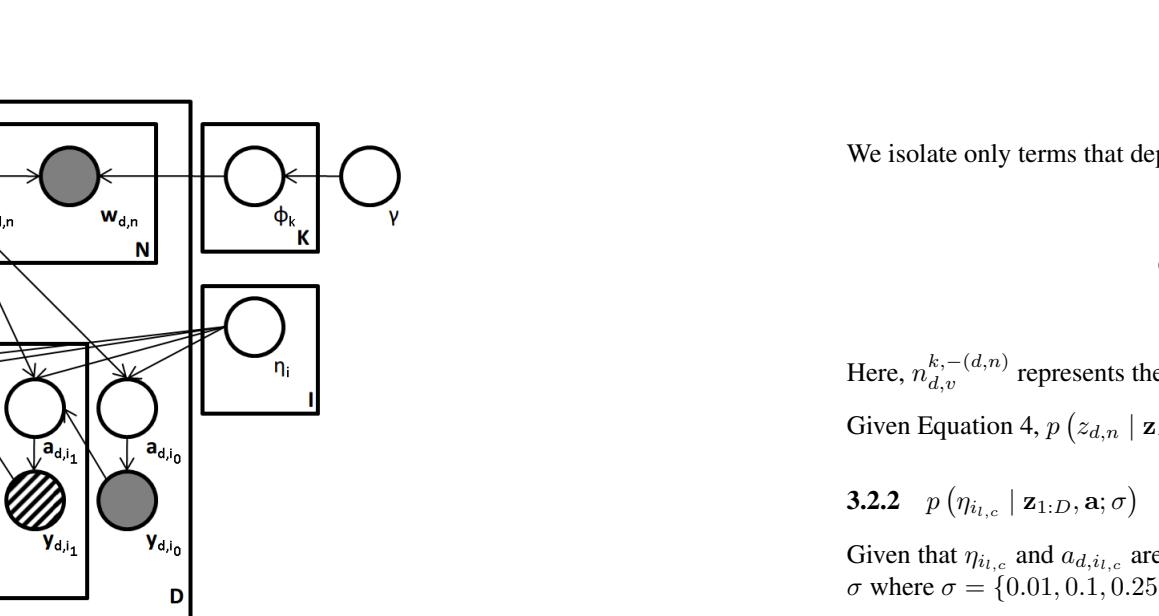


Figure 1: adapted sLDA model

A. Draw a latent variable

$$a_{d,i,c} \sim \begin{cases} N(z^T \eta_{i,c} + 1), & y_{parent,c} = 1 \\ truncN(z^T \eta_{i,c} + \varepsilon), & y_{parent,c} = -1 \end{cases}$$

where $\varepsilon = N^{-1} \sum_{n=1}^N z_{n,c}$ and $I = \{i_1, i_2, \dots, i_{|I|}\}$ and $i = \{i_1, i_2, \dots, i_{|I|,c}\}$

B. Draw a response variable

$$y_{d,i,c} | a_{d,i,c} \sim \begin{cases} 1, & a_{d,i,c} > 0 \\ -1, & otherwise \end{cases}$$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation 1.

$$\begin{aligned} p(\theta, z_{1:N}, \phi_{1:N}, \eta_{i,c}, a_{d,i,c}, \varepsilon, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}, \varepsilon; \sigma, \lambda) \\ = \frac{p(\theta, z_{1:N}, \phi_{1:N}, \eta_{i,c}, a_{d,i,c}, \varepsilon, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{1:N}, \varepsilon; \sigma, \lambda)}{\sum_{\theta, \phi, \eta, a, \varepsilon, \beta, \alpha, \alpha', \gamma} p(\theta, z_{1:N}, \phi_{1:N}, \eta_{i,c}, a_{d,i,c}, \varepsilon, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{1:N}, \varepsilon; \sigma, \lambda)} \end{aligned} \quad (1)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

3.2.1 $p(z_{m,n} | \mathbf{z}_{-(m,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z_m for a word instance, y_m , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

$$p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$(3) \quad p(\theta, z_{1:N}, \phi_{1:N}, \eta_{i,c}, a_{d,i,c}, \varepsilon, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}, \varepsilon; \sigma, \lambda) = \frac{p(\theta | \alpha) (\prod_{n=1}^N p(a_{n,c} | \theta) p(y_{n,c} | \eta_{n,c}, \beta)) (\prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) p(\phi_{k,n} | \phi_{k,n+1}, \beta, \alpha_k))}{\int_0^\infty p(\theta) \sum_{k=1}^K (\prod_{n=1}^N p(z_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta) d\theta} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) d\theta d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(4) \quad p(\mathbf{Y}, \mathbf{w}, \mathbf{z}, \mathbf{a}, \phi, \beta, \alpha, \varepsilon | \mathbf{x}, \mathbf{y}, \mathbf{z}_{-(d,n)}, \mathbf{a}_{-(d,n)}, \mathbf{w}_{-(d,n)}, \eta, \alpha, \beta, \gamma) = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(5) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(6) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(7) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(8) \quad = p(y_{i,c} | \mathbf{a}, \mathbf{y}_{-(i,c)}) = \delta(sign(a_{d,i,c}) = y_{i,c}) p(y_{i,c} | y_{parent,c}) \prod_{i,j \in children,c} p(y_{i,j} | y_{i,c})$$

$$(9) \quad p(y_{i,c} = -1 | y_{parent,c}) = \begin{cases} 1, & y_{parent,c} = -1 \\ 0.5, & y_{parent,c} = 1 \end{cases}$$

$$(10) \quad p(y_{i,c} = 1 | y_{parent,c}) = \begin{cases} 1, & y_{parent,c} = 1, \forall y_{i,k} \in y_{children,c}, y_{i,k} = -1 \\ 0.5, & y_{parent,c} = 1, \exists y_{i,k} \in y_{children,c}, y_{i,k} = 1 \\ 0, & otherwise \end{cases}$$

$$(11) \quad p(\mathbf{Y}, \mathbf{w}, \mathbf{z}, \mathbf{a}, \phi, \beta, \alpha, \varepsilon | \mathbf{x}, \mathbf{y}, \mathbf{z}_{-(d,n)}, \mathbf{a}_{-(d,n)}, \mathbf{w}_{-(d,n)}, \eta, \alpha, \beta, \gamma) = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(12) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(13) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(14) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(15) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(16) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(17) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(18) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(19) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(20) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(21) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(22) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(23) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(24) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod_{n=1}^N p(a_{n,c} | \theta, \eta_{n,c}, \beta, \alpha_k) \delta \theta d\theta \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) d\phi_{k,n} \prod_{n=1}^N p(y_{n,c} | \eta_{n,c}, \beta, \alpha_k) dy_{n,c}$$

$$(25) \quad = \int_{\phi_{1:K}} \prod_{k=1}^K p(\phi_{k,n} | \phi_{k,n-1}, \beta, \alpha_k) \prod$$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

The benefits of supervision in topic modeling

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured clinical data that has been at least in part manually categorized. Examples include patient records that are not linked to diagnoses and structured hierarchical directories of the same [2]. Product descriptions and catalogs (e.g. [2, 1]) as available from [2] and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records and the International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned to them[1]). In this work we show how to combine these two sources of information using a single model that allows one, among other things, to automatically categorize new text documents, suggest labels that might be inaccurate, and compute improved similarities between documents for information retrieval purposes. The models and techniques that we develop in this paper are applicable in other domains as well, for instance, unstructured representations of data that have been hierarchically categorized.

In this work we extend supervised latent Dirichlet allocation (sLDA) to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) augmented with per document “supervision” often taking the form of a single numerical or categorical “label.” More generally this “supervision” can be seen as extra data generated about a document; for instance its quality or relevance (e.g. online reviews), marks given to written work (e.g. graded essays), or the number of times a web document is linked. These labels are usually modeled as having been generated conditioned on the topic of topics found in a document. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. As one concrete example consider web retailers. They often have both a browsable hierarchy and free-text descriptions of all products they sell. The situation of each product in the product hierarchy (often multiply classified) can be seen as a form of multiple labeling and as a similar products based on the similarity of their textual description is an *u*. An equivalent challenge, particularly for larger retailers, is to situate the merchandise in as many categories as possible.

In this paper we describe the use of a topic model based on supervised latent Dirichlet allocation (sLDA) to identify topics within narrative discharge summaries and to automate the assignment of diagnostic codes, specifically International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes. There are a number of advantages to this approach. First, manually coding diagnoses is a time-consuming and notoriously unreliable process. Many diagnoses are omitted in the final record, and a high error rate is found even in the principal diagnoses [14].

• Benefits of combining human categorization information into “topic models”
 • LDA solved free text
 • supervised LDA improves LDA (extra info) and allows new inference (predict links, etc.)
 • amazon, freshdirect, netflix, dmoz, pandora (music genome)

2 Background

3 Methods

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

- For each topic, k :
 (a) Draw a distribution over words $\phi_k \sim Dir(\alpha, \alpha')$

- For each ICD9 code, c , at all levels in the tree, t :
 (a) Draw regression coefficient $\eta_{t,c} | \alpha' \sim N_K(-1, \sigma)$
 (b) Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$

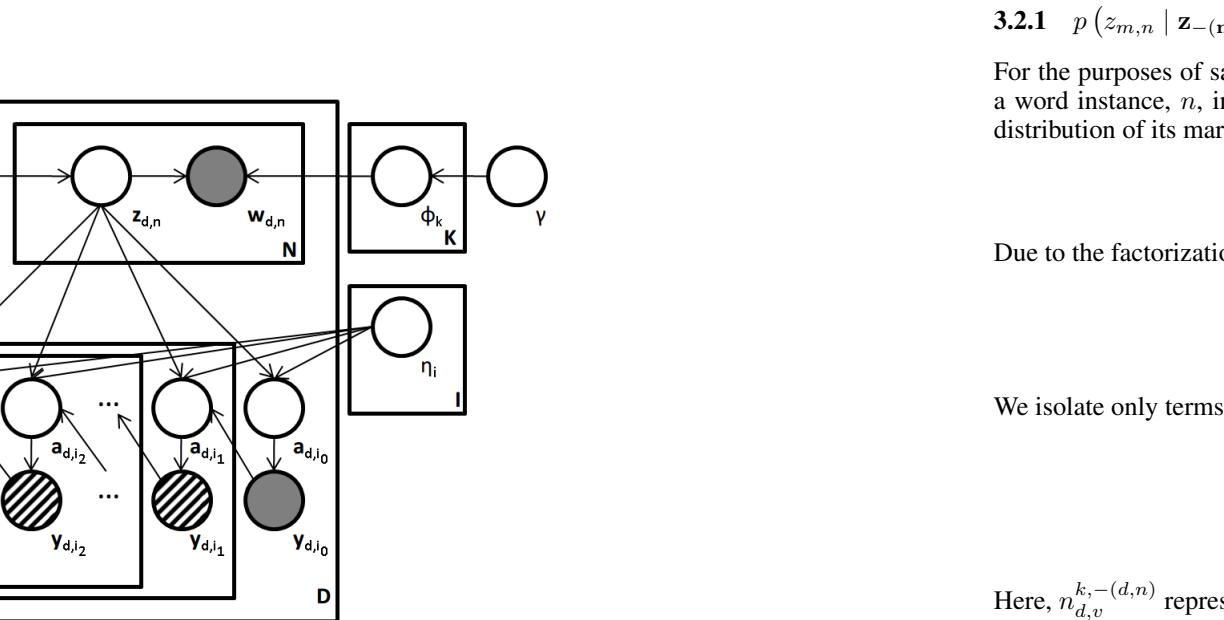


Figure 1: adapted sLDA model

4. For each document, d :

- Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$
 - For each word, n :
 - Draw topic assignment $z_{d,n} | \theta_d \sim Mult_K(\theta_d)$
 - Draw word $w_{d,n} | z_{d,n}, \beta_{1,K} \sim Mult(\beta_{z_n})$
 - For each level of the ICD-9 code tree, t :
 - For each ICD-9 code at this level, c :
 - Draw a latent variable $a_{d,t,c} \sim \begin{cases} \mathcal{N}(\bar{z}\eta_{t,c}, 1), & y_{parent,c} = 1 \\ trunc\mathcal{N}(\bar{z}\eta_{t,c}, 1), & y_{parent,c} = -1 \end{cases}$
- where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_C\}$ and $i_t = \{i_{t,0}, i_{t,1}, \dots, i_{t,C_t}\}$
- B. Draw a response variable $y_{d,i_{t,c}} | a_{d,i_{t,c}} \sim \begin{cases} 1, & a_{d,i_{t,c}} > 0 \\ -1, & otherwise \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_{t,c}}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{i_{t,c}}; \sigma, \lambda) \quad (1)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_{t,c}}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{i_{t,c}}; \sigma, \lambda)}{\int_{\theta, z_{1:N}, \phi_{1:K}, \eta_{i_{t,c}}, \alpha, \beta, \alpha', \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_{t,c}}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{i_{t,c}}; \sigma, \lambda)}$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular we will derive a collapsed Gibbs sampler.

3.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

3.2.1 $p(z_{m,n} | \mathbf{z}_{-(m,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

(2)

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant:

$$p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_{t,c} \in \mathcal{Z}} p(a_{d,i_{t,c}} | \mathbf{z}, \eta_{i_{t,c}}) p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [9].

$$\propto \prod_{i_{t,c} \in \mathcal{Z}} \exp \left\{ -\frac{(\bar{z}^T \eta_{i_{t,c}} - a_{d,i_{t,c}})^2}{2} \right\} \left(n_{d,i_{t,c}}^{k_{i_{t,c}}(d,n)} + \alpha \right) \frac{n_{i_{t,c},w_{d,n}}^{k_{i_{t,c}}(d,n)} + \gamma}{\sum_{i_{t,c} \in \mathcal{Z}} (n_{i_{t,c},w_{d,n}}^{k_{i_{t,c}}(d,n)} + \gamma)} \quad (4)$$

Here, $n_{d,i_{t,c}}^{k_{i_{t,c}}(d,n)}$ represents the count of word v in document d assigned to topic k omitting the $(d,n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

$$3.2.2 p(\eta_{i_{t,c}} | \mathbf{z}_{1:D}, \mathbf{a}; \sigma) \quad (5)$$

Given that $\eta_{i_{t,c}}$ and $a_{d,i_{t,c}}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$p(\eta_{i_{t,c}} | \mathbf{z}_{1:D}, \mathbf{a}; \sigma) = \mathcal{N}(\eta_{i_{t,c}} | \bar{\mu}_i, \hat{\Sigma}_i) \quad (5)$$

$$\bar{\mu}_i = \hat{\Sigma}_i (-1\sigma^{-1} + \bar{z}^T \mathbf{a}_{(i),i_{t,c}}) \quad (6)$$

$$\hat{\Sigma}_i^{-1} = \mathbf{I}\sigma^{-1} + \bar{z}^T \bar{z} \quad (7)$$

$$3.2.3 p(y_{d,i_{t,c}} | \mathbf{z}, \mathbf{Y}, \eta) \text{ and } p(y_{m,i} | \mathbf{a}) \quad (8)$$

In the augmented probit regression model, the posterior distribution of $a_{d,i_{t,c}}$ is distributed according to a truncated normal distribution where the response variable is observed.

$$p(a_{d,i_{t,c}} | \mathbf{z}, \mathbf{Y}, \eta) = \begin{cases} \text{trunc}\mathcal{N}^+(a_{d,i_{t,c}} | \eta_{i_{t,c}}^T \bar{z}, 1), & if \quad y_{d,i_{t,c}} = 1 \\ \text{trunc}\mathcal{N}^-(a_{d,i_{t,c}} | \eta_{i_{t,c}}^T \bar{z}, 1), & if \quad y_{d,i_{t,c}} = -1 \end{cases} \quad (9)$$

However, if $y_{d,i_{t,c}}$ is unobserved then $a_{d,i_{t,c}}$ must be sampled jointly with $y_{d,i_{t,c}}$ to ensure that the Markov chain is ergodic. Suppose that $a_{d,i_{t,c}}$ is sampled to have a negative value and $y_{d,i_{t,c}}$ is appropriately sampled at -1. Although there exist valid states where $a_{d,i_{t,c}} > 0$ and $y_{d,i_{t,c}} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_{t,c}} < 0 | y_{d,i_{t,c}} = -1) = 1$ and $p(y_{d,i_{t,c}} = -1 | a_{d,i_{t,c}} < 0) = 1$. Therefore, to ensure ergodicity, $a_{d,i_{t,c}}$ and $y_{d,i_{t,c}}$ must be sampled from the joint distribution as shown in Equation ??.

$$p(a_{d,i_{t,c}}, y_{d,i_{t,c}} | \mathbf{z}, \mathbf{Y}, \eta) \propto p(a_{d,i_{t,c}} | \mathbf{z}, \mathbf{Y}, \eta) p(y_{d,i_{t,c}} | a_{d,i_{t,c}}, \eta) \quad (10)$$

$$= \begin{cases} \mathcal{N}(a_{d,i_{t,c}} | \bar{z}^T \eta_{i_{t,c}}, 1) p(y_{d,i_{t,c}} | a_{d,i_{t,c}}, \eta_{i_{t,c}}), & y_{parent,i_{t,c}} = 1, \forall y_{i_{t,c}} \in y_{children,i_{t,c}}, y_{i_{t,c}} = -1 \\ \text{trunc}\mathcal{N}^-(a_{d,i_{t,c}} | \bar{z}^T \eta_{i_{t,c}}, 1) \delta(y_{d,i_{t,c}} = -1), & y_{parent,i_{t,c}} = -1 \\ \text{trunc}\mathcal{N}^+(a_{d,i_{t,c}} | \bar{z}^T \eta_{i_{t,c}}, 1) \delta(y_{d,i_{t,c}} = 1), & \forall y_{i_{t,c}} \in y_{children,i_{t,c}} \setminus y_{i_{t,c}}, y_{i_{t,c}} = 1 \\ 0, & otherwise \end{cases} \quad (11)$$

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was assigned to a patient but no malignant hypertension, it could be assumed that they did not have malignant hypertension. Unfortunately, ICD-9 codes observations never include observations of disease absence. ICD-9 codes are only documented when the condition is observed to be present. Additionally, ICD-9 codes are known to have relatively low sensitivity: conditions that are present are often not documented in a set of ICD-9 codes (Surjan, 1999). Given these facts, we made the following assumptions regarding each visit: recorded diagnoses and their ancestors were labeled as true; diagnoses that were observed at some time for a patient but not at the current visit were labeled as unobserved; and diagnoses that had never been listed for a patient were labeled as false for all of that patient's visits. This last assumption captures the belief that parts of the ICD-9 hierarchy that are never observed for a particular patient are almost certain to be false. Additionally, for computational purposes, we decided not to include an ICD-9 code at all if neither it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

4.3 ICD-9 Code Hierarchy

4.3.1 Rao-Blackwellization

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an is-a hierarchy [1]. An is-a hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors

$$p(y_{i_{t,c}} | \mathbf{a}, \mathbf{y}_{-(t,c)}) = \delta(sign(a_{d,i_{t,c}}) = y_{i_{t,c}}) p(y_{i_{t,c}} | y_{parent,i_{t,c}}) \prod_{i_{t,c} \in \mathcal{Z}} p(y_{i_{t,c}} | y_{parent,i_{t,c}}) \quad (12)$$

$$p(y_{i_{t,c}} = -1 | y_{parent,i_{t,c}}) = \begin{cases} 1, & y_{parent,i_{t,c}} = -1 \\ 0.5, & y_{parent,i_{t,c}} = 1 \end{cases} \quad (13)$$

$$p(a_{d,i_{t,c}}, y_{i_{t,c}} | \mathbf{z}, \mathbf{Y}, \eta) \quad (14)$$

$$= \begin{cases} \mathcal{N}(a_{d,i_{t,c}} | \bar{z}^T \eta_{i_{t,c}}, 1) p(y_{i_{t,c}} | a_{d,i_{t,c}}, \eta_{i_{t,c}}), & y_{parent,i_{t,c}} = 1, \forall y_{i_{t,c}} \in y_{children,i_{t,c}}, y_{i_{t,c}} = -1 \\ \text{trunc}\mathcal{N}^-(a_{d,i_{t,c}} | \bar{z}^T \eta_{i_{t,c}}, 1) \delta(y_{i_{t,c}} = -1), & y_{parent,i_{t,c}} = -1 \\ \text{trunc}\mathcal{N}^+(a_{d,i_{t,c}} | \bar{z}^T \eta_{i_{t,c}}, 1) \delta(y_{i_{t,c}} = 1), & \forall y_{i_{t,c}} \in y_{children,i_{t,c}} \setminus y_{i_{t,c}}, y_{i_{t,c}} = 1 \\ 0, & otherwise \end{cases} \quad (15)$$

$$p(\mathbf{Y}, \mathbf{w}, \mathbf{z}, \mathbf{y}, \mathbf{a}, \mathbf{phi}, \mathbf{beta}, \mathbf{xi}) = \prod_{i=1}^I \left[\prod_{j=1}^J p(y_{i,j} | \mathbf{z}_{i,j}, \mathbf{w}_{i,j}, \mathbf{phi}_{i,j}, \mathbf{beta}_{i,j}, \mathbf{xi}_{i,j}) \right] \prod_{i=1}^M \left[\prod_{j=1}^N p(w_{i,j} | \mathbf{z}_{i,j}, \mathbf{y}_{i,j}, \mathbf{phi}_{i,j}, \mathbf{beta}_{i,j}, \mathbf{xi}_{i,j}) \right] \prod_{i=1}^K \left[\prod_{j=1}^L p(z_{i,j} | \mathbf{y}_{i,j}, \mathbf{phi}_{i,j}, \mathbf{beta}_{i,j}, \mathbf{xi}_{i,j}) \right] \prod_{i=1}^L \left[\prod_{j=1}^M p(y_{i,j} | \mathbf{z}_{i,j}, \mathbf{w}_{i,j}, \mathbf{phi}_{i,j}, \mathbf{beta}_{i,j}, \mathbf{xi}_{i,j}) \right] \quad (16)$$

$$= \prod_{i=1}^I \left[\prod_{j=1}^J \left[\prod_{k=1}^K \Gamma(n_{i,j,k} + \beta_{i,j,k}) \Gamma(\sum_{k=1}^K n_{i,j,k} + \alpha_{i,j,k}) \right] \prod_{k=1}^K \Gamma(n_{i,j,k} + \beta_{i,j,k}) \Gamma(\sum_{k=1}^K n_{i,j,k} + \alpha_{i,j,k}) \right] \prod_{i=1}^M \left[\prod_{j=1}^N \left[\prod_{k=1}^L \Gamma(n_{i,j,k} + \beta_{i,j,k}) \Gamma(\sum_{k=1}^L n_{i,j,k} + \alpha_{i,j,k}) \right] \prod_{k=1}^L \Gamma(n_{i,j,k} + \beta_{i,j,k}) \Gamma(\sum_{k=1}^L n_{i,j,k} + \alpha_{i,j,k}) \right] \prod_{i=1}^K \left[\prod_{j=1}^L \left[\prod_{k=1}^M \Gamma(n_{i,j,k} + \beta_{i,j,k}) \Gamma(\sum_{k=1}^M n_{i,j,k} + \alpha_{i,j,k}) \right] \prod_{k=1}^M \Gamma(n_{i,j,k} + \beta_{i,j,k}) \Gamma(\sum_{k=1}^M n_{i,j,k} + \alpha_{i,j,k}) \right] \prod_{i=1}^L \left[\prod_{j=1}^M \left[\prod_{k=1}^K \Gamma(n_{i,j,k} + \beta_{i,j,k}) \Gamma(\sum_{k=1}^K n_{i,j,k} + \alpha_{i,j,k}) \right] \prod_{k=1}^K \Gamma(n_{i,j,k} + \beta_{i,j,k$$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA) in the same manner as was done in supervised LDA (sLDA) prior art. We find that the additional supervision signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in medical document labeling and product categorization tasks. Additionally, held-out likelihood

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [?], product descriptions and catalogs (e.g. [?] as available from [?]), and patient hospital treatment records and codes applied to a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA) in the same manner as was done in supervised LDA (sLDA) prior art. We find that the additional supervision signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in medical document labeling and product categorization tasks. Additionally, held-out likelihood

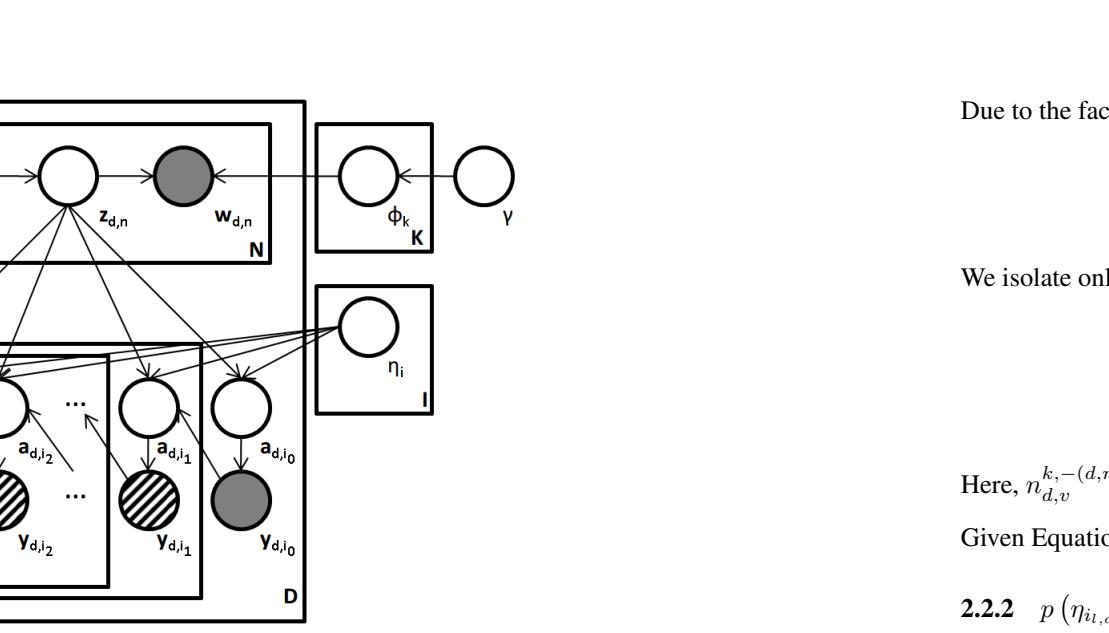


Figure 1: adapted sLDA model

In this work we extend supervised latent Dirichlet allocation (sLDA) [2] to take advantage of hierarchical supervision. LDA is latent Dirichlet allocation (LDA) [3] augmented with per-document “supervision” often taking the form of a single numerical or categorical “label”. More generally this “supervision” can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been drawn from an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [3].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiple situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical and web retail data suggests that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section ?? we introduce hierarchically supervised LDA (HSLDA), in Section ?? we review related work, and in Section ?? we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Methods

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, k:
 - (a) Draw a distribution over words $\phi_k \sim Dir(\gamma, 1 - \gamma)$
2. For each ICD9 code, c, at all levels in the tree, l:
 - (a) Draw regression coefficient $\eta_{i_l,c} \sim \mathcal{N}_K(-1, \sigma)$
3. Draw a prior over topic proportions $\beta \mid \alpha' \sim Dir_K(1, \alpha')$
4. For each document, d:
 - (a) Draw topic proportions $\theta_d \mid \beta \sim Dir_K(\beta, \alpha)$
 - (b) For each word, n:
 1. Draw topic assignment $z_{n,d} \mid \theta_d \sim Mult(\theta_d)$

1

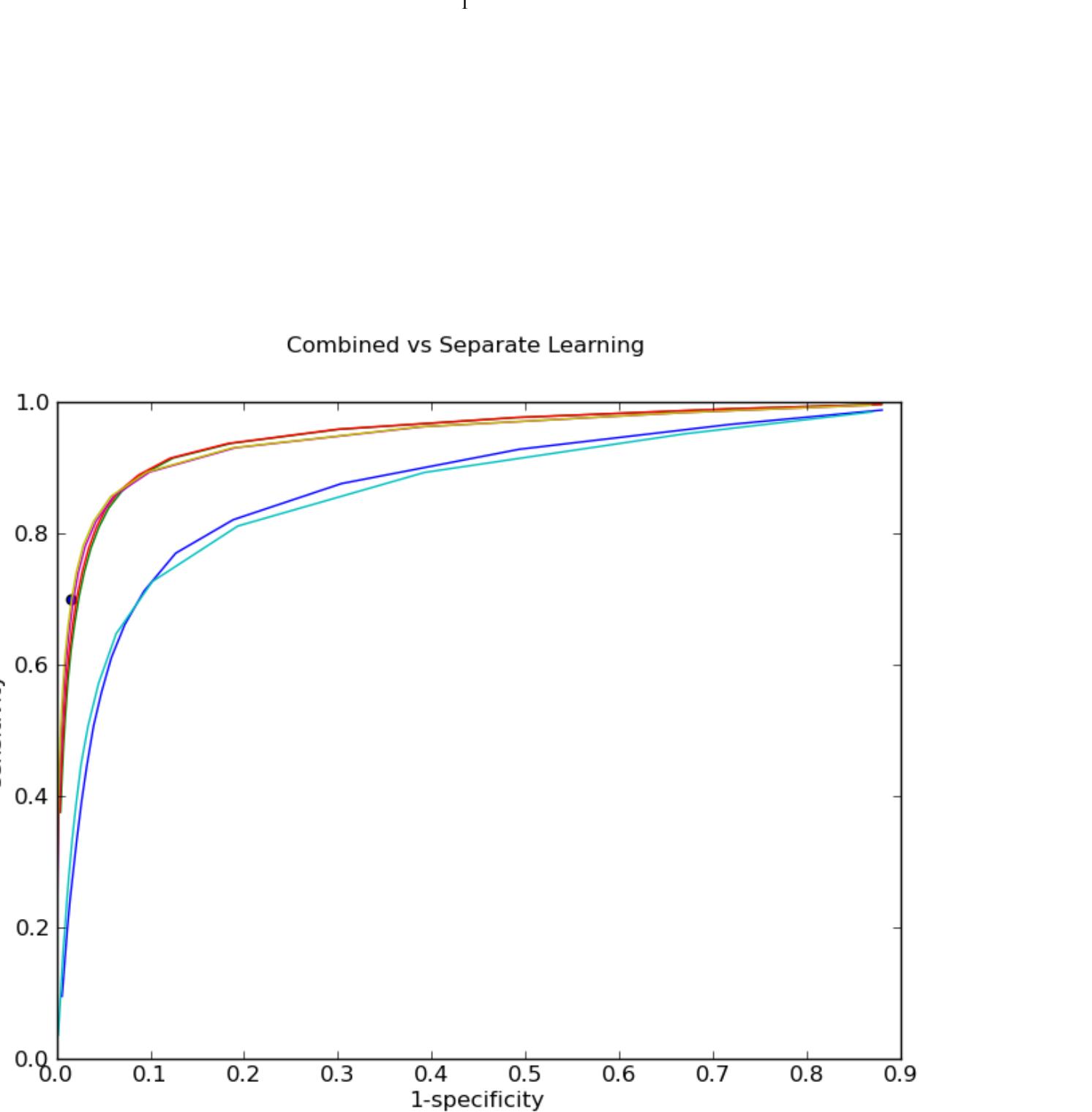


Figure 2: default

- i. Draw word $w_{n,d} \mid z_{n,d}, \beta_{1:K} \sim Mult(\beta_{z_{n,d}})$
- ii. For each level of the ICD-9 code tree, l:
 - A. For each ICD-9 code at that level, c:
 - A. Draw a latent variable $a_{d,i_l,c} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{i_l,c}, 1), & y_{parent_{i_l,c}} = 1 \\ trunc\mathcal{N}^+(\bar{z}^T \eta_{i_l,c}, 1), & y_{parent_{i_l,c}} = -1 \end{cases}$
 - where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_C\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$
 - B. Draw a response variable $y_{d,i_l,c} \mid a_{d,i_l,c} \sim \begin{cases} 1, & a_{d,i_l,c} > 0 \\ -1, & otherwise \end{cases}$

2.2.1 $p(z_{m,n} \mid \mathbf{z}_{-(m,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

$$p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_l \in \mathcal{I}} p(a_{d,i_l,c}) p(z_{d,n}, \eta_{i_l,c}, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [9].

$$\propto \prod_{i_l \in \mathcal{I}} \exp \left\{ -\frac{(\bar{z}^T \eta_{i_l,c} - a_{d,i_l,c})^2}{2} \right\} \binom{n_{d,(c)}}{n_{d,(c)} + \alpha_{\beta_k}} \frac{\Gamma(n_{d,(c)} + \gamma)}{\sum_{v=1}^V \binom{n_{d,(v)} + \alpha_{\beta_k}}{n_{d,(v)} + \gamma}} \quad (4)$$

Here, $n_{d,v}^{k,-(d,n)}$ represents the count of word v in document d assigned to topic k omitting the $(d,n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

2.2.2 $p(\eta_{i_l,c} \mid \mathbf{z}_{1:D}, \mathbf{a}, \sigma)$

Given that $\eta_{i_l,c}$ and $a_{d,i_l,c}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$p(\eta_{i_l,c} \mid \mathbf{z}_{1:D}, \mathbf{a}, \sigma) = \mathcal{N}(\eta_{i_l,c} \mid \hat{\mu}_i, \hat{\Sigma}_i) \quad (5)$$

$$\hat{\mu}_i = \hat{\Sigma}_i (-1\sigma^{-1} + \bar{Z}^T \eta_{i_l,c}) \quad (6)$$

$$\hat{\Sigma}_i^{-1} = \mathbf{I}\sigma^{-1} + \bar{Z}^T \bar{Z} \quad (7)$$

2.2.3 $p(a_{d,i_l,c} \mid \mathbf{z}, \mathbf{Y}, \eta)$ and $p(y_{m,i} \mid \mathbf{a})$

In the augmented probit regression model, the posterior distribution of $a_{d,i_l,c}$ is distributed according to a truncated normal distribution where the response variable is observed.

$$p(a_{d,i_l,c} \mid \mathbf{z}, \mathbf{Y}, \eta) = \begin{cases} \text{trunc}\mathcal{N}^+(a_{d,i_l,c} \mid \eta_i^T \bar{z}, 1, y_{d,i_l,c}) & if \quad y_{d,i_l,c} = 1 \\ \text{trunc}\mathcal{N}^-(a_{d,i_l,c} \mid \eta_i^T \bar{z}, 1, y_{d,i_l,c}) & if \quad y_{d,i_l,c} = -1 \end{cases} \quad (8)$$

$$p(y_{d,i_l,c} \mid \mathbf{z}, \mathbf{Y}_{-(d,i_l,c)}, \eta) \propto p(y_{d,i_l,c} \mid \mathbf{a}, \mathbf{y}_{-(d,i_l,c)}) p(a_{d,i_l,c} \mid \mathbf{z}, \mathbf{Y}, \eta) \quad (9)$$

$$p(y_{d,i_l,c} = 1 \mid y_{parent_{i_l,c}}) = \begin{cases} 1, & y_{parent_{i_l,c}} = 1 \\ 0.5, & y_{parent_{i_l,c}} = -1 \end{cases} \quad (10)$$

$$p(a_{d,i_l,c}, y_{d,i_l,c} \mid \mathbf{z}, \mathbf{Y}_{-(d,i_l,c)}) = \begin{cases} \mathcal{N}(a_{d,i_l,c} \mid \bar{z}^T \eta_{i_l,c}, 1) p(y_{d,i_l,c} \mid a_{d,i_l,c}), & y_{parent_{i_l,c}} = 1, \forall y_{i_l,c} \in children_{i_l,c}, y_{i_l,c} = -1 \\ \text{trunc}\mathcal{N}^-(a_{d,i_l,c} \mid \bar{z}^T \eta_{i_l,c}, 1) \delta(y_{d,i_l,c} = -1), & y_{parent_{i_l,c}} = -1 \\ \text{trunc}\mathcal{N}^+(a_{d,i_l,c} \mid \bar{z}^T \eta_{i_l,c}, 1) \delta(y_{d,i_l,c} = 1), & \exists y_{i_l,c} \in children_{i_l,c} \setminus y_{parent_{i_l,c}}, y_{i_l,c} = 1 \\ 0, & otherwise \end{cases} \quad (11)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

2.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_l,c}, \mathbf{a}_{i_l,c}, \mathbf{z}_d, \beta, \alpha, \alpha', \gamma \mid w_{1:N}, y_{1:N}, \eta, \sigma, \lambda) \quad (12)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_l,c}, \mathbf{a}_{i_l,c}, \mathbf{z}_d, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{1:N}, \eta, \sigma, \lambda)}{\int_{\theta, z_{1:N}, \phi_{1:K}, \eta_{i_l,c}, \mathbf{a}_{i_l,c}, \mathbf{z}_d, \beta, \alpha, \alpha', \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_l,c}, \mathbf{a}_{i_l,c}, \mathbf{z}_d, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{1:N}, \eta, \sigma, \lambda) \, d\theta \, dz_{1:N} \, d\phi_{1:K} \, d\eta_{i_l,c} \, da_{i_l,c} \, dz_d \, d\beta \, d\alpha \, d\alpha' \, d\gamma} \quad (13)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

2.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

2.2.1 $p(z_{m,n} \mid \mathbf{z}_{-(m,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

$$p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_l \in \mathcal{I}} p(a_{d,i_l,c}) p(z_{d,n}, \eta_{i_l,c}, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [9].

$$\propto \prod_{i_l \in \mathcal{I}} \exp \left\{ -\frac{(\bar{z}^T \eta_{i_l,c} - a_{d,i_l,c})^2}{2} \right\} \binom{n_{d,(c)}}{n_{d,(c)} + \alpha_{\beta_k}} \frac{\Gamma(n_{d,(c)} + \gamma)}{\sum_{v=1}^V \binom{n_{d,(v)} + \alpha_{\beta_k}}{n_{d,(v)} + \gamma}} \quad (4)$$

Here, $n_{d,v}^{k,-(d,n)}$ represents the count of word v in document d assigned to topic k omitting the $(d,n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

2.2.2 $p(\eta_{i_l,c} \mid \mathbf{z}_{1:D}, \mathbf{a}, \sigma)$

Given that $\eta_{i_l,c}$ and $a_{d,i_l,c}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$p(\eta_{i_l,c} \mid \mathbf{z}_{1:D}, \mathbf{a}, \sigma) = \mathcal{N}(\eta_{i_l,c} \mid \hat{\mu}_i, \hat{\Sigma}_i) \quad (5)$$

$$\hat{\mu}_i = \hat{\Sigma}_i (-1\sigma^{-1} + \bar{Z}^T \eta_{i_l,c}) \quad (6)$$

$$\hat{\Sigma}_i^{-1} = \mathbf{I}\sigma^{-1} + \bar{Z}^T \bar{Z} \quad (7)$$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents

$$p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_l \in \mathcal{I}} p(a_{d,i_l,c}) p(z_{d,n}, \eta_{i_l,c}, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [9].

$$\propto \prod_{i_l \in \mathcal{I}} \exp \left\{ -\frac{(\bar{z}^T \eta_{i_l,c} - a_{d,i_l,c})^2}{2} \right\} \binom{n_{d,(c)}}{n_{d,(c)} + \alpha_{\beta_k}} \frac{\Gamma(n_{d,(c)} + \gamma)}{\sum_{v=1}^V \binom{n_{d,(v)} + \alpha_{\beta_k}}{n_{d,(v)} + \gamma}} \quad (4)$$

Here, $n_{d,v}^{k,-(d,n)}$ represents the count of word v in document d assigned to topic k omitting the $(d,n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

2.2.3 $p(a_{d,i_l,c} \mid \mathbf{z}, \mathbf{Y}, \eta)$ and $p(y_{m,i} \mid \mathbf{a})$

Given that $\eta_{i_l,c}$ and <math

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA) in the same manner as was done in supervised LDA (sLDA) prior art. We find that the additional supervision signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in medical document labeling and product categorization tasks. Additionally, held-out likelihood

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [?], product descriptions and catalogs (e.g. [?] as available from [?]), and patient hospital treatment records and codes applied to them [?]. product descriptions and catalogs (e.g. hospital discharge records with ICD-9 codes assigned to them). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [2] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [3] augmented with per-document "supervision": often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been drawn from an inferred document topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiple situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesis that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 2 we introduce hierarchically supervised LDA (HSLDA), in Section 3.4 we review related work, and in Section ?? we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Methods

Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:

1. For each topic, k:
 - (a) Draw a distribution over words $\phi_k \sim Dir(\gamma, 1 - \gamma)$
2. For each ICD9 code, c, at all levels in the tree, l:
 - (a) Draw regression coefficient $\eta_{l,c} | \sigma \sim N_K(-1, \sigma)$
3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
4. For each document, d:
 - (a) Draw topic proportions $\theta_d | \beta \sim Dir_K(\beta, \alpha)$
 - (b) For each word, n:
 - (a) Draw topic assignment $z_{n,d} | \theta_d \sim Mult(\theta_d)$

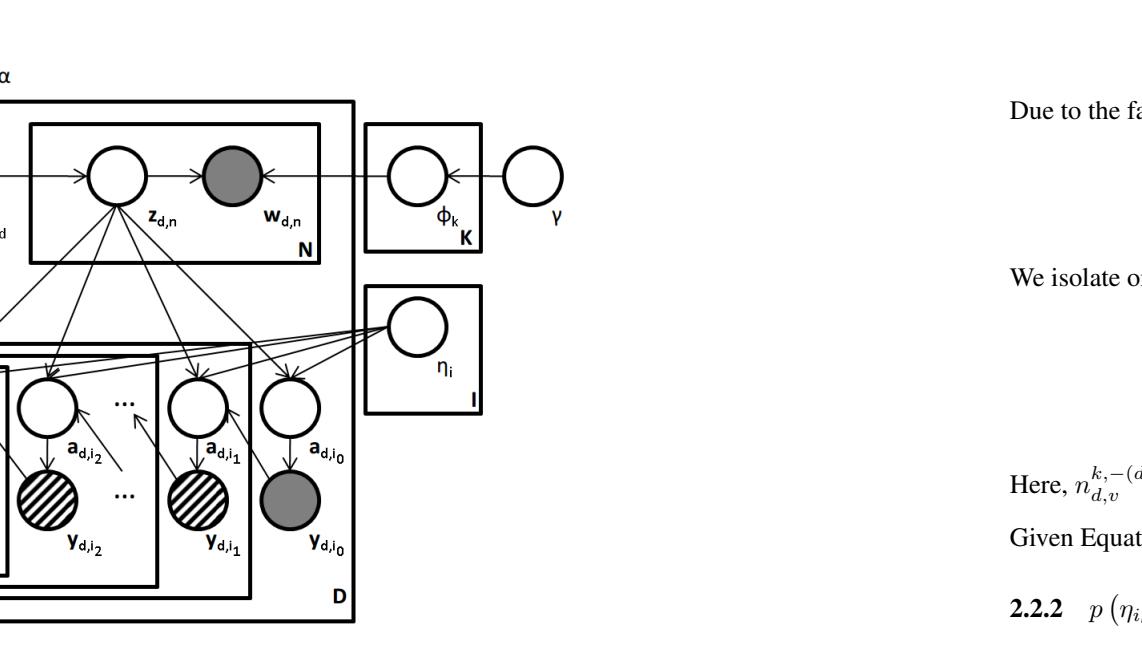


Figure 1: adapted sLDA model

$$p(z_{d,n} | z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) = p(z_{d,n}, z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) / \sum_{i_{l,c}} p(a_{i_{l,c}} | z_{l,c}) p(z_{d,n}, z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_{l,c}} p(a_{i_{l,c}} | z_{l,c}) p(z_{d,n}, z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [9].

$$\propto \prod_{i_{l,c}} \exp \left\{ -\frac{(\bar{z}^T \eta_{i_{l,c}} - a_{i_{l,c}})^2}{2} \right\} \binom{n_{d,(c)}^{k_{-(d,n)}} + \alpha_{\beta_k}}{\binom{n_{d,(c)}^{k_{-(d,n)}}}{n_{d,(c)}^{k_{-(d,n)}} + \alpha_{\beta_k}}} \frac{\Gamma(n_{d,(c)}^{k_{-(d,n)}} + \gamma)}{\sum_{v=1}^V \binom{n_{d,(c)}^{k_{-(d,n)}}}{n_{d,(c)}^{k_{-(d,n)}} + \gamma}} \quad (4)$$

Here, $n_{d,(c)}^{k_{-(d,n)}}$ represents the count of word v in document d assigned to topic k omitting the $(d, n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} | z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

$$p(\eta_{l,c} | z_{1:D}, \mathbf{a}, \sigma) \quad (2.2.4)$$

Given that $\eta_{l,c}$ and $a_{i_{l,c}}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$ii. \text{ Draw word } w_{n,d} | z_{n,d}, \beta_{1:K} \sim Mult(\beta_{z_n})$$

iii. For each level of the ICD-9 code tree, l:

A. For each ICD-9 code at the level, l:

A. Draw a latent variable

$$a_{d,i_{l,c}} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{i_{l,c}}, 1), & y_{parent_{l,c}} = 1 \\ \text{truncN}^+(\bar{z}^T \eta_{i_{l,c}}, 1), & y_{parent_{l,c}} = -1 \end{cases}$$

where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_C\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$

B. Draw a response variable $y_{d,i_{l,c}} | a_{d,i_{l,c}} \sim \begin{cases} 1, & a_{d,i_{l,c}} > 0 \\ -1, & \text{otherwise} \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents

that have been being drawn as an auxiliary variable.

2.1 Posterior Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,c}, \epsilon_{l,c}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{l,c}; \sigma, \lambda) \quad (1)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,c}, \epsilon_{l,c}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{l,c}; \sigma, \lambda)}{\int_{\theta, z_{1:N}, \phi_{1:K}, \eta_{l,c}, \epsilon_{l,c}, \alpha, \beta, \alpha', \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,c}, \epsilon_{l,c}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{l,c}; \sigma, \lambda) \quad (7)}$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

2.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [9]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

$$p(z_{m,n} | z_{-(m,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2.2.1)$$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

2.3 Prediction

3 Experiments

3.1 Data

Our data set was gathered from the clinical data warehouse of New York - Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The notes outline the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a controlled vocabulary which is the international standard diagnostic classification for epidemiological, statistical and administrative purposes (the WHO International Classification of Diseases, Ninth Revision, Clinical Modification (ICD-9-CM) codes). The codes are classified in a rooted tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. In each patient, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary. For the purposes of sLDA, ICD-9 codes will be used as labels for discharge summaries.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. Before beginning data processing, we generated a PHI-free dataset (see Data Pre-processing below).

3.2 Pre-Processing

Patient discharge summaries and their associated ICD-9 diagnoses are stored in two different places in the New York - Presbyterian data warehouse, and so had to be linked together before being fed into the sLDA algorithm. Each discharge note and set of diagnoses were assigned a patient unique identifier (PUIID), allowing the two types of data to be linked.

Natural Language Processing (NLP) techniques were used to process the free-text discharge summaries. First, the Natural Language Toolkit (<http://www.nltk.org/>) was used to tokenize the text. Next, feature selection was performed using a term frequency - inverse document frequency (TF-IDF) algorithm on the entire document set and sorting the words by their TF-IDF values. The top 10,000 words were manually evaluated to eliminate all potentially identifying information. Finally, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was observed to be absent, all of its descendants could be assumed to be absent as well (e.g., if a patient did not have malignant hypertension, it could be assumed that they did not have essential hypertension). Unfortunately, ICD-9 code observations never include observations of disease absence. ICD-9 codes are only documented when the condition is observed to be present. Additionally, ICD-9 codes are known to have relatively low sensitivity; conditions that are present are often not documented in a set of ICD-9 codes (Surjan, 1999). Given these facts, we made the following assumptions regarding each visit: recorded diagnoses and their ancestors were labeled as true; diagnoses that were observed at some time for a patient but not at the current visit were labeled as unobserved; and diagnoses that had never been listed for a patient were labeled as false for all of that patient's visits. This last assumption captures the belief that parts of the ICD-9 hierarchy that are never observed for a particular patient are almost certain to be false. Additionally, for computational purposes, we decided not to include an ICD-9 code at all if neither it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

3.3 ICD-9 Code Hierarchy

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an is-a hierarchy [1]. An is-a hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors

2.4 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [3].

This model, supervised latent dirichlet allocation (sLDA), builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [2].

Other models - predicting document links, other supervised latent variable models

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [13, 7, 12, 4], fully automatic assignment of ICD-9 codes to medical records, based on a few specific diseases [11] but the most recent years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [11] but the most recent years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [11] but the most recent years.

3.5 Evaluation

4 Results

5 Discussion

References

- [1] Amazon, Inc. <http://www.amazon.com/>, 2011.
- [2] DMoz open directory project, <http://www.dmoz.org/>, 2002.
- [3] International classification of disease, <http://bioprotocol.biomedontology.org/ontologies/35686>, May 2008.
- [4] Stanford network analysis platform, <http://snap.stanford.edu/>, 2004.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] P. Brown, D. Gochman, and J.R. Alber. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:903–1022, March 2003. ISSN 1532-4435.
- [7] P. Brown, D. Gochman, and J.R. Alber. The ngram classifier: A novel method of automatically creating text classifiers based on icd9 groupings. *Advances in Disease Surveillance*, 1:36, 2006.
- [8] K. Crammer, M. Dredze, K. Ganchev, P.P. Talukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007 Medical, Biomedical, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [9] R. Farkas and C. Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [10] H.R FreitasJunior, B.RibeiroNeto, R.F Vale, A.H.F.Laender, and L.Rs Lima. Categorization-driven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.
- [11] R.B Rao, S.Sandilya, R.S Niculescu, C.Germond, and H.Rao. Large scale diagnostic code classification for medical patient records. *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 416–425, 2003.
- [12] B.RibeiroNeto, A.H.F.Laender, and L.Rs Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for Information science and technology*, 52(2):391–401, 2001.
- [13] P. Bhattacharya, S. Bhattacharyya, and A. Ghosh. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [14] Y.W. Teh, M.I. Jordan, M.J. Beal, and D.M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [15] Hana Wallach, David Mimno, and Andrew McCallum. Rethinking id: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.

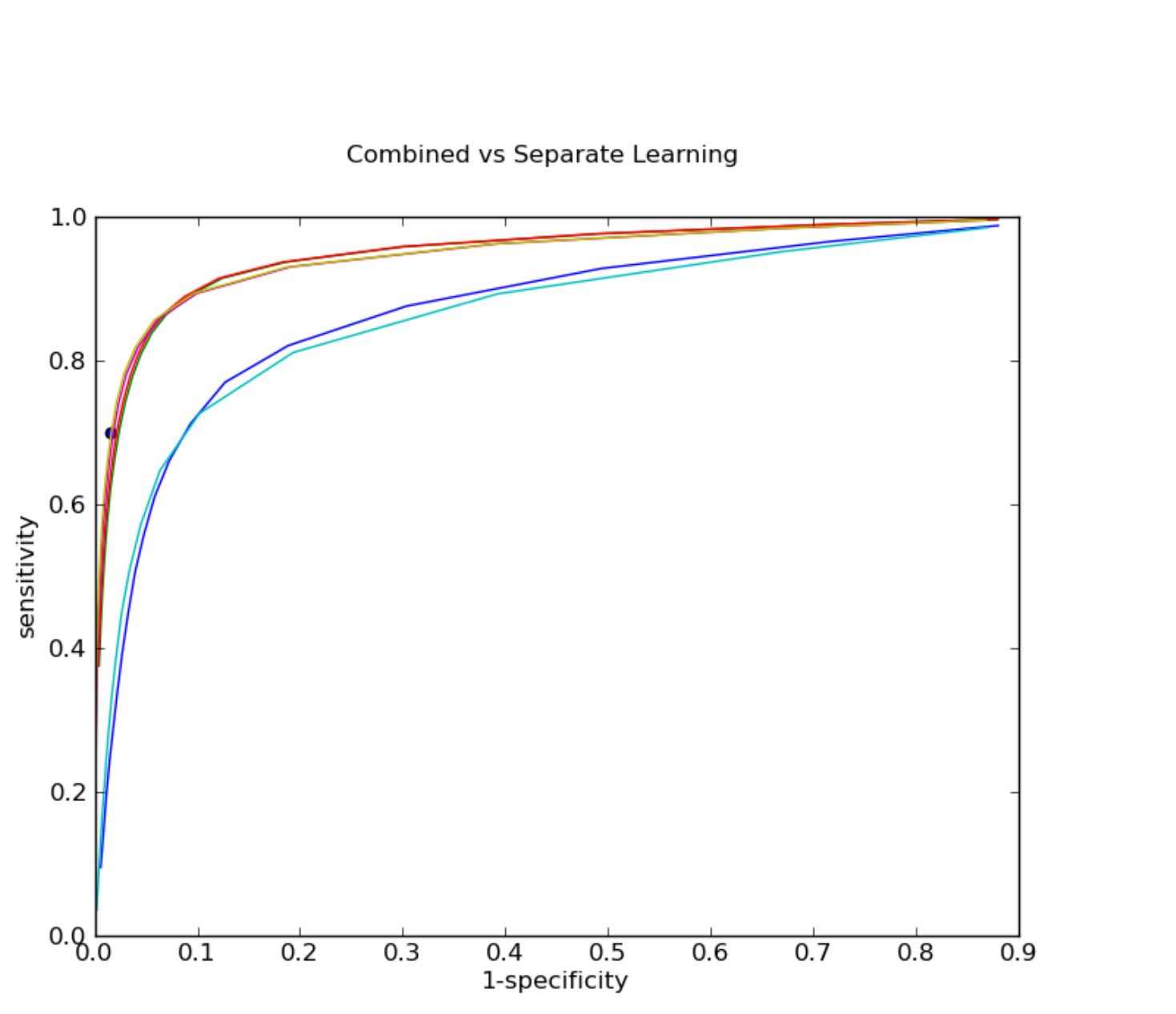


Figure 2: default

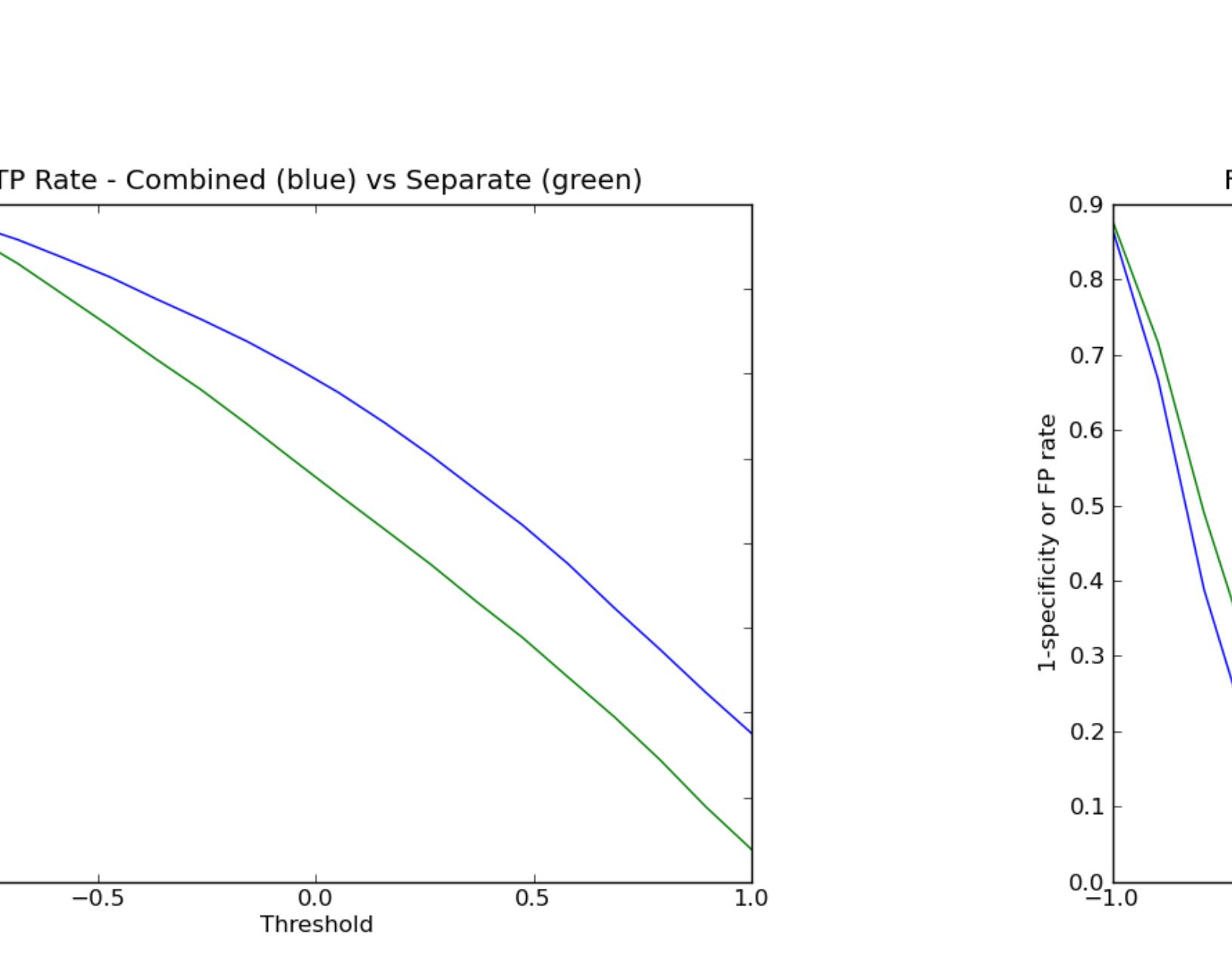
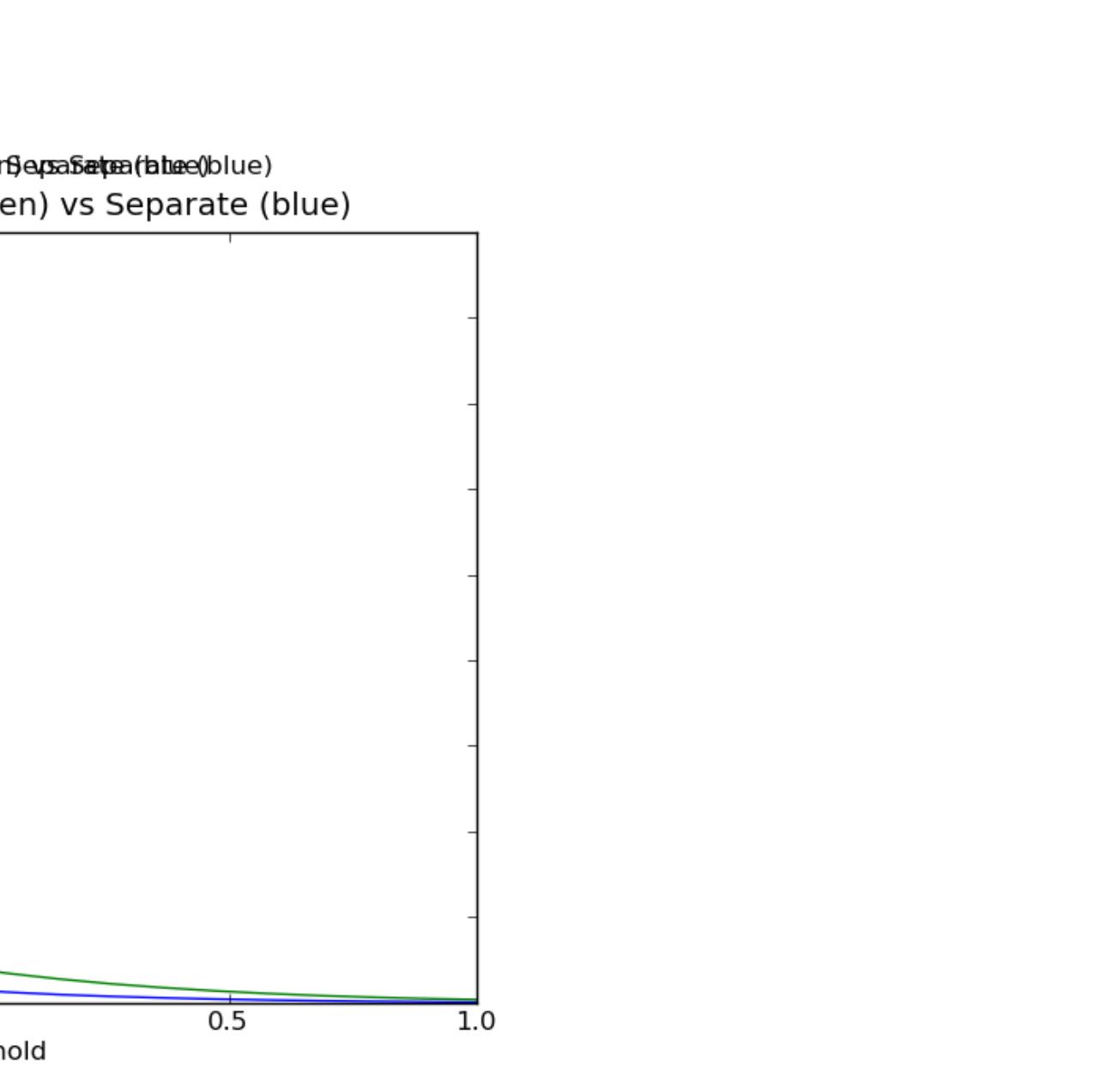


Figure 3: default



Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA) in the same manner as was done in supervised LDA (sLDA) prior art. We find that the additional supervision signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in medical document labeling and product categorization tasks. Additionally, held-out likelihood

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we work on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient records. Classification of these types of documents is often based on a combination of ICD-9-CM codes assigned to them. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image categories with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision": often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been drawn from an inferred document topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiple situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 2 we introduce hierarchically supervised LDA (HSLDA), in Section 3.4 we review related work, and in Section ?? we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Methods

- Given the number of topics, K, and broad gamma priors on hyperparameters, the generative process is as follows:
1. For each topic, k:
 - (a) Draw a distribution over words $\phi_k \sim Dir(\gamma, 1 - \gamma)$
 2. For each ICD-9 code, c, at all levels in the tree, l:
 - (a) Draw regression coefficient $\eta_{l,c} | \sigma \sim N_K(-1, \sigma)$
 3. Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
 4. For each document, d:
 - (a) Draw topic proportions $\theta_d | \beta \sim Dir_K(\beta, \alpha)$
 - (b) For each word, n:
 1. Draw topic assignment $z_{n,d} | \theta_d \sim Mult(\theta_d)$

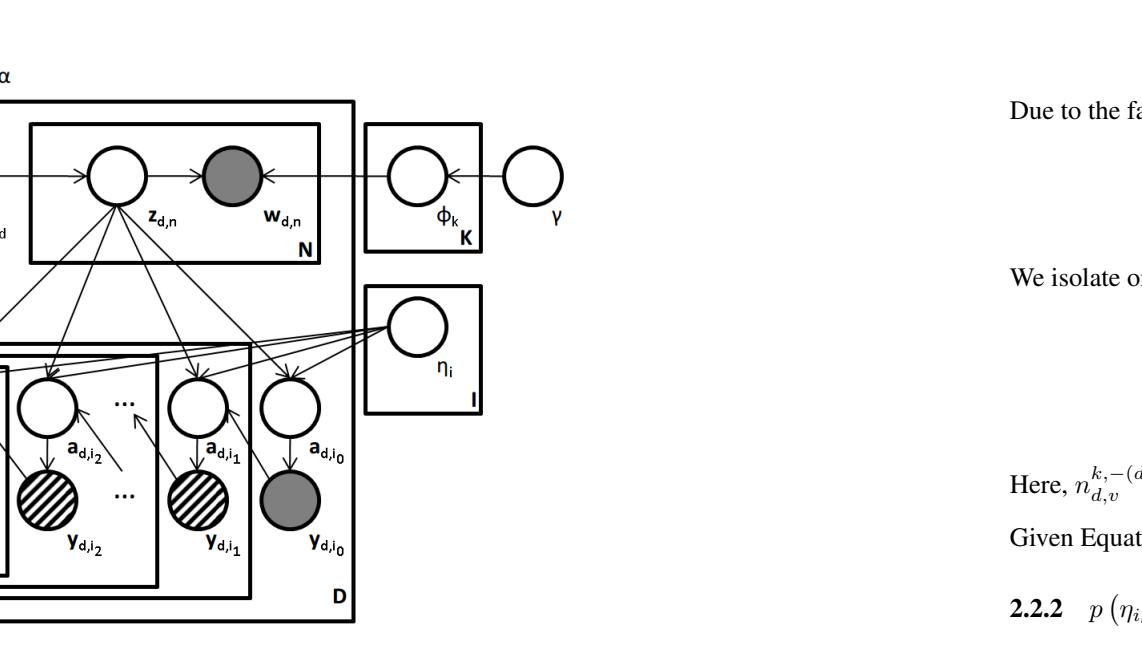


Figure 1: adapted sLDA model

ii. Draw word $w_{n,d} | z_{n,d}, \beta_{1:K} \sim Mult(\beta_{z_{n,d}})$

c) For each level of the ICD-9 code tree, l:

i. For each ICD-9 code at the level, c:

A. Draw a latent variable

$$a_{d,i_{l,c}} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{l,c}, 1), & y_{parent_{l,c}} = 1 \\ \text{truncN}^+(\bar{z}^T \eta_{l,c}, 1), & y_{parent_{l,c}} = -1 \end{cases}$$

where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_C\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$

B. Draw a response variable $y_{d,i_{l,c}} | a_{d,i_{l,c}} \sim \begin{cases} 1, & a_{d,i_{l,c}} > 0 \\ -1, & \text{otherwise} \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents

as have been done in an inferred document topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,c} | \mathbf{z}, \mathbf{y}, \alpha, \beta, \alpha', \gamma | w_{1:N}, y_{1:N}, \eta_{l,c} | \sigma, \lambda) \quad (1)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,c} | \mathbf{z}, \mathbf{y}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{1:N}, \eta_{l,c} | \sigma, \lambda)}{\int_{\theta, z_{1:N}, \phi_{1:K}, \eta_{l,c} | \mathbf{z}, \mathbf{y}, \alpha, \beta, \alpha', \gamma} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{l,c} | \mathbf{z}, \mathbf{y}, \alpha, \beta, \alpha', \gamma, w_{1:N}, y_{1:N}, \eta_{l,c} | \sigma, \lambda) d\theta dz_{1:N} d\phi_{1:K} d\eta_{l,c}}$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

2.2 Gibbs Sampling

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [12]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

$$p(z_{m,n} | \mathbf{z}_{-(m,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant.

where $\mathbf{Y}_{-(d,i_{l,c})}$ denotes all of the response variables excluding the response variable being sampled.

2.3 ICD-9 Code Hierarchy

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an is-a hierarchy [3].

An is-a hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors

$$p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_{l,c} \in \mathcal{Z}} p(a_{d,i_{l,c}} | \mathbf{z}, \eta_{i_{l,c}}) p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{m,n}$ and absorb all other constant terms into the normalization constant [12].

$$\propto \prod_{i_{l,c} \in \mathcal{Z}} \exp \left\{ -\frac{(\bar{z}^T \eta_{i_{l,c}} - a_{d,i_{l,c}})^2}{2} \right\} \binom{n_{d,(l,c)}}{n_{d,(l,c)} + \alpha_{i_{l,c}}} \frac{n_{d,(l,c)}^{n_{d,(l,c)}}}{\sum_{v=1}^V \binom{n_{d,(l,c)}}{v} n_{d,(l,c)} + \gamma} \quad (4)$$

Here, $n_{d,(l,c)}$ represents the count of word v in document d assigned to topic k omitting the $(d, n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

$$p(\eta_{i_{l,c}} | \mathbf{z}_{1:D}, \mathbf{a}, \sigma) \quad (2.2.4)$$

Given that $\eta_{i_{l,c}}$ and $a_{d,i_{l,c}}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$p(\eta_{i_{l,c}} | \mathbf{z}_{1:D}, \mathbf{a}, \sigma) = \mathcal{N}(\eta_{i_{l,c}} | \hat{\mu}_i, \hat{\Sigma}_i) \quad (5)$$

$$\hat{\mu}_i = \hat{\Sigma}_i (-1\sigma^{-1} + \bar{Z}^T \eta_{i_{l,c}}) \quad (6)$$

$$\hat{\Sigma}_i^{-1} = \mathbf{I}\sigma^{-1} + \bar{Z}^T \bar{Z} \quad (7)$$

$$p(a_{d,i_{l,c}} | \mathbf{z}, \mathbf{Y}, \eta) \text{ and } p(y_{m,i} | \mathbf{a}) \quad (2.2.3)$$

In the augmented probit regression model, the posterior distribution of $a_{d,i_{l,c}}$ is distributed according to a truncated normal distribution where the response variable is observed.

$$p(a_{d,i_{l,c}} | \mathbf{z}, \mathbf{Y}, \eta) = \begin{cases} \text{truncN}^+(a_{d,i_{l,c}}, \eta_i^T \bar{z}, 1, y_{d,i_{l,c}}) & \text{if } y_{d,i_{l,c}} = 1 \\ \text{truncN}^-(a_{d,i_{l,c}}, \eta_i^T \bar{z}, 1, y_{d,i_{l,c}}) & \text{if } y_{d,i_{l,c}} = -1 \end{cases} \quad (8)$$

$$p(y_{d,i_{l,c}} | \mathbf{a}, \mathbf{y}_{-(l,c)}) = \delta(\text{sign}(a_{d,i_{l,c}}) = y_{d,i_{l,c}}) p(y_{d,i_{l,c}} | y_{parent_{l,c}}) \prod_{i_{j,c} \in children_{l,c}} p(y_{i_{j,c}} | y_{d,i_{l,c}}) \quad (9)$$

$$p(y_{d,i_{l,c}} = -1 | y_{parent_{l,c}}) = \begin{cases} 1, & y_{parent_{l,c}} = -1 \\ 0.5, & y_{parent_{l,c}} = 1 \end{cases} \quad (10)$$

$$p(a_{d,i_{l,c}}, y_{d,i_{l,c}} | \mathbf{z}, \mathbf{Y}_{-(d,i_{l,c})}, \eta) = \begin{cases} \text{truncN}^-(a_{d,i_{l,c}} | \bar{z}^T \eta_{l,c}, 1) \delta(y_{d,i_{l,c}} = -1), & y_{parent_{l,c}} = 1, y_{d,i_{l,c}} \in children_{l,c}, y_{l,c} = -1 \\ \text{truncN}^+(a_{d,i_{l,c}} | \bar{z}^T \eta_{l,c}, 1) \delta(y_{d,i_{l,c}} = 1), & y_{parent_{l,c}} = -1, y_{d,i_{l,c}} \in children_{l,c}, y_{l,c} = 1 \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was unobserved, it must be sampled jointly with $y_{d,i_{l,c}}$ to ensure that the Markov chain is ergodic. Suppose that $a_{d,i_{l,c}}$ is sampled to have a negative value and $y_{d,i_{l,c}}$ is appropriately sampled. Although there exist valid states where $a_{d,i_{l,c}} > 0$ and $y_{d,i_{l,c}} = 1$, they will never be reached by such a Markov chain since $p(a_{d,i_{l,c}} < 0 | y_{d,i_{l,c}} = -1) = p(y_{d,i_{l,c}} = 1 | a_{d,i_{l,c}} < 0) = 1$. Therefore, to ensure ergodicity, $a_{d,i_{l,c}}$ and $y_{d,i_{l,c}}$ must be sampled from the joint distribution as shown in Equation ??.

$$p(a_{d,i_{l,c}}, y_{d,i_{l,c}} | \mathbf{z}, \mathbf{Y}_{-(d,i_{l,c})}, \eta) \propto p(a_{d,i_{l,c}} | y_{parent_{l,c}}) p(y_{d,i_{l,c}} | a_{d,i_{l,c}}) \quad (12)$$

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was unobserved, it must be sampled jointly with $y_{d,i_{l,c}}$ to ensure that the Markov chain is ergodic. Suppose that $a_{d,i_{l,c}}$ is sampled to have a negative value and $y_{d,i_{l,c}}$ is appropriately sampled. Although there exist valid states where $a_{d,i_{l,c}} > 0$ and $y_{d,i_{l,c}} = 1$, they will never be reached by such a Markov chain since $p(a_{d,i_{l,c}} < 0 | y_{d,i_{l,c}} = -1) = p(y_{d,i_{l,c}} = 1 | a_{d,i_{l,c}} < 0) = 1$. Therefore, to ensure ergodicity, $a_{d,i_{l,c}}$ and $y_{d,i_{l,c}}$ must be sampled from the joint distribution as shown in Equation ??.

$$p(a_{d,i_{l,c}}, y_{d,i_{l,c}} | \mathbf{z}, \mathbf{Y}_{-(d,i_{l,c})}, \eta) \propto p(a_{d,i_{l,c}} | y_{parent_{l,c}}) p(y_{d,i_{l,c}} | a_{d,i_{l,c}}) \quad (12)$$

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was unobserved, it must be sampled jointly with $y_{d,i_{l,c}}$ to ensure that the Markov chain is ergodic. Suppose that $a_{d,i_{l,c}}$ is sampled to have a negative value and $y_{d,i_{l,c}}$ is appropriately sampled. Although there exist valid states where $a_{d,i_{l,c}} > 0$ and $y_{d,i_{l,c}} = 1$, they will never be reached by such a Markov chain since $p(a_{d,i_{l,c}} < 0 | y_{d,i_{l,c}} = -1) = p(y_{d,i_{l,c}} = 1 | a_{d,i_{l,c}} < 0) = 1$. Therefore, to ensure ergodicity, $a_{d,i_{l,c}}$ and $y_{d,i_{l,c}}$ must be sampled from the joint distribution as shown in Equation ??.

$$p(a_{d,i_{l,c}}, y_{d,i_{l,c}} | \mathbf{z}, \mathbf{Y}_{-(d,i_{l,c})}, \eta) \propto p(a_{d,i_{l,c}} | y_{parent_{l,c}}) p(y_{d,i_{l,c}} | a_{d,i_{l,c}}) \quad (12)$$

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was unobserved, it must be sampled jointly with $y_{d,i_{l,c}}$ to ensure that the Markov chain is ergodic. Suppose that $a_{d,i_{l,c}}$ is sampled to have a negative value and $y_{d,i_{l,c}}$ is appropriately sampled. Although there exist valid states where $a_{d,i_{l,c}} > 0$ and $y_{d,i_{l,c}} = 1$, they will never be reached by such a Markov chain since $p(a_{d,i_{l,c}} < 0 | y_{d,i_{l,c}} = -1) = p(y_{d,i_{l,c}} = 1 | a_{d,i_{l,c}} < 0) = 1$. Therefore, to ensure ergodicity, $a_{d,i_{l,c}}$ and $y_{d,i_{l,c}}$ must be sampled from the joint distribution as shown in Equation ??.

$$p(a_{d,i_{l,c}}, y_{d,i_{l,c}} | \mathbf{z}, \mathbf{Y}_{-(d,i_{l,c})}, \eta) \propto p(a_{d,i_{l,c}} | y_{parent_{l,c}}) p(y_{d,i_{l,c}} | a_{d,i_{l,c}}) \quad (12)$$

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was unobserved, it must be sampled jointly with $y_{d,i_{l,c}}$ to ensure that the Markov chain is ergodic. Suppose that $a_{d,i_{l,c}}$ is sampled to have a negative value and $y_{d,i_{l,c}}$ is appropriately sampled. Although there exist valid states where $a_{d,i_{l,c}} > 0$ and $y_{d,i_{l,c}} = 1$, they will never be reached by such a Markov chain since $p(a_{d,i_{l,c}} < 0 | y_{d,i_{l,c}} = -1) = p(y_{d,i_{l,c}} = 1 | a_{d,i_{l,c}} < 0) =$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction, making LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks and show both improved label prediction performance and show evidence that the learned topic model improves as a result of using this signal too.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [2], product descriptions and catalogs (e.g., as found in [4]), medical discharge records via International Classification of Diseases Ninth Revision, Clinical Modification (ICD-9-CM) codes assigned [1]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision": often taking the form of a single numerical or categorical "label." More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditional drawn as an inferred document-specific topic mixture. It will have been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 2 we introduce hierarchically supervised LDA (HSLDA), in Section 3 we review related work, and in Section ?? we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Model

Given the number of topics, K , and broad gamma priors on hyperparameters, the generative process is as follows:

- For each topic, k :
 - Draw a distribution over words $\phi_k \sim Dir_V(1, \gamma)$
- For each ICD-9 code, i , at all levels in the tree:
 - Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
 - Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
- For each document, d :
 - Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$

1

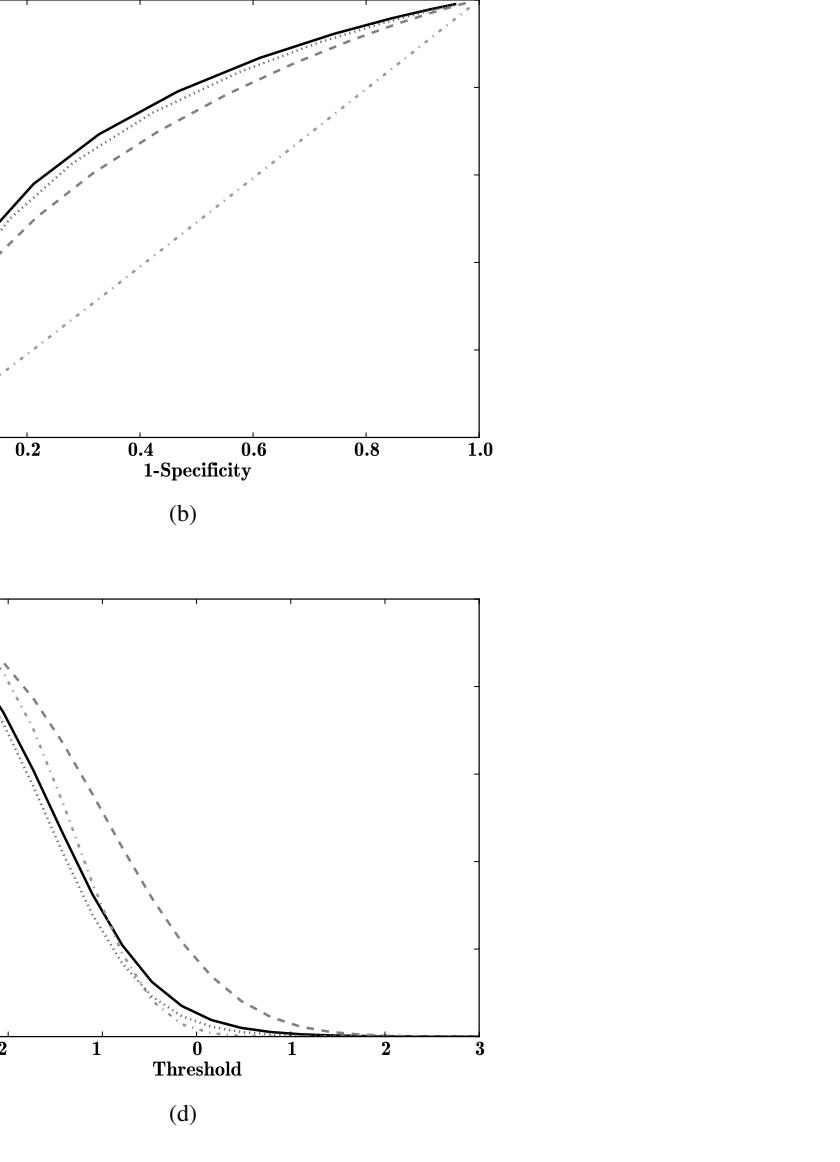


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed is ADLER, dotted is ADLER, and dot-dashed is ADLER.

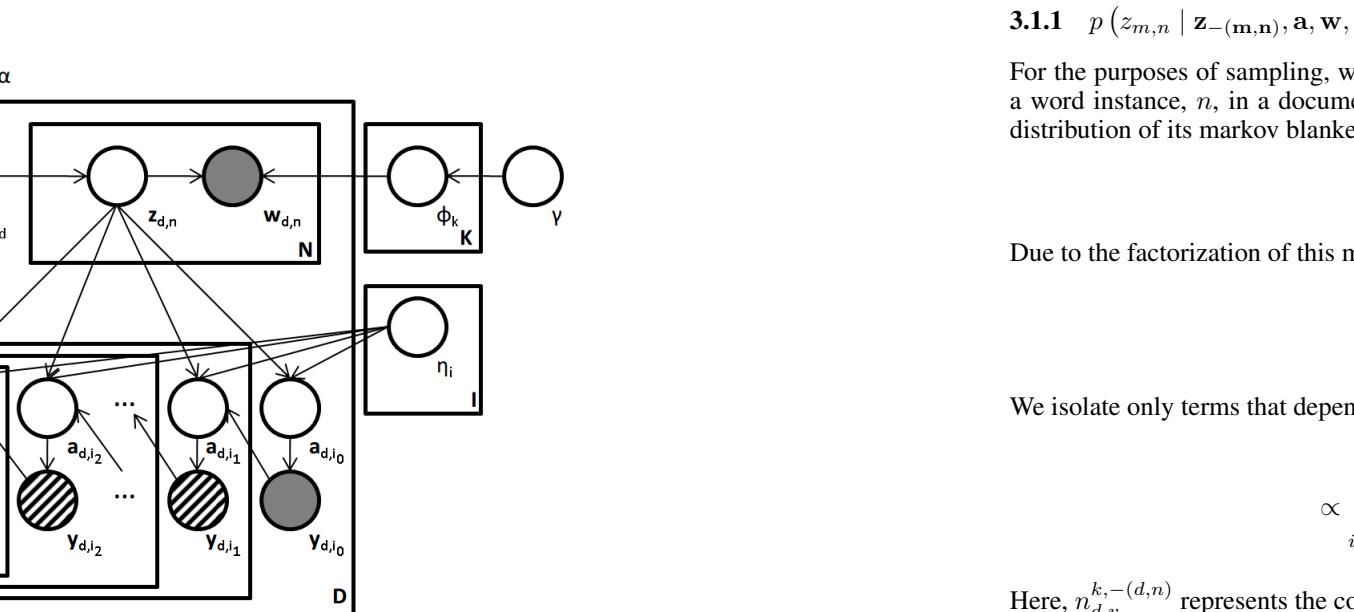


Figure 1: adapted sLDA model

(b) For each word, n :

- Draw topic assignment $z_{n,d} | \theta_d \sim Mult_K(\theta_d)$
- Draw word $w_{n,d} | z_{n,d}, \beta_{1,K} \sim Mult_V(\beta_{1,n})$
- For each level of the ICD-9 code tree, i :
 - For each ICD-9 code at this level, c :
 - Draw a latent variable $a_{d,i,c} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_{i,c}, 1), & y_{d,i,c} = 1 \\ \text{truncN}^+(\bar{z}^T \eta_{i,c}, 1), & y_{d,i,c} = -1 \end{cases}$

$$\text{where } \bar{z} = N^{-1} \sum_{n=1}^N z_n \text{ and } \mathcal{I} = \{i_0, i_1, \dots, i_C\} \text{ and } i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$$

- Draw a response variable $y_{d,i,c} | a_{d,i,c} \sim \begin{cases} 1, & a_{d,i,c} > 0 \\ -1, & \text{otherwise} \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3 Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$\begin{aligned} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathcal{I}, a_{d,i,c} \in \mathcal{I}, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{d,i,c} \in \mathcal{I}; \sigma, \lambda) \\ = \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathcal{I}, a_{d,i,c} \in \mathcal{I}, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{d,i,c} \in \mathcal{I}; \sigma, \lambda)}{\int_{\theta} \int_{\phi} \int_{\eta} \int_{a} \int_{\beta} \int_{\alpha} \int_{\alpha'} \int_{\gamma} \sum_{\mathcal{I}} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathcal{I}, a_{d,i,c} \in \mathcal{I}, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{d,i,c} \in \mathcal{I}; \sigma, \lambda)} \end{aligned} \quad (1)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [12]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

2

- [16] P Ruch, J Gobell, I Thairii, and A Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [17] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [18] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking lda: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, pages 1973–1981, 2009.

3

- ### 6.1 Rao-Blackwellization
- To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

$$\begin{aligned} p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} \in \mathcal{I}, a_{d,i,c} \in \mathcal{I}, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{d,i,c} \in \mathcal{I}; \sigma, \lambda) &= \frac{p(\theta | \alpha) \prod_{n=1}^N p(z_{n,d}) p(w_{n,d} | z_n, \phi_{1:K}) \prod_{k=1}^K p(\phi_k | \beta) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c})}{\int_{\theta} \int_{\phi} p(\theta | \alpha) \sum_{\mathcal{I}} p(\phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c})} \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \\ &= \int_{\theta} \int_{\phi} \int_{\eta} \int_{a} \int_{\beta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \\ &= \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C \int_{\theta} \int_{\phi} \int_{\eta} \prod_{n=1}^N p(z_{n,d} | \theta, z_{n,d-1}, \phi_{1:K}) \prod_{i=1}^I p(\eta_{i,c} | \beta) \prod_{d=1}^D p(a_{d,i,c} | \beta) \prod_{i=1}^I p(y_{d,i,c} | a_{d,i,c}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(w_{n,d} | z_{n,d}, \phi_{1:K}) \\ &\quad \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(a_{d,i,c} | \theta, z_{n,d}, \phi_{1:K}, \eta_{i,c}) \prod_{d=1}^D \prod_{i=1}^I \prod_{c=1}^C p(y_{d,i,c} | a_{d,i,c}, \theta, z_{n,d}, \$$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction. Our approach compares favorably to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks and show both improved label prediction performance and show evidence that the learned topic model improves as a result of using this signal too.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [2], product descriptions and catalogs (e.g., as found in [4]), free-text clinical discharge records via International Classification of Diseases 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [1]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision": often taking the form of a single numerical or categorical "label." More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditional drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have a browserable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section ?? we introduce hierarchically supervised LDA (HSLDA), in Section 4.4 we review related work, and in Section ?? we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Model

Given the number of topics, K , and broad gamma priors on hyperparameters, the generative process is as follows:

- For each topic, k :
(a) Draw a distribution over words $\phi_k \sim Dir_V(1, \gamma)$
- For each ICD-9 code, i , at all levels in the tree:
(a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
- Draw a prior over topic proportions $\beta | \alpha' \sim Dir_K(1, \alpha')$
- For each document, d :
(a) Draw topic proportions $\theta_d | \beta \sim Dir_K(\beta, \alpha)$

1

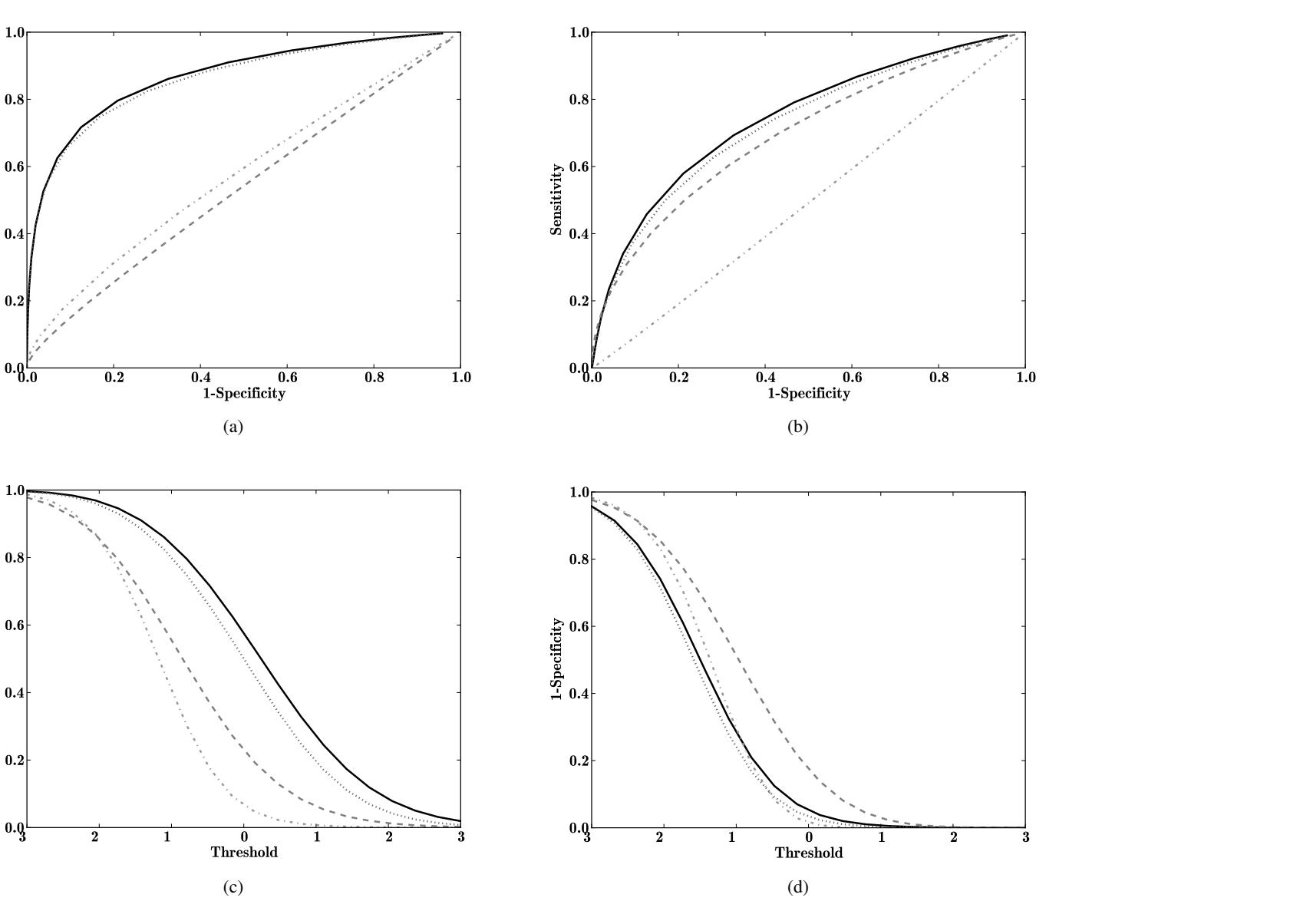


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed is ADLER, dotted is ADLER, and dot-dashed is ADLER.

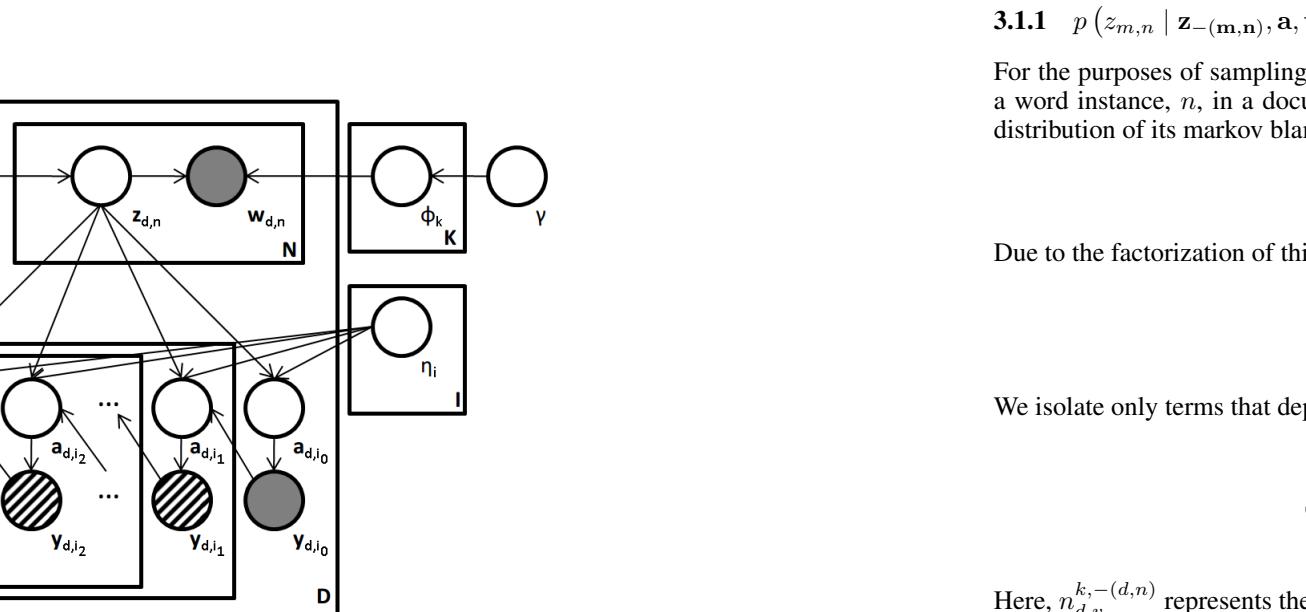


Figure 1: adapted sLDA model

3.1.1 $p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z_d , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant:

$$p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_c \in \mathcal{I}} \prod_{i_c \in \mathcal{I}} p(a_{d,i_c} | \mathbf{z}, \eta_{i_c}) p(z_{d,n}, \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{d,n}$ and absorb all other constant terms into the normalization constant [12]:

$$\propto \prod_{i_c \in \mathcal{I}} \exp \left\{ -\frac{(\bar{z}^T \eta_{i_c} - a_{i_c})^2}{2} \right\} \left(n_{d,i_c}^{k_{i_c}(d,n)} + \alpha \right) \sum_{v=1}^V \binom{n_{d,i_c}(d,n)}{n_{d,i_c} v} + \gamma \quad (4)$$

Here, $n_{d,i_c}^{k_{i_c}(d,n)}$ represents the count of word v in document d assigned to topic k omitting the $(d,n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} | \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\eta_{i_c} | \mathbf{z}_{1:D}, \mathbf{a}; \sigma)$

Given that η_{i_c} and a_{d,i_c} are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$:

$$p(\eta_{i_c} | \mathbf{z}_{1:D}, \mathbf{a}; \sigma) = \mathcal{N}(\eta_{i_c} | \bar{\mu}_i, \hat{\Sigma}_i) \quad (5)$$

$$\bar{\mu}_i = \hat{\Sigma}_i (-1\sigma^{-1} + \bar{Z}^T \mathbf{a}_{(i)}) \quad (6)$$

$$\hat{\Sigma}_i^{-1} = \mathbf{I}\sigma^{-1} + \bar{Z}^T \bar{Z} \quad (7)$$

3.1.3 $p(a_{d,i_c} | \mathbf{z}, \mathbf{Y}, \eta)$ and $p(y_{m,i} | \mathbf{a})$

In the augmented probit regression model, the posterior distribution of a_{d,i_c} is distributed according to a truncated normal distribution where the response variable is observed.

$$p(a_{d,i_c} | \mathbf{z}, \mathbf{Y}, \eta) = \begin{cases} \text{truncN}^+(a_{d,i_c} | \eta_i^T \bar{z}, 1), & \text{if } y_{d,i_c} = 1 \\ \text{truncN}^-(a_{d,i_c} | \eta_i^T \bar{z}, 1), & \text{if } y_{d,i_c} = -1 \end{cases} \quad (8)$$

$$p(y_{m,i} | \mathbf{a}) = \delta(\text{sign}(a_{d,i_c}) = y_{m,i}) p(y_{m,i} | \text{parents}_{m,i}) \prod_{i_c \in \text{children}_m} p(y_{i_c} | y_{m,i}) \quad (9)$$

$$p(y_{m,i} = -1 | \text{parent}_{m,i}) = \begin{cases} 1, & \text{parent}_{m,i} = -1 \\ 0.5, & \text{parent}_{m,i} = 1 \end{cases} \quad (10)$$

$$p(a_{d,i_c}, y_{d,i_c} | \mathbf{z}, \mathbf{Y}, \eta) = \begin{cases} \mathcal{N}(a_{d,i_c} | \bar{z}^T \eta_{i_c}, 1) p(y_{d,i_c} | a_{d,i_c}), & y_{parent_{i_c}} = 1, \forall y_{i_c} \in \text{children}_{i_c}, y_{i_c} = -1 \\ \text{truncN}^-(a_{d,i_c} | \bar{z}^T \eta_{i_c}, 1) \delta(y_{d,i_c} = -1), & y_{parent_{i_c}} = -1 \\ \text{truncN}^+(a_{d,i_c} | \bar{z}^T \eta_{i_c}, 1) \delta(y_{d,i_c} = 1), & \exists y_{i_c} \in \text{children}_{i_c} \setminus y_{i_c} = 1 \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

where $\mathbf{Y}_{-(d,i_c)}$ denotes all of the response variables excluding the response variable being sampled.

4.3 ICD-9 Code Hierarchy

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an *is-a* hierarchy [3]. An *is-a* hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors.

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c} \in \mathcal{I}, a_{i_c} \in \mathcal{A}, \beta, \alpha', \gamma | w_{1:N}, y_{i_c} \in \mathcal{I}; \sigma, \lambda) \quad (12)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c} \in \mathcal{I}, a_{i_c} \in \mathcal{A}, \beta, \alpha', \gamma | w_{1:N}, y_{i_c} \in \mathcal{I}; \sigma, \lambda)}{\int_0^\infty p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a})} \quad (13)$$

$$= \int_0^\infty \int_{\phi_{1:K}} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (14)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (15)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (16)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (17)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (18)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} \int_{\gamma} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (19)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} \int_{\gamma} \int_{\lambda} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (20)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} \int_{\gamma} \int_{\lambda} \int_{\sigma} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (21)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} \int_{\gamma} \int_{\lambda} \int_{\sigma} \int_{\phi_{1:K}} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (22)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} \int_{\gamma} \int_{\lambda} \int_{\sigma} \int_{\phi_{1:K}} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (23)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} \int_{\gamma} \int_{\lambda} \int_{\sigma} \int_{\phi_{1:K}} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (24)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} \int_{\gamma} \int_{\lambda} \int_{\sigma} \int_{\phi_{1:K}} p(\theta | \sigma) \sum_{k=1}^K p(\phi_{1:k} | \beta) \prod_{i_c=1}^N p(a_{i_c} | \mathbf{z}, \eta_{i_c}, \alpha', \beta, \gamma) \prod_{i_c=1}^N p(z_{i_c} | \mathbf{z}_{-(i_c)}, \mathbf{a}, \phi_{1:i_c}, \beta, \gamma) \prod_{i_c=1}^N p(y_{i_c} | \mathbf{a}) \quad (25)$$

$$= \int_0^\infty \int_{\phi_{1:K}} \int_{\eta_{i_c} \in \mathcal{I}} \int_{a_{i_c} \in \mathcal{A}} \int_{\beta} \int_{\alpha'} \int_{\gamma} \int_{\lambda} \int_{\sigma} \int_{\phi_{1:K}} p(\theta | \sigma) \sum_{k=1$$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially outperforms out-of-sample label prediction compared to simple LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks and show both improved label prediction performance and show evidence that the learned topic model improves as a result of using this signal too.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [2], product descriptions and catalogs (e.g., as found in [4]), medical discharge summaries and International Classification of Diseases 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [1]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision": often taking the form of a single numerical or categorical "label." More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g., online review scores), marks given to written work (e.g., essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditional drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiple situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and Web retail data suggests that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section ?? we introduce hierarchically supervised LDA (HSLDA), in Section 4.4 we review related work, and in Section ?? we apply HSLDA to health care and Web retail data, showing predictive performance and improved topic generation.

2 Model

Given the number of topics, K , and broad gamma priors on hyperparameters, the generative process is as follows:

- For each topic, k :
 (a) Draw a distribution over words $\phi_k \sim Dir(\gamma, 1)$
- For each ICD9 code, c , at all levels in the tree, I :
 (a) Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
- Draw a prior over topic proportions $\beta | \alpha' \sim Dir(1, \alpha')$
- For each document, d :
 (a) Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$

Figure 1: adapted sLDA model

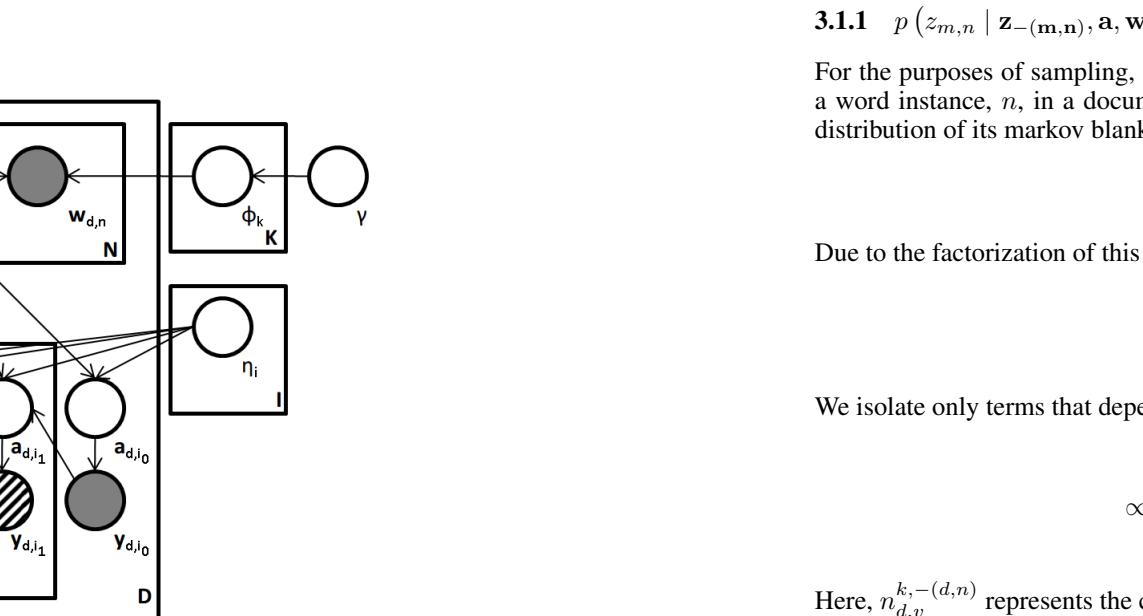


Figure 1: adapted sLDA model

(b) For each word, n :

- Draw topic assignment $z_{d,n} | \theta_d \sim Mult_K(\theta_d)$
- Draw word $w_{n,d} | z_{d,n}, \beta_{1,K} \sim Mult_V(\beta_{1,n})$
- For each level of the ICD-9 code tree, I :
 i. For each ICD-9 code at this level, c :
 A. Draw a latent variable $a_{d,i,c} \sim \begin{cases} N(z^T \eta_{i,c}, 1), & y_{d,i,c} = 1 \\ truncN(z^T \eta_{i,c}, 1), & y_{d,i,c} = -1 \end{cases}$
 where $z = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_C\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$
- B. Draw a response variable $y_{d,i,c} | a_{d,i,c} \sim \begin{cases} 1, & a_{d,i,c} > 0 \\ -1, & otherwise \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3 Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} | \tau_{i,c} \in \mathcal{I}, a_{i,c} \in \mathcal{A}, \beta, \alpha, \alpha', \gamma | w_{1:N}, y_{d,i,c} | \sigma, \lambda) \quad (1)$$

$$= \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} | \tau_{i,c} \in \mathcal{I}, a_{i,c} \in \mathcal{A}, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{d,i,c} | \sigma, \lambda)}{\int_0^\infty \int_0^\infty \dots \int_0^\infty p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i,c} | \tau_{i,c} \in \mathcal{I}, a_{i,c} \in \mathcal{A}, \beta, \alpha, \alpha', \gamma, w_{1:N}, y_{d,i,c} | \sigma, \lambda) d\eta_{i,c} da_{i,c} \dots d\theta} \quad (2)$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [12]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

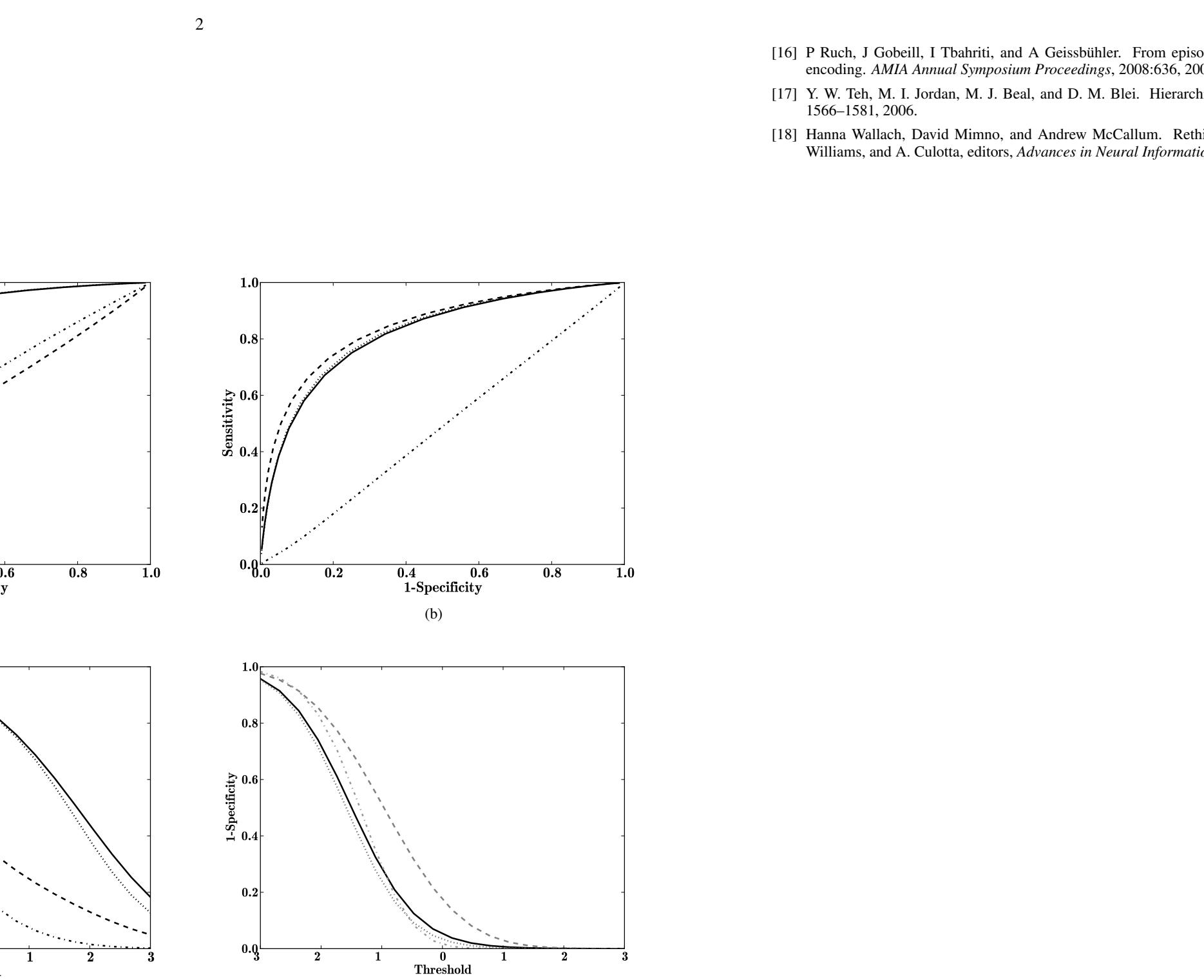


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: ?? includes ancestor prediction performance. ?? results are for given (leaf) labels alone. Bottom row: ?? are the sensitivity curves from ?? aligned on threshold value, ?? are the 1-specificity curves from ?? aligned on threshold value.

3.1.4 $p(z_{d,n} | z_{-(d,n)}, \alpha, w, \eta, \alpha, \beta, \gamma)$

For the purposes of sampling, we will be able to derive a representation of the joint distribution isolating a particular latent variable, z , for a word instance, n , in a document instance, d . The conditional probability with respect to this latent variable is proportional to the joint distribution of its markov blanket up to a constant:

$$p(z_{d,n} | z_{-(d,n)}, \alpha, w, \eta, \alpha, \beta, \gamma) \propto p(z_{d,n}, z_{-(d,n)}, \alpha, w, \eta, \alpha, \beta, \gamma) \quad (2)$$

Due to the factorization of this model, we can rewrite the joint distribution as the following:

$$\propto \prod_{i_l \in \mathcal{I}} p(a_{i_l,c} | z, \eta_{i_l,c}) p(z_{d,n}, z_{-(d,n)}, \alpha, w, \eta, \alpha, \beta, \gamma) \quad (3)$$

We isolate only terms that depend on $z_{d,n}$ and absorb all other constant terms into the normalization constant [12].

$$\propto \prod_{i_l \in \mathcal{I}} \exp \left\{ -\frac{(z^T \eta_{i_l,c} - a_{i_l,c})^2}{2} \right\} \left(n_{d,c}^{k_{d,c}(d,n)} + \alpha \right) \sum_{i_l=1}^V \binom{n_{i_l,c}(d,n)}{n_{i_l,c}} + \gamma \quad (4)$$

Here, $n_{d,c}^{k_{d,c}}$ represents the count of word v in document d assigned to topic k omitting the $(d, n)^{th}$ word count.

Given Equation 4, $p(z_{d,n} | z_{-(d,n)}, \alpha, w, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\eta_{i_l,c} | z_{1:D}, \alpha, \sigma)$

Given that $\eta_{i_l,c}$ and $a_{d,i_l,c}$ are distributed normally, this posterior distribution is also normal. We evaluated the model over various values of σ where $\sigma = \{0.01, 0.1, 0.25, 1, 2\}$.

$$p(\eta_{i_l,c} | z_{1:D}, \alpha, \sigma) = \mathcal{N}(\bar{\eta}_{i_l,c}, \hat{\Sigma}_{i_l,c}) \quad (5)$$

$$\bar{\eta}_{i_l,c} = \hat{\Sigma}_{i_l,c}^{-1} (-1\sigma^{-1} + \bar{Z}^T \mathbf{a}_{i_l,c}) \quad (6)$$

$$\hat{\Sigma}_{i_l,c}^{-1} = \mathbf{I}\sigma^{-1} + \bar{Z}^T \bar{Z} \quad (7)$$

3.1.3 $p(a_{d,i_l,c} | z, \mathbf{Y}, \eta)$ and $p(y_{m,i} | \mathbf{a})$

In the augmented probit regression model, the posterior distribution of $a_{d,i_l,c}$ is distributed according to a truncated normal distribution where the response variable is observed.

$$p(a_{d,i_l,c} | z, \mathbf{Y}, \eta) = \begin{cases} truncN^+(a_{d,i_l,c} | \eta_{i_l,c}^T \bar{z}_1, y_{d,i_l,c}) & if \ y_{d,i_l,c} = 1 \\ truncN^-(a_{d,i_l,c} | \eta_{i_l,c}^T \bar{z}_1, y_{d,i_l,c}) & if \ y_{d,i_l,c} = -1 \end{cases} \quad (8)$$

However, if $y_{d,i_l,c}$ is unobserved then $a_{d,i_l,c}$ must be sampled jointly with $y_{d,i_l,c}$ to ensure that the Markov chain is ergodic. Suppose that $a_{d,i_l,c}$ is sampled to have a negative value and $y_{d,i_l,c}$ is appropriately sampled at -1. Although there exist valid states where $a_{d,i_l,c} > 0$ and $y_{d,i_l,c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_l,c} < 0 | y_{d,i_l,c} = -1) = 1$ and $p(y_{d,i_l,c} = -1 | a_{d,i_l,c} < 0) = 1$. Therefore, to ensure ergodicity, $a_{d,i_l,c}$ and $y_{d,i_l,c}$ must be sampled from the joint distribution as shown in Equation ??.

$$p(a_{d,i_l,c}, y_{d,i_l,c} | z, \mathbf{Y}, \eta) \propto p(a_{d,i_l,c} | z, \mathbf{Y}, \eta) p(y_{d,i_l,c} | a_{d,i_l,c}) \quad (9)$$

$$p(y_{d,i_l,c} = -1 | y_{d,i_l,c}) = \begin{cases} 1, & y_{d,i_l,c} = -1 \\ 0, & y_{d,i_l,c} = 1 \end{cases} \quad (10)$$

$$= \begin{cases} \mathcal{N}(a_{d,i_l,c} | \bar{z}^T \eta_{i_l,c}, 1) p(y_{d,i_l,c} | a_{d,i_l,c}), & y_{d,i_l,c} = 1, \forall y_{i_l,c} \in y_{children_{i_l,c}}, y_{i_l,c} = -1 \\ truncN^-(a_{d,i_l,c} | \bar{z}^T \eta_{i_l,c}, 1) \delta(y_{d,i_l,c} = -1), & y_{d,i_l,c} = -1 \\ truncN^+(a_{d,i_l,c} | \bar{z}^T \eta_{i_l,c}, 1) \delta(y_{d,i_l,c} = 1), & \exists y_{i_l,c} \in y_{children_{i_l,c}} \setminus y_{d,i_l,c} = 1 \\ 0, & otherwise \end{cases} \quad (11)$$

where $\mathbf{Y}_{-(d,i_l,c)}$ denotes all of the response variables excluding the response variable being sampled.

3.4 ICD-9 Code Hierarchy

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an $is-a$ hierarchy [3]. An $is-a$ hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors.

4.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

$$[16] P. Ruch, J. Gobell, I. Thivierge, and A. Geissbühler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.$$

$$[17] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.$$

$$[18] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking lda: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, pages 1973–1981, 2009.$$

4.2 Pre-Processing

Patient discharge summaries and their associated ICD-9 diagnoses are stored in two different places in the New York - Presbyterian data warehouse, and so had to be linked together before being fed into the sLDA algorithm. Each discharge note and set of diagnoses were assigned a unique identifier (PVID), allowing the two types of data to be linked.

Natural Language Processing (NLP) techniques were used to process the free-text discharge summaries. First, the Natural Language Toolkit (<http://www.nltk.org/>) was used to tokenize the text. Next, feature selection was performed using a term frequency - inverse document frequency (TF-IDF) algorithm on the entire document set and sorting the words by their TF-IDF values. The top 10,000 words were manually evaluated to eliminate all potentially identifying information. Finally, each discharge summary was converted to a bag-of-words,

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The $is-a$ relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a diagnosis was observed to be absent, all of its descendants could also be absent (e.g., if a patient did not have malignant hypertension, it could be assumed that they did not have essential hypertension). Unfortunately, ICD-9 code observations never include observations of disease absence. ICD-9 codes are only documented when the condition is observed to be present. Additionally, ICD-9 codes are known to have relatively low sensitivity: conditions that are present are often not documented in a set of ICD-9 codes (Surjan, 1999). Given these facts, we made the following assumptions regarding each visit: recorded diagnoses and their ancestors were labeled as true; diagnoses that were observed at some time for a patient but not at the current visit were labeled as unobserved; and diagnoses that had never been listed for a patient were labeled as false for all of that patient's visits. This last assumption captures the belief that parts of the ICD-9 hierarchy that are never observed for a particular patient are almost certain to be false. Additionally, for computational purposes, we decided not to include an ICD-9 code at all if neither it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

ICD-9 code at all if neither it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

4.3 Evaluation

Here, we report the accuracy of the HSLDA model. The HSLDA model is evaluated using the F1 score, precision, recall, and AUC.

4.3.1 HSLDA vs. LDA

Table 1 compares the performance of HSLDA against LDA on two datasets: a medical dataset and a retail dataset.

4.3.2 HSLDA vs. sLDA

Table 2 compares the performance of HSLDA against sLDA on two datasets: a medical dataset and a retail dataset.

4.3.3 HSLDA vs. HSLDA

Table 3 compares the performance of HSLDA against HSLDA on two datasets: a medical dataset and a retail dataset.

4.3.4 HSLDA vs. HSLDA

Table 4 compares the performance of HSLDA against HSLDA on two datasets: a medical dataset and a retail dataset.

4.3.5 HSLDA vs. HSLDA

Table 5 compares the performance of HSLDA against HSLDA on two datasets: a medical dataset and a retail dataset.

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction, compared to standard LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks and show both improved label prediction performance and show evidence that the learned topic model improves as a result of using this signal too.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [2], product descriptions and catalogs (e.g., as found in [4]), medical discharge summaries, International Classification of Diseases 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [1]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-word image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per document "supervision": often taking the form of a single numerical or categorical "label." More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditional drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multi-level) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section ?? we introduce hierarchically supervised LDA (HSLDA), in Section 4.4 we review related work, and in Section ?? we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Model

Given the number of topics, K , and broad gamma priors on hyperparameters, the generative process is as follows:

- For each topic, k :
 - Draw a distribution over words $\phi_k \sim Dir(\gamma, 1 - \gamma)$
- For each ICD9 code, c , at all levels in the tree, I :
 - Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
 - Draw a prior over topic proportions $\beta | \alpha' \sim Dir(1, \alpha')$
- For each document, d :
 - Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$

Figure 1: adapted sLDA model

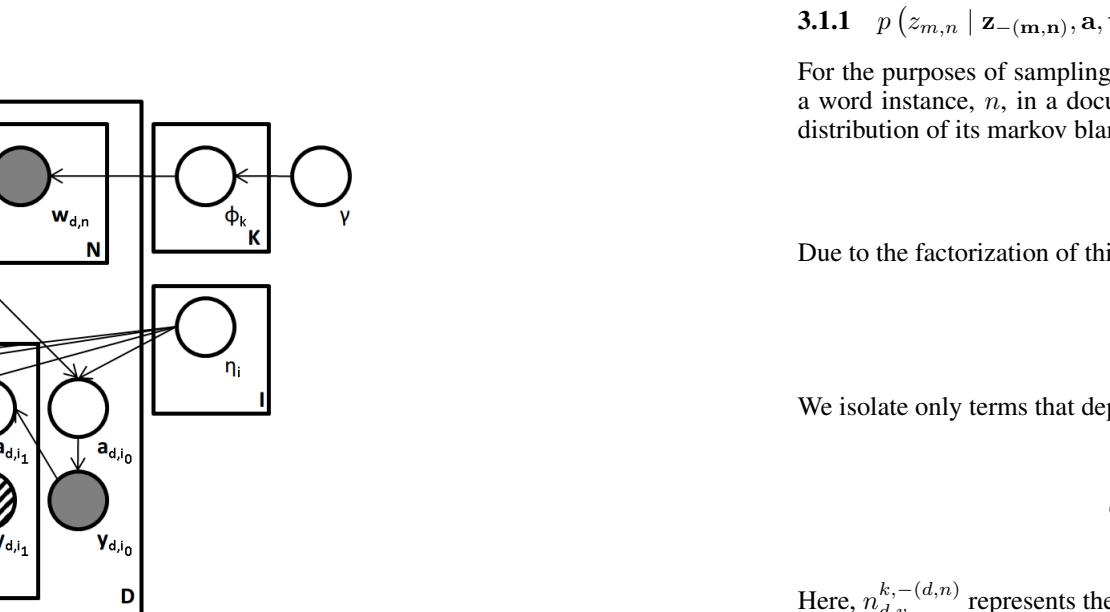


Figure 1: adapted sLDA model

(b) For each word, n :

- Draw topic assignment $z_{n,d} | \theta_d \sim Mult_K(\theta_d)$
- Draw word $w_{n,d} | z_{n,d}, \beta_{i,K} \sim Mult_V(\beta_{i,K})$
- For each level of the ICD-9 code tree, I :
 - For each ICD-9 code at this level, c :
 - Draw a latent variable $a_{d,i_c} \sim \begin{cases} N(z^T \eta_{i_c}, 1), & y_{parent_{i_c}} = 1 \\ truncN(z^T \eta_{i_c}, 1), & y_{parent_{i_c}} = -1 \end{cases}$
- where $\bar{z} = N^{-1} \sum_{n=1}^N z_n$ and $I = \{i_0, i_1, \dots, i_L\}$ and $i_l = \{i_{l,0}, i_{l,1}, \dots, i_{l,C_l}\}$
- B. Draw a response variable $y_{d,i_c} | a_{d,i_c} \sim \begin{cases} 1, & a_{d,i_c} > 0 \\ -1, & otherwise \end{cases}$

The generative model for the ICD-9 codes is equivalent to a probit regression model. In our case, the regression is conditional on the parents according to the constraints of the ICD-9 code tree. The latent variable utilized here is also known as an auxiliary variable.

3 Inference

Given an observation of a set of ICD-9 codes and a document, the posterior distribution for the latent variables is given by Equation .

$$p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}; \sigma, \lambda) = \frac{p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}; \sigma, \lambda)}{\int_0^1 \int_{0, \eta_{i_c}, \alpha, \alpha', \beta, \gamma} \sum_{i=1}^N p(\theta, z_{1:N}, \phi_{1:K}, \eta_{i_c}, \alpha, \alpha', \gamma | w_{1:N}, y_{1:N}; \sigma, \lambda)}$$

The denominator for this distribution is the marginal probability of the data and cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model, marginalizing $\theta_{1:D}$ and $\phi_{1:K}$. For details regarding collapsing in LDA models see Griffiths and Steyvers [12]. To derive the Gibbs sampler we evaluate the individual conditional probability distributions for all unobserved variables.

2 Model

Given the number of topics, K , and broad gamma priors on hyperparameters, the generative process is as follows:

- For each topic, k :
 - Draw a distribution over words $\phi_k \sim Dir(\gamma, 1 - \gamma)$
- For each ICD9 code, c , at all levels in the tree, I :
 - Draw regression coefficient $\eta_{i,c} | \sigma \sim N_K(-1, \sigma)$
 - Draw a prior over topic proportions $\beta | \alpha' \sim Dir(1, \alpha')$
- For each document, d :
 - Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\beta, \alpha)$

Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

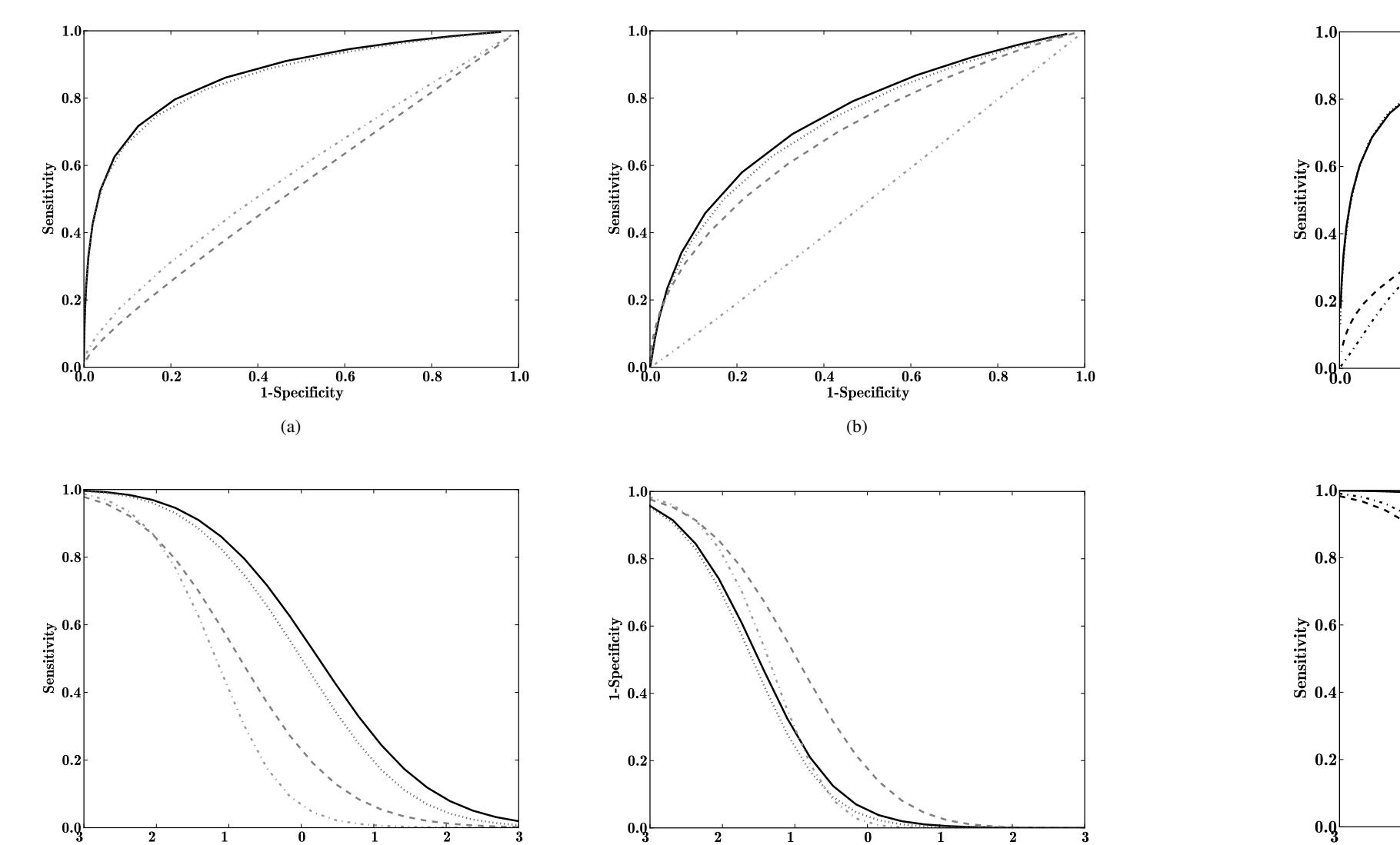


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.2 Prediction

4 Experiments

4.1 Data

Our data set was gathered from the clinical data warehouse of New York - Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The notes outline the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a controlled vocabulary which is the international standard diagnostic classification for epidemiological, administrative and research purposes (http://www.who.int/classifications/icd/en/). The codes are classified in a rooted tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary. For the purposes of sLDA, ICD-9 codes will be used as labels for discharge summaries.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. Before beginning data processing, we generated a PHI-free dataset (see Data Pre-processing below).

4.2 Pre-Processing

Patient discharge summaries and their associated ICD-9 diagnoses are stored in two different places in the New York - Presbyterian data warehouse, and so had to be linked together before being fed into the sLDA algorithm. Each discharge note and set of diagnoses were assigned a patient unique identifier (PUIID), allowing the two types of data to be linked.

However, if y_{d,i_c} is unobserved then a_{d,i_c} must be sampled jointly with y_{d,i_c} to ensure that the Markov chain is ergodic. Suppose that a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, to ensure ergodicity, a_{d,i_c} and y_{d,i_c} must be sampled from the joint distribution as shown in Equation ??.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1$, it will never be reached by such a Markov chain since $p(a_{d,i_c} < 0 | y_{d,i_c} = -1) = 1$ and $p(y_{d,i_c} = -1 | a_{d,i_c} < 0) = 1$. Therefore, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of protected health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. The is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present, all of its ancestors could be assumed to be present as well (e.g., if a patient had malignant hypertension, it could be assumed that they also had essential hypertension). Second, if a_{d,i_c} is sampled to have a negative value and y_{d,i_c} is appropriately sampled at -1. Although there exist valid states where $a_{d,i_c} > 0$ and $y_{d,i_c} = 1</$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model where probit regressors are conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the addition of the hierarchical constraint dramatically constrained labels, substantially improving out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks and show both improved label prediction performance and show evidence that the learned topic model improves as a result of using this signal too.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper, we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [1]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document “supervision,” often taking the form of a single numerical or categorical “label.” More generally this “supervision” can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been drawn from an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision is critical in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and/or multiple labelings in a similar manner. Consider web retail products. A user may purchase a few items and produce a few labels they sell. The signal of each product in a product hierarchy (often multiple) contains a specific, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simple unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section ?? we introduce hierarchically supervised LDA (HSLDA), in Section 4.4 we review related work, and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Model

We define here a hierarchically supervised LDA model. We assume a pre-specified set of labels \mathcal{L} . Each document is assigned a response of either -1 or 1, but potentially many labels in \mathcal{L} . A response of 1 or -1 to label l indicates if the document respectively is or is not l . The label l for a document d will be used interchangeably to refer the observed response of document d to label l . The label set is assumed to be an “is-a” hierarchy. This means that if a label l_1 is a parent of label l_2 in the hierarchy and document d has a positive response to the label l_2 then document d also has a positive response to the label l_1 . Seen in a negative light, if document d has a negative response to the label l_1 then document d also has a negative response to the l_2 label. To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models. Approximate inference is performed using Gibbs sampling.

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , and the standard deviation σ used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$. We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

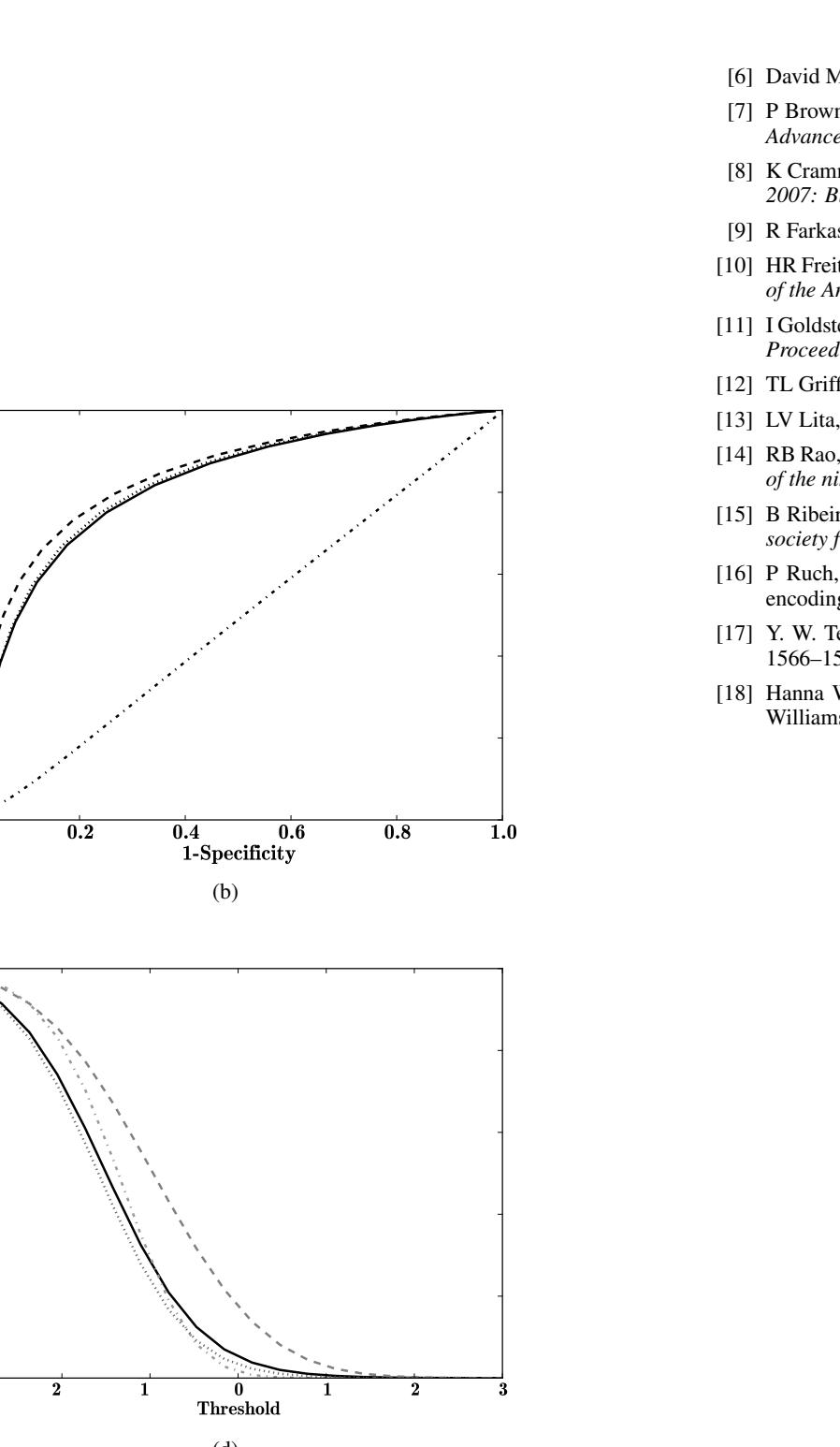


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

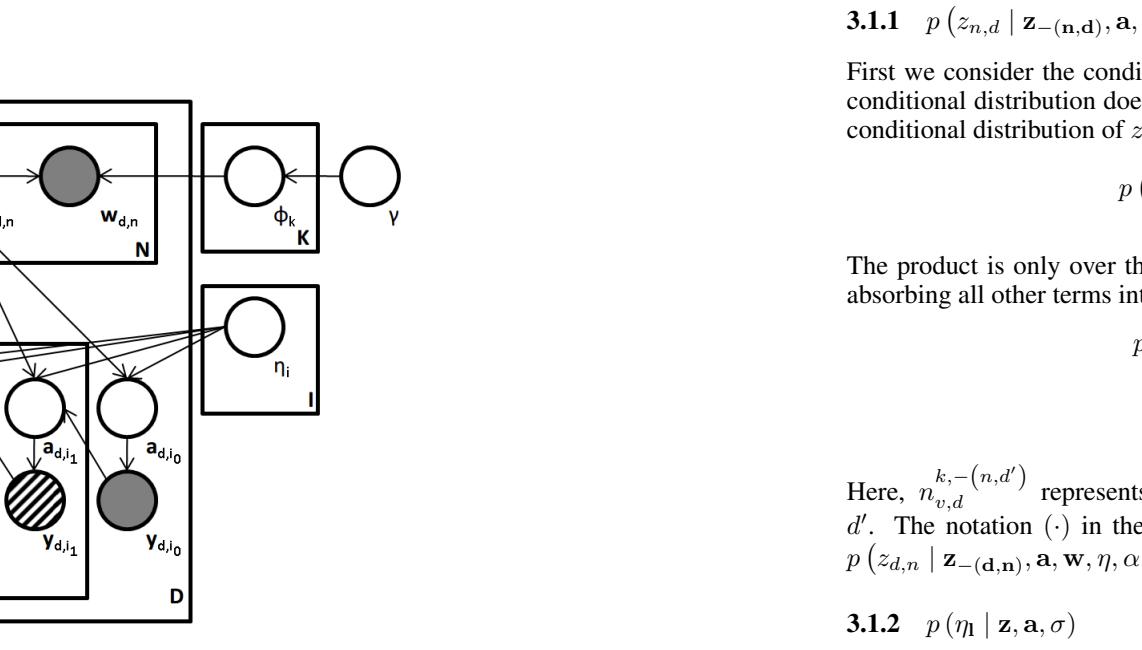


Figure 1: adapted sLDA model

3.1.1 $p(z_{n,d} | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include $\theta_{1:D}$ and $\phi_{1:K}$ because they have been integrated out as in the collapsed Gibbs sampler [12]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket, explicitly this means

$$p(z_{n,d} | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant [12] we find

$$p(z_{n,d} = k | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \binom{n_{k-(n,d)} + \gamma}{n_{k-(n,d)} + \alpha + \beta_k} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(\bar{z}_d^T \eta - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $n_{k-(n,d)}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 4, $p(z_{d,n} | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\eta | \mathbf{z}, \alpha, \sigma)$

We now consider the conditional distribution of the regression coefficients $a_{l,d}$ for $l \in \mathcal{L}$. Given that η_l and $a_{l,d}$ are distributed normally, the posterior distribution of η_l is normally distributed with mean $\bar{\mu}$ and covariance Σ such that

$$\Sigma^{-1} = \mathbf{I}\sigma^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\bar{\mu} = \Sigma_l (-\mathbf{I}\sigma^{-1} + \mathbf{Z}^T \mathbf{a}_l). \quad (4)$$

This is a standard result from normal Bayesian linear regression [2]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \bar{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

3.1.3 $p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \eta)$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^T \bar{z}_d) \right\} I(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta | \mathbf{z}, \alpha', \alpha)$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion. This flexible distribution allows for an asymmetric prior over document level distributions over topics [18].

Posterior inference is performed using the “direct assignment” method of Teh et al. [17].

$$\beta \sim \text{Dir}(m_{(1,1)}, m_{(1,2)}, \dots, m_{(1,K)}) \quad (6)$$

$$p(m_{(d,k)} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha_k)}{\Gamma(\alpha_k + n_{d,k})} s(n_{d,k}, m) (\alpha_k \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α , α' , and γ are given broad Gamma(2, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

3.2 Prediction

4 Experiments

4.1 Data

Our data set was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient’s chief complaint, diagnostic findings, therapy administered, patient’s response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a controlled terminology which is the international standard diagnostic classification for diseases and procedures in medicine [1].

This model, supervised latent dirichlet allocation (sLDA), builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

Other models - predicting document links, other supervised latent variable models

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [16, 10, 15, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [14] but the most recent and promising work on the subject was inspired by the 2007 Medical NLP Challenge: “International Challenge: Classifying Clinical Free Text Using Natural Language Processing” (website). Most of the classification strategies included word matching and rule-based algorithms. [11, 8, 9]. The dataset given to the participants consisted of documents that were 1-2 lines each and all of the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which attempted to work with a document to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

6 Discussion

• what about the nonparametric version of this?

• discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

5 Results

epidemiological, health management, and clinical purposes (<http://www.who.int/classifications/icd/en>). The codes are classified in a root-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subserves the child diagnosis. In each edge, ICD-9 codes are generated mainly by trained medical coders, who review all the information in the discharge summary. For purposes of sLDA, ICD-9 codes will be used as labels for discharge summaries.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. Before beginning data processing, we generated a PHI-free dataset (see Data Pre-processing below).

4.2 Pre-Processing

Patient discharge summaries and their associated ICD-9 diagnoses are stored in two different places in the NewYork-Presbyterian data warehouse, and so had to be linked together before being fed into the sLDA algorithm. Each discharge note and set of diagnoses were assigned a patient unique identifier (PUID) and a visit unique identifier (UUID), allowing the two types of data to be linked.

Natural Language Processing (NLP) techniques were used to process the free-text discharge summaries. First, the Natural Language Toolkit (<http://www.nltk.org/>) was used to tokenize the text. Next, feature selection was performed using a term frequency - inverse document frequency (TF-IDF) algorithm on the entire document set and sorting the words by their TF-IDF values. The top 10,000 words were manually evaluated to eliminate all potentially identifying information. Finally, each discharge summary was converted to a bag-of-words, listing the frequencies of the remaining, free of potential health information, top 10,000 words.

Preparation of the diagnostic codes involved inference over the ICD-9 hierarchy. An is-a relationships of the hierarchy allowed us to make two important assumptions. First, if a diagnosis was observed to be present in a document, it was assumed to be present as well (e.g. if a patient had malignant hypertension, it was assumed that all of its descendants were also present). Second, if a diagnosis was observed to be absent, it could be assumed that all of its descendants were also absent (e.g. if a patient did not have malignant hypertension). Unfortunately, ICD-9 code observations never include observations of disease absence. ICD-9 codes are only documented when the condition is observed to be present. Additionally, ICD-9 codes are known to have relatively low sensitivity: conditions that are present are often not documented in a set of ICD-9 codes (Surjan, 1999). Given these facts, we made the following assumptions regarding each visit: recorded diagnoses and their ancestors were labeled as true; diagnoses that were observed at some time for a patient but not at the current visit were labeled as unobserved; and diagnoses that had never been listed for a patient were labeled as false for all of that patient’s visits. This last assumption captures the belief that parts of the ICD-9 hierarchy that are never observed for a particular patient are almost certain to be false. Additionally, for computational purposes, we decided not to include an ICD-9 code at all if neither it nor one of its descendants had been assigned to a patient in any of the records in our dataset.

4.3 ICD-9 Code Hierarchy

Here, we augment the sLDA model such that the supervised signal is distribution over the ICD-9 code tree, which is an is-a hierarchy. An is-a hierarchy is represented by the tree data structure where each node has only a single parent and nodes cannot be parents of ancestors (i.e. there are no loops). In this particular case, the ICD-9 code hierarchy is also partially a prefix tree where the labels for certain nodes are prefixes for child nodes. Given this rule does not apply to all nodes in the hierarchy, we did not use this feature to determine the structure of the hierarchy. Instead we acquired a dataset that explicitly defined the relationships between the nodes of the hierarchy [3]. In documentation of ICD-9 codes for billing purposes, only a subset of the nodes can be used, however the nodes higher in the hierarchy contain semantic information about the categories of codes that are their descendants. For this reason, we included these nodes in our model.

4.4 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

This model, supervised latent dirichlet allocation (sLDA), builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

4.5 Evaluation

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

6.1 Rao-Blackwellization

To derive the Gibbs sampler in general, we integrate over the parameters θ and $\phi_{1:K}$ resulting in the collapsed joint distribution shown in Equations 5-8.

$$p(\theta, z_{1:N} | w_{1:N}, y_{1:N}, \phi_{1:K}, \alpha, \beta, \mu, \xi) = \frac{p(\theta | \alpha)}{\int_0^\infty p(\theta | \alpha) d\theta} \frac{p(z_{1:N} | \theta, \phi_{1:K}, \xi)}{\int_0^\infty p(z_{1:N} | \theta, \phi_{1:K}, \xi) dz_{1:N}} \frac{p(w_{1:N} | z_{1:N}, \phi_{1:K}, \xi)}{\int_0^\infty p(w_{1:N} | z_{1:N}, \phi_{1:K}, \xi) dw_{1:N}} \frac{p(y_{1:N} | w_{1:N}, \phi_{1:K}, \xi)}{\int_0^\infty p(y_{1:N} | w_{1:N}, \phi_{1:K}, \xi) dy_{1:N}} \frac{\prod_{l=1}^L p(\phi_{l:K} | \theta, \alpha_l, \beta_l, \mu_l)}{\prod_{l=1}^L \int_0^\infty p(\phi_{l:K} | \theta, \alpha_l, \beta_l$$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We describe a family of probit regressors for conditionally dependent label hierarchy and to train Dirichlet allocation (LDA). We find that the addition of labels to the document frequency matrix, however, only constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks and show both improved label prediction performance and show evidence that the learned topic model improves as a result of using this signal too.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [1]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document “supervision,” often taking the form of a single numerical or categorical “label.” More generally this “supervision” can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generated/provided as inputs to an inferred topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [1].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retailers who sell multiple products under a single category and frequently sell products they sell. The structure of each product in a product hierarchy often contains multiple, hierarchical labelings of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simple unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section ?? we introduce hierarchically supervised LDA (HSLDA). In Section 4.4 we review related work, and in Section ?? we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Model

We define here a hierarchically supervised LDA model. We assume a pre-specified set of labels \mathcal{L} . Each document is assigned a response of either -1 or 1 for at least one, but potentially many, labels in \mathcal{L} . A response of 1 or -1 to label l indicates if the document respectively is or is not l . The label l for a document d will be used interchangeably to refer the observed response of document d to label l . The label set is assumed to be an “is-a” hierarchy. This means that if a label l_1 is a parent of label l_2 in the hierarchy and document d has a positive response to the label l_2 then document d also has a positive response to the label l_1 . Seen in a negative light, if document d has a negative response to the l_1 label then document d also has a negative response to the l_2 label. To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models. Approximate inference is performed using Gibbs sampling.

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , and the standard deviation σ used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

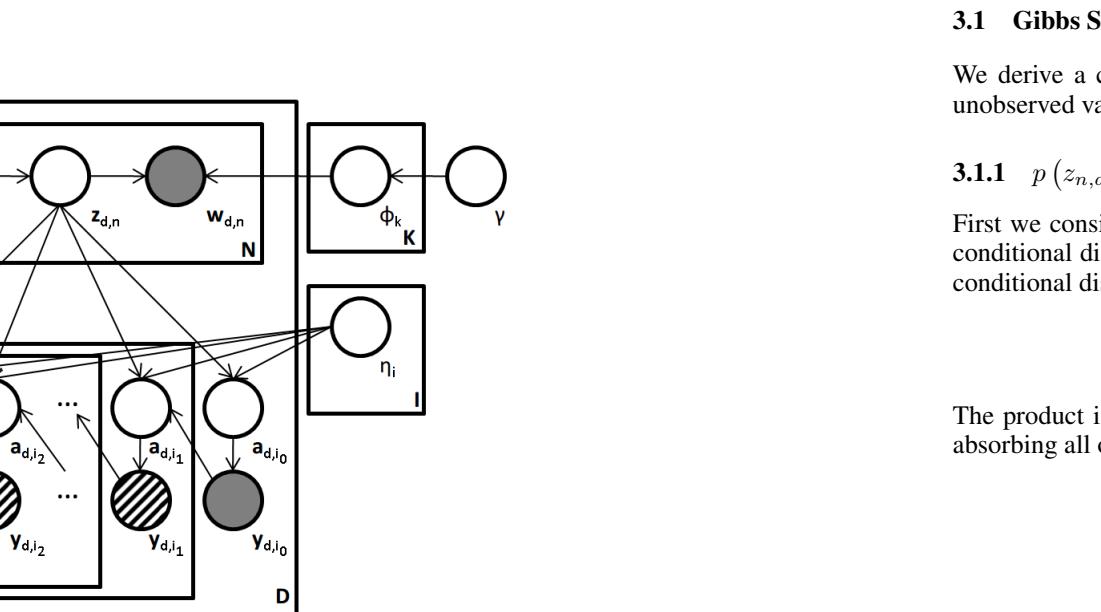


Figure 1: adapted sLDA model

1. For each topic $k = 1, \dots, K$:

- Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1})$, where $\mathbf{1}$ is a vector of ones of length V

2. For each label $l \in \mathcal{L}$:

- Draw the regression coefficients $\eta_l | \sigma \sim \mathcal{N}_K(-1, \sigma I_K)$, where I_K is the K dimensional identity matrix
- 3. Draw a prior over topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha')$

4. For each document $d = 1, \dots, D$:

- Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
- For $n = 1, \dots, N_d$:
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{1,K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
- For each label $l \in \mathcal{L}$:
 - Draw $a_{l,d} | z_{1:N_d, d}, \eta_l, y_{l, \text{parent}(l), d} \sim \begin{cases} \mathcal{N}(\bar{z}^T \eta_l, 1), & y_{l, \text{parent}(l), d} = 1 \\ \mathcal{N}(\bar{z}^T \eta_l, 1) I(a_{l,d} < 0), & y_{l, \text{parent}(l), d} = -1 \end{cases}$ where $\bar{z}_d = N_d^{-1} \sum_{n=1}^{N_d} z_{n,d}$
 - Set the response variable $y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{l, \text{parent}(l), d} = 1 \\ -1 & \text{otherwise} \end{cases}$

$$y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{l, \text{parent}(l), d} = 1 \\ -1 & \text{otherwise} \end{cases}$$

This type of generative model is known as a probit regression model. Probit regression is like logistic regression except instead of modeling the logit of $P(y=1)$ using a linear form we model $\Phi^{-1}(P(y=1))$ using a linear form, where $\Phi(\cdot)$ is the CDF for a standard normal distribution. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because the are introduced to make exact Gibbs sampling possible and are not of primary interest.

3 Inference

In the Bayesian approach to statistical modeling the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of \mathcal{L} that is fully observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{l,d}$ for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\beta_{1,K}$ and θ_d in $\text{Dir}_K(\cdot)$ and $\text{Dir}_V(\cdot)$ respectively. The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of tractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior. In particular, we will derive a collapsed Gibbs sampler.

1

2

4.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a controlled terminology which is the international standard diagnostic classification for epidemiological, health management, and clinical purposes. The ICD-9 codes are organized in a root-to-leaf structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing “Pneumonia due to adenovirus” is a child of the code representing “Viral pneumonia” where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifier) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [citation]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free texts. Aside from topic, one of the goals is to compare the sensitivity of predictions from the HSLDA model in comparison to the codes in a case where a test closer to ground truth is available. For this we will compare whether predictions for the ICD-9 code associated with anemia are better predicted by HSLDA or LDA given the ICD-9 codes. Anemia was chosen because hemoglobin values are readily available and the definition of anemia according the World Health Organization is approximately 12.5, with a threshold of 12 for women and 13 for men [1].

Here, $n_{v,d}^{k,-(n,d')}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d' . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{n,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

4.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset and partially directly from the Amazon.com website. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, “DVD / Games / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine”. Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords. This is a standard result from normal Bayesian linear regression [?]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \bar{z}_d , and $\mathbf{a}_t = [a_{t,1}, a_{t,2}, \dots, a_{t,D}]^T$.

3.1.3 $p(\eta_l | \mathbf{z}, \mathbf{Y}, \eta)$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} | \mathbf{z}_{1:N_d, d}, \eta_l, y_{l, \text{parent}(l), d}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ \frac{-1}{2} (a_{l,d} - \eta_l^T \bar{z}_d) \right\} I(a_{l,d} y_{l, \text{parent}(l), d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta | \mathbf{z}, \mathbf{a}, \alpha)$

The two main methods of evaluation for the model are prediction and topic quality. We compare model performance against 3 similar models to demonstrate that each component of the model is important for performance. Specifically, we evaluate models including independent regressors + sLDA (hierarchical constraints on labels ignored), HSLDA fit by running LDA first then running tree-conditional regressions, and HSLDA fit with fixed random regression parameters.

4.3 Evaluation

The two main methods of evaluation for the model are prediction and topic quality. We compare model performance against 3 similar models to demonstrate that each component of the model is important for performance. Specifically, we evaluate models including independent regressors + sLDA (hierarchical constraints on labels ignored), HSLDA fit by running LDA first then running tree-conditional regressions, and HSLDA fit with fixed random regression parameters.

4.3.1 Prediction

The two measures for predictive performance used here include the true positive rate and the false positive rate. We evaluate model performance on held out data. A more ideal evaluation of performance would include a manually labeled hierarchy since it is well known that ICD-9 codes have a relatively low sensitivity.

Posterior inference is performed using the “direct assignment” method of Teh et al. [17].

$$\beta \sim \text{Dir}(m_{(1,1)}, m_{(1,2)}, \dots, m_{(1,K)}) \quad (6)$$

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α , α' , and γ are given broad Gamma(2, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Experiments

In the section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of someone who has been hospitalized. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

3

4.4 Related Work

Latent Dirichlet Allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In

other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with each document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA based on least squares regression [5].

There have been many models that incorporate both latent models of text and some form of supervision [citations]. One set of models that are particularly relevant to HSLDA are Chang and Blei's [citation] hierarchical models for document networks. In that family of models, they encountered a similar scenario where the lack of a link did not truly indicate absence. Therefore, as in the work of Chang and Blei, we employ regularization in the form of a negative prior on the regression parameters to provide a prior that indicates a bias towards being truly negative in the absence of a code.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [16, 10, 15, 17], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few specific diseases [14] but most recent and promising work on the subject was inspired by the 2007 Medical NLP Challenge: “Identifying Clinical Terms Classifying Clinical Free Text Using Natural Language Processing” [18]. The data set given to the participants consisted of documents that were 1-2 lines each and all of the documents were radiology reports – clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Li et al publication [13]. Li proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

5 Results

6 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...

• what about the nonparametric version of this?

• discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

References

[1] Amazon, Inc. <http://www.amazon.com/>, 2011.

[2] DMOZ open directory project, <http://www.dmoz.org/>, 2002.

[3] Stanford network analysis platform, <http://snap.stanford.edu/>, 2004.

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer hierarchical product catalogs [1] as well as medical records [3] and patient discharge summaries [4]. Clinical Modification Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [5]). In this work we show how to combine two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [4] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [5] augmented with document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 4.4 we review related work. In Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section ?? we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [5].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [4].

1

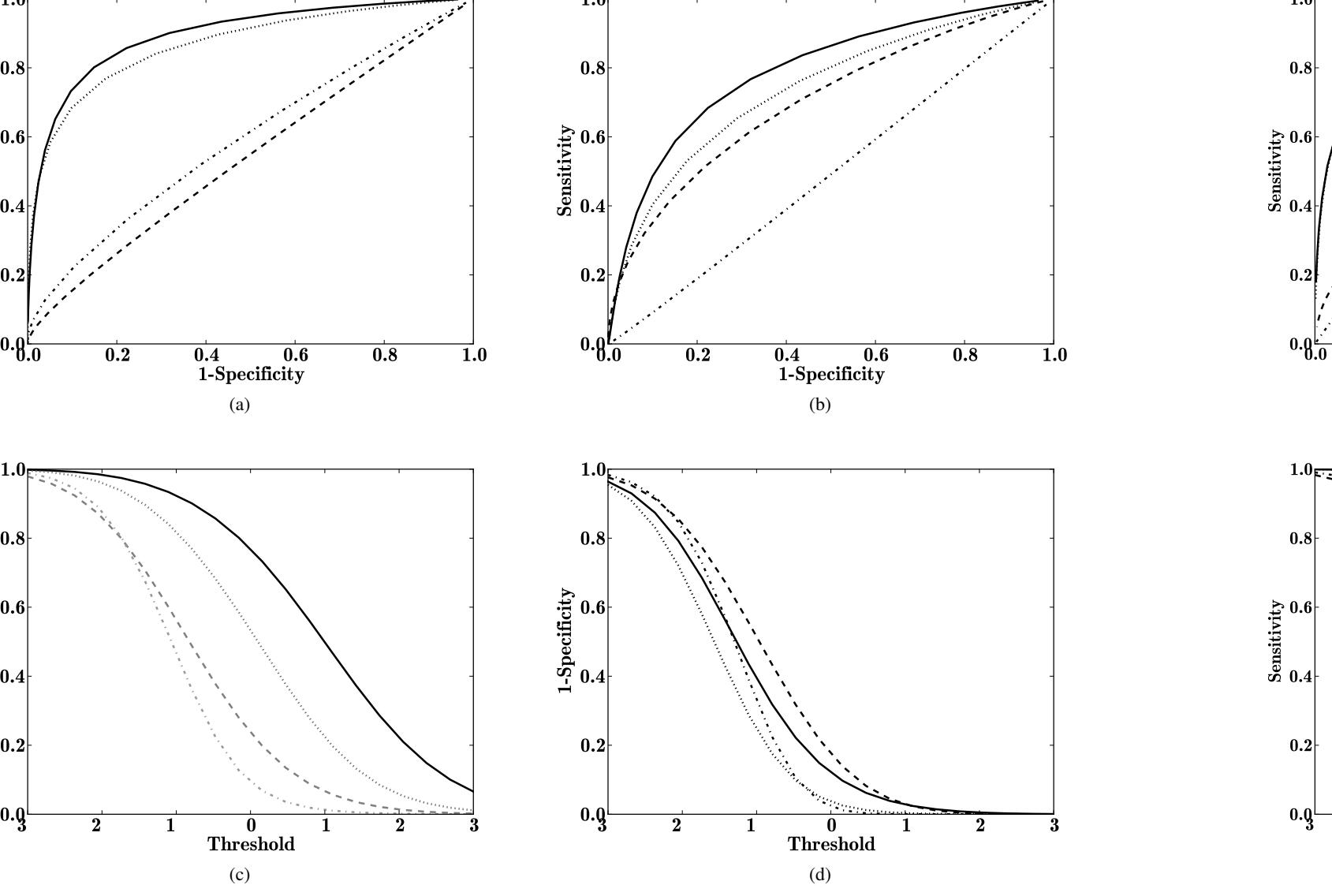


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [2 ? ?]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models, they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 0 or 1 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_1 then it will also have a positive response to label l_2 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

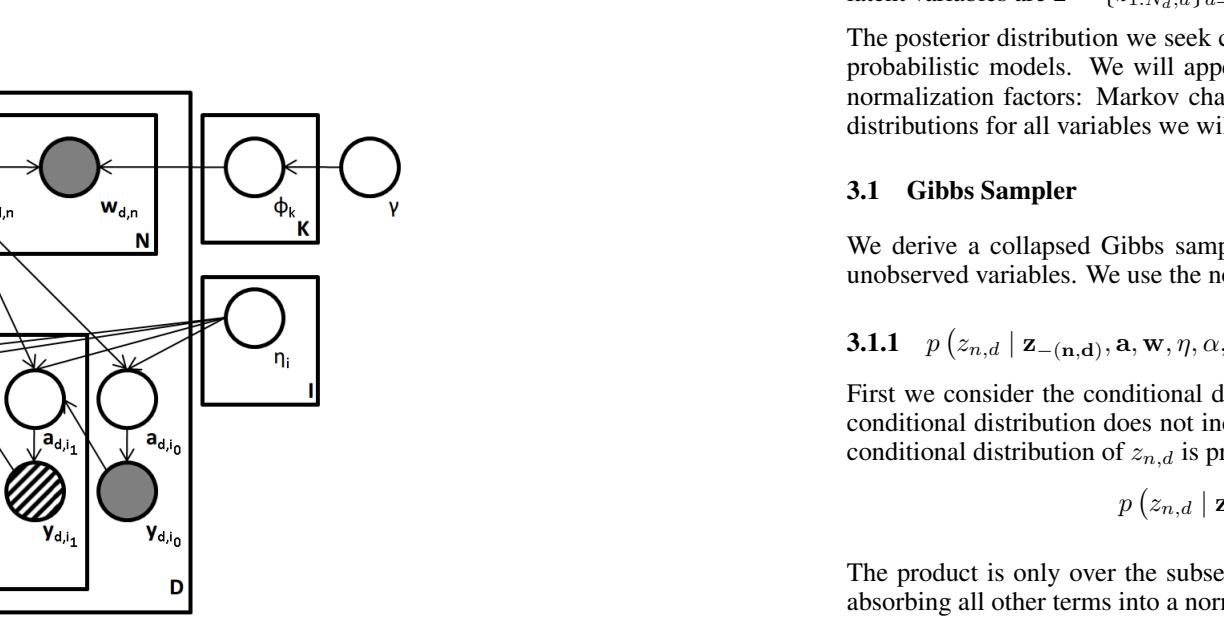


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$:
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where γ is a vector of ones of length V
2. For each label $l \in \mathcal{L}$:
 - Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_I, \sigma I_K)$, where I_K is the K dimensional identity matrix
 - 3. Draw a prior over topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha')$
 - 4. For each document $d = 1, \dots, D$:
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

2

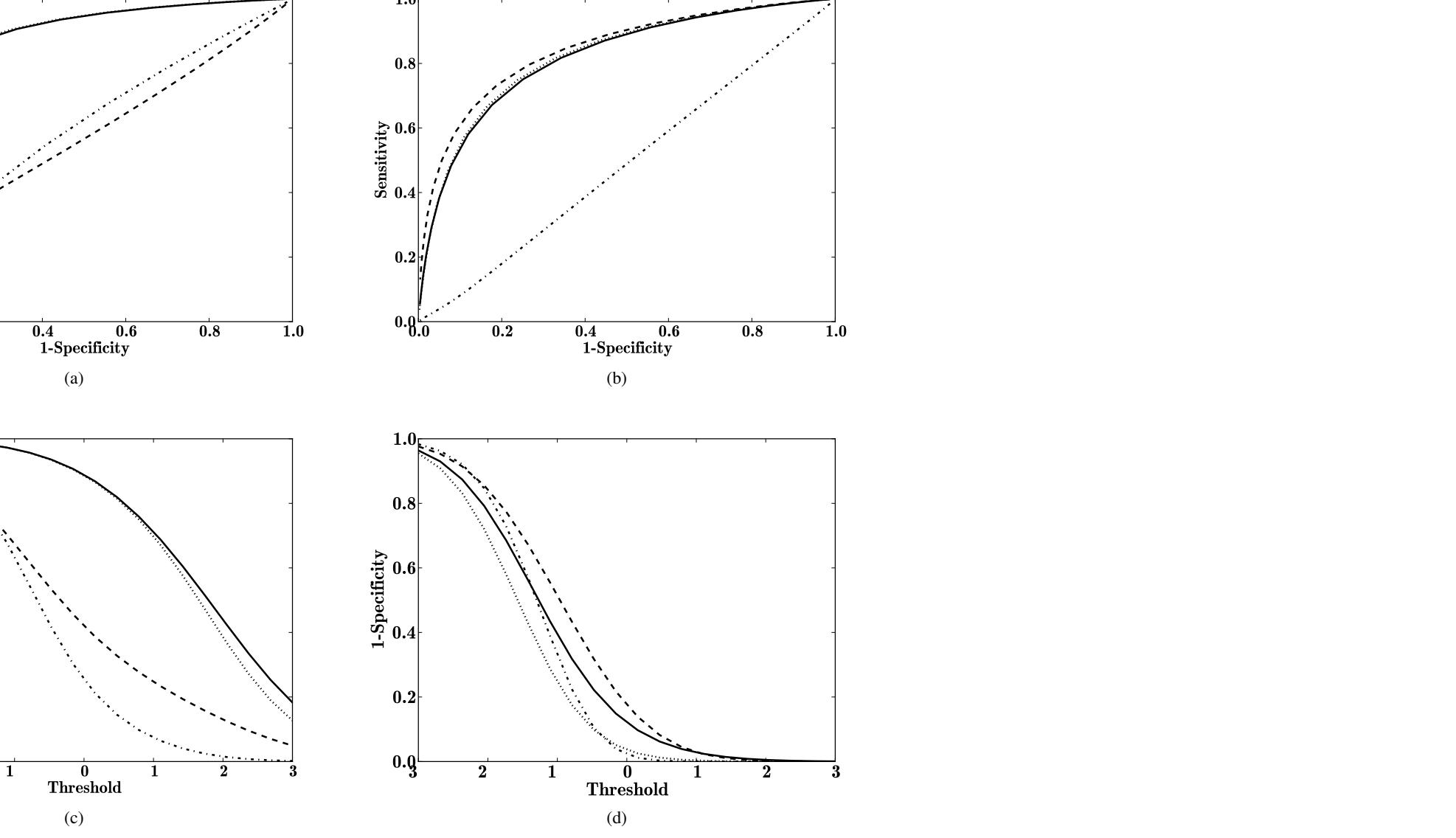


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.1.3 $p(a_{l,d} \mid a_l, \eta_l, \mathbf{z}, \mathbf{Y}, \eta)$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^T \mathbf{z}_{l,d}) \right\} I(a_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta \mid \mathbf{z}, \alpha', \alpha)$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [17]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [16].

$$\beta \sim \text{Dir}(m_{(j),1} + \alpha', m_{(j),2} + \alpha', \dots, m_{(j),K} + \alpha') \quad (6)$$

$$p(m_{d,k} \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m_k) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Experiments

In the section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

4.1 Data and Pre-Processing

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z} \setminus z_{n,d}$.

3.1.1 $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include $\theta_{1,D}$ and $\theta_{1,k}$ because they have been integrated out as in the collapsed Gibbs sampler [11]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \eta) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [11] we find

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} \frac{n_{k,(n,d)}^{a_{l,d} - (n_{k,(n,d)} + \gamma)}}{(n_{k,(n,d)} + \alpha \beta_k)^{\frac{n_{k,(n,d)} + \gamma}{2}}} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(z_{n,d}^T \eta - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $n_{k,(n,d)}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, we know that $p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\eta \mid \mathbf{z}, \mathbf{a}, \sigma)$

We now consider the conditional distribution of the regression coefficients η_l for $l \in \mathcal{L}$. Given that η_l and $a_{l,d}$ are distributed normally, the posterior distribution of η_l is normally distributed with mean μ and covariance Σ such that

$$\Sigma^{-1} = \mathbf{I} \sigma^{-2} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\mu = \hat{\Sigma} \left(\mathbf{I} \frac{\mu}{\sigma^2} + \mathbf{Z}^T \mathbf{a}_l \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [2]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is $\mathbf{z}_{l,d}$, and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

3

[15] P. Ruch, J. Gobell, I. Thirumalai, and A. Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.

[16] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.

[17] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking lda: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems* 22, pages 1973–1981, 2009.

4

[18] T.L. Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.

[19] LV Liita, S. Nicaise, and J. Bi. Large scale diagnostic code classification for medical patient records, 2008.

[20] RB Rao, S Sandhya, RS Niculescu, C Ganesan, and H Rao. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 416–425, 2003.

[21] T.L. Griffiths and M. Steyvers. Automatic text categorization for medical reports. *AMIA Annual Symposium Proceedings*, 2007:279–280, 2007.

[22] K. Crummer, M. Dredze, P. Gauchi, P.F. Talukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.

[23] R. Faria and G. Szarvas. Automatic construction of rule-based icd-9-cm coding system. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.

[24] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1523-4433.

[25] P. Brown, DG Cockayne, JR Allegre, and T. Madala. The ngram cc classifier: A novel method of automatically creating cc classifiers based on icd9 groupings. *Advances in Disease Surveillance*, 1:30, 2006.

[26] H. Heffernan, B. Ribeiro-Neto, RF Vale, AHF Laender, and LRS De Lima. Category-driven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.

[27] I. Galustyan, A. Aramyan, and Ö.Uzuner. Their approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279–280, 2007.

[28] TL Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.

[29] LV Liita, S. Nicaise, and J. Bi. Large scale diagnostic code classification for medical patient records, 2008.

[30] RB Rao, S Sandhya, RS Niculescu, C Ganesan, and H

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions and catalogs (e.g. [1]) as well as from [3] product manual discharge summaries and ICD-9-CM codes applied to them. In particular, we compare discharge records vs. International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned ([2]). In this work we show how to combine two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [4] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [5] augmented with a form of document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [5].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [4].

1

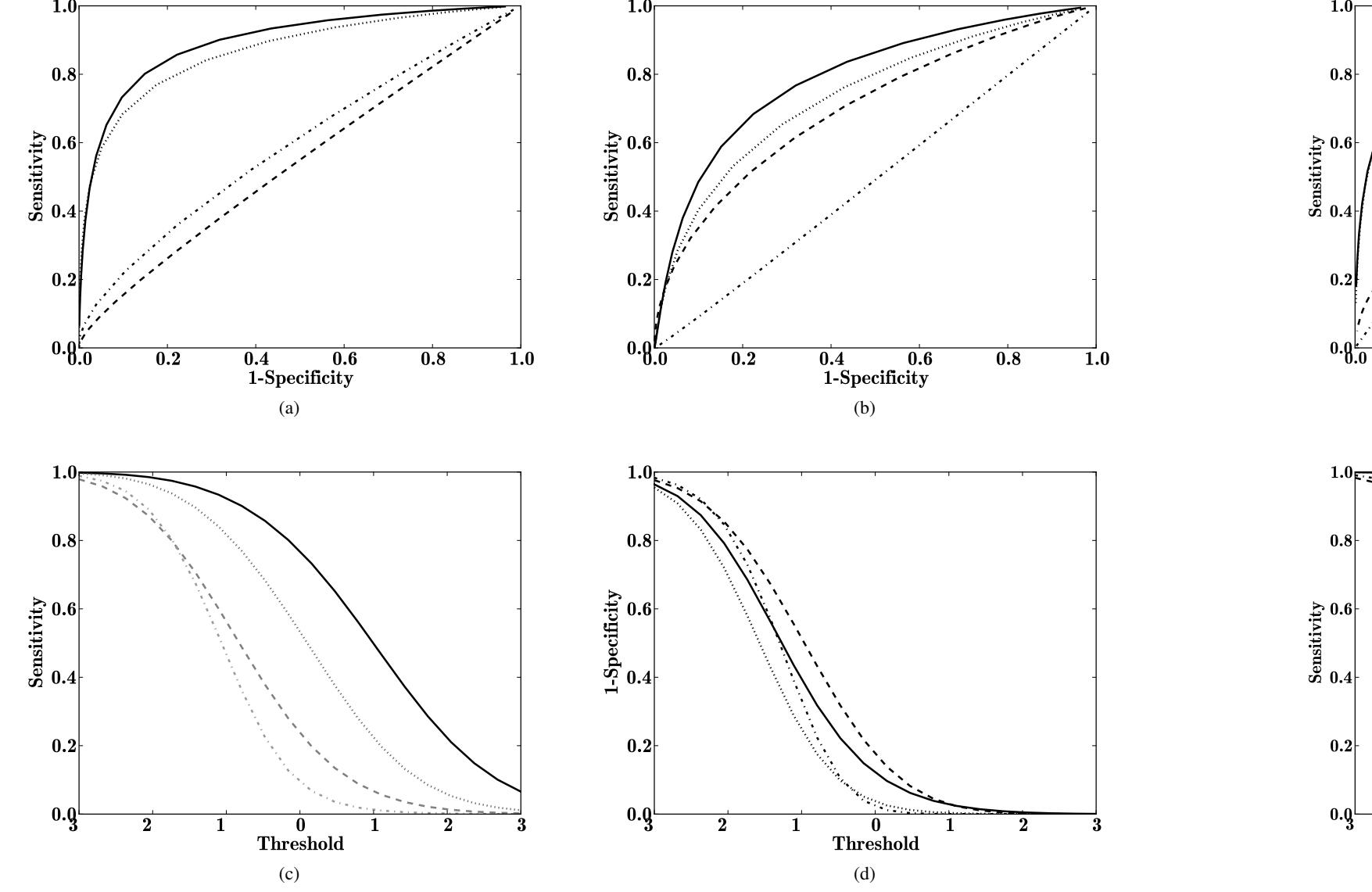


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [2, 3, 7, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models, they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [15, 9, 14, 6], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [13] but the most recent and promising work is done on large-scale datasets, such as the Stanford Network Analysis Platform (SNAP) [1] and the Amazon.com product "Natural Language Processing" website. Most of the classification strategies included word matching and rule-based algorithms [10, 7, 8]. The data set given to the participants consisted only of documents that were 1-2 lines each and all of the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes that could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [12]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or 0 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

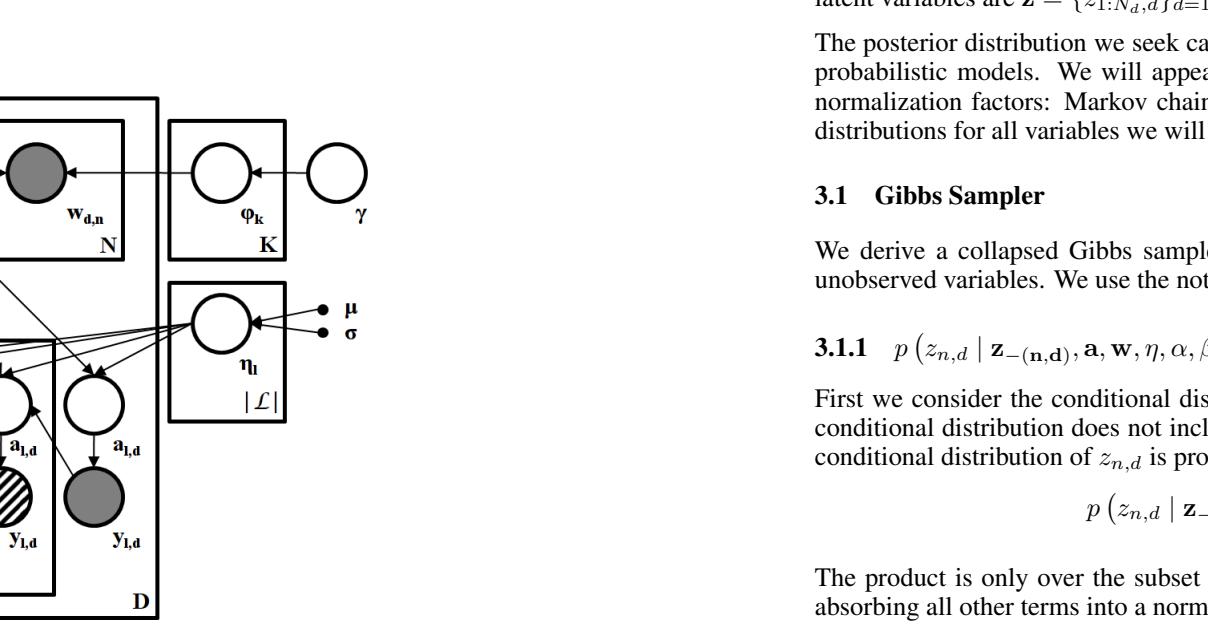


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

- For each topic $k = 1, \dots, K$:
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V
- For each label $l \in \mathcal{L}$:
 - Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix
 - Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$
 - For each document $d = 1, \dots, D$:
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

• For each document $d = 1, \dots, D$:

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha)$

• Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V

• Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_l, \sigma I_K)$, where I_K is the K dimensional identity matrix

• Draw a prior over topic proportions $\beta \mid \sigma \sim \text{Dir}_K(\alpha')$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
a jp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions in catalogs [11] as well as from [4], patient medical records versus International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with a form of document "supervision"; often taking the form of a single numerical or categorical "label." More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generated model as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

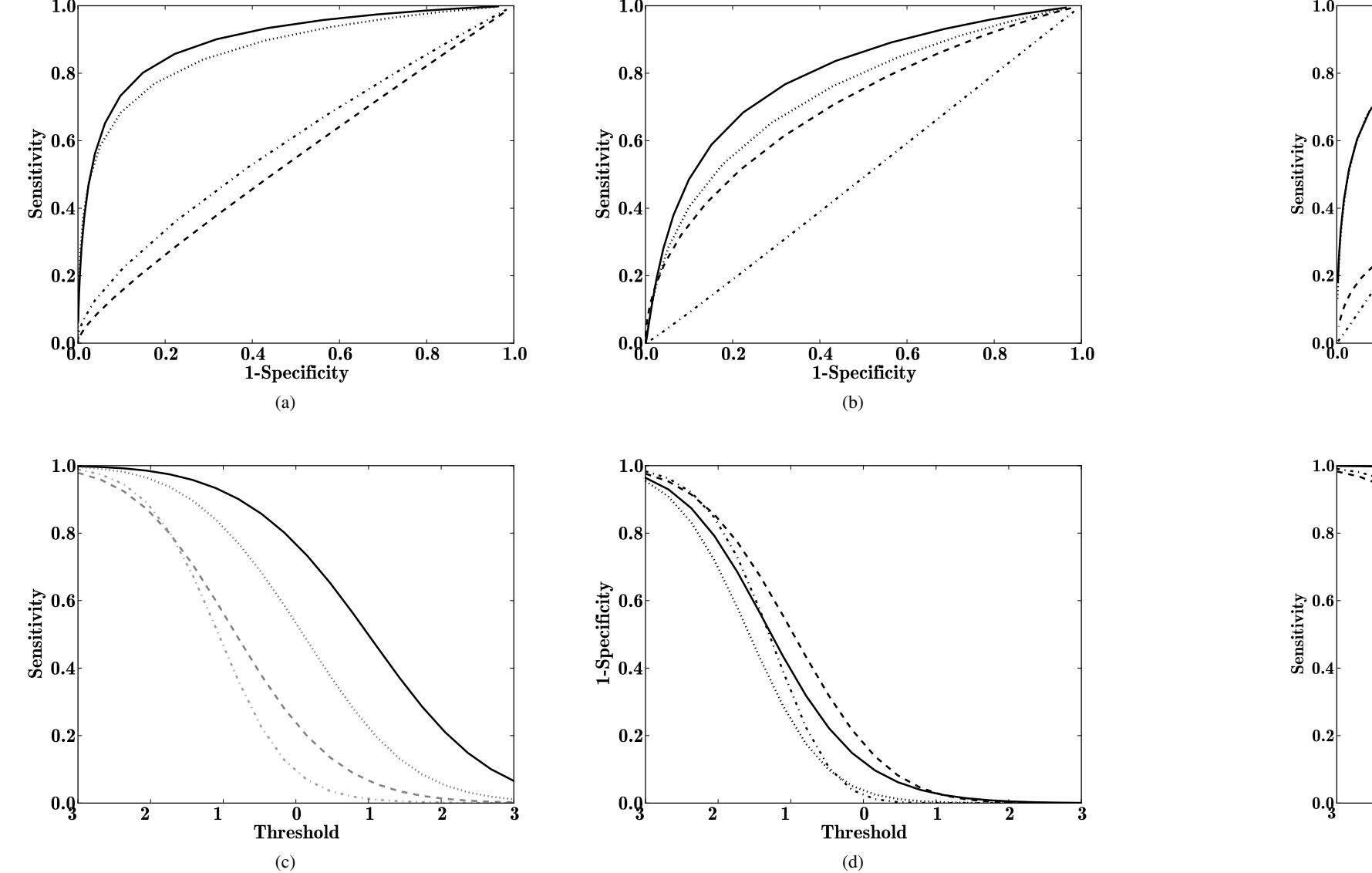


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 23, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 12, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work has come from the medical domain [1, 2, 11, 12, 13, 14, 15, 16]. Chang and Blei's work on "Hierarchical Topic Modeling Using Natural Language Processing" (website). Most of the classification strategies included word matching and rule-based algorithms. [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all of the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or 0 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_1 then it will also have a positive response to label l_2 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

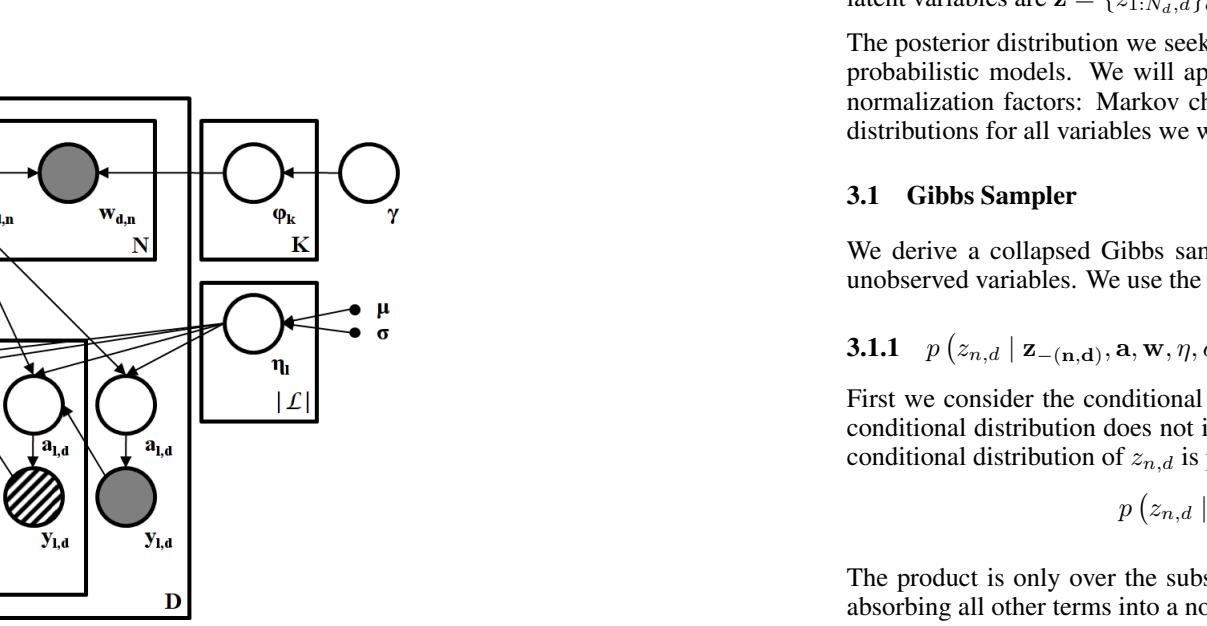


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$:
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V
2. For each label $l \in \mathcal{L}$:
 - Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_I, \sigma I_K)$, where I_K is the K dimensional identity matrix
 - 3. Draw a prior over topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha')$
 - 4. For each document $d = 1, \dots, D$:
 - Draw topic proportions $\theta_d \mid \beta, \alpha' \sim \text{Dir}_K(\alpha')$

1. For each topic $k = 1, \dots, K$:
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma)$, where 1 is a vector of ones of length V
2. For each label $l \in \mathcal{L}$:
 - Draw the regression coefficients $\eta_l \mid \sigma \sim \mathcal{N}_K(\mu_I, \sigma I_K)$, where I_K is the K dimensional identity matrix
 - 3. Draw a prior over topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha')$
 - 4. For each document $d = 1, \dots, D$:
 - Draw topic proportions $\theta_d \mid \beta, \alpha' \sim \text{Dir}_K(\alpha')$

We will now consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include θ_{d,N_d} and η_{l,N_d} because they have been integrated out as in the collapsed Gibbs sampler [14]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \eta) p(z_{n,d} \mid z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [14] we find

$$p(z_{n,d} = k \mid z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \frac{n_{k,n}^{k-n_{n,d}} + \gamma}{(n_{\cdot,n}^{k-n_{n,d}} + \alpha \beta_k)} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(z_{n,d}^T \eta_l - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $n_{v,n}^{k-n_{n,d}}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{d,n} \mid z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.1 $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include θ_{d,N_d} and η_{l,N_d} because they have been integrated out as in the collapsed Gibbs sampler [14]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \eta) p(z_{n,d} \mid z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [14] we find

$$p(z_{n,d} = k \mid z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \frac{n_{k,n}^{k-n_{n,d}} + \gamma}{(n_{\cdot,n}^{k-n_{n,d}} + \alpha \beta_k)} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(z_{n,d}^T \eta_l - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $n_{v,n}^{k-n_{n,d}}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{d,n} \mid z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\eta \mid \mathbf{z}, \alpha, \sigma)$

We now consider the conditional distribution of the regression coefficients η_l for $l \in \mathcal{L}$. Given that η_l and $a_{l,d}$ are distributed normally, the posterior distribution of η_l is normally distributed with mean μ and covariance Σ such that

$$\Sigma^{-1} = \mathbf{I} \sigma^{-2} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\mu = \hat{\Sigma} \left(\frac{\mu}{\sigma^2} \mathbf{Z}^T \mathbf{a}_d \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [2]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is z_d , and $\mathbf{a}_d = [a_{1,d}, a_{2,d}, \dots, a_{K,d}]^T$.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

4.1.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partially directly from the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine". Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

4

[17] Daniel Ramage, David Hall, Rameesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL http://www.aclweb.org/anthology/cn1909_248.pdf.

[18] R.F. Bell, S.Sandhu, B.S.Narasimha, C.Garroway, and H.R. Kottur. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on knowledge discovery and data mining*, pages 416–425, 2003.

[19] B.RibeiroNeto, AHF.Lacerda, and LRS.Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and technology*, 52(3):391–401, 2001.

[20] P.Bach, J.Gobbi, I.Thakriti, and A.Gruenbacher. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636–638, 2008.

[21] Y.W.Teh, M.I.Jordan, M.J.Beal, and D.M.Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.

[22] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.

[23] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

5

[11] R.Farkas and G.Szavars. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008. doi: 10.1186/1471-2105-9-S10-309.

[12] R.Freitas,Junior, B.RibeiroNeto, RF.Vale, AHF.Lacerda, and LRS.Lima. Categorizationdriven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.

[13] I.Goldstein, A.Arzamasyan, and O.Uzuner. Their approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.

[14] TL.Griffiths and M.Steyers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions and categories [11] as well as from [4], patient medical records versus International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with document "supervision"; often taking the form of a single numerical or categorical "label." More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generated modelled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured and supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

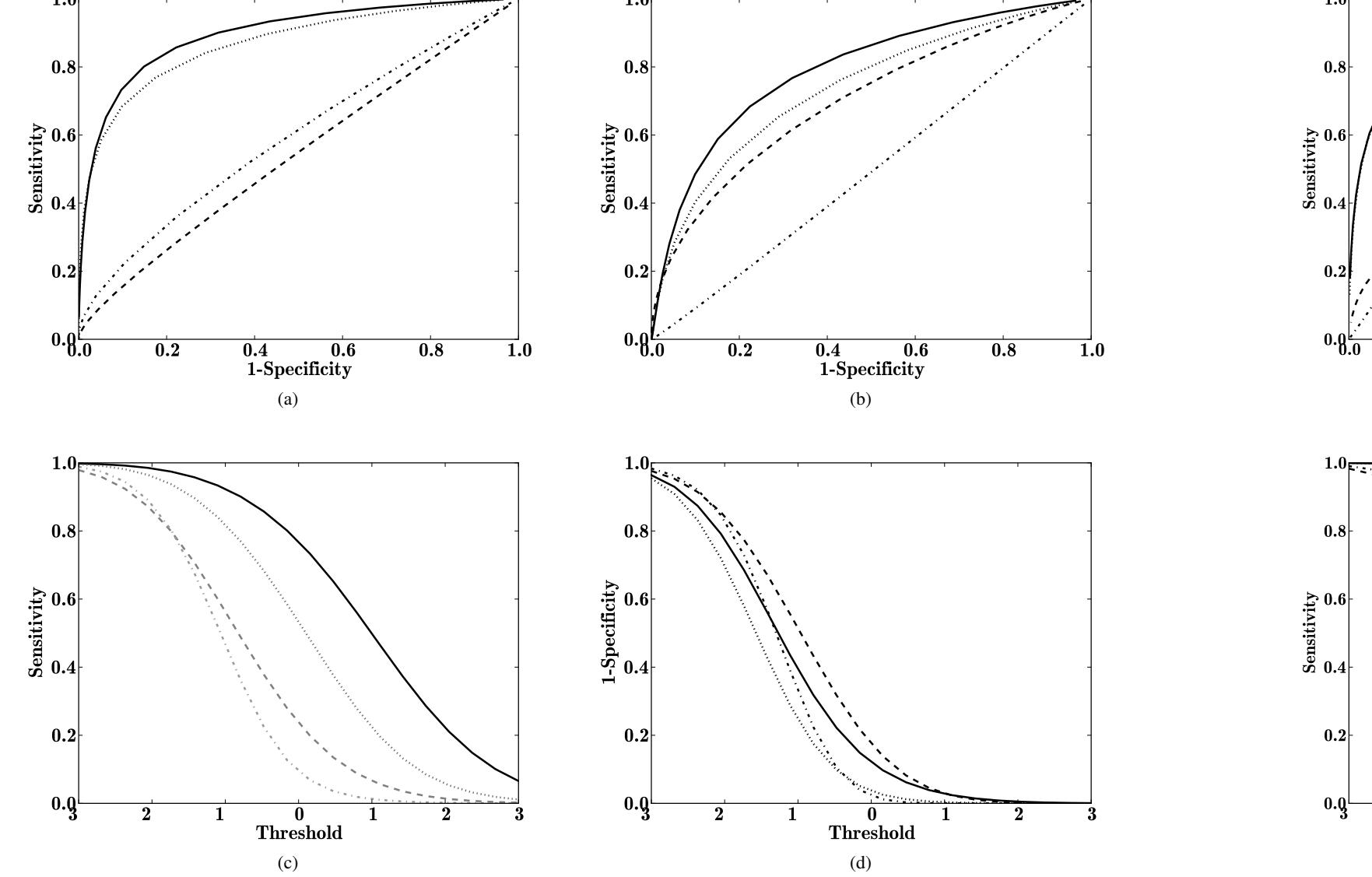


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 23, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 12, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work has come from the medical domain [19, 20, 21, 22, 23]. Chang and Blei's work is available on their website "Clinical Natural Language Processing" (website). Most of the classification strategies included word matching and rule-based algorithms. [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all of the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or 0 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

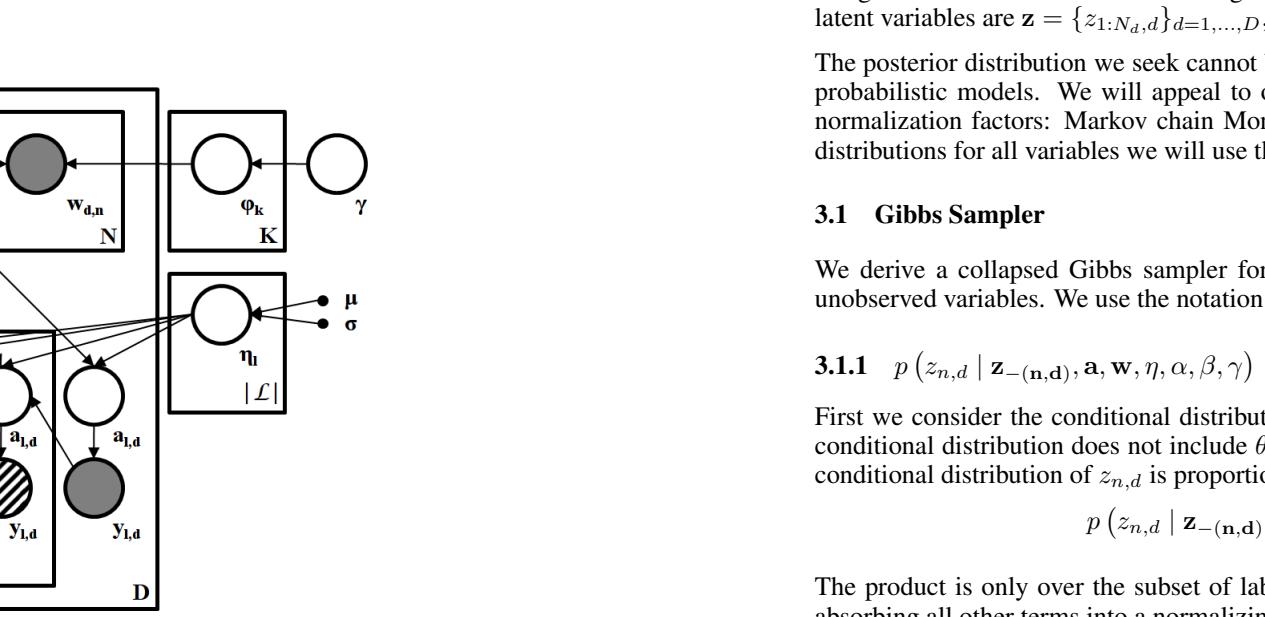


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$

2

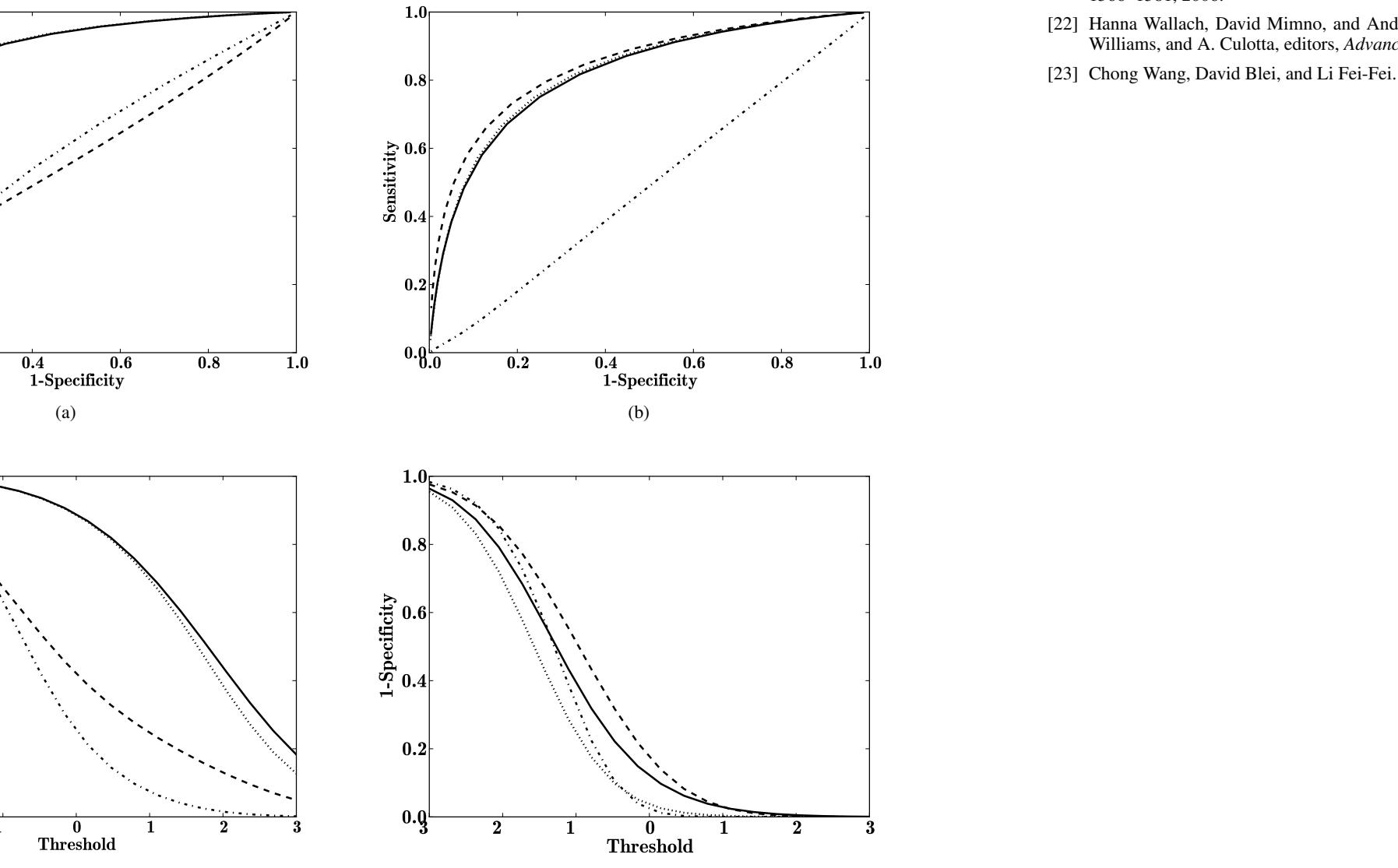


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

• For $n = 1, \dots, N_d$

- Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
- Draw word $w_{n,d} \mid z_{n,d}, \beta_{1:K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
- For each label $l \in \mathcal{L}$
 - Draw $a_{l,d} \mid \bar{\mathbf{z}}_d, \beta_l, y_{ps(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \beta_l, 1), & y_{ps(l)} = 1 \\ \mathcal{N}(\mathbf{z}^T \beta_l, 1) \mathbb{I}(a_{l,d} < 0), & y_{ps(l)} = -1 \end{cases}$
 - Set the response variable $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{ps(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

3.1.3 $p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta)$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta^T \bar{\mathbf{z}}_d)^2 \right\} I(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta \mid \mathbf{z}, \alpha', \alpha)$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [22]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [21].

$$\beta \sim \text{Dir}(m_{(.),1} + \alpha', m_{(.),2} + \alpha', \dots, m_{(.),K} + \alpha') \quad (6)$$

$$p(m_{d,k} \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m_{(.,k)})(\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Experiments

In the section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

4.1 Data and Pre-Processing

4.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary are part of a standard nomenclature which is the international standard diagnostic classification for epidemiological, statistical and administrative purposes. The ICD-9 codes are assigned in a hierarchical manner, where each code representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing "pneumonia due to adenovirus" is a child of the code representing "viral pneumonia" where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

4.1.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partially directly from the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine". Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

4.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including sLDA with independent regressions (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. We will refer to these three comparison models as the sLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen as to highlight performance in absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter for regression parameters.

4.3 Evaluation

For each dataset, a held out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held out set with their ancestors and considers all other non-existent labels to be negative. This uniform treatment of labels allows for a fair comparison of the models.

5 Results

6 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

References

- [1] Amazon, Inc. <http://www.amazon.com/2011>, 2011.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1533-4758.
- [7] P Brown, DG Cockrane, and JR Allegre. The ngram cc classifier: A novel method of automatically creating cc classifiers based on icd9 groupings. *Advances in Disease Surveillance*, 1:36, 2006.
- [8] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/AOS/1277944690.
- [9] K Cramer, M Dredze, K Ganesh, PF Tufekci, and S Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biomedical, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [10] Yan Yan David S. Nilasena Martha J. Radford Brian F. Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [11] R Farkas and G Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC Bioinformatics*, 9(Suppl 3):S

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions and categories [11] as well as from [4] medical product descriptions and International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggest that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

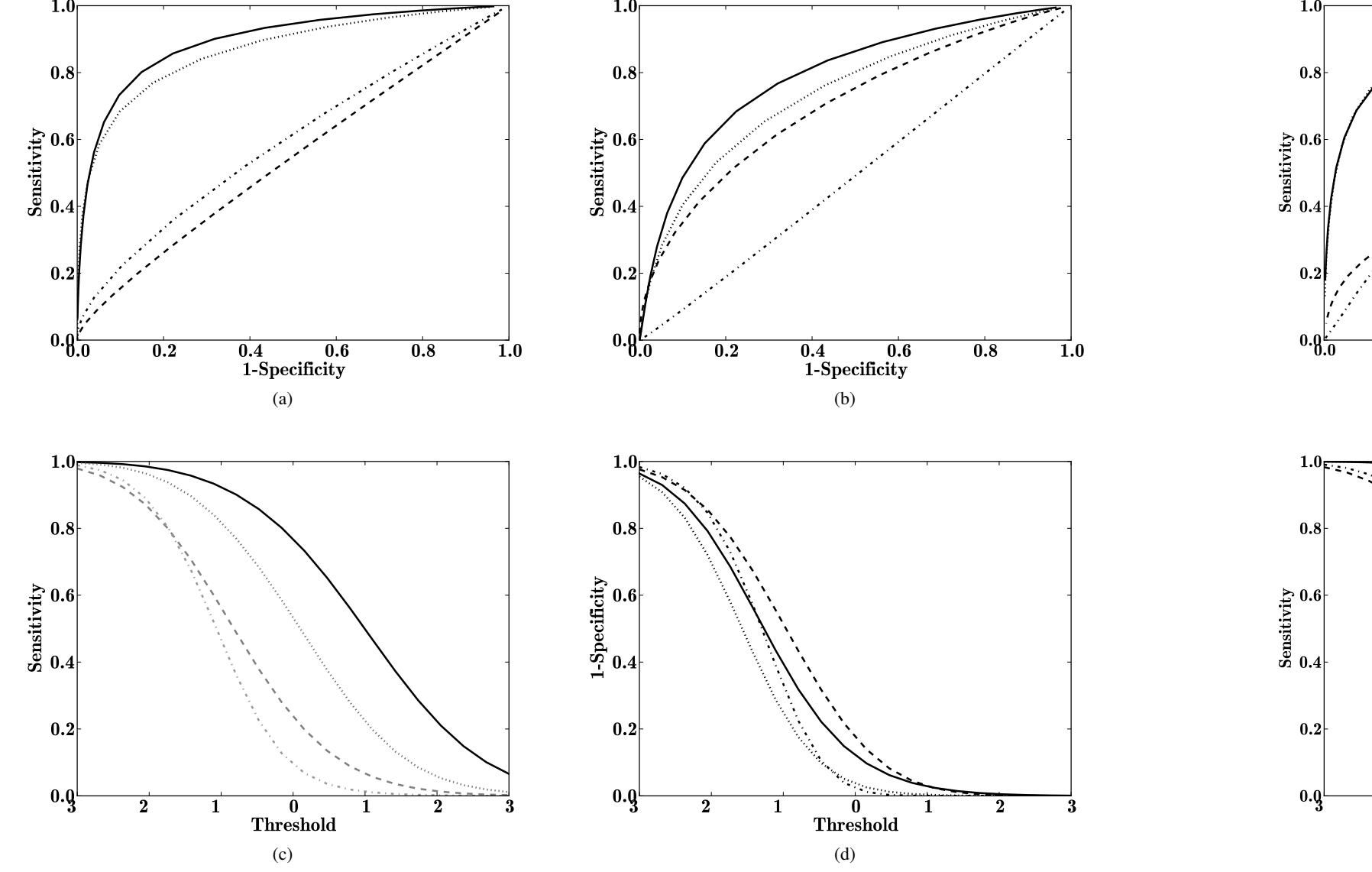


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

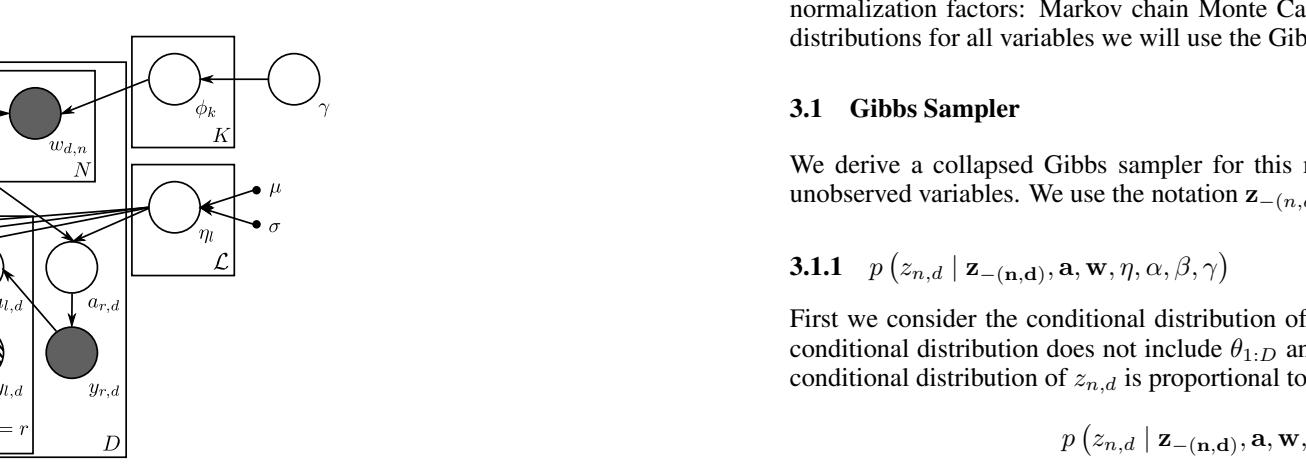
There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 23, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 12, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work has come from the medical domain. Chang and Blei [8] propose a hierarchical topic model for clinical text using "Natural Language Processing" (website). Most of the classification strategies included word matching and rule-based algorithms. [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.



Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions in catalogs [11] as well as from [4], patient medical records vs International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generated modelled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

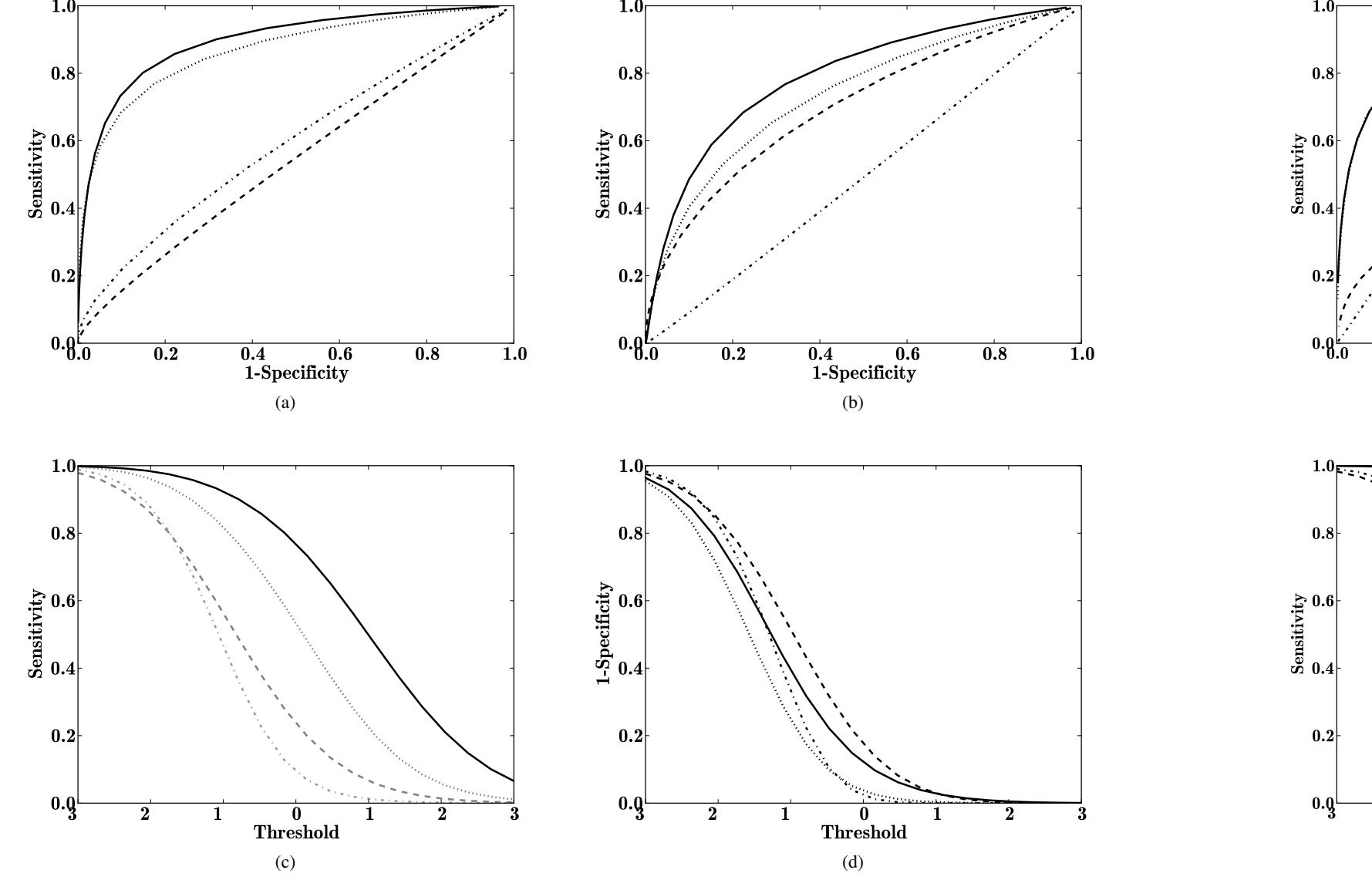


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 23, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 12, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work has come from larger datasets, such as the Stanford Network Analysis Platform (SNAP) [11] and the Clinical Record Interoperability Language Processing (website). Most of the classification strategies included word matching and rule-based algorithms, [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all of the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or 0 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

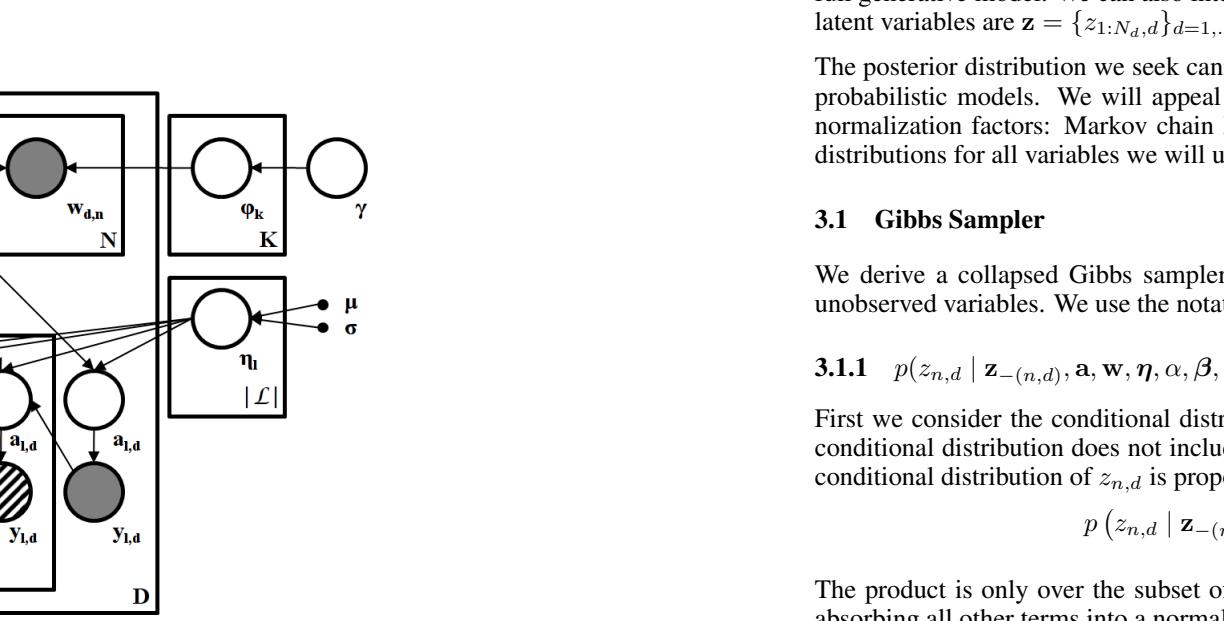


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$

2

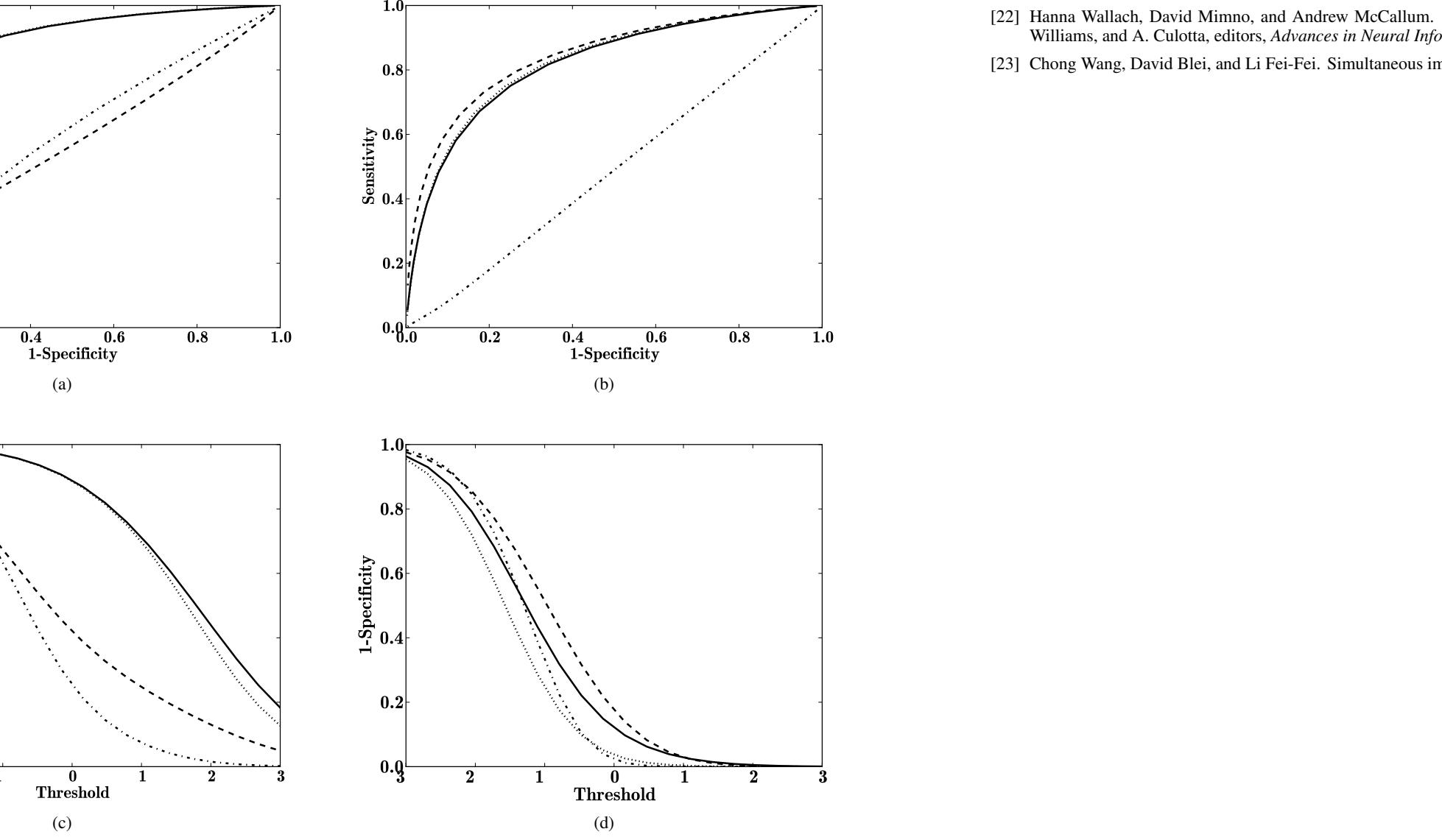


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.1.3 $p(a_{l,d} \mid a_l, \eta_l, \mathbf{z}_d, \mathbf{Y}, \eta)$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}_d, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^T \mathbf{z}_d) \right\} I(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta \mid \mathbf{z}, \alpha', \alpha)$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [22]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [21].

$$\beta \sim \text{Dir}(m_{(.,1)} + \alpha', m_{(.,2)} + \alpha', \dots, m_{(.,K)} + \alpha') \quad (6)$$

$$p(m_{d,k} \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Experiments

In the section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

4.1 Data and Pre-Processing

4.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary are part of a nomenclature terminology which is the international standard diagnostic classification for epidemiological, administrative, and clinical purposes. The ICD-9 codes are often assigned in a hierarchical manner, where each code representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing "pneumonia due to adenovirus" is a child of the code representing "viral pneumonia" where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

4.1.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partially directly from the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine". Each product was labeled with multiple categories.

This is a standard result from normal Bayesian linear regression [?]. Here, Z is a $D \times K$ matrix such that row d of Z is \mathbf{z}_d , and $\mathbf{a}_i = [a_{i,1}, a_{i,2}, \dots, a_{i,D}]^T$.

3

- [17] Daniel Ramage, David Hall, Ramesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL http://www.aclweb.org/anthology/cn1699510_1699543.pdf.

- [18] R.F. Bae, S.Sandhu, B.S. Balakrishnan, C.Gerrard, and H.R. Kothiyal. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on knowledge discovery and data mining*, pages 416–425, 2003.

- [19] B.BerissoNeto, AHF.Lacerda, and LRS.Da.Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for Information science and Technology*, 52(3):391–401, 2001.

- [20] P. Bush, J. Gobbi, I. Tahiri, and A. Grisolia. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636–638, 2008.

- [21] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.

- [22] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.

- [23] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

4

- [11] R Farkas and G Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008. doi: 10.1186/1471-2105-9-S10-S10.

- [12] DMOZ open directory project. <http://www.dmoz.org/>, 2002.

- [13] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.

- [14] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.

- [15] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.

- [16] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/AOS/1277940809.

- [17] K.Cramer, M.Dredze, K.Ganchev, PP.Talukdar, and S.Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.

- [18] I.Goldstein, A.Azurmanyan, and O.Uzuner. Their approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:27

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions and categories [11] as well as from [49], patient medical records vs International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applied in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label." More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggest that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

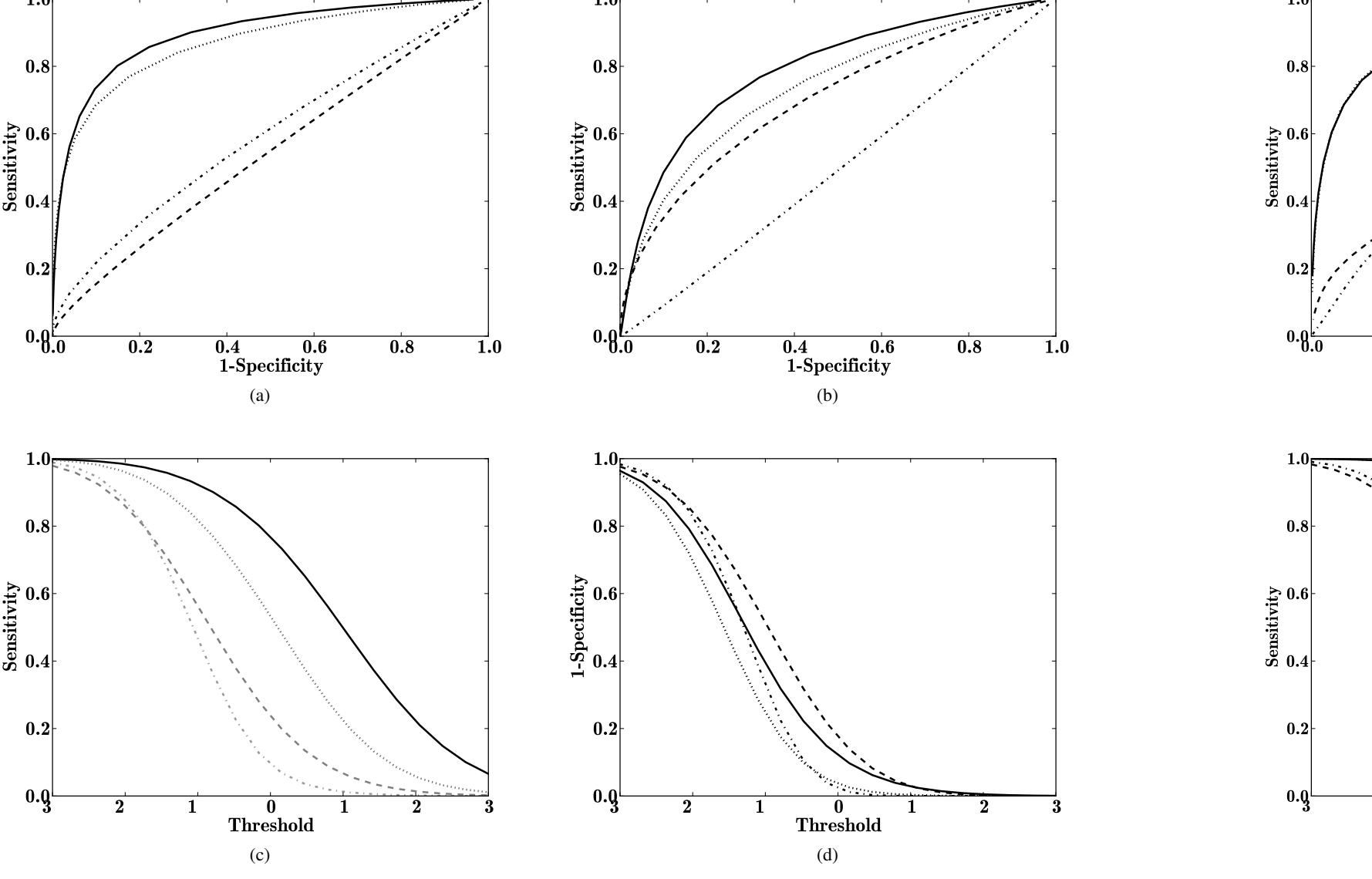


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 13, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work has come from the medical domain [11, 13, 49]. Chang and Blei's work [8] is one of the most prominent papers in this field, along with Natural Language Processing (website). Most of the classification strategies included word matching and rule-based algorithms. [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or at least one, but potentially many labels in \mathcal{L} . The label, l_i , for a document, d , will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

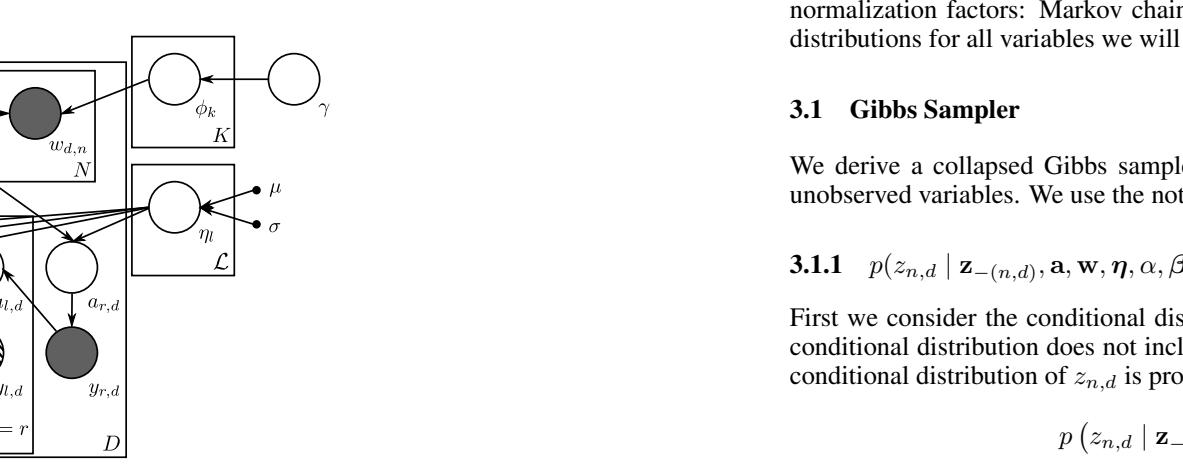


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{l,K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$:
 - Draw a distribution over words $\phi_l \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
 - Draw topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{l,K} \sim \text{Multinomial}(\beta_{z_{n,d}})$

2

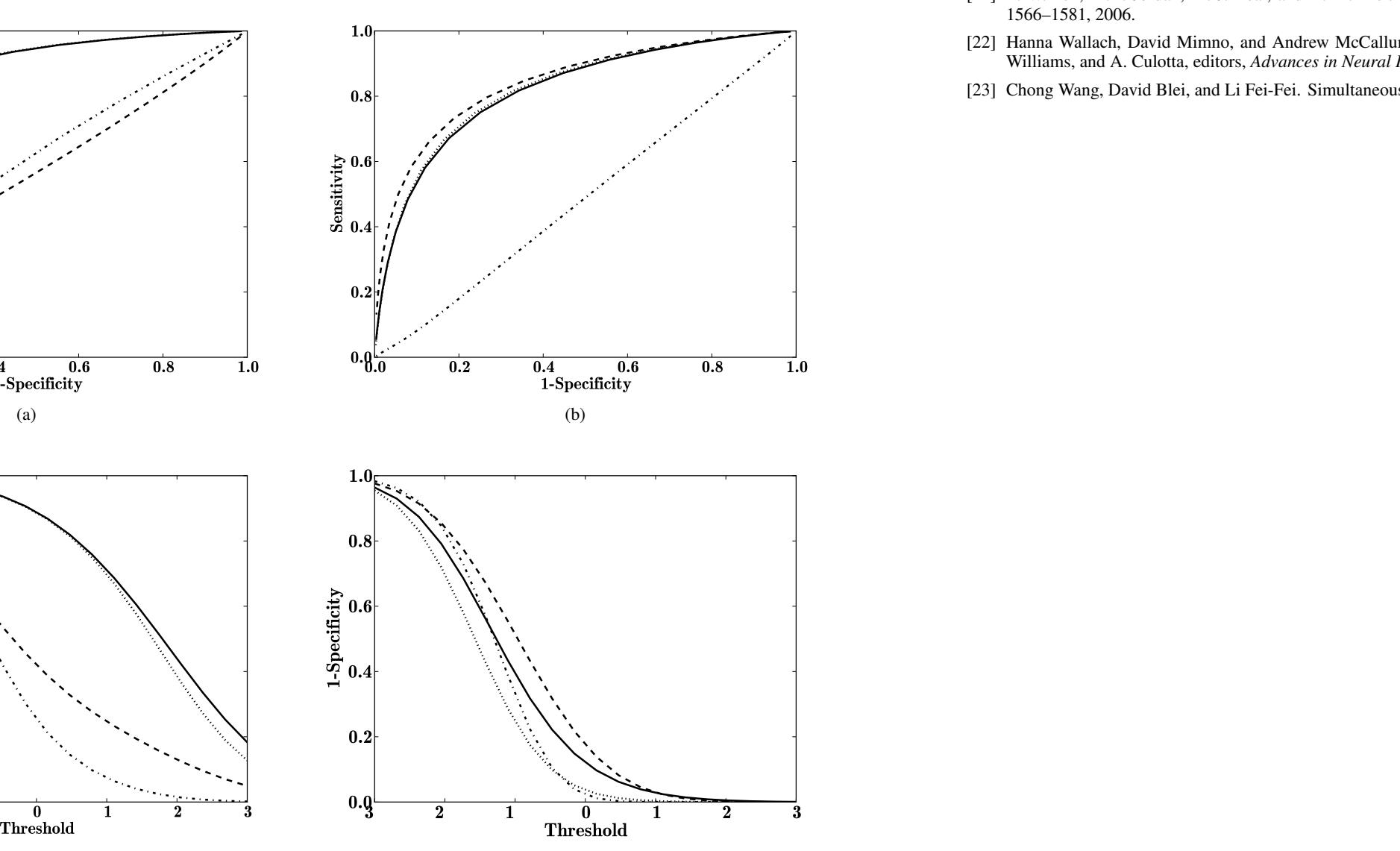


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

- Draw $a_{t,d} | \bar{\mathbf{z}}_d, \beta_t, y_{ps(t),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \beta_t, 1), & y_{ps(t)} = 1 \\ \mathcal{N}(\mathbf{z}^T \beta_t, 1) \mathbb{I}(a_{t,d} < 0), & y_{ps(t)} = -1 \end{cases}$
where $\bar{\mathbf{z}}_d = N_d^{-1} \sum_{i=1}^{N_d} \mathbf{z}_{i,d}$
- Set the response variable

$$y_{t,d} | a_{t,d} = \begin{cases} 1 & \text{if } a_{t,d} > 0 \text{ and } y_{ps(t),d} = 1 \\ -1 & \text{otherwise} \end{cases}$$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{t,d}$ utilized here are also known as auxiliary variables because the are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, β_L , in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3.1.4 $p(\beta | \mathbf{z}, \alpha')$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [22]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [21].

$$\beta \sim \text{Dir}(m_{(.),1} + \alpha', m_{(.),2} + \alpha', \dots, m_{(.),K} + \alpha') \quad (6)$$

$$p(m_{d,k} | \mathbf{z}, \mathbf{m}_{(.,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m_{(.,k)}) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

$$p(\alpha), p(\alpha'), p(\gamma)$$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Experiments

4.1 Data and Pre-Processing

4.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a standard nomenclature which is the International Standard Diagnostic Classification for Epidemiology and Statistics (ICD-9-CM). The ICD-9-CM codes used in this study were generated in a hierarchical manner, where each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing "pneumonia due to adenovirus" is a child of the code representing "viral pneumonia" where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

4.1.2 Product Category Prediction

Data for these experiments were obtained partly from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partly directly from the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine". Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

3

3.1.3 $p(a_{t,d} | \mathbf{z}, \mathbf{Y}, \eta)$

The auxiliary variables $a_{t,d}$ must be sampled for documents $d = 1, \dots, D$ and $t \in \mathcal{L}_d$. The conditional posterior distribution of $a_{t,d}$ is the truncated normal distribution

$$p(a_{t,d} | \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{t,d} - \eta^T \bar{\mathbf{z}}_d)^2 \right\} I(a_{t,d} y_{t,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta | \mathbf{z}, \alpha')$

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter

4.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including sLDA with independent regressions (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. We will refer to these comparison models as the sLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen as to highlight performance in absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter

4.3 Evaluation

For each dataset, a held out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held out set and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held out set and considers all other non-existent labels to be negative. This uniform treatment of labels allows for a fair comparison of the models.

The two measures for predictive performance used here include the true positive rate and the false positive rate evaluated based on $p(y_{t,d} | w_{t,d})$ for each label in each model.

5 Results

6 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions and categories [11] as well as from [49], medical record descriptions and International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are assigned in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label." More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggest that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

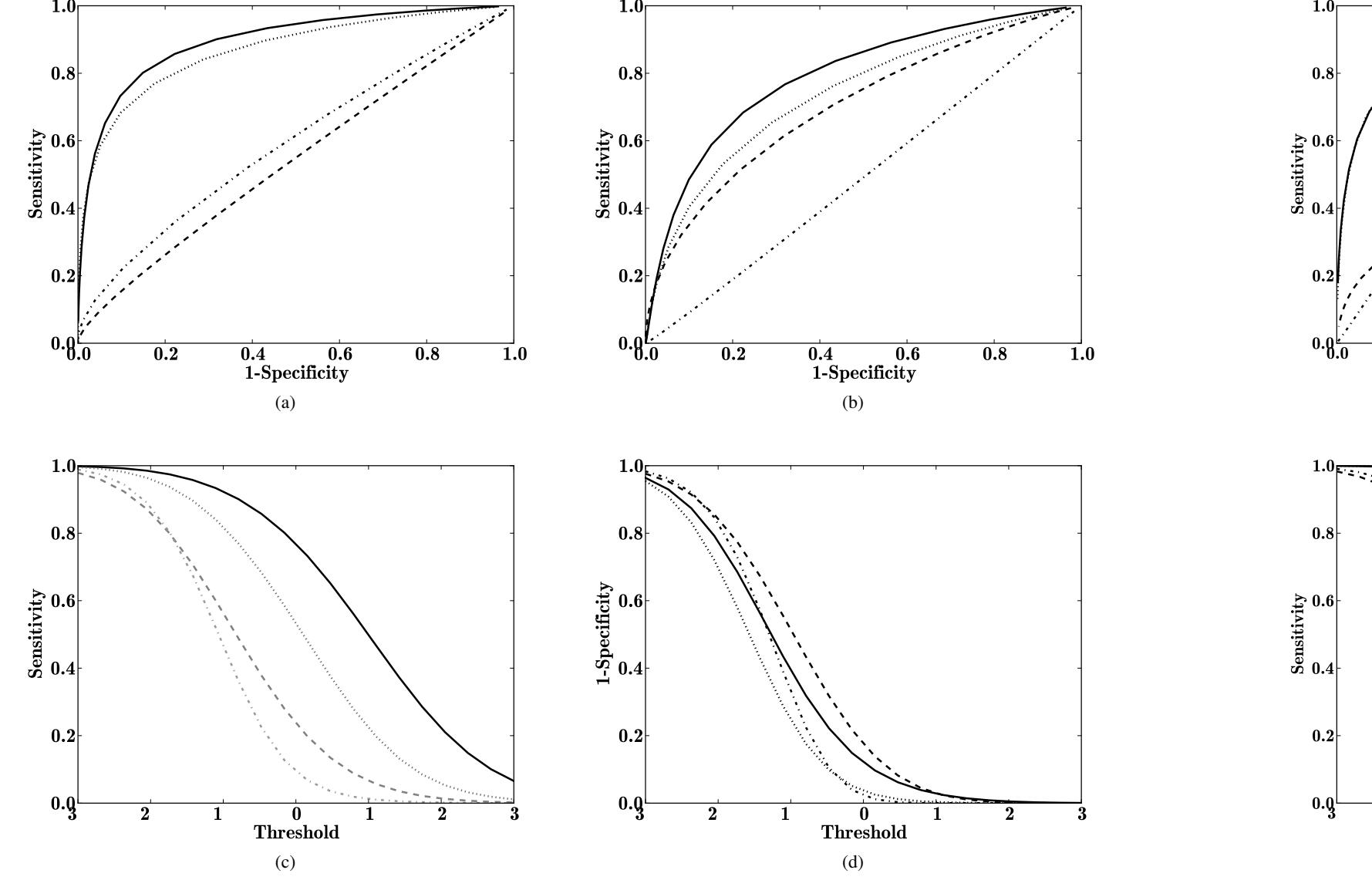


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 13, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work has come from Chang and Blei's work [13] and Griffiths and Steyvers' work [14]. Chang and Blei's work is based on "Clinical Natural Language Processing" (website). Most of the classification strategies included word matching and rule-based algorithms. [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

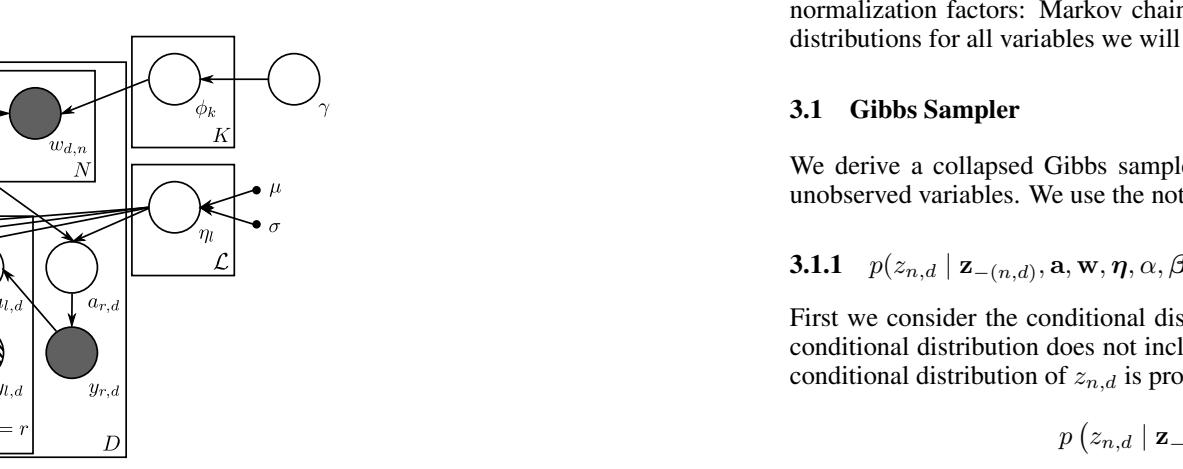


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α_0 , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{1:K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$:
 - Draw a distribution over words $\phi_l \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
 - Draw topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{1:K} \sim \text{Multinomial}(\beta_{z_{n,d}})$

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions and categories [11] as well as from [49], patient medical records vs International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are assigned in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggest that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

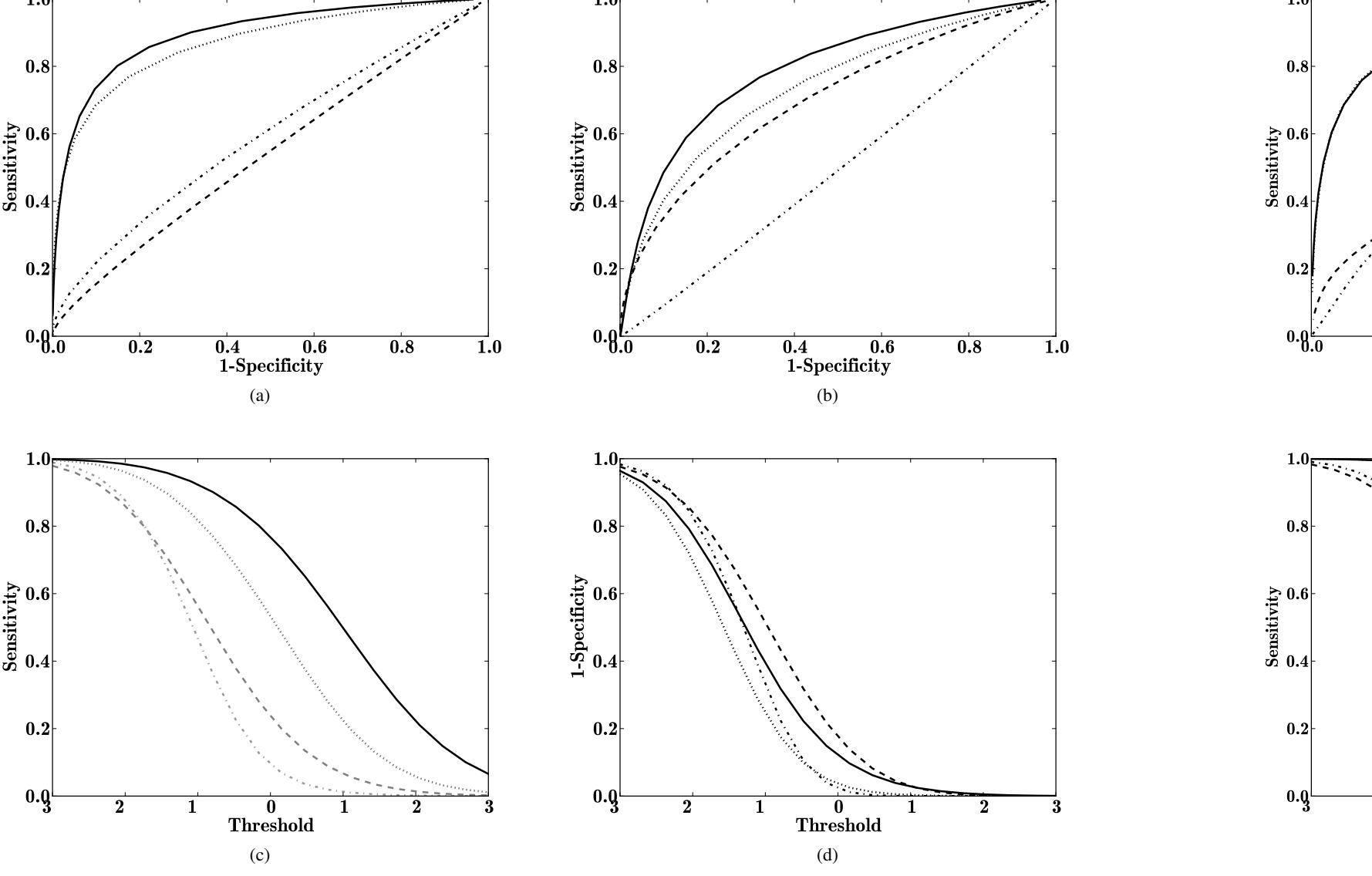


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 23, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work has come from Chang and Blei [8] and Griffiths and Steyvers [13]. Chang and Blei's work is based on their "Collaborative Natural Language Processing" (website). Most of the classification strategies included word matching and rule-based algorithms. [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to l_2 then it will also have a positive response to l_1 . Conversely, if document d has a negative response to l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

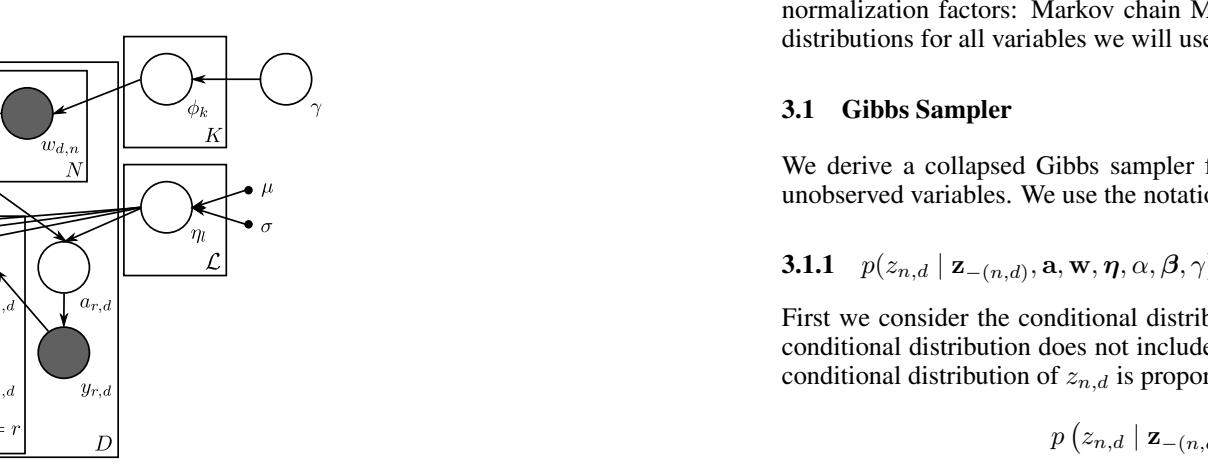


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{1:K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$:
 - Draw a distribution over words $\phi_l \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
 - Draw topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{1:K} \sim \text{Multinomial}(\beta_{z_{n,d}})$

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and customer reviews [1], product descriptions and categories [11] as well as from [49] product manual transcripts and insurance policies (e.g., hospital discharge records vs International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applied in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggest that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

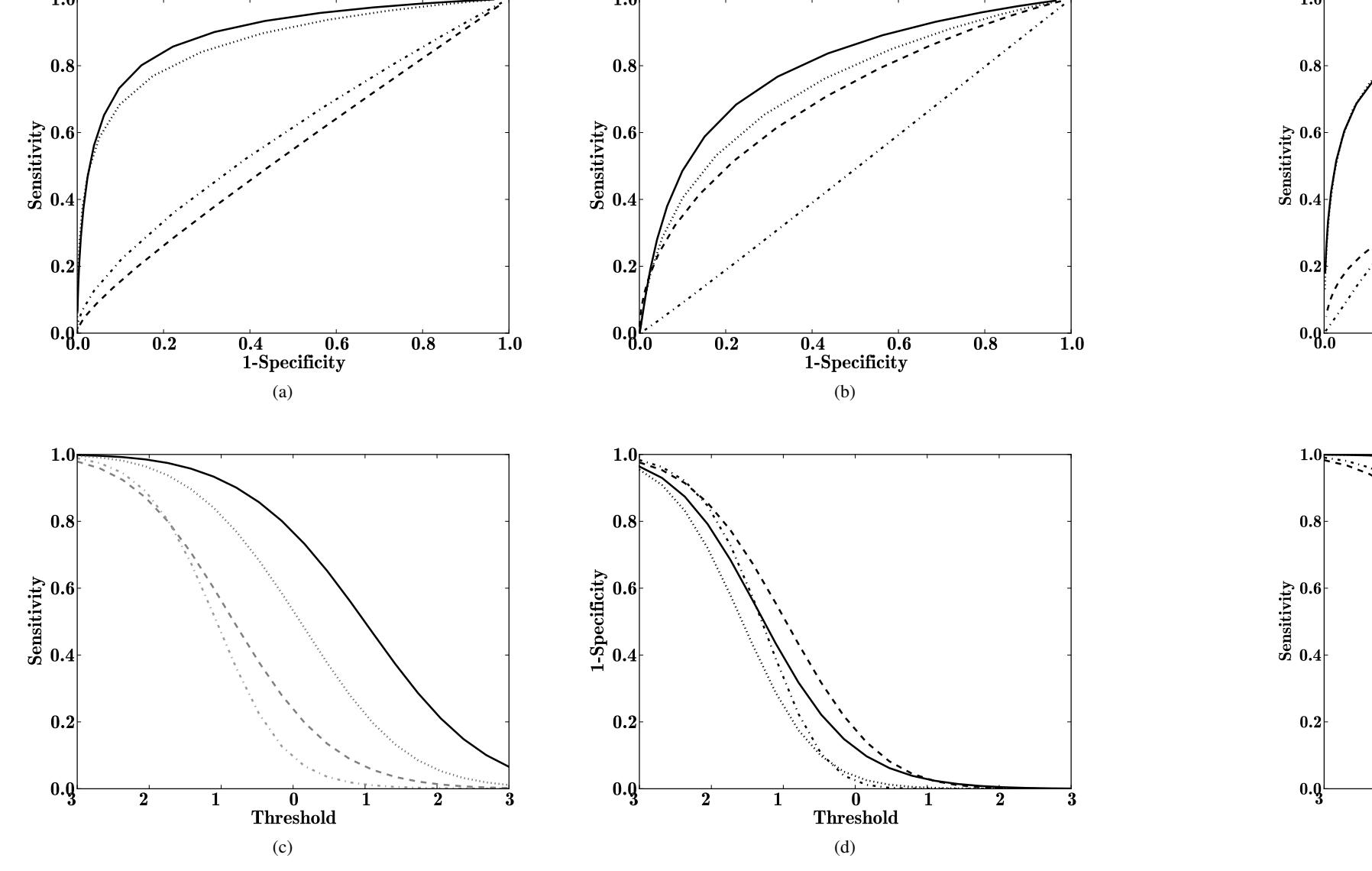


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 13, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work has come from Chang and Blei's work [13] and Griffiths and Steyvers' work [14]. Our work is built on the work of Chang and Blei's work [13] and extends it to incorporate hierarchical supervision. Most of the classification strategies included word matching and rule-based algorithms. [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all of the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

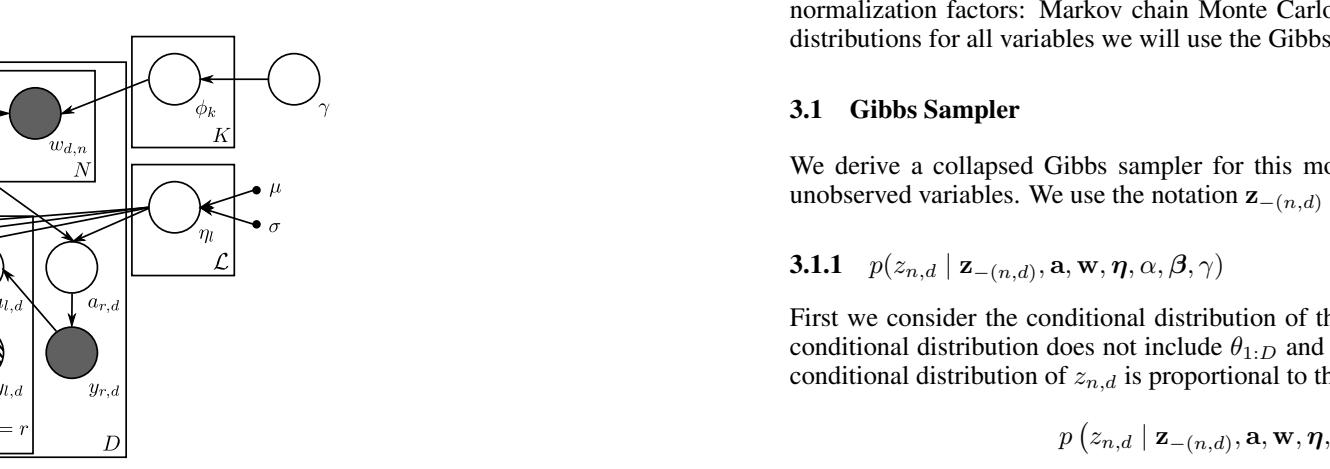


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma_l)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu\mathbf{1}_K, \sigma\mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha')$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha\beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{1:K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$:

2

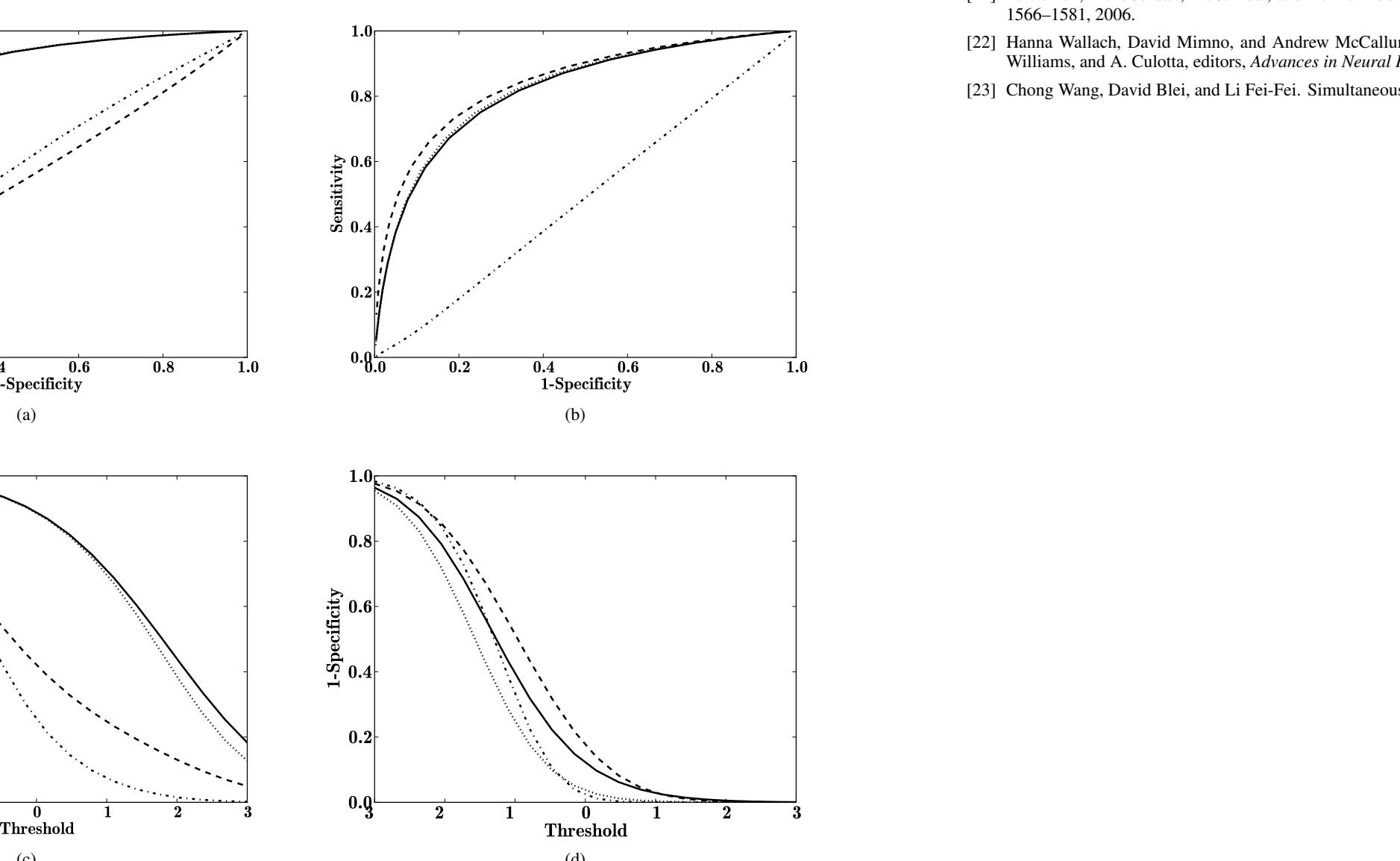


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.1.3 $p(a_{t,d} | \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta})$
The auxiliary variables $a_{t,d}$ must be sampled for documents $d = 1, \dots, D$ and $t \in \mathcal{L}_d$. The conditional posterior distribution of $a_{t,d}$ is the truncated normal distribution

$$p(a_{t,d} | \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{t,d} - \boldsymbol{\eta}^T \mathbf{z}_d) \right\} \mathbb{I}(a_{t,d} y_{ps(t),d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta | \mathbf{z}, \alpha', \alpha)$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [22]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

Posterior inference is performed using the "direct assignment" method of Teh et al. [21].

$$\beta \sim \text{Dir}(m_{(.),1} + \alpha', m_{(.),2} + \alpha', \dots, m_{(.),K} + \alpha') \quad (6)$$

$$p(m_{d,k} | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m_{(d,k)}) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Experiments

In the section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

4.1 Data and Pre-Processing

4.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a standard nomenclature which is the International Standard Diagnostic Classification for Epidemiology and Statistics (ICD-9-CM). The ICD-9-CM codes used in this study were generated in a hierarchical manner, with a child-parent relationship. The ICD-9-CM codes were grouped into categories representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing "pneumonia due to adenovirus" is a child of the code representing "viral pneumonia" where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text. We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

4.1.2 Product Category Prediction

Data for these experiments were obtained partly from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partly directly from the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine". Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

3

[17] Daniel Ramage, David Hall, Rameesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL <http://www.aclweb.org/anthology/cn09-1069.pdf>.

[18] RB Fung, S Sankar, BS Nallapati, CC Aggarwal, and HR Li. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on knowledge discovery and data mining*, pages 416–425, 2003.

[19] B Ribeiro-Neho, AHF Laender, and LRS Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for Information science and Technology*, 52(3):391–401, 2001.

[20] P Bush, J Gobbi, I Thathri, and A Griswold. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636–638, 2008.

[21] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.

[22] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.

[23] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

4

[11] R Farkas and G Szarvas. Automatic construction of rule-based icd-9-cm coders. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.

[12] HR Freitas-Júnior, B Ribeiro-Neto, RF Vale, AHF Laender, and LRS Lima. Categorizationdriven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.

[13] I Goldstein, A Azarbayyan, and O Üzuner. Their approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.

[14] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.

[15] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904, 2009.

[16] LV Lita, S Yu, S Niculescu, and J Bi. Large scale diagnostic code classification for medical patient records, 2008.

4.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how these two types of data can be used together to improve label prediction. Specifically, we show that, if one uses the hierarchical structure of the labels, one can obtain improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

In this work we extend the supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction out-of-sample [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work. In Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_m\}$. Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be interchangeably referred to the observed response of document d to label l_i . The label set is assumed to be structured as an "IS-A" hierarchy. To understand this, consider a hierarchy where label l_i is a parent of label l_j . If document d has a positive response to label l_i then it will also have a positive response to label l_j . Conversely, if document d has a negative response to label l_i then it will also have a negative response to l_j . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ weight



Figure 1: adapted sLDA model

parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $v = 1, \dots, V$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{v,d} \mid z_{n,d}, \beta_{1,k} \sim \text{Multinomial}(\beta_{1,k})$
 - For each label $l \in \mathcal{L}$
 - Draw topic assignment $z_{n,d} \mid \theta_d, y_{p(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \eta_l, 1), & y_{p(l),d} = 1 \\ \mathcal{N}(\mathbf{z}^T \eta_l, 1) | (a_{l,d} < 0), & y_{p(l),d} = -1 \end{cases}$
where $\mathbf{z}_d = \sum_{k=1}^K z_{n,d} \phi_k$
 - Set the response variable
 $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{p(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution – the inverse of which is known as the probit. In this case, the regression is conditional on the parents of the auxiliary variables $a_{l,d}$ utilized here are also known as an auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing. In our model, we collapse over the labels \mathcal{L} not fully observed for a document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ for $l \in \mathcal{L} \setminus \mathcal{L}_d$ in the full generative model. We can also integrate out the parameters $\phi_{l,d}$ and θ_d for $l \in \mathcal{L} \setminus \mathcal{L}_d$. Therefore, in our model the latent variables are $\mathbf{z} = (z_{n,d})_{d=1}^D, \eta = (\eta_l)_{l \in \mathcal{L}}, \theta = (\theta_d)_{d=1}^D, \beta = (\beta_{1,k})_{k=1}^K, \alpha, \alpha', \gamma$.

The posterior distribution we seek cannot be solved in closed form. This is often the case in calculating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MC3) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text. this, to our knowledge, is the first principled approach to doing so

References

- [1] Amazon, Inc. <http://www.amazon.com/>, 2011.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [7] P Brown, DG Cochran, and JR Allegri. The ngram cc classifier: A novel method of automatically creating cc classifiers based on ic9 groupings. *Advances in Disease Surveillance*, 1:30, 2006.
- [8] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [9] K Cramer, M Dredze, K Ganchev, PP Talukdar, and S Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [10] Yan David S, Nilasena Martha J, Radford Brian F, Gage Elena Birman-Deych, Amy D Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [11] R Farkas and G Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [12] HR FreitasJunior, B RibeiroNeto, RF Vale, AHF Laender, and LRS Lima. Categorization-driven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.
- [13] I Goldstein, A Arzumanyan, and Ö Uzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [14] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [15] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [16] LV Lita, S Yu, S Niculescu, and J Bi. Large scale diagnostic code classification for medical patient records. 2008.
- [17] Danijar Ramage, David Hall, Ramesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 – Volume 1*, pages 248–256, Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL <https://aclanthology.org/E09-1030.pdf>.
- [18] RB Rao, S Sandilya, RS Niculescu, C Germond, and HRao. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 416–425, 2003.
- [19] B RibeiroNeto, AHF Laender, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for Information science and Technology*, 52(5):391–401, 2001.
- [20] P Rueh, J Gobell, I Thashit, and A Grivethshuber. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [21] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [22] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, C. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems* 22, pages 1973–1981, 2009.
- [23] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus z_{n,d}$. We are interested in the predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 23, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models, they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [14] we find

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \left(c_{(k,-(n,d))}^{k,-(n,d)} + \alpha_k \right) \frac{\eta_{n,d}^{w_{n,d}}}{\sum_{l \in \mathcal{L}_d} \eta_l^{w_{n,l}}} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(\bar{z}_d \cdot \boldsymbol{\eta}_l - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $c_{(k,-(n,d))}^{k,-(n,d)}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.2 p($\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \mathbf{w}, \boldsymbol{\gamma}$)

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean $\bar{\boldsymbol{\eta}}$ and covariance $\bar{\Sigma}$ such that

$$\Sigma^{-1} = \mathbf{I}^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\bar{\boldsymbol{\eta}} = \bar{\Sigma} \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a} \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [2]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

3.3.1 p($\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \mathbf{w}, \boldsymbol{\gamma}$)

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \boldsymbol{\eta}_l^T \bar{\mathbf{z}}_d)^2 \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.3.4 p($\boldsymbol{\beta} \mid \mathbf{z}, \mathbf{a}, \alpha'$)

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [22]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [21].

$$\boldsymbol{\beta} \sim \text{Dir}(m_{(j),1} + \alpha', m_{(j),2} + \alpha', \dots, m_{(j),K} + \alpha') \quad (6)$$

$$p(m_{(d,k)} = m \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \boldsymbol{\beta}) = \frac{\Gamma(\alpha_k)}{\Gamma(\alpha_k + n_{d,k})} s(n_{d,k}, m) (\alpha_k \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.4.5 p(α), p(α'), p(γ)

The hyperparameters $\alpha</$

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi,bartlett@stat,noemie@dbmi,fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and current hierarchical taxonomies of products [2], product descriptions for categories [14] and patient discharge summaries and International Classification of Disease 9th Revision, Clinical Modification (ICD-9 CM) codes assigned to them [18]. Medical discharge records vs International Classification of Disease 9th Revision, Clinical Modification (ICD-9 CM) codes assigned to them [13]. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggest that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 4 we review related work, in Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 5 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

1.1 Related Work

Latent Dirichlet Allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

1

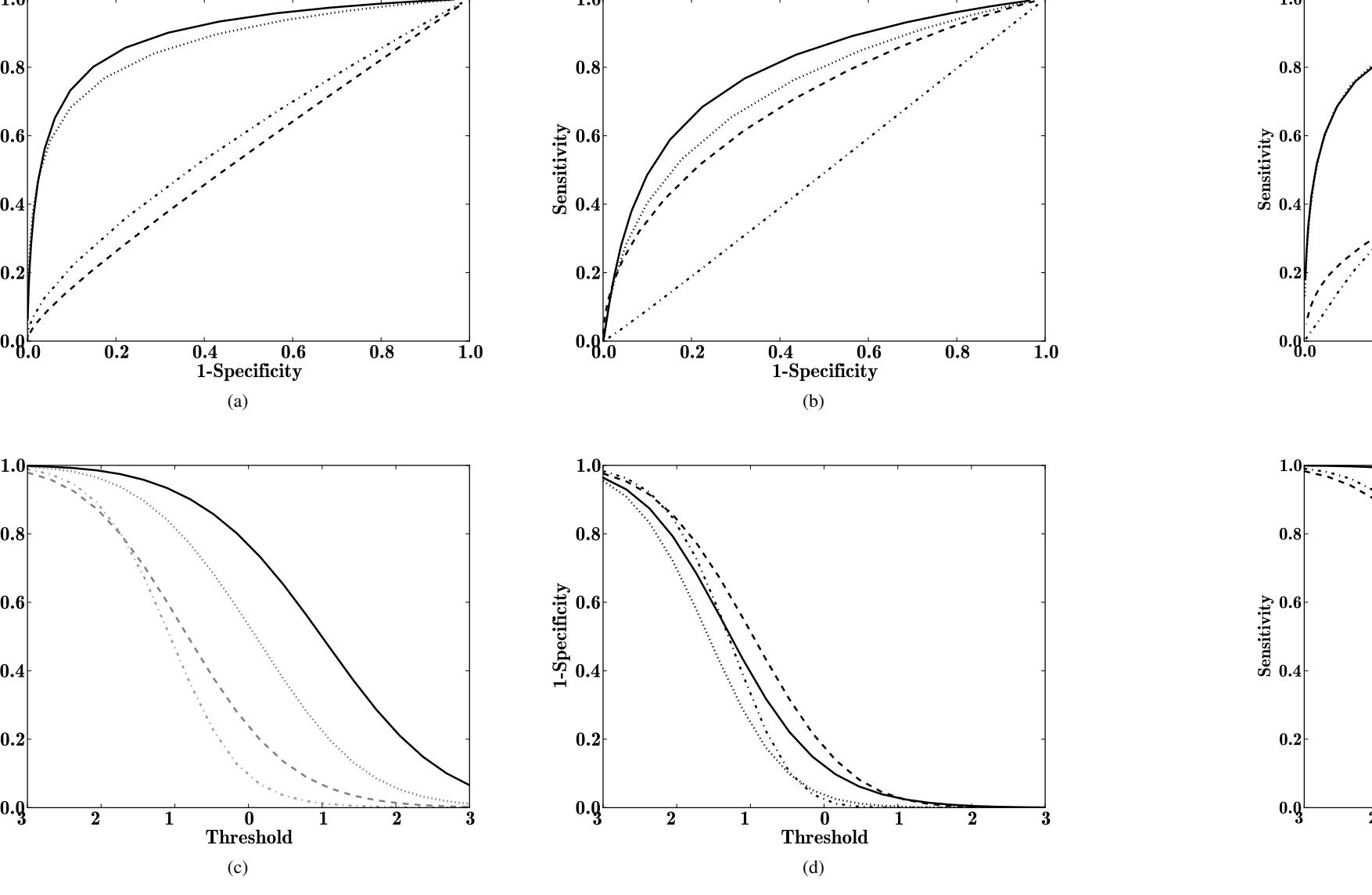


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 23, 8]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

While there have been some attempts to automatically classify groups of patients as a potential preliminary step to ICD-9 code assignment [20, 19, 7], fully automatic assignment of ICD-9 codes to medical text became a more prevalent research topic only in the last few years. A subset of earlier work proposed various methods on small corpora, based on a few special diseases [18] but the most recent and promising work is from Chang et al. [15], Chen et al. [8] and Gao et al. [11]. Other work in this space includes Cui et al. [10] and Li et al. [12] among Natural Language Processing (website). Most of the classification strategies included word matching and rule-based algorithms. [13, 9, 11]. The data set given to the participants consisted only of documents that were 1-2 lines each and all the documents were radiology reports - clearly limiting the scope of potential ICD-9 codes which could be assigned. The only paper which has attempted to work with a document scope as large as ours was the 2008 Lita et al. publication [16]. Lita proposed support vector machine and bayesian ridge regression methods to assign appropriate labels to the documents but did not utilize the ICD-9 hierarchy to leverage more comprehensive predictions.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either 1 or 0 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_1 then it will also have a positive response to label l_2 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models:

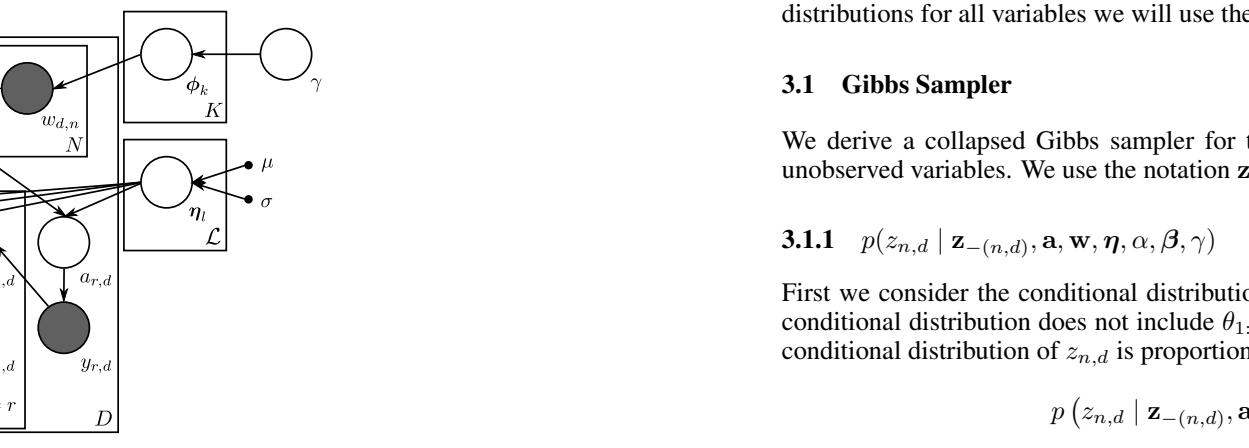


Figure 1: adapted sLDA model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α , α' , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \beta_{1,K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$:
 - Draw topic assignment $z_{n,d,l} | z_{n,d}, \beta_{1,K} \sim \text{Multinomial}(z_{n,d} \cdot \eta_l)$
 - Draw label response $y_{n,d,l} | z_{n,d,l}, \theta_d \sim \text{Probit}(z_{n,d,l} \cdot \mu + \sigma \epsilon_{n,d,l})$

2

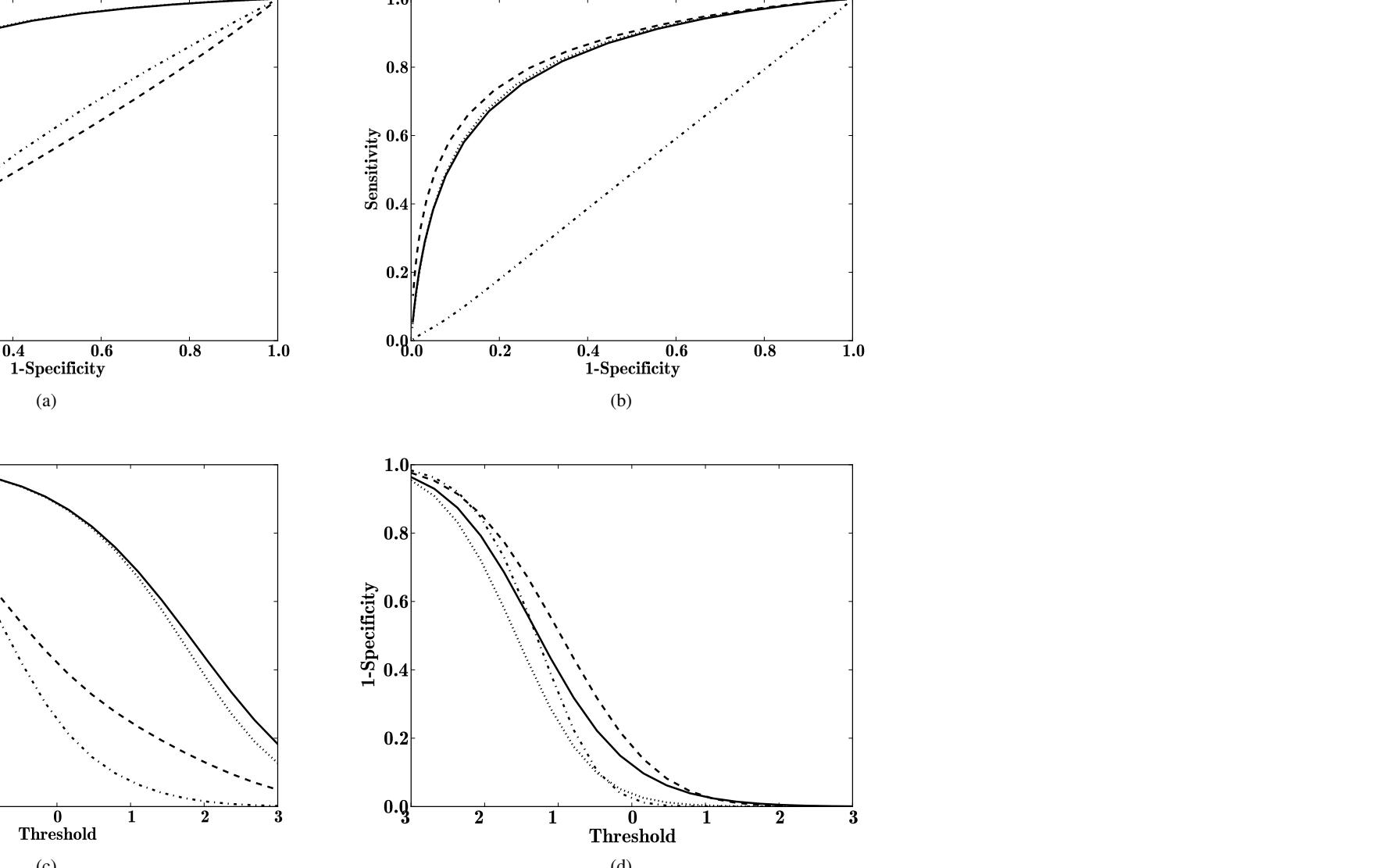


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

- Draw $a_{t,d} | \bar{\mathbf{z}}_d, \beta_t, y_{pa(t),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \boldsymbol{\eta}_t, 1), & y_{pa(t)} = 1 \\ \mathcal{N}(\mathbf{z}^T \boldsymbol{\eta}_t, 1) \mathbb{I}(a_{t,d} < 0), & y_{pa(t)} = -1 \end{cases}$
 where $\bar{\mathbf{z}}_d = \sum_{i=1}^N z_{i,d} \mathbb{I}(y_{pa(i),d} > 0)$
 - Set the response variable

$$y_{t,d} | a_{t,d} = \begin{cases} 1 & \text{if } a_{t,d} > 0 \text{ and } y_{pa(t),d} = 1 \\ -1 & \text{otherwise} \end{cases}$$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{t,d}$ utilized here are also known as auxiliary variables because the are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, β_L , in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

Posterior inference is performed using the "direct assignment" method of Teh et al. [21].

$$\beta \sim \text{Dir}(m_{(.),1} + \alpha', m_{(.),2} + \alpha', \dots, m_{(.),K} + \alpha')$$

$$p(m_{d,k} | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha\beta_k)}{\Gamma(\alpha\beta_k + n_{d,k})} s(n_{d,k}, m_{(d,k),k}) (\alpha\beta_k)^m$$

where $s(n, m)$ represents stirings numbers of the first kind.

$$p(\alpha, \gamma) p(\alpha', \gamma)$$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Experiments

In the section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

4.1 Data and Pre-Processing

4.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary are parts of a multi-level ontology which is the international standard diagnostic classification for epidemiology, statistics and clinical purposes. The ICD-9-CM is a subset of a much larger ontology called SNOMED representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing "pneumonia due to adenovirus" is a child of the code representing "viral pneumonia" where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

4.1.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partially directly from the the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine". Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

3

4

4.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including LDA with independent regressions (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. We will refer to the three comparison models as the sLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen as to highlight performance in absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter for regression parameters.

4.3 Evaluation

For each dataset, a held out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

The two measures for predictive performance used here include the true positive rate and the false positive rate evaluated based on $p(y_{t,d} | w_{1:N,d})$ for each label in each model.

5 Results

6 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...
 • what about the nonparametric version of this?
 • discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

References

- [1] Amazon, Inc. <http://www.amazon.com/>, 2011.
- [2] DMoz open directory project, <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM, international classification of diseases, clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform, <http://snap.stanford.edu/>, 2004.
- [5] D. Blei, and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing Systems*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, 2003. ISSN 1533-4432.
- [7] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOS743.
- [8] Yan Yan David S. Vilasina Martha J. Radford Brian F. Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–51, 2005.
- [9] T.L. Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [10] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DissLDL: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [11] Daniel Klein, David Hall, Ranesh Narayan, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-label settings. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume*

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noe@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how these two types of structured data can be used together to improve label prediction. Specifically, we will show how these labels that are not inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well; namely, any unstructured representations of data that have been hierarchically classified (e.g., images/tweets with bag-of-feature representations).

In this work we extend the supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document “supervision”; often taking the form of a single numerical or categorical “label”. More generally this “supervision” can be seen as extra data about a document, for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction out-of-sample [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and Web retail data suggests that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 1.1 we review related work, in Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 4 we apply HSLDA to health care and Web retail data, showing predictive performance and improved topic generation.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_m\}$. Each document is assigned a response of either 1 or 0 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be interchangeably referred to the observed response of document d to label l_i . The label set is assumed to be structured as an “Is-A” hierarchy. To understand this, consider a hierarchy where label l_j is a parent of label l_i . If document d has a positive response to label l_j then it will also have a positive response to label l_i . Conversely, if document d has a negative response to label l_i then it will also have a negative response to l_j . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ weight

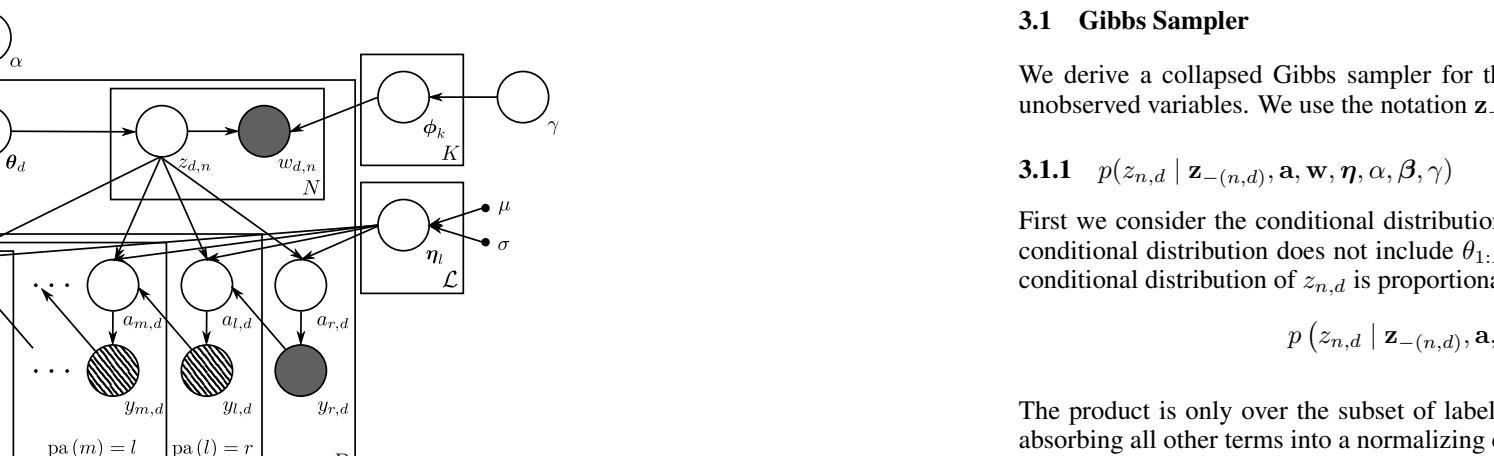


Figure 1: adapted sLDA model

parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $v = 1, \dots, N$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{v,d} \mid z_{n,d}, \beta_{1,k} \sim \text{Multinomial}(\beta_{1,k})$
 - For each label $l \in \mathcal{L}$
 - Draw topic assignment $z_{n,d} \mid \theta_d, y_{l,d} \sim \text{Multinomial}(\eta_l)$, $y_{l,d} = 1$
where $\bar{z}_d = \sum_{v=1}^N z_{n,d}$
- Set the response variable

$$y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{p(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution – the inverse of which is known as the probit. In this case, the regression is conditional on the parents of the observed variable $a_{l,d}$. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing. In our case, we are interested in the labels \mathcal{L} , not fully observed for all documents. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{l,d}$ for $l \in \mathcal{L} \setminus \mathcal{L}_d$ in the full generative model. We can also integrate out the parameters η_l and the variables $a_{l,d}$ as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\{z_{n,d}\}_{d=1,\dots,D}, \eta = \{\eta_l\}_{l \in \mathcal{L}}, \theta = \{\theta_d\}_{d \in \mathcal{L}}, \beta = \{\beta_{1,k}\}_{k=1,\dots,K}, \beta_{\mathcal{L}}, \alpha, \alpha', \gamma$.

The posterior distribution we seek cannot be solved in closed form. This is often the case in calculating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior:

$$\beta \sim \text{Dir}(m_{(.),1} + \alpha', m_{(.),2} + \alpha', \dots, m_{(.),K} + \alpha') \quad (6)$$

$$p(m_{d,k} = m \mid \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{\Gamma(\alpha_k)}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} (a_{d,k} - \eta_l^T \bar{z}_d)^2\right\} I(a_{d,k} > 0). \quad (5)$$

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_{d} \setminus z_{n,d}$. We need to consider the conditional distributions for each of the words in the document. There are many models for making predictions based on text free, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constraints the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on text free, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constraints the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

3.2 Results

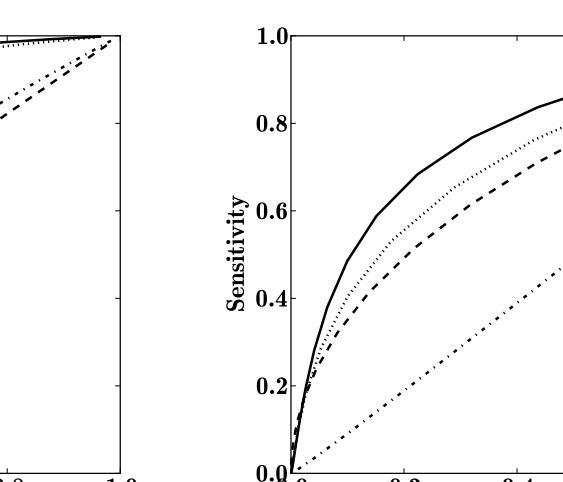


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

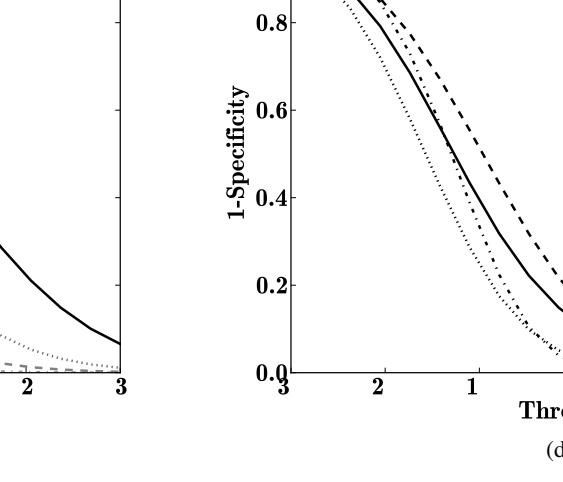


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.



Figure 4

Figure 4: Comparison of three models. We applied HSLDA along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. We will refer to the three comparison models as the sLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen to highlight performance in absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

7 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

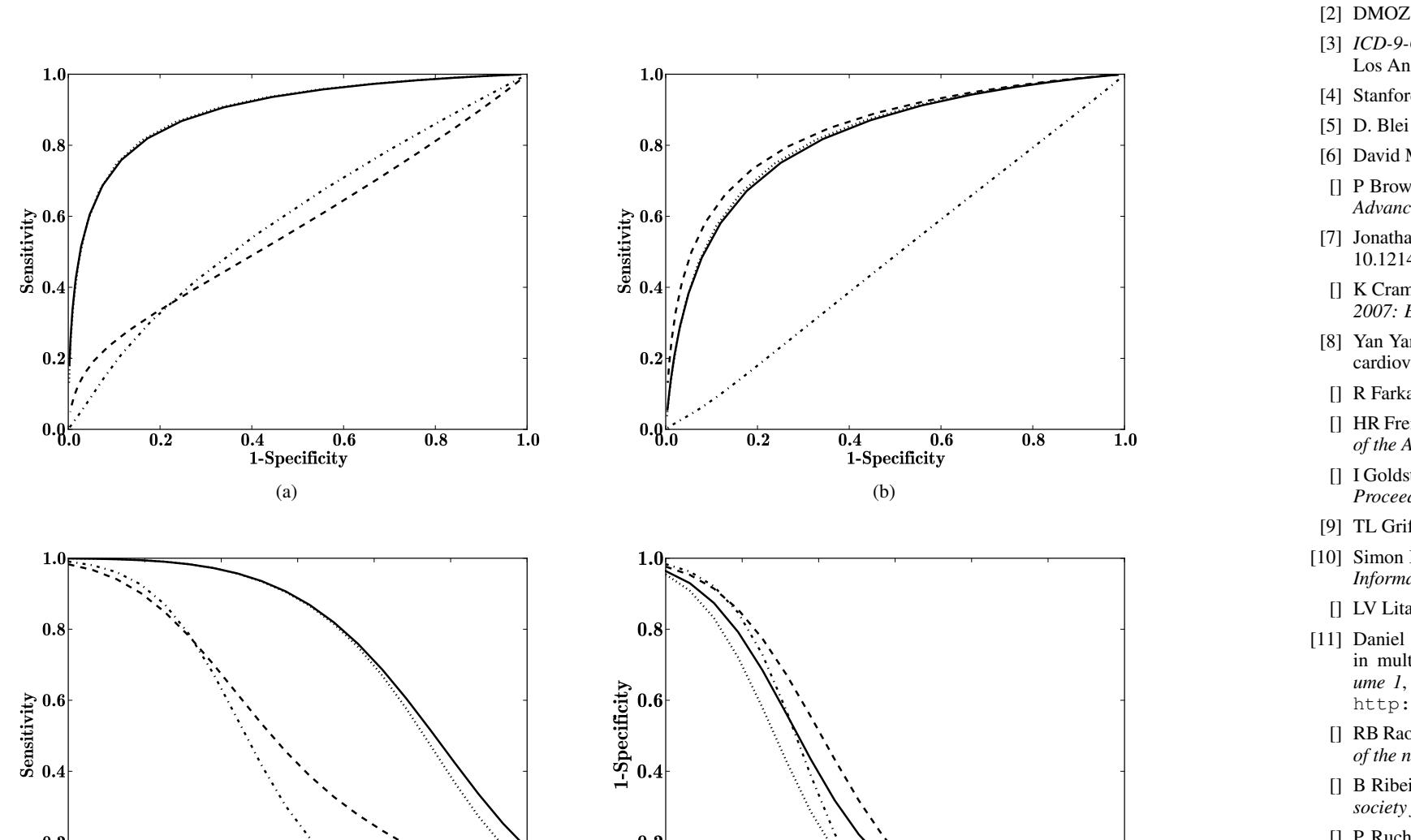


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text. this, to our knowledge, is the first principled approach to doing so

References

- [1] Amazon, Inc. <http://www.amazon.com/>, 2011.
- [2] DMZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [7] P Brown, DG Cochairman, and JR Allegri. The ngram cc classifier: A novel method of automatically creating cc classifiers based on ic9r groupings. *Advances in Disease Surveillance*, 1:30, 2006.
- [8] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/AOAS309.
- [9] K Crammer, M Dredze, K Ganchev, PP Talukdar, and S Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [10] Yan Yan David S. Nilasena Martha J. Radford Brian F. Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [11] R Farkas and G Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [12] HR FreitasJunior, RB RibeiroNeto, RF Vale, AHF Laender, and LRS Lima. Categorizationdriven crosslanguage retrieval of medical information. *Journal of the American Society for Information Science and Technology*, 57(4):501–510, 2006.
- [13] I Goldstein, A Azumetsyan, and Ö Üzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [14] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [15] Simon Lacoste-julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. *In Neural Information Processing Systems*, pages 897–904.
- [16] LV Lita, S Yu, S Niculescu, and J Bi. Large scale diagnostic code classification for medical patient records. 2008.
- [17] Danijar Ramage, David Hall, Rameesh Nullapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, pages 248–256, Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL http://www.aclweb.org/EMNLP2009/pdf/poster_cfm14-1699510_1.pdf.
- [18] RB Rao, S Sandilya, RS Niculescu, C Germond, and H Rao. Clinical and financial outcomes analysis with existing hospital patient records. *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 416–425, 2003.
- [19] B RibeiroNeto, AHF Laender, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and technology*, 52(5):591–401, 2001.
- [20] P Rueh, J Gobell, I Thashir, and A Grivethshuber. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [21] W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [22] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.
- [23] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for booking and insurance purposes (e.g., hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how these two types of data can be used together to improve label prediction. In particular, we find that using the hierarchical structure of the labels can lead to more accurate, compute-improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well; namely, any unstructured representations of data that have been hierarchically categorized (e.g., image captions with bag-of-feature models, news stories with topic models).

In this work we extend the supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document “supervision”; often taking the form of a single numerical or categorical “label”. More generally this “supervision” can be seen as extra data about a document, for instance its quality or relevance (e.g., online review scores), marks given to written work (e.g., essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observe big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 4 we review related work, in Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 5 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_m\}$. Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be interchangeably referred to the observed response of document d to label l_i . The label set is assumed to be structured as an “is-a” hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ weight

1
...
2

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text. this, to our knowledge, is the first principled approach to doing so

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing. In our model we often integrate over labels \mathcal{L} not fully observed for a document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{l,d}$ for $l \in \mathcal{L} \setminus \mathcal{L}_d$ inside the full generative model. We can also integrate out the parameters $\sigma_{l,D}$ and $\theta_{l,D}$ as in Griffiths and Steyvers [14]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{(n,d)}\}_{d=1,\dots,D}, \boldsymbol{\eta} = \{\eta_l\}_{l \in \mathcal{L}}, \mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}, d=1,\dots,D}, \beta, \alpha, \alpha', \gamma$.

The posterior distribution we seek cannot be solved in closed form. This is often the case in calculating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior:

- the first few steps are calculating $\eta_l | \mathbf{z}, \mathbf{a}, \beta, \alpha, \alpha', \gamma$
- the remaining steps are calculating $a_{l,d} | \mathbf{z}, \mathbf{y}, \mathbf{a}, \beta, \alpha, \alpha', \gamma$
- the final step is calculating $y_{l,d} | \mathbf{z}, \mathbf{a}, \beta, \alpha, \alpha', \gamma$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [22]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the “direct assignment” method of Teh et al. [21]:

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_{d \setminus z_{(n,d)}}$ in predictor of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on a free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

There have been many models that incorporate both latent models of text and some form of supervision [17, 15, 23, 8]. One set of models that are particularly relevant to sLDA are Chang & Blei’s hierarchical models for document networks (Relational Topic Models). In that family of models, they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

The two measures for predictive performance used here include the true positive rate and the false positive rate evaluated based on $p(y_{l,d} | w_{1:N,d})$ for each label in each model.

5 Evaluation

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{(n,d)}$ and absorbing all other terms into a normalizing constant as in [14] we find

$$p(z_{(n,d)} = k | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} | \mathbf{z}, \boldsymbol{\eta}_l) p(z_{(n,d)} | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

Here, $\sum_{v \in \mathcal{V}} n_{v,d}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{(n,d)} | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

5 Experiments

In the section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

5.1 Data and Pre-Processing

5.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of New-York-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The note outlines the patient’s chief complaint, diagnostic findings, therapy administered, patient’s response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a controlled terminology which is the international standard diagnostic classification for epidemiological, health management, and clinical purposes. The ICD-9 codes are organized in a rooted-tree structure, where each edge representing an is-a relationship between parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code representing “Pneumonia due to adenovirus” is a child of the code representing “viral pneumonia” where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term-frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [10]. Regardless of that fact, this is one of the only sources for information on patient diagnoses outside the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

5.1.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partially directly from the Amazon.com website [11]. The product IDs and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, “DVD / Games / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine”. Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

5.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3
 ...
 2

3.1.1 $p(z_{(n,d)} | \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean $\bar{\boldsymbol{\eta}}_l$ and covariance $\bar{\Sigma}_l$ such that

$$\bar{\Sigma}_l^{-1} = \mathbf{I}_{|\mathcal{L}|} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\bar{\boldsymbol{\eta}}_l = \bar{\Sigma}_l^{-1} \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_{l,d} \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [2]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is $\mathbf{z}_{(n,d)}$, and $\mathbf{a}_{l,d} = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

3.1.2 $p(\boldsymbol{\eta}_l | \mathbf{z}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ \frac{1}{2} (a_{l,d} - \boldsymbol{\eta}_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.3 $p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [22].

The prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the “direct assignment” method of Teh et al. [21]:

$$\beta \sim \text{Dir}(m_{(\cdot),1} + \alpha', m_{(\cdot),2} + \alpha', \dots, m_{(\cdot),K} + \alpha') \quad (6)$$

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha_k)}{\Gamma(\alpha_k + n_{d,k})} s(n_{d,k}, m) (\alpha_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.1.4 $p(\beta | \mathbf{z}, \alpha', \gamma)$

The applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate

models including sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3
 ...
 2

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters $\alpha, \alpha',$ and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Related Work

Latent Dirichlet Allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter for regression parameters.

4 Evaluation

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how these two types of data can be used together to improve label prediction. Specifically, we show that using hierarchical constraints on labels that are inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically categorized (e.g., image-classifications with bag-of-features representations).

In this work we extend the supervised latent Dirichlet allocation (sLDA) [6] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction out-of-sample [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and Web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 4 we review related work, in Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 5 we apply HSLDA to health care and Web retail data, showing predictive performance and improved topic generation.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_m\}$. Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be interchangeably referred to the observed response of document d to label l_i . The label set is assumed to be structured as an "IS-A" hierarchy. To understand this, consider a hierarchy where label l_j is a parent of label l_i . If document d has a positive response to label l_j then it will also have a positive response to label l_i . Conversely, if document d has a negative response to label l_i then it will also have a negative response to l_j . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ weight

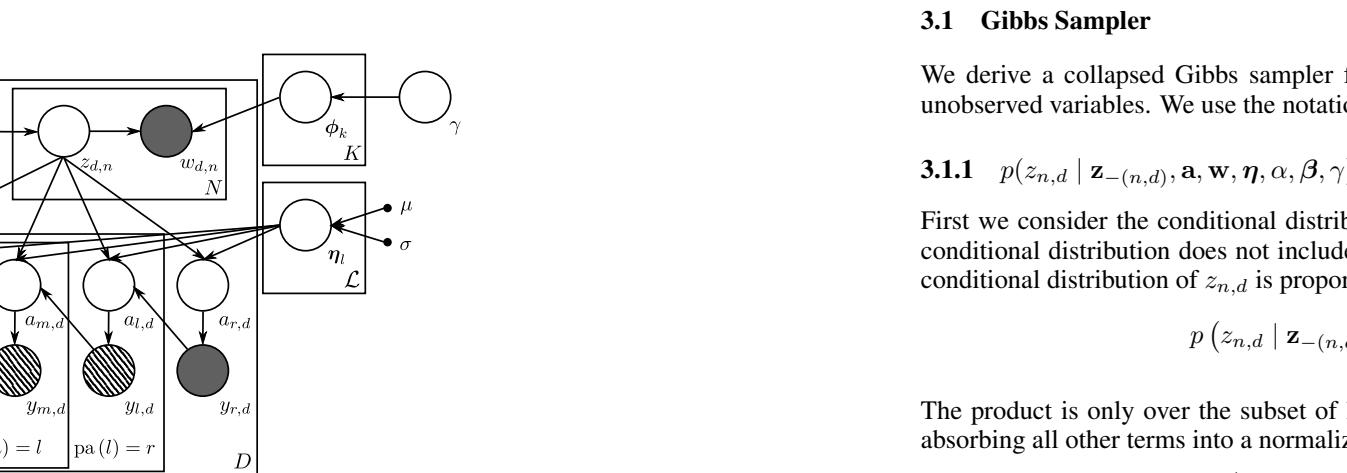


Figure 1: adapted sLDA model

parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $v = 1, \dots, V_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \beta_{1,k} \sim \text{Multinomial}(\beta_{1,k})$
 - For each label $l \in \mathcal{L}$
 - Draw topic assignment $z_{n,d} \mid \theta_d, y_{pa(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \eta_l, 1), & y_{pa(l)} = 1 \\ \mathcal{N}(\mathbf{z}^T \eta_l, 1) \mid (a_{l,d} < 0), & y_{pa(l)} = -1 \end{cases}$
 - Where $\bar{z}_d = \sum_{n=1}^N z_{n,d}$
 - Set the response variable

$$y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{pa(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution – the inverse of which is known as the probit. In this case, the regression is conditional on the parents of each label $l \in \mathcal{L}$.

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$(5)$$

$$p(a_{l,d} \mid z_d, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ \frac{1}{2} (a_{l,d} - \eta_l^T \bar{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0).$$

This conditional distribution can be sampled using an inverse CDF method.

$$(5.1.4)$$

$$p(\beta \mid \mathbf{z}, \alpha', \alpha)$$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [13]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [12].

$$\beta \sim \text{Dir}(m_{(j),1} + \alpha', m_{(j),2} + \alpha', \dots, m_{(j),K} + \alpha')$$

$$p(m_{d,k} = m \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha \beta_k)^m$$

where $s(n, m)$ represents stirling numbers of the first kind.

$$(7)$$

$$p(\alpha, p(\alpha'), p(\gamma))$$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Related Work

Latent Dirichlet Allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

5 Evaluation

5.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_{d \setminus n,d}$.

$$(3.1.1)$$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include $\theta_{1,D}$ and $\phi_{1,V}$ because they have been integrated out as in the collapsed Gibbs sampler [9]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$(1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$(2)$$

$$\text{Here, } c_{v,(n,d')}^{k,-(n,d')} \text{ represents the number of words of type } v \text{ in document } d \text{ assigned to topic } k \text{ omitting the } n^{\text{th}} \text{ word of document } d'. \text{ The notation } (\cdot) \text{ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, }$$

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma).$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$(2)$$

$$\text{Here, } c_{v,(n,d')}^{k,-(n,d')} \text{ represents the number of words of type } v \text{ in document } d \text{ assigned to topic } k \text{ omitting the } n^{\text{th}} \text{ word of document } d'. \text{ The notation } (\cdot) \text{ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, }$$

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma).$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$(2)$$

$$\text{Here, } c_{v,(n,d')}^{k,-(n,d')} \text{ represents the number of words of type } v \text{ in document } d \text{ assigned to topic } k \text{ omitting the } n^{\text{th}} \text{ word of document } d'. \text{ The notation } (\cdot) \text{ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, }$$

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma).$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$(2)$$

$$\text{Here, } c_{v,(n,d')}^{k,-(n,d')} \text{ represents the number of words of type } v \text{ in document } d \text{ assigned to topic } k \text{ omitting the } n^{\text{th}} \text{ word of document } d'. \text{ The notation } (\cdot) \text{ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, }$$

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma).$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$(2)$$

$$\text{Here, } c_{v,(n,d')}^{k,-(n,d')} \text{ represents the number of words of type } v \text{ in document } d \text{ assigned to topic } k \text{ omitting the } n^{\text{th}} \text{ word of document } d'. \text{ The notation } (\cdot) \text{ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, }$$

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma).$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$(2)$$

$$\text{Here, } c_{v,(n,d')}^{k,-(n,d')} \text{ represents the number of words of type } v \text{ in document } d \text{ assigned to topic } k \text{ omitting the } n^{\text{th}} \text{ word of document } d'. \text{ The notation } (\cdot) \text{ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, }$$

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma).$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$(2)$$

$$\text{Here, } c_{v,(n,d')}^{k,-(n,d')} \text{ represents the number of words of type } v \text{ in document } d \text{ assigned to topic } k \text{ omitting the } n^{\text{th}} \text{ word of document } d'. \text{ The notation } (\cdot) \text{ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, }$$

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma).$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how these two types of data can be used together to build more accurate and informative models. Specifically, we show that these two types of data are not only useful for learning topics, but can also be used to improve label prediction. The models and techniques that we develop in this paper are applicable in other domains as well; namely, any unstructured representations of data that have been hierarchically categorized (e.g. images/tweets with bag-of-feature representations).

In this work we extend the supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document “supervision”; often taking the form of a single numerical or categorical “label”. More generally this “supervision” can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction out-of-sample [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 4 we review related work, in Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 5 we apply HSLDA to health care and web retail data, showing predictive performance and improved topic generation.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be interchangeably referred to the observed response of document d to label l_i . The label set is assumed to be structured as an “is-a” hierarchy. To understand this, consider a hierarchy where label l_i is a parent of label l_j . If document d has a positive response to label l_i then it will also have a positive response to label l_j . Conversely, if document d has a negative response to label l_i then it will also have a negative response to l_j . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ weight

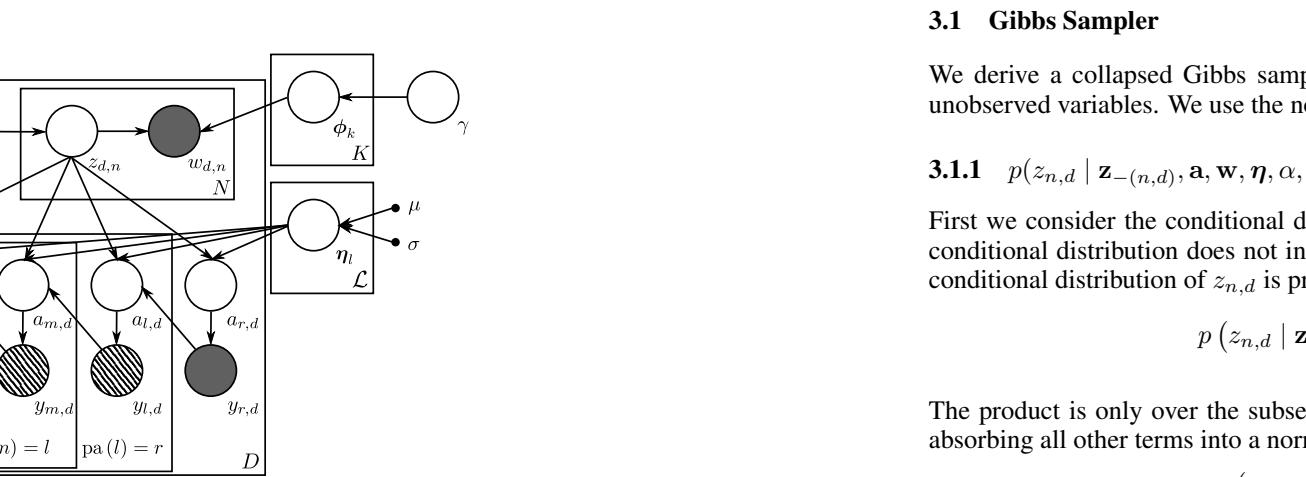


Figure 1: HSLDA graphical model

parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $i = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \beta_{1,K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$
 - Draw topic assignment $z_{l,d} \mid \theta_d, y_{\text{pa}(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \eta_l, 1), & y_{\text{pa}(l),d} = 1 \\ \mathcal{N}(\mathbf{z}^T \eta_l, 1) \mid (a_{l,d} < 0), & y_{\text{pa}(l),d} = -1 \end{cases}$
 - Where $\mathbf{z}_d = \sum_{n=1}^N z_{n,d}$
 - Set the response variable $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{\text{pa}(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution – the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as an auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing. In our case, the labels in \mathcal{L} are not fully observed for all documents. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{\text{pa}(l),d}$ for $l \in \mathcal{L} \setminus \mathcal{L}_d$ in the full generative model. We can also integrate out the parameters ϕ_k and θ_d as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\{z_{i,N_d,d}\}_{d=1,\dots,D}, \eta = \{\eta_l\}_{l \in \mathcal{L}}, \alpha = \{\alpha_l\}_{l \in \mathcal{L}}, \beta_{\mathcal{L}}, \alpha', \alpha$, and γ .

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text. this, to our knowledge, is the first principled approach to doing so

References

- [1] Amazon, Inc. <http://www.amazon.com/>, 2011.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [7] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [8] Yan Yan David S. Nilasena Martha J. Radford Brian F. Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [9] T.L. Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [10] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [11] Daniel Ramage, David Hall, Ramesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL <http://aclweb.org/citation.cfm?id=1699510.1699543>.
- [12] Y.W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [13] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, C. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.
- [14] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

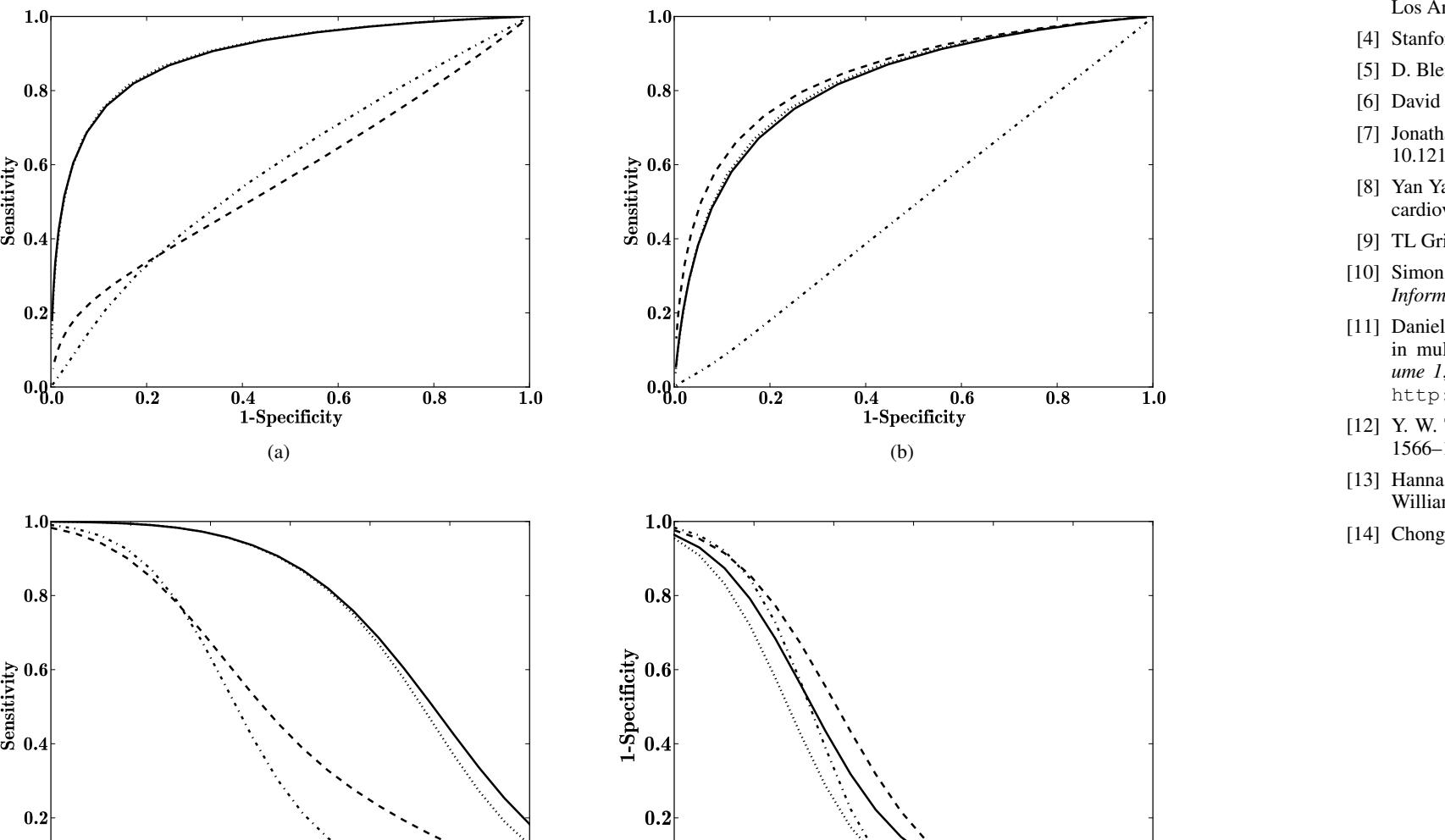


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus \mathbf{z}_{n,d}$.

3.1.1 $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include $\theta_{1,D}$ and $\phi_{1,K}$ because they have been integrated out as in the collapsed Gibbs sampler [9]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \left(c_{(k),d}^{k, -(n,d)} + \alpha \boldsymbol{\beta}_k \right) \frac{w_{n,d}^{k, -(n,d) + \gamma}}{c_{(k),d}^{k, -(n,d)} + \alpha \boldsymbol{\beta}_k + V} \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(\mathbf{z}_d^T \boldsymbol{\eta}_l - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $c_{(k),d}^{k, -(n,d)}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \alpha)$

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean $\bar{\mu}$ and covariance $\bar{\Sigma}$ such that

$$\Sigma^{-1} = \mathbf{I}^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\bar{\mu} = \bar{\Sigma} \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [2]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

3.1.3 $p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta})$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ - \frac{1}{2} (a_{l,d} - \boldsymbol{\eta}_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta \mid \mathbf{z}, \alpha')$

In our model we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [13]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the “direct assignment” method of Teh et al. [12].

$$\beta \sim \text{Dir}(m_{(j),1} + \alpha', m_{(j),2} + \alpha', \dots, m_{(j),K} + \alpha') \quad (6)$$

$$p(m_{(d,k)} = m \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α' , α , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that the use of hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on structured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and medical discharge summaries. ICD-9-CM codes [1] are available from the National Center for Health Statistics (NCHS) [4] and medical treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., medical discharge summaries paired with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per document “supervision”; often taking the form of a single numerical or categorical “label”. More generally this supervision is just extra per document data; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that specifically leverage hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply labeled, bag-of-words data. We will refer to the grouped bag-of-word data as a document. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathcal{W}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the ordered set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured in an “is-a” hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

1

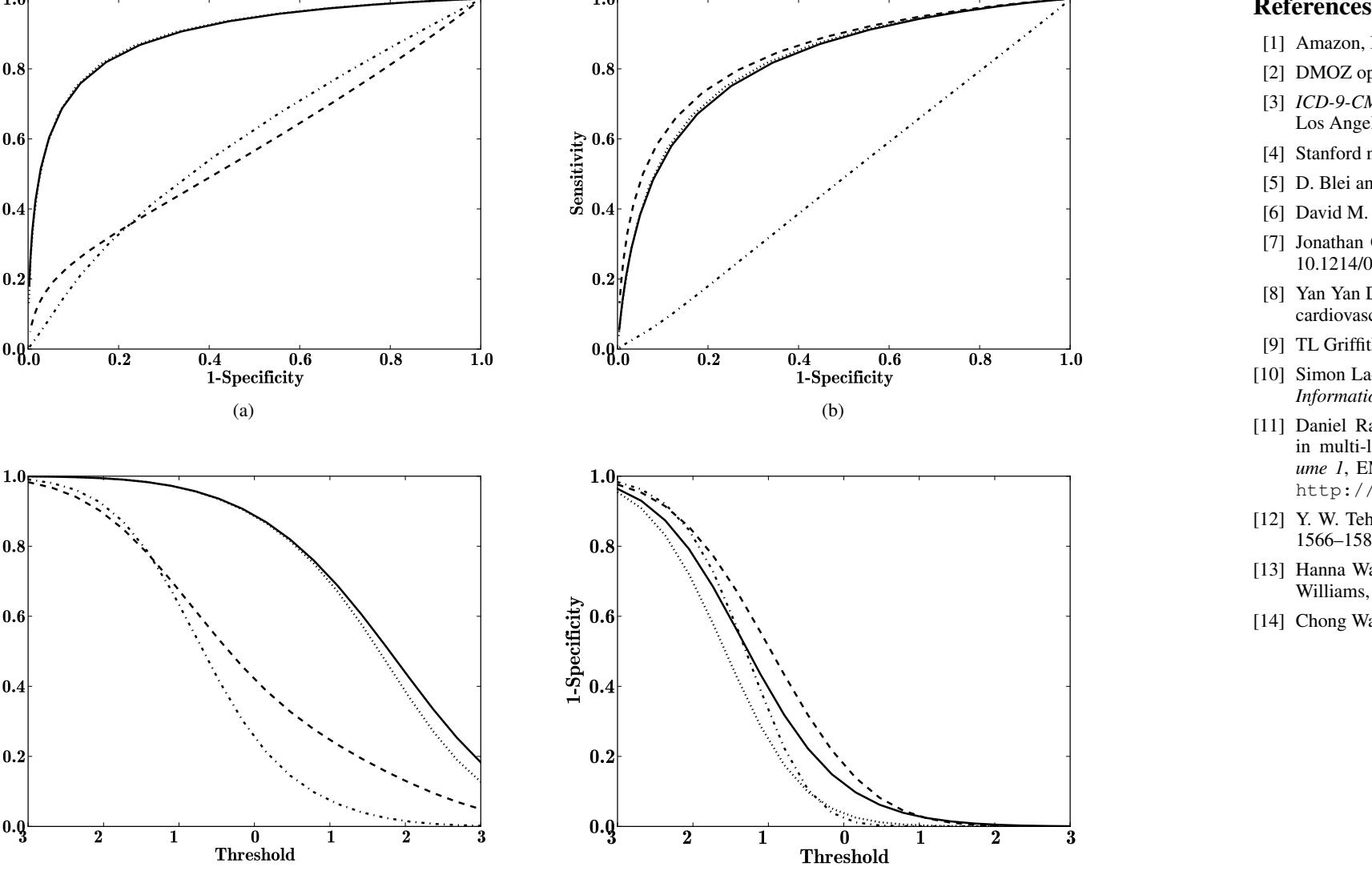


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

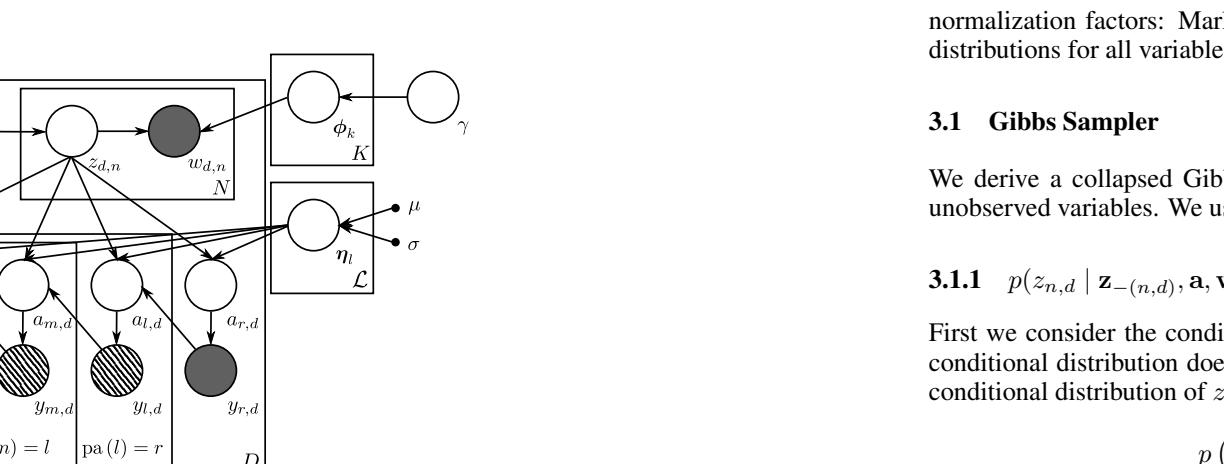


Figure 1: HSLDA graphical model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{I}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K -dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha')$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha' \sim \text{Dir}_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \beta_{1,K} \sim \text{Multinomial}(\beta_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$:
 - Draw $a_{l,d} \mid \mathbf{z}_d, \beta_l, y_{\text{pa}(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \eta_l, 1), & y_{\text{pa}(l),d} = 1 \\ \mathcal{N}(\mathbf{z}^T \eta_l, 1) \mathbb{I}(a_{l,d} < 0), & y_{\text{pa}(l),d} = -1 \end{cases}$
 - Set the response variable $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{\text{pa}(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logit function as the link function, the probit regression model uses the CDF for a standard normal distribution – the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because the arc introduced to make Gibbs sampling possible are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, β_l , in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing.

In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{l,d}$ for $l \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\phi_{k,l}$ and $\theta_{l,D}$ as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\mathbf{z} = \{\mathbf{z}_{1:N_d,d}\}_{d=1,\dots,D}$, $\eta = \{\eta_l\}_{l \in \mathcal{L}}$, $\mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}, d=1,\dots,D}$, $\beta = \{\beta_l\}_{l \in \mathcal{L}}$, $\alpha = \{\alpha_l\}_{l \in \mathcal{L}}$, and γ .

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable

normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

4 Related Work

Latent Dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. As is known to topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of topics, of documents in a corpus. The conditional distribution does not include $\theta_{l,D}$ and $a_{l,D}$ because they have been integrated out as in the collapsed Gibbs sampler [9]. The conditional distribution of $a_{l,d}$ is proportional to the joint distribution of its markov blanket, and LDA followed by least squares regression [5].

There have been many models that incorporate bag-of-words and some form of supervision [11, 10, 14, 7]. One set of models that have been relevant to HSLDA is Chang and Blei's hierarchical models for document networks (referred to as Topic Models). In that family of models, they have considered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

5 Experiments

In this section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

5.1 Data and Pre-Processing

5.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of New York-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The notes outline the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a controlled terminology which is the international standard diagnostic classification for epidemiological, health management, and clinical purposes. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing “Pneumonia due to adenovirus” is a child of the code representing “viral pneumonia” where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [8]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, such as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

5.1.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partially directly from the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, “DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine”. Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

5.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by least squares regression, and HSLDA fit with fixed random regression parameters. We will refer to the three comparison models

as the sLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen as to highlight performance in absence of hierarchical constraints.

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter for regression parameters.

5 Evaluation

For each dataset, a held out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

The two measures for predictive performance used here include the true positive rate and the false positive rate evaluated based on $p(y_{l,d} \mid w_{1:N_d,d})$ for each label in each model.

6 Results

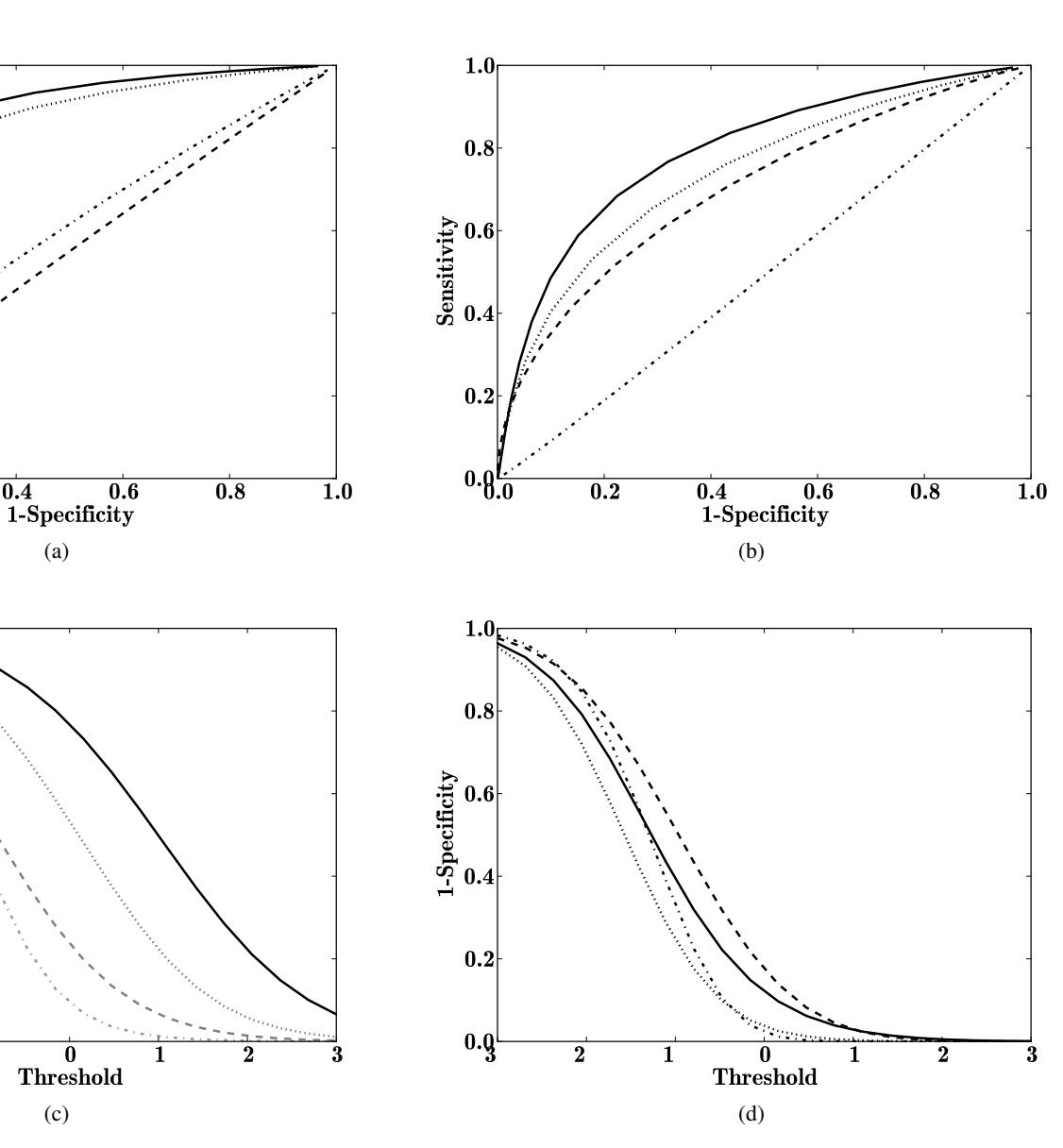


Figure 2: Out-of-sample ICD-9 code prediction. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3

7 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

References

- [1] Amazon, Inc. <http://www.amazon.com/>, 2011.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [7] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOS309.
- [8] Yan David S, Nilasena Martha J, Radford Brian F, Gage Elena Birman-Deych, Amy D, Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–495, 2005.
- [9] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [10] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [11] Daniel Ramage, David Hall, Rameesh Narlapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, Stroudsburg, PA, USA, 2009

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-words data. Out-of-sample label prediction is the primary goal of this work; however, improved dimensionality reduction is also of interest. We define a model that uses probit regressors on a conditionally dependent label hierarchy tied to latent Dirichlet allocation (LDA). We find that the additional signal that comes from multiple, hierarchically constrained labels substantially improves out-of-sample label prediction in comparison to supervised LDA approaches that don't utilize information derived from the structure of the label space. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topic model also improves as a result of using this signal.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how these two types of data can be used together to improve label prediction. Specifically, we show that, if one uses the hierarchical structure of the labels, one can build a more accurate model for label prediction. This will lead to better label prediction, and, in turn, to better models that are more accurate.

In this work we extend the supervised latent Dirichlet allocation (sLDA) [6] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label". More generally this "supervision" can be seen as extra data about a document; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn given an inferred document-specific topic mixture. It has been demonstrated that the signal provided by this supervision can result in better, task-specific document models and can also lead to good label prediction out-of-sample [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply-situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that hierarchical labels should, at least in theory, provide better signal than the simpler unstructured supervision previously considered. Results from applying our model to medical record and Web retail data suggests that this is likely to be the case. In particular, we observed big gains in our primary goal of out-of-sample label prediction when using hierarchical supervision.

The remainder of this paper is structured as follows. In Section 4 we review related work, in Section 2 we introduce hierarchically supervised LDA (HSLDA), and in Section 5 we apply HSLDA to health care and Web retail data, showing predictive performance and improved topic generation.

2 Model

We define here a hierarchically supervised LDA model. Although we will focus on document modeling in our description and experiments, this model applies equally well to other collections of discrete data with hierarchically constrained labels.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be interchangeably referred to as the observed response of document d to label l_i . The label set is assumed to be structured as an "is-a" hierarchy. To understand this, consider a hierarchy where label l_i is a parent of label l_j . If document d has a positive response to label l_i then it will also have a positive response to label l_j . Conversely, if document d has a negative response to label l_i then it will also have a negative response to l_j . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ weight

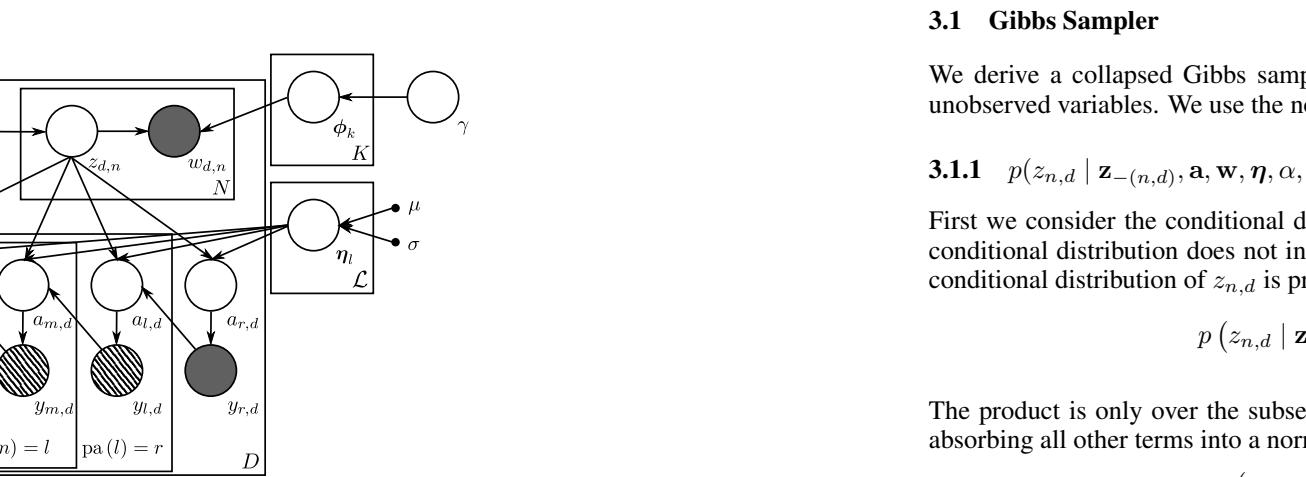


Figure 1: adapted sLDA model

parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $v = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_k \sim \text{Multinomial}(\phi_k z_{n,d})$
 - For each label $l \in \mathcal{L}$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{l,K} \sim \text{Multinomial}(\phi_{l,K} z_{n,d})$
 - Set the response variable

$$y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{\text{pa}(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents of each label l . The latent variables $a_{l,d}$ utilized here are also known as an auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing. In our case, the labels in \mathcal{L} are not fully observed for each document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ for $l \in \mathcal{L} \setminus \mathcal{L}_d$ in the full generative model. We can also integrate out the parameters $\phi_{l,K}$ and η_l as Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $= \{z_{i,N_d,d}\}_{d=1,\dots,D}, \eta = \{\eta_l\}_{l \in \mathcal{L}}, \theta = \{\theta_d\}_{d \in \mathcal{L}}, a = \{a_{l,d}\}_{l \in \mathcal{L}_d, d=1,\dots,D}, \beta, \alpha, \alpha', \gamma$.

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text. this, to our knowledge, is the first principled approach to doing so

References

- [1] Amazon, Inc. <http://www.amazon.com/>, 2011.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [7] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [8] Yan Yan David S. Nilasena Martha J. Radford Brian F. Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [9] T.L. Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [10] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [11] Daniel Ramage, David Hall, Ramesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL <http://aclweb.org/portal.acm.org/citation.cfm?id=1699510.1699543>.
- [12] Y.W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1586, 2006.
- [13] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, C. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, 22, pages 1973–1981, 2009.
- [14] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

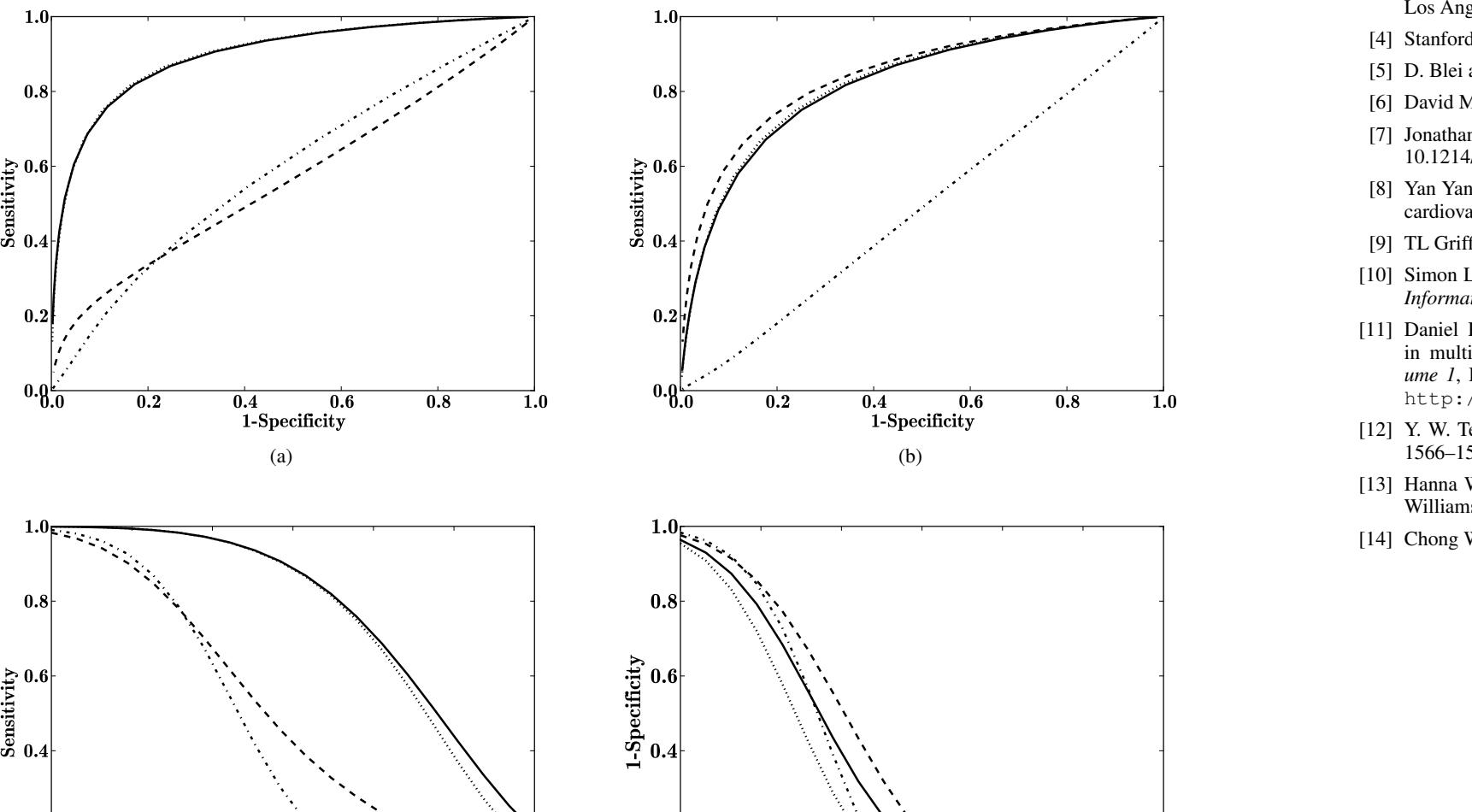


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus z_{n,d}$. We note that $\mathbf{z}_{-(n,d)}$ is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

There have been many models that incorporate both latent models of text and some form of supervision [11, 10, 14, 7]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models, they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

3.2 Evaluation

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

3.3 Evaluation

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

For each dataset, a held-out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held-out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held-out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placements in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data is also of interest. We find that the use of hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and medical discharge summaries. ICD-9-CM codes [1] are available for thousands of medical terms and diseases, including clinical treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., medical discharge summaries paired with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per document “supervision”; often taking the form of a single numerical or categorical “label”. More generally this supervision is just extra per document data; for instance its quality or relevance (e.g., online review scores), marks given to written work (e.g., essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that specifically leverage hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply labeled, bag-of-words data. We will refer to the grouped bag-of-word data as a document. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathcal{W}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the ordered set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$.

We assume a pre-specified set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label l_i for a document d , will be used interchangeably to refer to the observed response of document d to label l_i . The label set is assumed to be structured in an “is-a” hierarchy. To understand this, consider a hierarchy where label l_1 is a parent of label l_2 . If document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

1

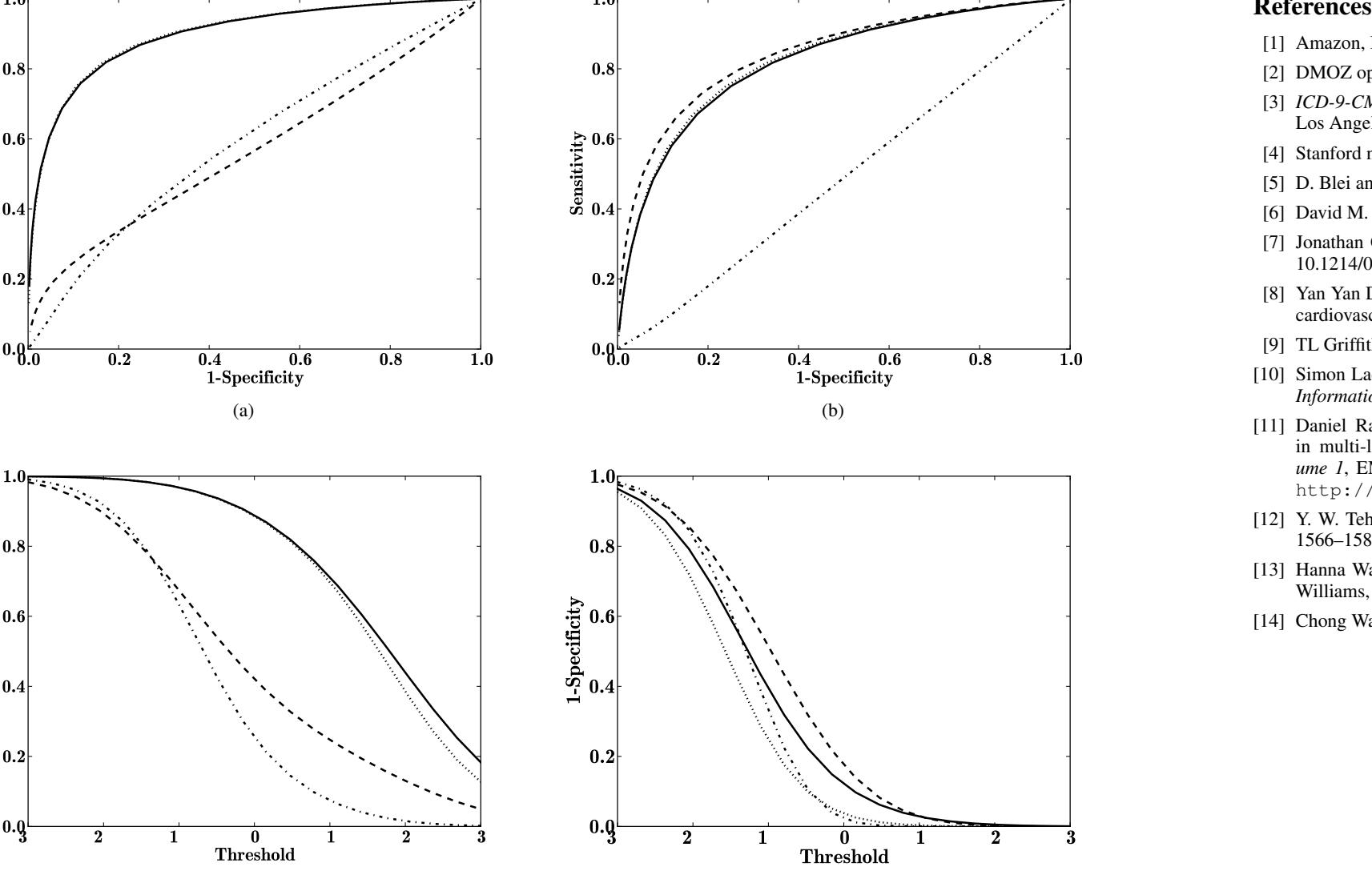


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

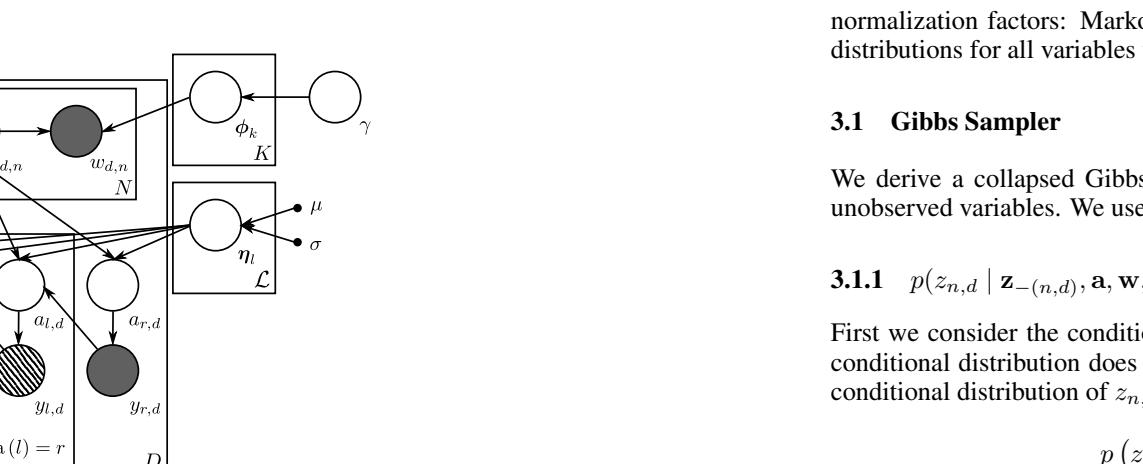


Figure 1: HSLDA graphical model

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(1/\mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu\mathbf{1}_K, \sigma\mathbf{I}_K)$, where \mathbf{I}_K is the K -dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha'\mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha' \sim \text{Dir}_K(\alpha'\beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{k,d} \sim \text{Multinomial}(\phi_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$:
 - Draw $a_{l,d} \mid \mathbf{z}_d, \eta_l, y_{\text{pa}(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}_d^\top \eta_l, 1), & y_{\text{pa}(l),d} = 1 \\ \mathcal{N}(\mathbf{z}_d^\top \eta_l, 1)\mathbb{I}(a_{l,d} < 0), & y_{\text{pa}(l),d} = -1 \end{cases}$
 - Set the response variable $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{\text{pa}(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logit function as the link function, the probit regression model uses the CDF for a standard normal distribution – the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because the arc introduced to make Gibbs sampling possible are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, β_L , in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing.

In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{l,d}$ for $l \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\phi_{k,d}$ and $\theta_{l,D}$ as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1:N_d,d}\}_{d=1,\dots,D}$, $\eta = \{\eta_l\}_{l \in \mathcal{L}}$, $\mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}, d=1,\dots,D}$, $\beta, \alpha, \alpha', \gamma$.

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable

normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

4 Related Work

Latent Dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. As is known to topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, sLDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of topics, of documents in a corpus. The conditional distribution does not include $\theta_{l,D}$ and $a_{l,D}$ because they have been integrated out as in the collapsed Gibbs sampler [9]. The conditional distribution of $a_{l,d}$ is proportional to the joint distribution of its markov blanket, and LDA followed by least squares regression [5].

There have been many models that incorporate bag-of-words and some form of supervision [11, 10, 14, 7]. One set of models that have been relevant to HSLDA is Chang and Blei's hierarchical models for document networks (referred to as Topic Models). In that family of models, they have focused on a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

5 Experiments

In this section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

5.1 Data and Pre-Processing

5.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of New York-Presbyterian Hospital. The data consisted of free-text discharge summaries and their respective ICD-9 codes. A discharge summary is a clinical report prepared by a physician or other health professional at the conclusion of a hospital stay or series of treatments. The notes outline the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of a controlled terminology which is the international standard diagnostic classification for epidemiological, health management, and clinical purposes. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing “Pneumonia due to adenovirus” is a child of the code representing “viral pneumonia” where the former is a type of the latter. The ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (a name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [8]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, such as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

5.1.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partially directly from the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, “DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine”. Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

5.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by least squares regression, and HSLDA fit with fixed random regression parameters. We will refer to the three comparison models

as the sLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen as to highlight performance in absence of hierarchical constraints.

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2, -8, 2, 6, 11\}$), the mean prior parameter for regression parameters.

5 Evaluation

For each dataset, a held out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

The two measures for predictive performance used here include the true positive rate and the false positive rate evaluated based on $p(y_{l,d} \mid w_{1:N_d,d})$ for each label in each model.

6 Results

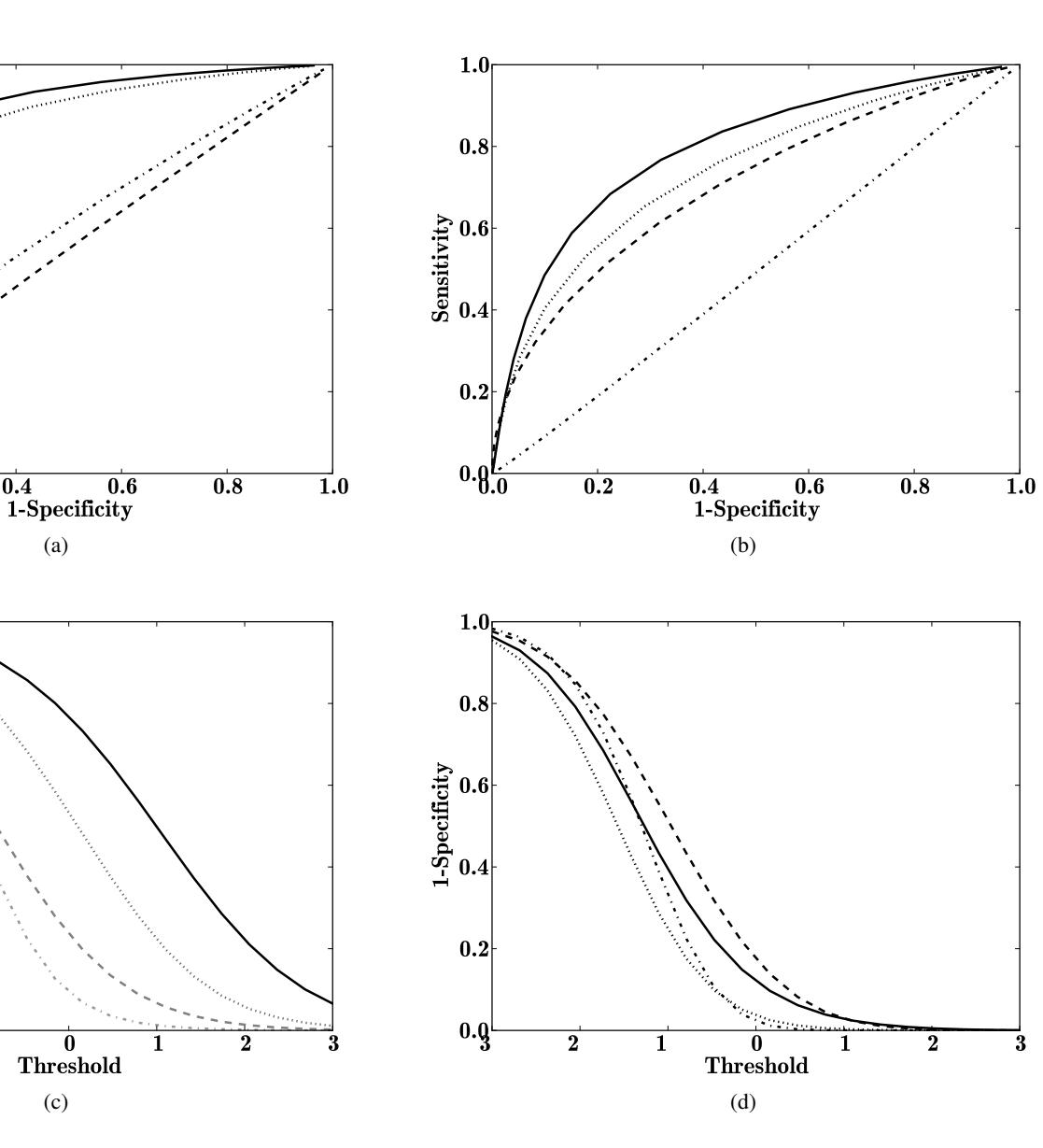


Figure 2: Out-of-sample ICD-9 code prediction. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3

7 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

References

- [1] Amazon, Inc. <http://www.amazon.com/>, 2011.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [7] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [8] Yan David S, Nilasena Martha J, Radford Brian F, Gage Elena Birman-Deych, Amy D, Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–495, 2005.
- [9] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [10] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [11] Daniel Ramage, David Hall, Rameesh Narlapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp909@dbm, bartlett@stat, noemiedbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs [e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in the paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work, we extend the Latent Dirichlet Allocation (LDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with a document-specific prior, taking the form of a single numerical or categorical "label". More generally this supervision is just extra per-document data, for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and Web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We note for exposition purposes that this label set has hard "is-a" parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a partially rooted tree. Without loss of generality we will consider a single root $r \in \mathcal{L}$. Each document has a response $y_{l,d} \in \{-1, 1\}$ to every label which indicates whether the label applies to document d or not.

1

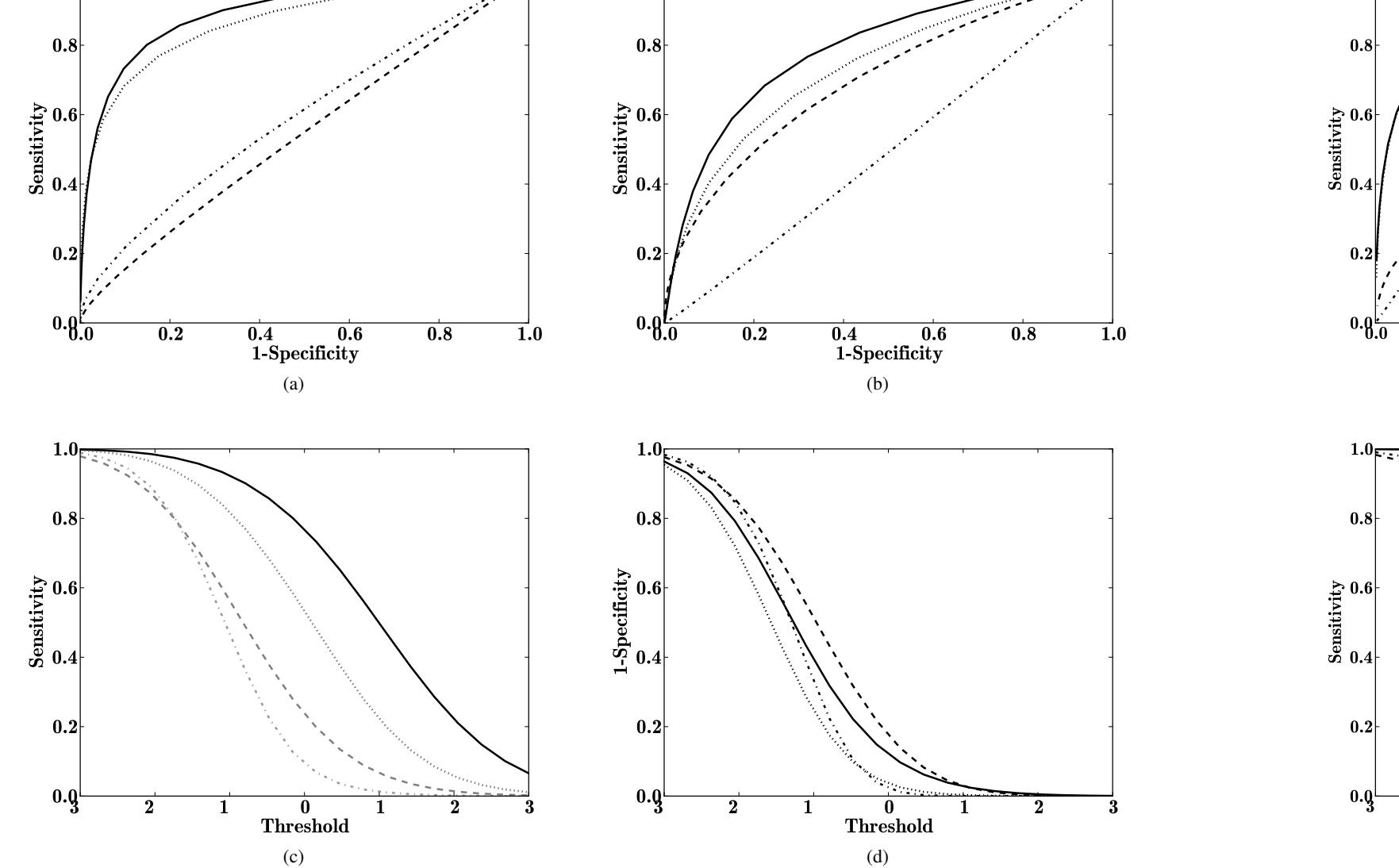


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

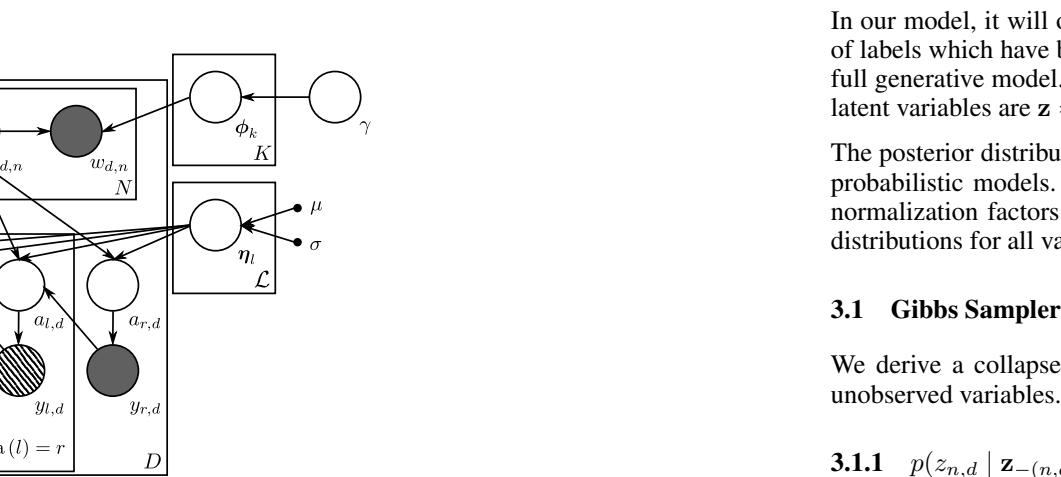


Figure 1: HSLDA graphical model

The is-a hierarchical constraint is a hard constraint that document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l .

The free parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(1/\mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu\mathbf{I}_K, \sigma\mathbf{I}_K)$, where \mathbf{I}_K is the K -dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha'\mathbf{I}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha\beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{1,K} \sim \text{Multinomial}(\phi_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$
 - Draw $a_{l,d} \mid z_d, \eta_l \sim \text{Beta}(\eta_l, 1 - \eta_l)$
 - Draw $y_{l,d} \mid z_d, \eta_l, pa(l), d \sim \begin{cases} \mathcal{N}(z_d^T \eta_l, 1), & y_{pa(l),d} = 1 \\ \mathcal{N}(z_d^T \eta_l, 1)^2, & y_{pa(l),d} = -1 \end{cases}$
 - where $\bar{z}_d = N_d^{-1} \sum_{n=1}^{N_d} z_{n,d}$
 - Set the response variable $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{pa(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing.

In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{l,d}$ for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\phi_{1,K}$ and θ_d as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1,N_d,d}\}_{d=1,\dots,D}$, $\boldsymbol{\eta} = \{\eta_l\}_{l \in \mathcal{L}_d}$, $\mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}_d, d=1,\dots,D}$, $\beta_{\mathcal{L}}$, α' , α , and γ .

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved labels. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus z_{n,d}$.

$$3.1.1 \quad p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include θ_d and $\phi_{1,K}$ because they have been integrated out in the collapsed Gibbs sampler [9]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \frac{(c_{v,n,d}^{k, -(n,d)} + \gamma)}{(c_{v,n,d}^{k, -(n,d)} + \alpha \beta_k)} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(z_{n,d}^T \boldsymbol{\eta}_l - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $c_{v,n,d}^{k, -(n,d)}$ represents the number of words of type v in document d assigned to topic k omitting the n th word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

$$3.1.2 \quad p(\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \sigma)$$

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean $\hat{\boldsymbol{\mu}}_l$ and covariance $\hat{\Sigma}_l$ such that

$$\hat{\Sigma}_l^{-1} = \mathbf{I}^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\hat{\boldsymbol{\mu}}_l = \hat{\Sigma}_l^{-1} \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [?]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \bar{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

$$3.1.3 \quad p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta)$$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \boldsymbol{\eta}^T \bar{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

$$3.1.4 \quad p(\beta \mid \mathbf{z}, \alpha', \alpha)$$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [13]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [12].

$$\beta \sim \text{Dir}(m_{(j),1} + \alpha', m_{(j),2} + \alpha', \dots, m_{(j),K} + \alpha') \quad (6)$$

$$p(m_{d,k} = m \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k,m})} s(n_{d,k,m}) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3

Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, LDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

There have been many models that incorporate both latent models of text and some form of supervision [11, 10, 14, 7]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models, they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

5 Experiments

In the section we describe the application of HSLDA in two hierarchically structured domains. Firstly, we evaluate using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

5.1 Data and Pre-Processing

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbm, bartlett@stat, noemiedbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs [e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-word image representations).

In this work, we extend the Latent Dirichlet Allocation (LDA) [5] to take advantage of hierarchical supervision. HSLDA is latent Dirichlet allocation (LDA) [6] augmented with hierarchical supervision, which takes the form of a simple numerical or categorical "label." More generally, this supervision is just extra per-document data, for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_d \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$.

Each label $l \in \mathcal{L}$ is a parent of l in \mathcal{L} also in the set of labels. We will for exposition purposes assume that this label set has hard "is-a" parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a response $y_{l,d} \in \{-1, 1\}$ to every label which indicates whether the label applies to document d or not.

1

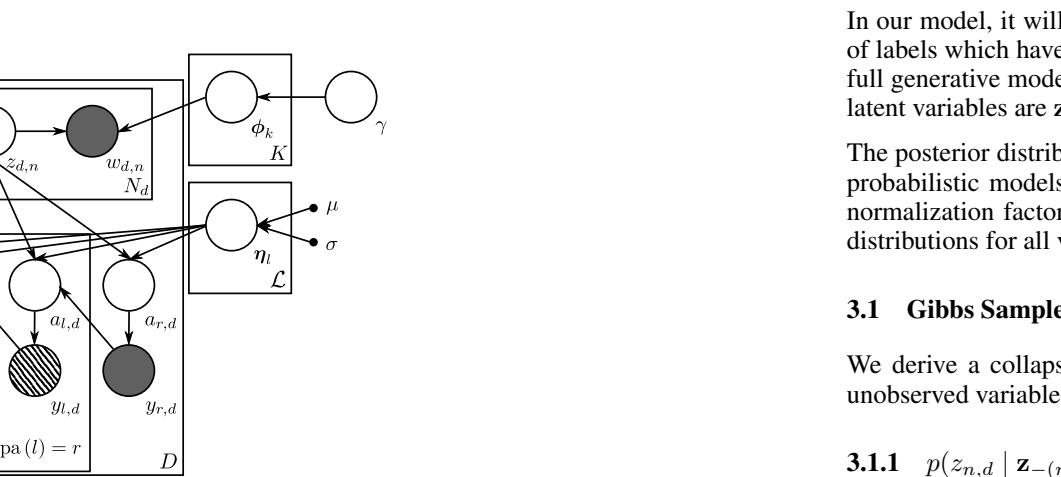


Figure 1: HSLDA graphical model

The is-a hierarchical constraint is a hard constraint that document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of the words using a generative cascade of conditional probit regression models.

Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l .

The free parameters of the model are the number of topics, K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K -dimensional identity matrix
3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{l,d} \sim \text{Multinomial}(\phi_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$
 - Draw $a_{l,d} \mid z_d, \eta_l, y_{\text{pa}(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \eta_l, 1), & y_{\text{pa}(l),d} = 1 \\ \mathcal{N}(\mathbf{z}^T \eta_l, 1) \mathbb{I}(y_{\text{pa}(l),d} < 0), & y_{\text{pa}(l),d} = -1 \end{cases}$
 - where $\mathbf{z}_d = N_d^{-1} \sum_{n=1}^{N_d} z_{n,d}$
 - Set the response variable

$$y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{\text{pa}(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing.

2

In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l',d}$ and $y_{l',d}$ for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\phi_{l,d}$ and θ_d as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1,N_d,d}\}_{d=1,\dots,D}, \boldsymbol{\eta} = \{\eta_l\}_{l \in \mathcal{L}}, \mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}_d, d=1,\dots,D}, \beta, \alpha, \gamma$.

7

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

3.1 Gibbs Sampler

8

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus z_{n,d}$.

3.1.1 $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$

9

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include $\theta_{l,d}$ and $\phi_{l,d}$ because they have been integrated out as in the collapsed Gibbs sampler [9]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

10

$p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma).$ (1)

11

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

12

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \frac{(c_{v,n,d}^{k, -(n,d)} + \gamma)}{(c_{v,n,d}^{k, -(n,d)} + \alpha \beta_k)} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(\mathbf{z}_d^T \boldsymbol{\eta}_l - a_{l,d})^2}{2} \right\}. (2)$$

13

Here, $c_{v,n,d}^{k, -(n,d)}$ represents the number of words of type v in document d assigned to topic k omitting the n th word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

14

3.1.2 $p(\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \sigma)$

15

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean $\hat{\boldsymbol{\mu}}_l$ and covariance $\hat{\Sigma}_l$ such that

16

$$\hat{\Sigma}_l^{-1} = \mathbf{I}^{-1} + \mathbf{Z}^T \mathbf{Z} (3)$$

$$\hat{\boldsymbol{\mu}}_l = \hat{\Sigma}_l^{-1} \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right). (4)$$

17

This is a standard result from normal Bayesian linear regression [?]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

3.1.3 $p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta)$

18

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

19

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \boldsymbol{\eta}_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). (5)$$

20

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta \mid \mathbf{z}, \alpha')$

21

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [13]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

22

Posterior inference is performed using the "direct assignment" method of Teh et al. [12].

23

$$\beta \sim \text{Dir} (m_{(.),1} + \alpha', m_{(.),2} + \alpha', \dots, m_{(.),K} + \alpha') (6)$$

24

$$p(m_{d,k} = m \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha \beta_k)^m (7)$$

25

where $s(n, m)$ represents stirling numbers of the first kind.

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

</div

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp909@dbm, bartlett@stat, noemiedbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work, we extend the Latent Dirichlet Allocation (LDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] extended with a document-specific prior, taking the form of a simple empirical or categorical "blended". More generally this supervision is just extra per-document data, for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and Web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$.

Each label $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard "is-a" parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a partially rooted tree. Without loss of generality we will consider a single root $r \in \mathcal{L}$. Each document has a

response $y_{l,d} \in \{-1, 1\}$ to every label which indicates whether the label applies to document d or not.

1

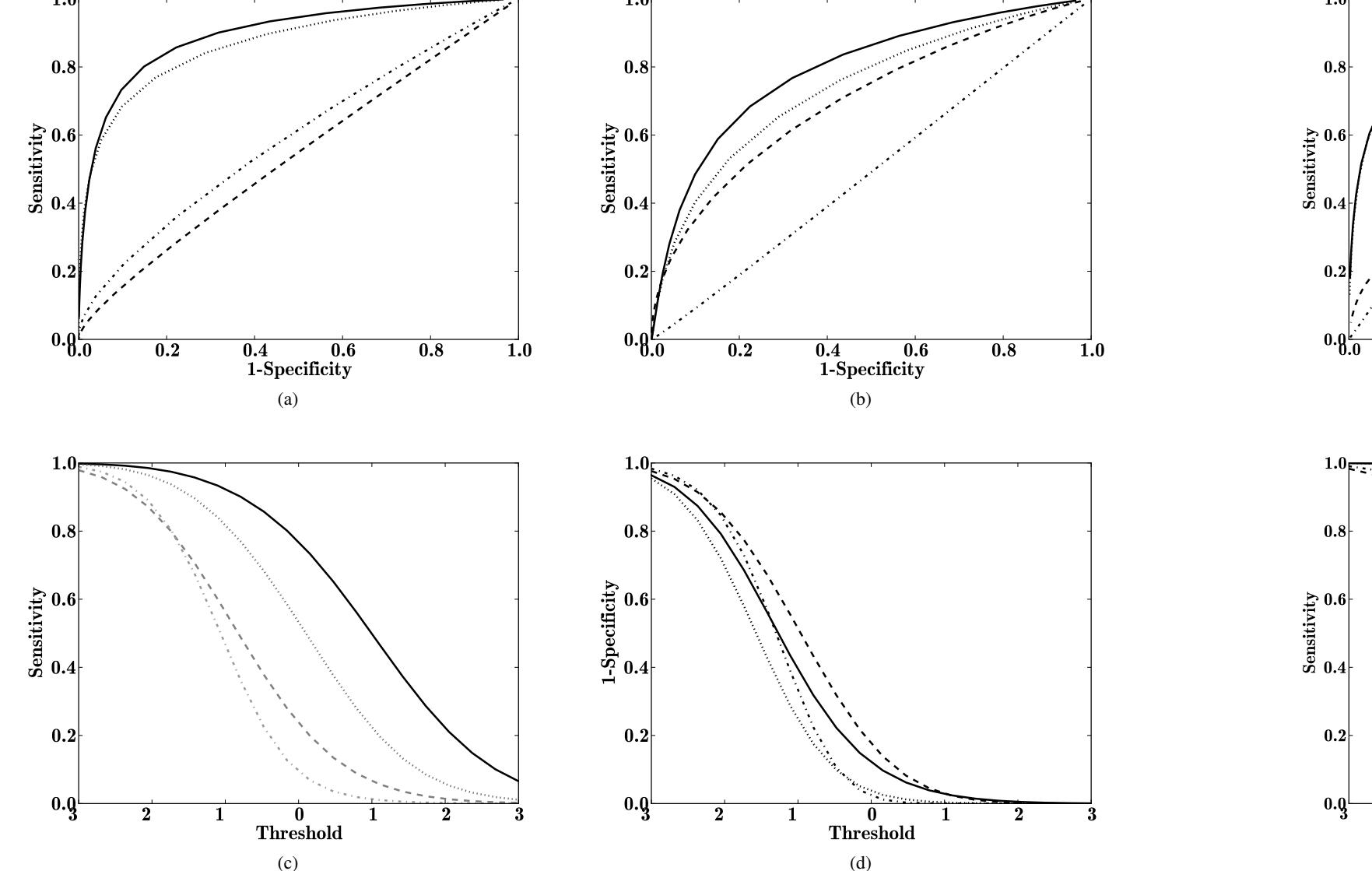


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

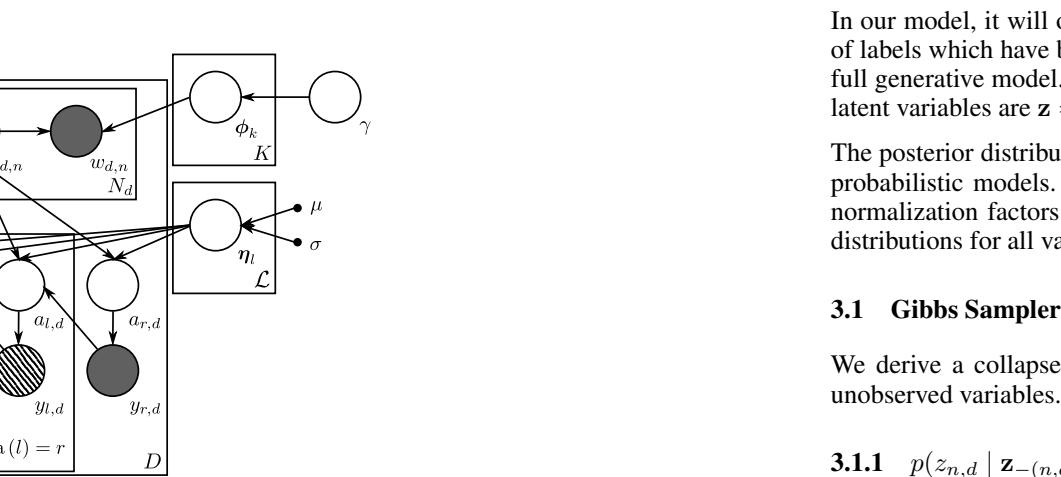


Figure 1: HSLDA graphical model

The is-a hierarchical constraint is a hard constraint that document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l .

The free parameters of the model are the number of topics, K , the number of unique words in the vocabulary, V , the number of documents, D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(1/\mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu\mathbf{1}_K, \sigma\mathbf{I}_K)$, where \mathbf{I}_K is the K -dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha'\mathbf{I}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha\beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{1,k} \sim \text{Multinomial}(\phi_{z_{n,d},k})$
 - For each label $l \in \mathcal{L}$
 - Draw $a_{l,d} \mid z_d, \eta_l \sim \text{Beta}(\eta_l^{z_d}, \eta_l^{1-z_d})$
 - Where $\bar{z}_d = N_d^{-1} \sum_{n=1}^N z_{n,d}$
 - Set the response variable $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{pa(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing.

In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l',d}$ and $y_{l',d}$ for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\phi_{1,k}$ and θ_d as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1,N_d,d}\}_{d=1,\dots,D}$, $\boldsymbol{\eta} = \{\eta_l\}_{l \in \mathcal{L}}$, $\mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}_d, d=1,\dots,D}$, β , α' , and γ .

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus z_{n,d}$.

$$3.1.1 \quad p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include θ_d and $\phi_{1,k}$ because they have been integrated out in the collapsed Gibbs sampler [9]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \frac{c_{v,n,d}^{k - z_{n,d}} + \gamma}{(c_{v,n,d} + \alpha\beta_k)^{\frac{k - z_{n,d}}{2}}} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(z_{n,d} - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $c_{v,n,d}^{k - z_{n,d}}$ represents the number of words of type v in document d assigned to topic k omitting the n th word of document d . The notation $\langle \cdot \rangle$ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

$$3.1.2 \quad p(\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \sigma)$$

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean $\hat{\boldsymbol{\mu}}_l$ and covariance $\hat{\Sigma}_l$ such that

$$\hat{\Sigma}_l^{-1} = \mathbf{I}^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\hat{\boldsymbol{\mu}}_l = \hat{\Sigma}_l^{-1} \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [?]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \bar{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

$$3.1.3 \quad p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta_l)$$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta_l) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^T \bar{z}_d) \right\} \mathbb{I}(a_{l,d} \eta_l^T \bar{z}_d > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

$$3.1.4 \quad p(\beta \mid \mathbf{z}, \alpha')$$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [13]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [12].

$$\beta \sim \text{Dir}(m_{(\cdot),1} + \alpha', m_{(\cdot),2} + \alpha', \dots, m_{(\cdot),K} + \alpha') \quad (6)$$

$$p(m_{d,k} = m \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha\beta_k)}{\Gamma(\alpha\beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha\beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3

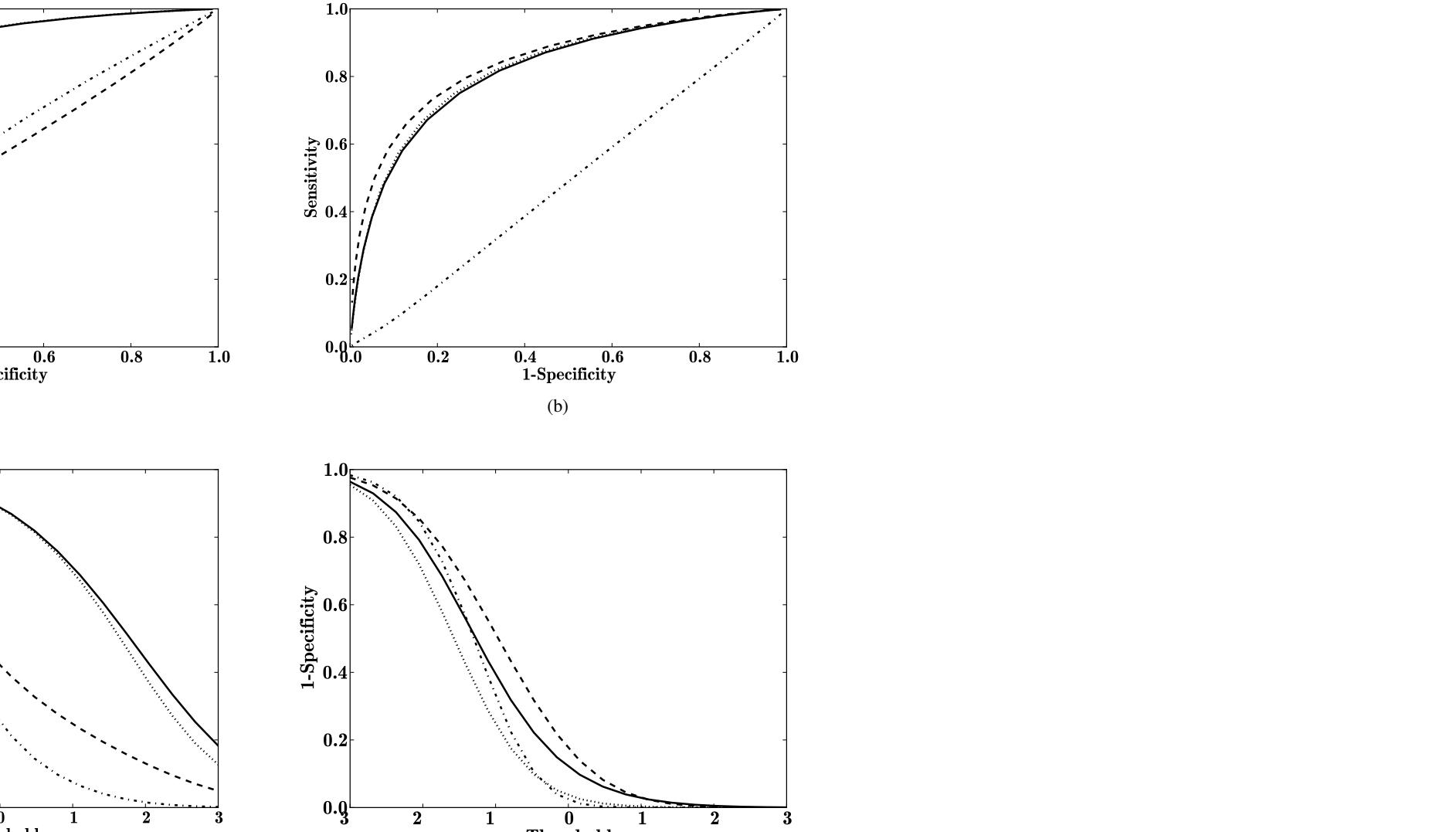


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp909@dbm, bartlett@stat, noemiedbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs [e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work, we extend the Latent Dirichlet Allocation (LDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with a document-specific prior, taking the form of a simple empirical or categorical "blended". More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and Web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We want for exposition purposes that this label set has hard "is-a" parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a single root $r \in \mathcal{L}$. Each document has a response $y_{l,d} \in \{-1, 1\}$ to every label which indicates whether the label applies to document d or not.

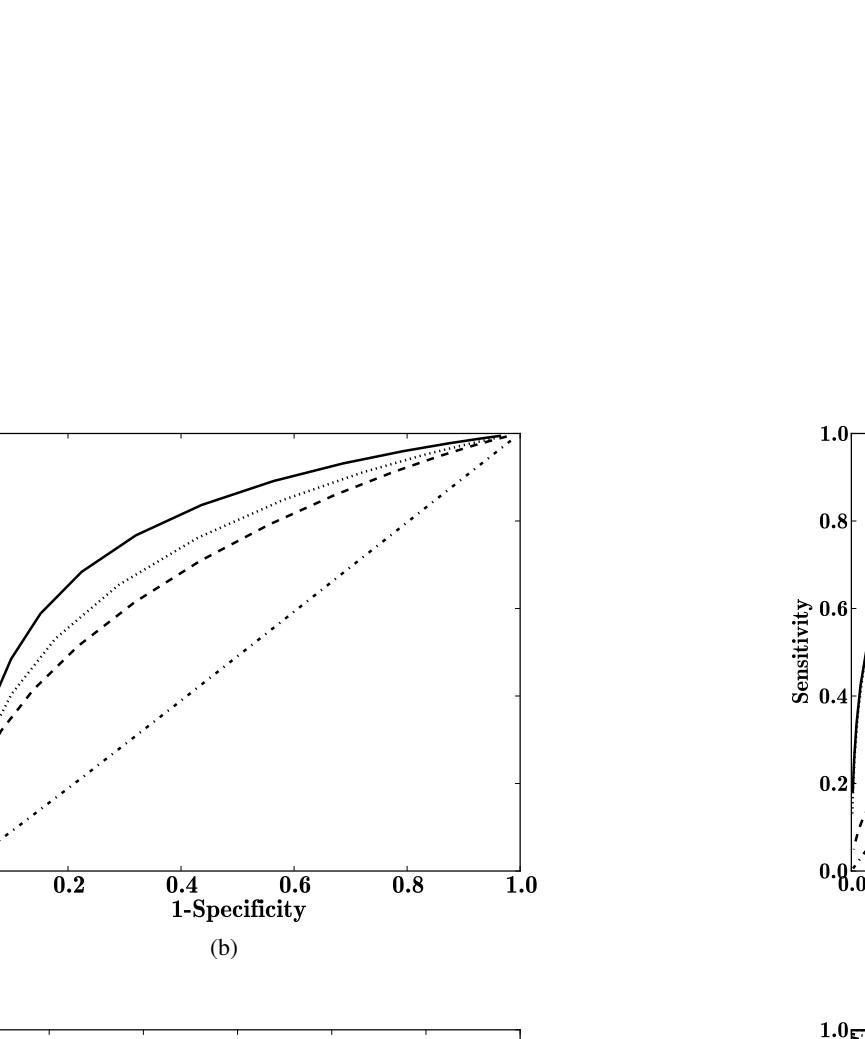


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

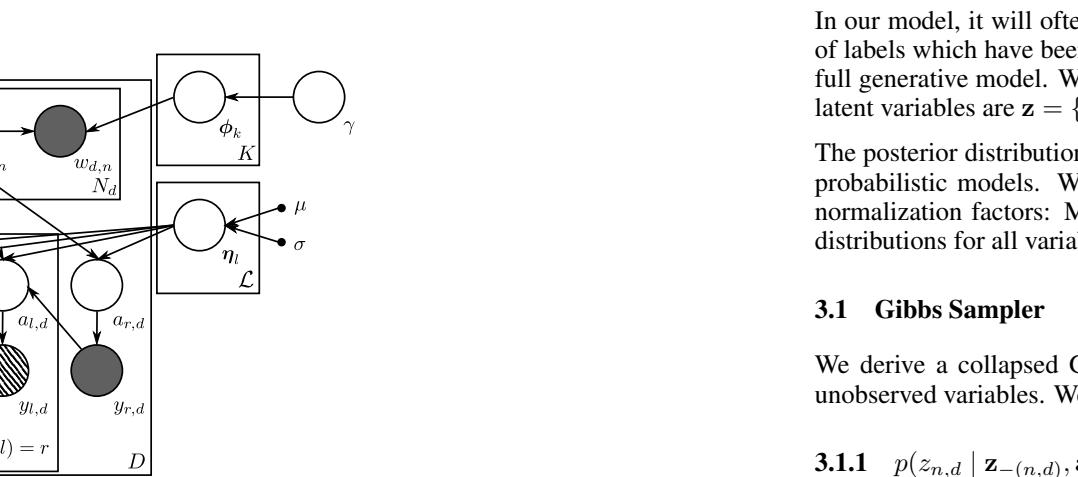


Figure 1: HSLDA graphical model

The is-a hierarchical constraint is a hard constraint that document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l .

The free parameters of the model are the number of topics, K , the number of unique words in the vocabulary, V , the number of documents, D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α^*, α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K -dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(1/\mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu\mathbf{1}_K, \sigma\mathbf{I}_K)$, where \mathbf{I}_K is the K -dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha^* \sim \text{Dir}_K(\alpha^*)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha\beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{l,K} \sim \text{Multinomial}(\phi_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$
 - Draw $a_{l,d} \mid z_d, \eta_l \sim \text{Beta}(\eta_l^{z_d}, 1 - \eta_l^{z_d})$
 - Draw $y_{l,d} \mid a_{l,d}, \eta_l \sim \text{Beta}(\eta_l a_{l,d}, 1 - \eta_l a_{l,d})$
 - Set the response variable $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{l,d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing.

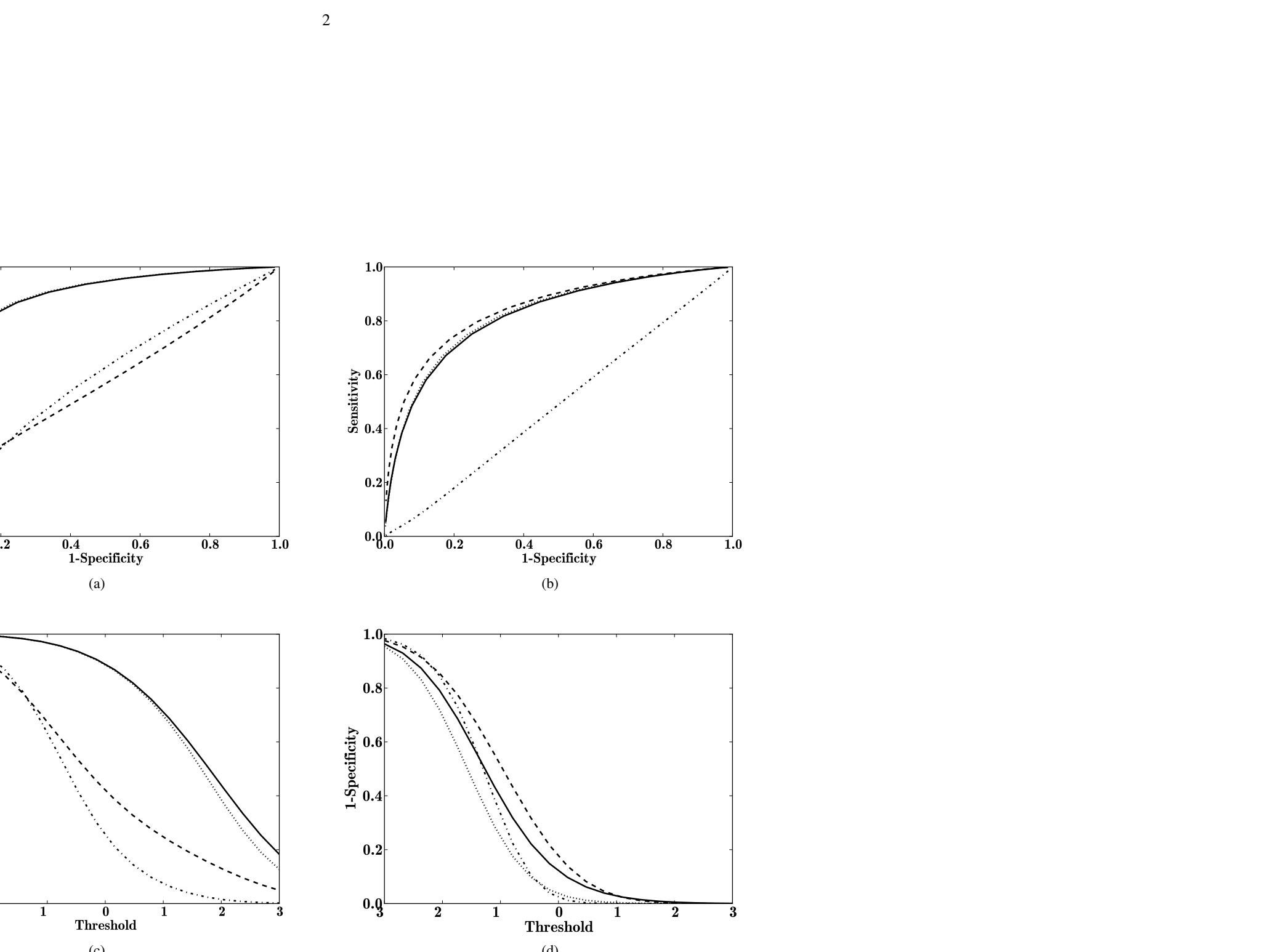


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l',d}$ and $y_{l',d}$ for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\phi_{l,K}$ and $\theta_{l,d}$ as in Griffiths and Steyvers [8]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1,N_d,d}\}_{d=1,\dots,D}$, $\boldsymbol{\eta} = \{\eta_l\}_{l \in \mathcal{L}}$, $\mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}, d \in \mathcal{L}_d, l=1,\dots,D}$, $\beta_{\mathcal{L}}$, α^* , and γ .

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus z_{n,d}$.

$$3.1.1 \quad p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include $\theta_{l,d}$ and $\phi_{l,K}$ because they have been integrated out in the collapsed Gibbs sampler [8]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [8] we find

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \frac{(c_{k,-(n,d)}^{k-n})^{a_{k,d}}}{(c_{k,-(n,d)}^{k-n})^{a_{k,d}} + \alpha \beta_k} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(\mathbf{z}_d^T \boldsymbol{\eta}_l - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $c_{v,-(n,d)}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{d,n} \mid \mathbf{z}_{-(d,n)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

$$3.1.2 \quad p(\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \sigma)$$

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean $\hat{\boldsymbol{\mu}}_l$ and covariance Σ such that

$$\Sigma^{-1} = \mathbf{I}^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\hat{\boldsymbol{\mu}}_l = \hat{\Sigma}^{-1} \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [?]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

$$3.1.3 \quad p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta_l)$$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta_l) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} \eta_l > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

$$3.1.4 \quad p(\beta \mid \mathbf{z}, \alpha^*, \alpha)$$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [12]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [11].

$$\beta \sim \text{Dir}(m_{(.,1} + \alpha^*, m_{(.,2} + \alpha^*, \dots, m_{(.,K} + \alpha^*) \quad (6)$$

$$p(m_{d,k} = m \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3

3.1.5 $p(\alpha), p(\alpha'), p(\gamma)$

The hyperparameters α , α' , and γ are given broad Gamma(1, 1000) prior distributions and sampled via the Metropolis-Hastings algorithm.

4 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating an exponential family response variable. Although there are many models for making predictions based on free text, LDA is unique in that it is a generative model, it represents documents as a mixed-membership, and constrains the inference of the latent structure of the documents by its predictability of the response variable. In other words, sLDA infers topics such that the model is capable of a high predictive likelihood for words in a document and the response variable associated with a document. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

There have been many models that incorporate both latent models of text and some form of supervision [10, 9, 13, 7]. One set of models that are particularly relevant to HSLDA are Chang and Blei's hierarchical models for document networks (Relational Topic Models). In that family of models, they encountered a similar scenario where the lack of a link did not truly indicate absence. In hierarchically labeled data, negative labels are uncommon and the lack of a label in the hierarchy is not equivalent to a negative label. Therefore, as in the work of Chang and Blei, we employ regularization to account for the lack of negative labels. This will be discussed further in 2.

5 Experiments

In the section we describe the application of HSLDA for prediction in two hierarchically structured domains. Firstly, we describe using discharge summaries to predict diagnoses, encoded as ICD-9 codes. Discharge summaries are documents that are authored by clinicians to summarize the course of a hospitalization. ICD-9 codes are used mainly for billing purposes to indicate the conditions for which a patient was treated. Secondly, we describe using Amazon.com product descriptions to predict product categories.

5.1 Data and Pre-Processing

5.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. The data consisted of free-text discharge summaries and corresponding ICD-9 codes. A discharge summary is a detailed report required by a physician or a hospital administrator at the conclusion of a hospital stay or series of treatments. The report outlines the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placements in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs [4], and patient records and their associated billing codes. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with tag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (LDA) [6] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per document "supervision"; often taking the form of a single numerical or categorical "label". More generally this supervision is just extra per document data; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. We call this Hierarchically Supervised Latent Dirichlet Allocation (HSLDA). We demonstrate HSLDA on two domains: medical discharge summaries and retail product categorization. We evaluate HSLDA's performance on two metrics: classification accuracy and topic quality.

The situation of each product in a product hierarchy (often multi-level) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $\mathbf{w}_{(n,d)} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{\mathbf{w}_{(1,d)}, \dots, \mathbf{w}_{(N_d,d)}\}$ be the set of N_d observations in document d . Let \mathcal{L}_d be the set of documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$ has a parent $\text{pa}(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set is hard ("is-a" parent child constraints (Explained Later)), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a response $y_{l,d} \in \{-1, 1\}$ to every label which indicates whether the label applies to document d or not.

The is-a-hierarchical constraint is a hard constraint that document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

1

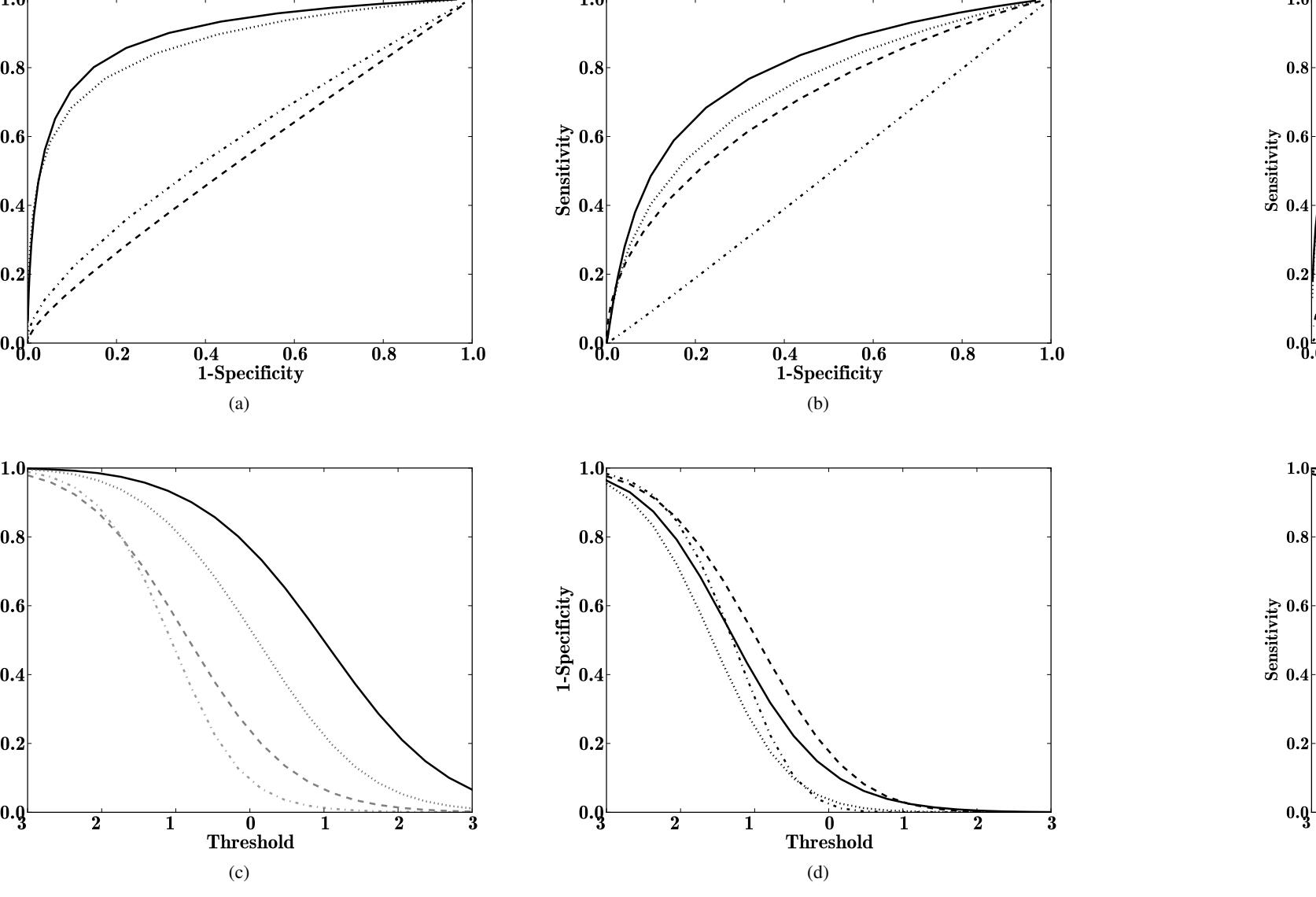


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

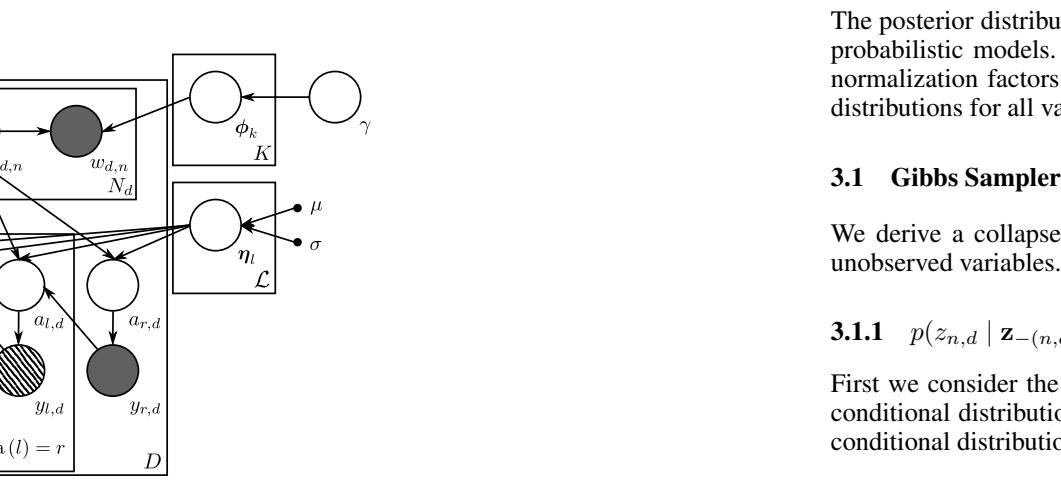


Figure 1: HSLDA graphical model

Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l .

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_{l,K} \sim \text{Multinomial}(\phi_{l,z_{n,d}})$
 - For each label $l \in \mathcal{L}$:
 - Draw $a_{l,d} | \mathbf{z}_d, \eta_l, y_{\text{pa}(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}^T \eta_l, 1), & y_{\text{pa}(l),d} = 1 \\ \mathcal{N}(\mathbf{z}^T \eta_l, 1) \mathbb{I}(a_{l,d} < 0), & y_{\text{pa}(l),d} = -1 \end{cases}$
 - where $\mathbf{z}_d = \sum_{n=1}^{N_d} \mathbf{z}_{n,d}$
 - Set the response variable $y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{\text{pa}(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as an auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing. In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{l',d}$ for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\phi_{l,K}$ and θ_d as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1:N_d,d}\}_{d=1,\dots,D}$, $\eta = \{\eta_l\}_{l \in \mathcal{L}}$, $\mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}_d, d=1,\dots,D}$, β , α' and γ .

2

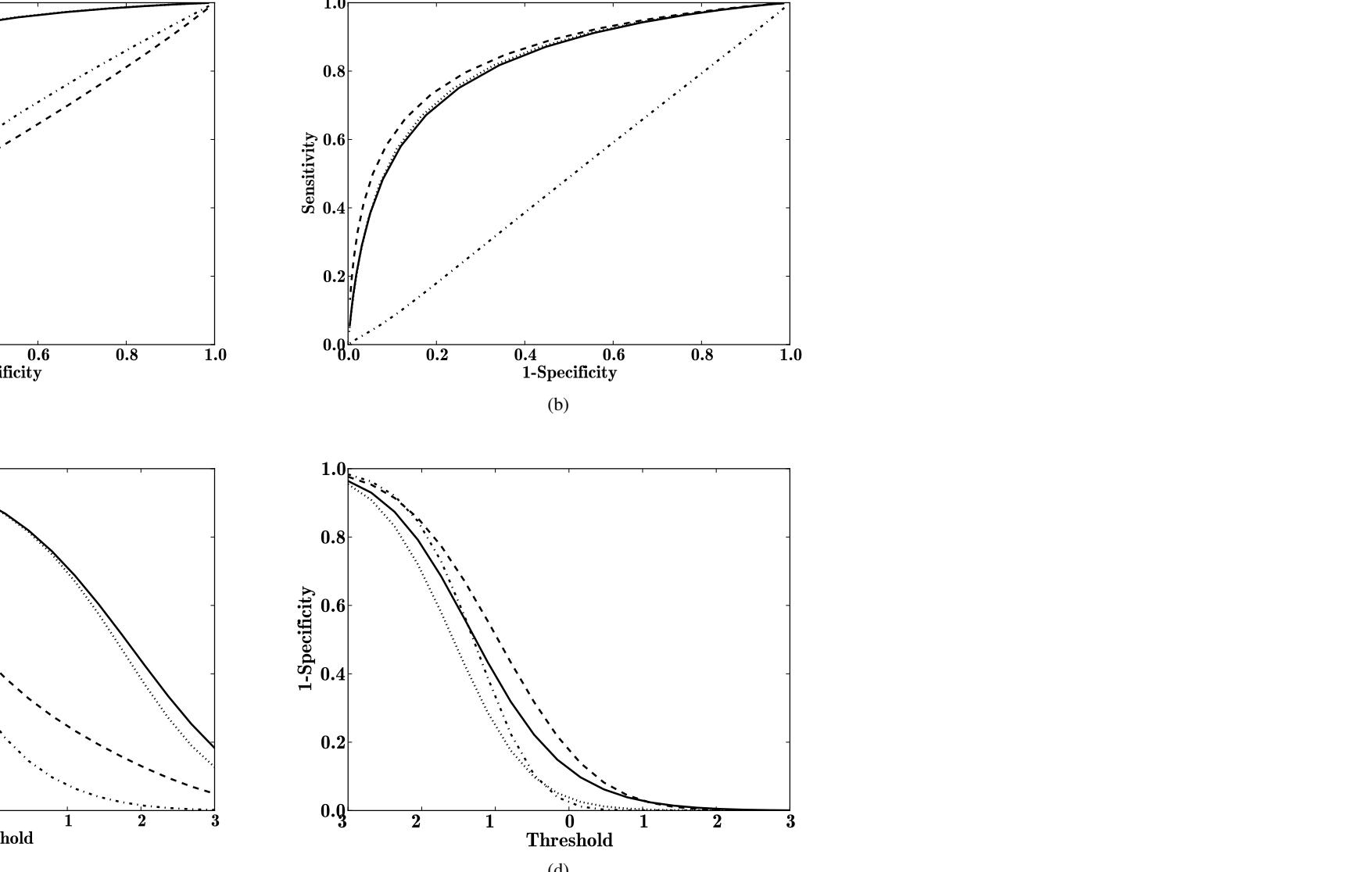


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

4 Related Work

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $z_{-(n,d)}$ to denote $z_d | z_{-(n,d)}$.

3.1.1 $p(z_{n,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$p(z_{n,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \frac{(c_{(n,d)}^{z_{n,d}})^{\alpha} \gamma^{\beta}}{(c_{(n,d)}^{z_{n,d}} + \alpha \beta_k) \left(\frac{c_{(n,d)}^{z_{n,d}}}{c_{(n,d)}^{z_{n,d}} + V} + V \right)} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(\mathbf{z}_d^T \boldsymbol{\eta}_l - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $c_{(n,d)}^{z_{n,d}}$ represents the number of words of type v in document d assigned to topic k omitting the n th word of document d . The notation $(\cdot)_d$ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{d,n} | z_{-(d,n)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\boldsymbol{\eta}_l | \mathbf{z}, \mathbf{a}, \sigma)$

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean μ_l and covariance Σ such that

$$\Sigma^{-1} = \mathbf{I}\sigma^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\mu_l = \Sigma \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [7]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

The task of learning ICD-9 codes has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge organized by [7]. The task in the challenge, however, differed from ours in scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [7, 2, 1]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [2, 2, 2, 1].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=309.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing. The text was tokenized with NLTK 2.0. Each discharge summary was tokenized with NLTK 2.0. Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

3.1.3 $p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \eta)$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta | \mathbf{z}, \alpha', \alpha)$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [13]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

Posterior inference is performed using the "direct assignment" method of Teh et al. [12].

$$\beta \sim \text{Dir}((m_{(.,1)} + \alpha', m_{(.,2)} + \alpha', \dots, m_{(.,K)} + \alpha') \quad (6)$$

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nlm.nih.gov>

5 Experiments

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noemie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp909@dbm, bartlett@stat, noemiel@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs [e.g. [1] as available from [4]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [3]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend Latent Dirichlet Allocation (LDA) [5] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [6] augmented with labels on documents taking the form of a single numerical or categorical "is-related-to". More generally this supervision is just extra per document data, for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been provided/drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [5].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider global web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$.

Each label $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We want to express our purpose that this label set has hard "is-a" parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a response $y_{l,d} \in \{-1, 1\}$ to every label which indicates whether the label applies to document d or not.

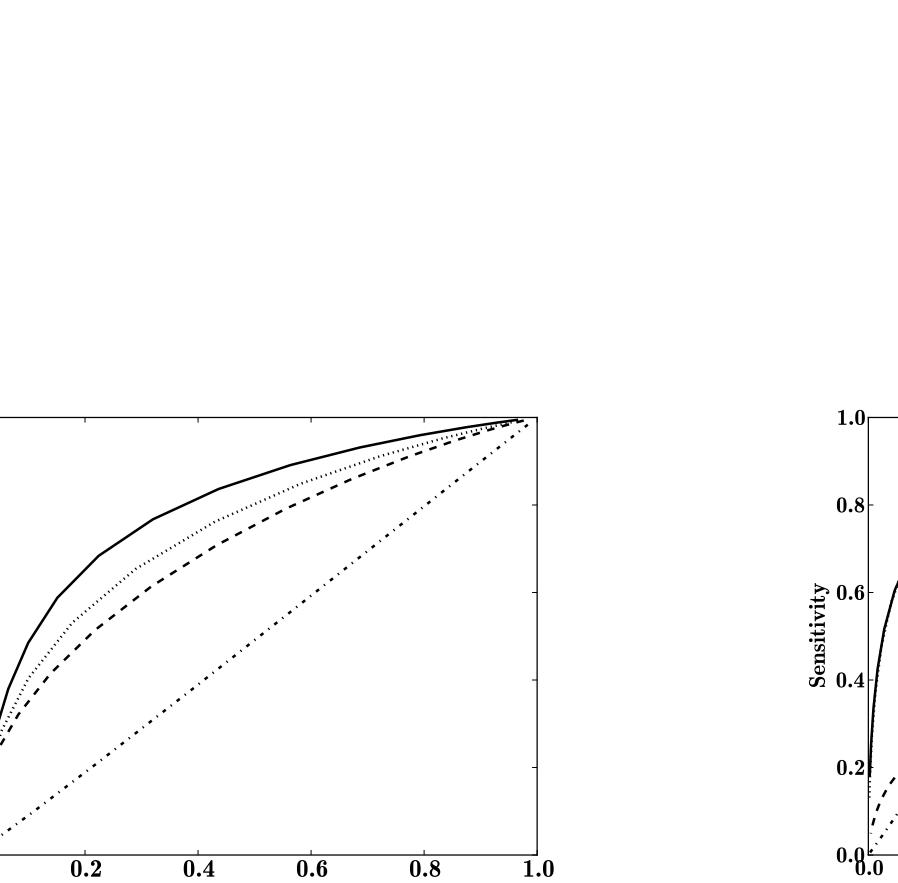


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.



Figure 1: HSLDA graphical model

The is-a hierarchical constraint is a hard constraint that document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l .

The free parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $Dir_K(\cdot)$ and the K dimensional normal distribution as $N_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is shown in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim Dir(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim N_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where $\mathbf{1}_K$ is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim Dir_K(\alpha' \mathbf{I}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim Dir_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim Multinomial(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{l,K} \sim Multinomial(\phi_{z_{n,d}})$
 - For each label $l \in \mathcal{L}$
 - Draw $a_{l,d} \mid z_d, \eta_l \sim p(a_{l,d} \mid \eta_l)$
 - Draw $y_{l,d} \mid z_d, \eta_l, pa(l), d \sim \begin{cases} \mathcal{N}(z_d^T \eta_l, 1), & y_{pa(l)} = 1 \\ \mathcal{N}(z_d^T \eta_l, 1) \mathbb{I}(a_{l,d} < 0), & y_{pa(l)} = -1 \end{cases}$
 - where $\bar{z}_d = N_d^{-1} \sum_{n=1}^{N_d} z_{n,d}$
 - Set the response variable $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{pa(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution - the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing.

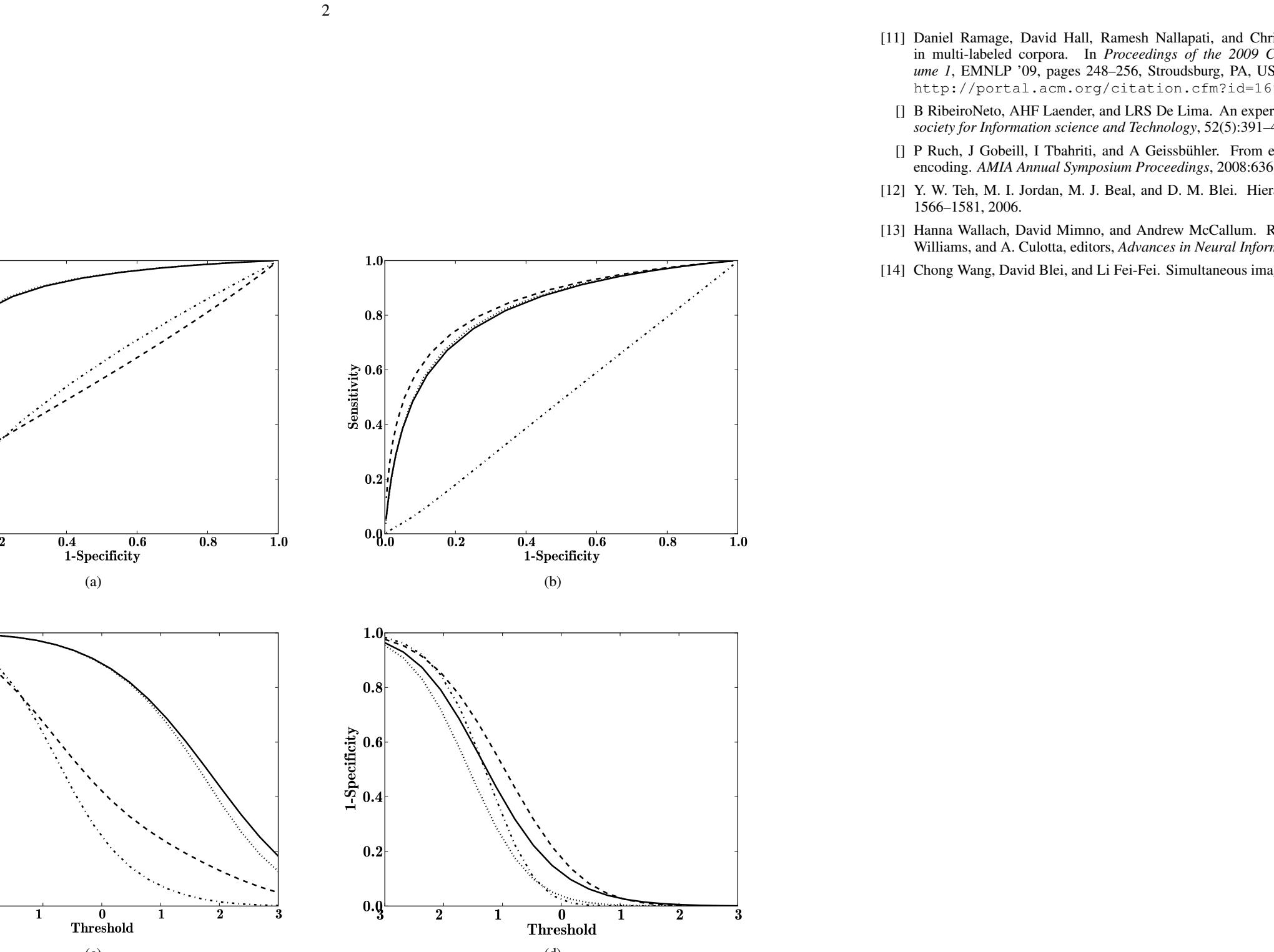


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l',d}$ and $y_{l',d}$ for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters $\phi_{l',K}$ and $\theta_{l',d}$ as in Griffiths and Steyvers [9]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1,N_d,d}\}_{d=1,\dots,D}$, $\boldsymbol{\eta} = \{\eta_l\}_{l \in \mathcal{L}}$, $\mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}, d \in \mathcal{L}_d, l=1,\dots,D}$, β, α, σ , and γ .

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus z_{n,d}$.

$$3.1.1 \quad p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \eta_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [9] we find

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \frac{(c_{(k,-(n,d))}^{k-n-(a_{l,d})} + \alpha \beta_k)}{(c_{(k,-(n,d))}^{k-n-(a_{l,d})} + \alpha \beta_k) \exp \left\{ -\frac{(z_d^T \eta_l - a_{l,d})^2}{2} \right\}}. \quad (2)$$

Here, $c_{(k,-(n,d))}^{k-n-(a_{l,d})}$ represents the number of words of type v in document d assigned to topic k omitting the n^{th} word of document d . The notation (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma)$ can be sampled through enumeration.

$$3.1.2 \quad p(\eta_l \mid \mathbf{z}, \mathbf{a}, \sigma)$$

We now consider the conditional distribution of the regression coefficients η_l for $l \in \mathcal{L}$. Given that η_l and $a_{l,d}$ are distributed normally, the posterior distribution of η_l is normally distributed with mean $\hat{\mu}_l$ and covariance Σ such that

$$\Sigma^{-1} = \mathbf{I}^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\hat{\mu}_l = \hat{\Sigma}^{-1} \left(\frac{l}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [?]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \bar{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

$$3.1.3 \quad p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta)$$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^T \bar{z}_d) \right\} \mathbb{I}(a_{l,d} \eta_l^T \bar{z}_d > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

$$3.1.4 \quad p(\beta \mid \mathbf{z}, \alpha')$$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [13]. This prior shares many features with the hierarchical Dirichlet process and inference over the introduction to make exact Gibbs sampling possible and are not of primary interest.

Posterior inference is performed using the "direct assignment" method of Teh et al. [12].

$$\beta \sim Dir((m_{(j,1)} + \alpha', m_{(j,2)} + \alpha', \dots, m_{(j,K)} + \alpha')) \quad (6)$$

$$p(m_{(d,k)} = m \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + n_{d,k})} s(n_{d,k}, m) (\alpha \beta_k)^m \quad (7)$$

where $s(n, m)$ represents stirling numbers of the first kind.

3 Experiments

3.1 Data and Pre-Processing

3.1.1 Diagnosis Prediction

Our data set was gathered from the clinical data warehouse of New York-Presbyterian Hospital. The data consisted of free-text discharge summaries and ICD-9 codes from 1999 to 2004. A discharge summary is a clinical report required by a physician or a medical professional at the conclusion of a hospital stay or series of treatments. The report outlines the patient's chief complaint, diagnostic findings, therapy administered, patient's response to the chosen therapy, the treatment plan and the recommendations upon discharge. The ICD-9 codes used to structure the discharge summary data are part of the international standard diagnostic classification for epidemiological, health management, and clinical purposes. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code representing "Pneumonia due to adenovirus" is a child of the code representing "Viral pneumonia" where the former is a type of the latter. In the hospital, ICD-9 codes are generated manually by trained medical coders, who review all the information in the discharge summary.

The text of the discharge summaries were pre-processed such that each document would be represented as counts over a 10,000 word vocabulary. The Natural Language Toolkit was used to tokenize the text. A vocabulary was identified by first sorting terms based on a global term-frequency-inverse document frequency measure. The top 10,000 words which were not identifying in some way (name, place, or identifying number) were selected for inclusion in the vocabulary.

For each hospitalization there are usually several ICD-9 codes assigned for billing purposes. These codes are known to be quite specific but not very sensitive [8]. Regardless of that fact, this is one of the only sources for information on patient diagnoses aside from the free text.

We worked within the guidelines of the Health Insurance Portability and Accountability Act (HIPAA), which protects patient privacy and the security of potentially identifying medical material, known as personal health information (PHI). HIPAA covers any information within a medical record that was created, used, or disclosed during the course of providing a health care service and that can be used to identify an individual. This study was approved by the Institutional Review Board.

3.1.2 Product Category Prediction

Data for these experiments were obtained partially from the Stanford Network Analysis Platform (SNAP) Amazon product metadata dataset [4] and partially directly from the Amazon.com website [1]. The product ID's and categorizations were obtained from the SNAP dataset and the product descriptions were obtained directly from the website.

We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine". Each product was labeled with multiple categories.

The vocabulary for this experiment was created by including the most frequent 30K words omitting stopwords.

4 Related Work

Latent dirichlet allocation (LDA) is a generative

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placements in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-words data is also of interest. We find that using the signal from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs [3], and patient records and their associated billing codes. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (LDA) [4] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [5] augmented with per document “supervision”; often taking the form of a single numerical or categorical “label.” More generally this supervision is just extra per document data; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. We call this Hierarchically Supervised Latent Dirichlet Allocation (HSLDA). We demonstrate this model on two domains: medical discharge summaries and retail product categorization. We evaluate model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter for regression parameters.

The situation of each product in a product hierarchy (often multi-level) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $W_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let \mathcal{L}_d be the set of d documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set is hard “leaf” parent child constraints (Explained Later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a response $y_{l,d} \in \{-1, 1\}$ to every label which indicates whether the label applies to document d or not.

The is-a-hierarchical constraint is a hard constraint that document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

1

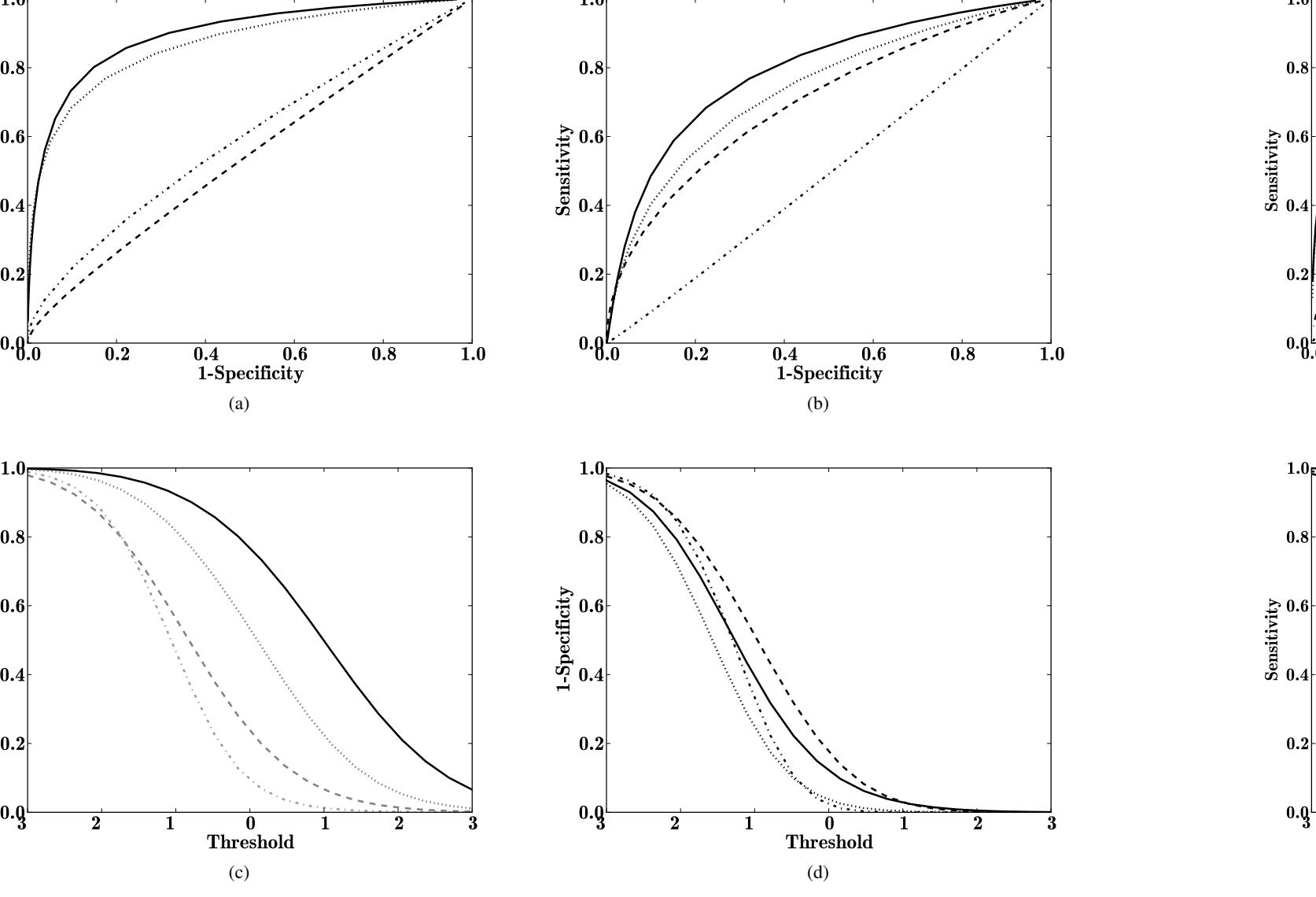


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

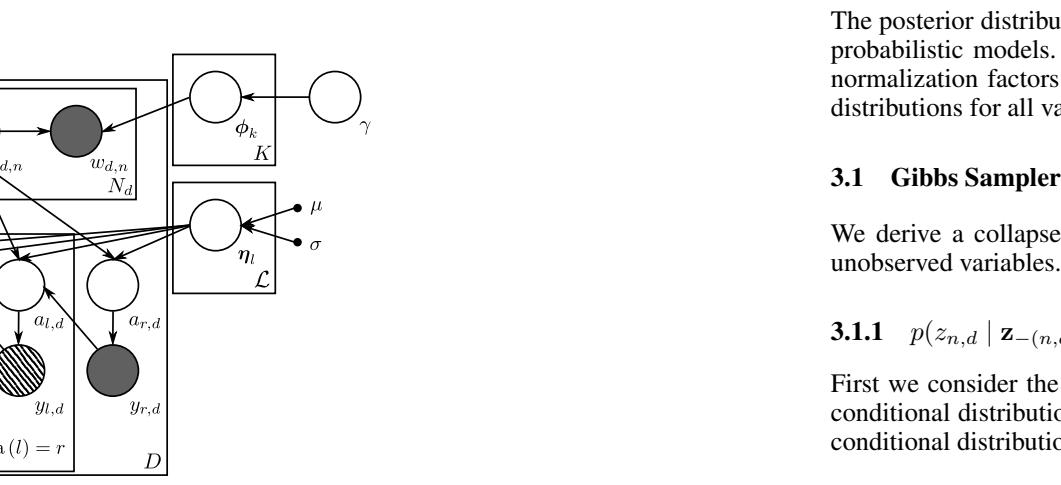


Figure 1: HSLDA graphical model

Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l .

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure 1.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(1/\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_{l,z} \sim \text{Multinomial}(\phi_{l,z})$
 - For each label $l \in \mathcal{L}$:
 - Draw $a_{l,d} | \mathbf{Z}_d, \eta_l, y_{\text{pred}(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{Z}^T \eta_l, 1), & y_{\text{pred}(l)} = 1 \\ \mathcal{N}(\mathbf{Z}^T \eta_l, 1) \mathbb{I}(a_{l,d} < 0), & y_{\text{pred}(l)} = -1 \end{cases}$
 - where $\mathbf{Z}_d = \sum_{n=1}^{N_d} \mathbf{z}_{n,d}$
 - Set the response variable $y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{\text{pred}(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution – the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as an auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

3 Inference

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing. In our model, it will often be the case that the set of labels \mathcal{L} is not fully observed for every document. We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ and $y_{l,d}$ for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ from the full generative model. We can also integrate out the parameters θ_d and η_l for $l' \in \mathcal{L} \setminus \mathcal{L}_d$ in Griffiths and Steyvers [12]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1:N_d,d}\}_{d=1,\dots,D}$, $\boldsymbol{\eta} = \{\eta_l\}_{l \in \mathcal{L}}$, $\mathbf{a} = \{a_{l,d}\}_{l \in \mathcal{L}_d, d=1,\dots,D}$, β , α , α' and γ .

2

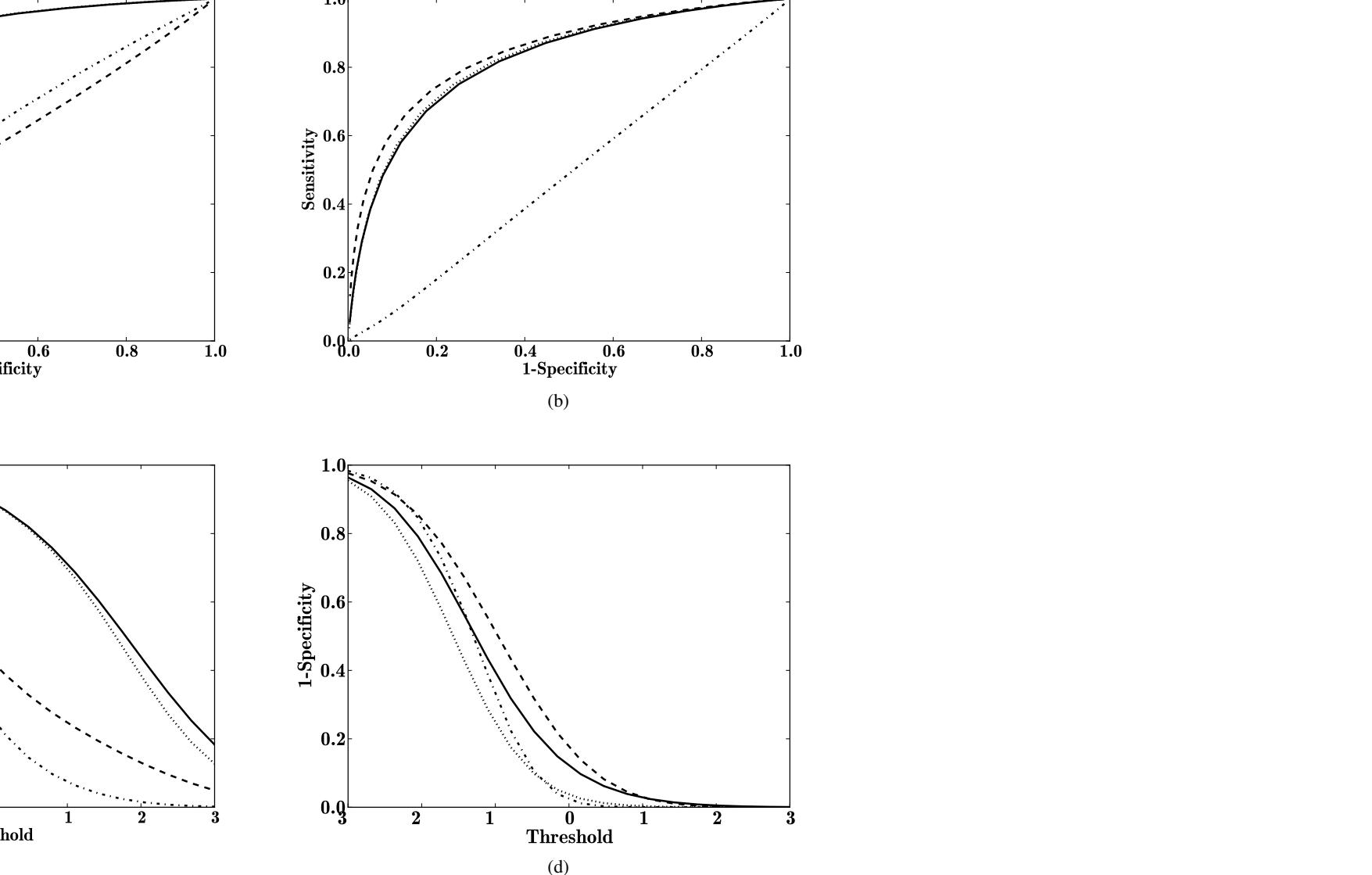


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

4 Related Work

The posterior distribution we seek cannot be solved in closed form. This is often the case in evaluating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $z_{-(n,d)}$ to denote $z_d | z_{<n,d}$.

3.1.1 $p(z_{n,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include $\theta_{l,d}$ and $\phi_{l,z}$ because they have been integrated out as in the collapsed Gibbs sampler [12]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [12] we find

$$p(z_{n,d} = k | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \frac{(c_{(.,d)}^{z_{n,d}-1} \gamma)^{\alpha}}{(c_{(.,d)}^{z_{n,d}} + \alpha \beta_k)} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(\mathbf{z}_d^T \boldsymbol{\eta}_l - a_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $c_{(.,d)}^{z_{n,d}}$ represents the number of words of type v in document d assigned to topic k omitting the n th word of document d . The notation $(.)$ in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 2, $p(y_{l,d} | z_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\boldsymbol{\eta}_l | \mathbf{z}, \mathbf{a}, \sigma)$

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean μ_l and covariance Σ such that

$$\Sigma^{-1} = \mathbf{I}\sigma^{-2} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\mu_l = \Sigma \left(\frac{\mu}{\sigma^2} + \frac{\mathbf{Z}^T \mathbf{a}}{\sigma^2} \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [?]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{a} = [a_{1,d}, a_{2,d}, \dots, a_{D,d}]^T$.

3.1.3 $p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \mathbf{Y}, \eta)$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} | z_{-(n,d)}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

3.1.4 $p(\beta | \mathbf{z}, \alpha', \alpha)$

In our model, we place a hierarchical Dirichlet prior over topic assignments. This flexible distribution allows for an asymmetric prior over document level distributions over topics [21]. This prior shares many features with the hierarchical Dirichlet process and inference over this distribution proceeds in a very similar fashion.

3.1.5 Product Category Prediction

In this experiment, we look at product descriptions and their categorizations according to a product hierarchy. Product ID's and categorizations were obtained from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were crawled from the Amazon.com website directly. We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, “DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine.”

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nltk.org>

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

6 Results

7 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

5.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. We will refer to the three comparison models as the sLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen as to highlight performance relative to the hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable to regression parameters.

5.3 Evaluation

For each dataset, a held out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but we show lower-dimensional representations of the bag-of-words data is also of interest and that using a size of size from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on hierarchical structures of the same [2], product descriptions and catalogs [3], and patient records and their associated billing codes. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-word image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [4] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [5] augmented with per document "supervision"; often taking the form of a single numerical or categorical "label." More generally this supervision is just extra per document data; for instance its quality or relevance (e.g. online review scores), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathcal{W}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label l except the root $l_0 \in \mathcal{L}$ is a parent node $\in \mathcal{L}$ also in the set of labels. We consider the problem of learning that this label set has "leaf" parent-child constraints (explained later) although this assumption is not relied on the cost of more complex models. Such a label hierarchy forms a multiply-rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d , i.e., $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

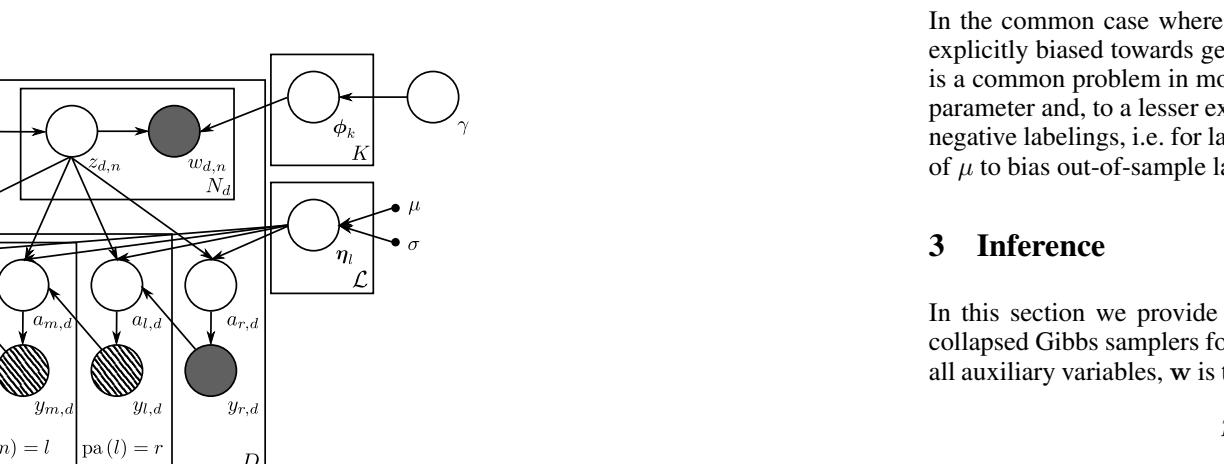


Figure 1: HSLDA graphical model

as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA "topics" (distributions over the elements of Σ). ϕ_d is a distribution over "words," θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_{l,z_{n,d}} \sim \text{Multinomial}(\phi_{l,z_{n,d}})$
 - Set $y_{r,d} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_d, \eta_l \sim \mathcal{N}(\eta_l^T z_d, 1)$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Apply label l to document d according to $a_{l,d}$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_d, \eta_l \sim \mathcal{N}(\eta_l^T z_d, 1)$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Apply label l to document d according to $a_{l,d}$

is also a standard probit regression result [?].

HSLDA departs from stock LDA in that we estimate a hierarchical Dirichlet prior over topic assignments (i.e., β is a parameter in our model). This was shown to ADLER [21]. Sampling β is done using the "direct assignment" method of Teh et al. [20]

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be slightly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing

5.1.1 Diagnosis Prediction

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient's complaints, findings and diagnoses, along with treatment and hospital course. For each admission trained medical coders review the information in the discharge summary and assign a series of diagnosis codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a labeling of a patient's diagnoses based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for "Pneumonia due to adenovirus" is a child of the code for "Viral pneumonia," where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [8], and sometimes make mistakes [10].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP community challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 test documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [7, 11, 19]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [14, 18, 16, 19, 15].

The dataset was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.3 average tokens per sentence on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=500.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Category Prediction

In this experiment, we look at product descriptions and their categorizations according to a product hierarchy. Product ID's and categorizations were obtained from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were crawled from the amazon.com website directly. We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine."

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

5.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by running LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. We will refer to the three comparison models as the sLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen to highlight performance in absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter for regression parameters.

5.3 Evaluation

For each dataset, a held out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (?) and the regression coefficients (?) which are analytic. This simplifies posterior inference substantially.

¹<http://www.cdc.gov/nchs/icd/icd9cm.htm>

²<http://www.nltk.org>

The two measures for predictive performance used here include the true positive rate and the false positive rate evaluated based on $p(y_{l,d} | w_{1:N,d})$ for each label in each model.

6 Results

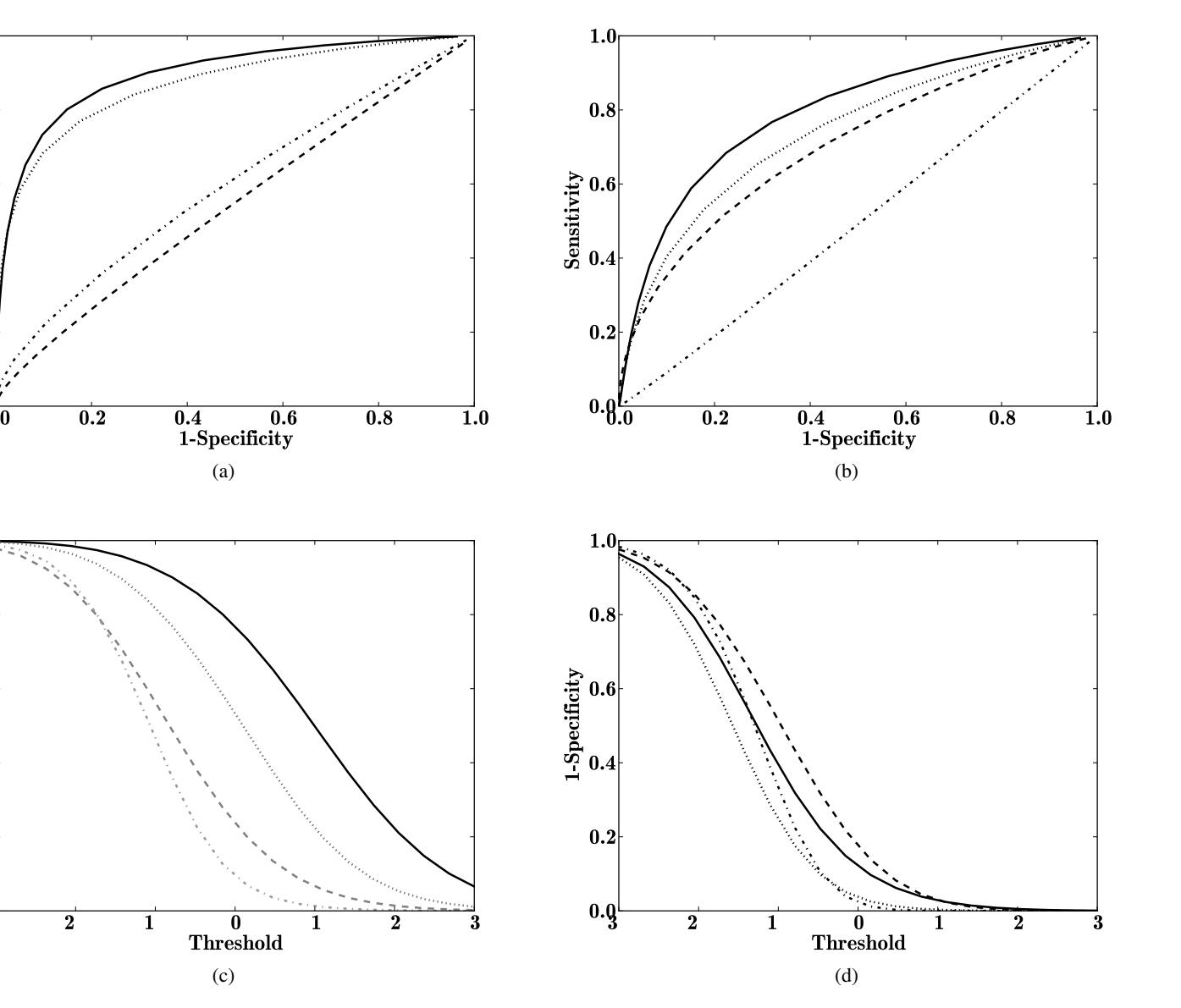


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

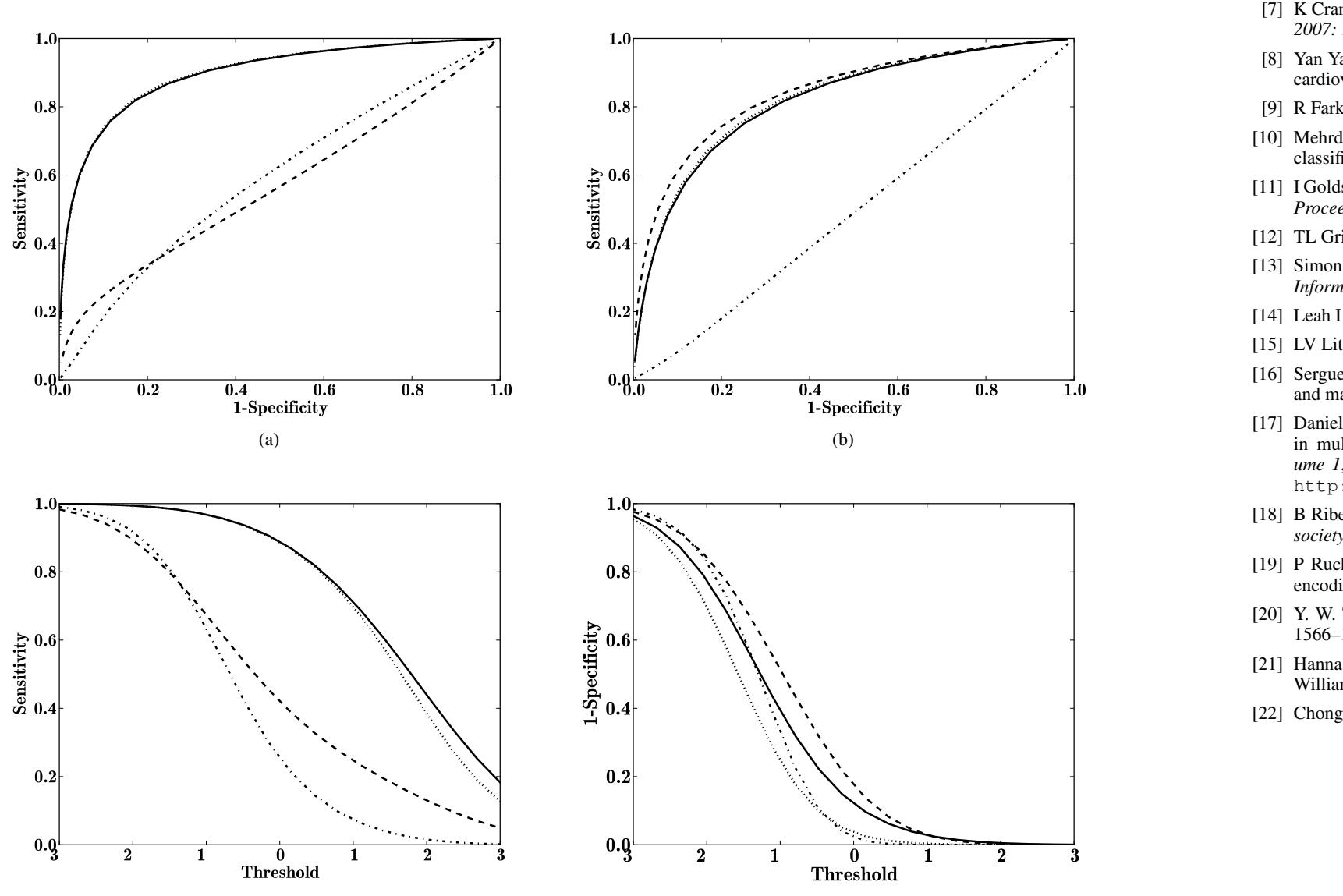


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noe@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but we show lower-dimensional representations of the bag-of-words data is also of interest and that using a size of sum from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured medical data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and catalogues of hierarchical directories of the same [2], product descriptions and catalogs [3], and patient records and their associated billing codes. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [4] to take advantage of hierarchical supervision, sLDA is latent Dirichlet allocation (LDA) [5] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label." More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiply labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label l is a parent node in \mathcal{L} , and a parent node l in \mathcal{L} also in the set of labels. We make no explicit assumption that this label set has a "root" element; parent-child constraints (explained later) are, although not required, at the cost of more complexity. Such a label hierarchy forms a multiply-rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an $is-a$ label hierarchy are that if the l th label is applied to document d , i.e. $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{pa(l),d} = 1, y_{pa(pa(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

1

2

6 Results

7 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

..

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
- [2] DM02 open directory project. <http://www.dmoz.org/>, 2002.
- [3] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [4] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [5] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [6] Jonathan Chang and David M. Blei. Hierarchical models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [7] K Crammer, M Dredze, K Ganchev, PP Talukdar, and S Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [8] Yan Yan David S. Nilesena Martha J. Radford Brian F. Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [9] R Farkas and G Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [10] Mehrdad Farzandipour, Abbas Sheikhtabari, and F. Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:78–84, 2010.
- [11] I Goldstein, A Arzumanyan, and O Uzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [12] T.L. Griffiths and M. Steyvers. Finding scientific topics. *PNAI*, 101(suppl. 1):S228–S235, 2004.
- [13] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [14] Leah Larkey and Bruce Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [15] LV Liita, S Yu, S Niculescu, and J. B. Large scale diagnostic code classification for medical patient records, 2008.
- [16] Serguei Pakhomov, James Buntrock, and Christopher Chute. Automating the assignment of diagnosis codes to patient encounters using example-based and machine learning techniques. *Journal of the American Medical Informatics Association (JAMIA)*, 13(5):516–525, 2006.
- [17] Daniel Ramage, David Hall, Ramesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL <http://portal.acm.org/citation.cfm?id=1699510.1699543>.
- [18] B RiberaNeto, AHF Laender, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and technology*, 52(3):391–401, 2001.
- [19] P Radchenko, J. Mitchell, I. Tshishirini, and A. Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic coding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [20] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [21] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.
- [22] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

6



Figure 1: HSLDA graphical model

as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plotted variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA "topics" (distributions over the elements of Σ), ϕ_k is a distribution over "words," θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $\mathcal{N}_K(\cdot)$ is the K -dimensional Normal distribution, I_K is the K -dimensional identity matrix, \mathbf{a}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $m = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_{v,d} \sim \text{Multinomial}(\phi_{v,d})$
 - Set $y_{p(l),d} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_{l,d}, \eta_l \sim \mathcal{N}(\mathbf{z}_{l,d}^T \eta_l, 1)$
 - If $a_{l,d} < 0$, then $y_{p(l),d} = -1$
 - Apply label l to document d according to $a_{l,d}$
 - Draw $a_{l,d} | z_{l,d}, \eta_l \sim \mathcal{N}(\mathbf{z}_{l,d}^T \eta_l, 1) \mathbb{I}(a_{l,d} < 0)$
 - Apply label l to document d according to $a_{l,d}$
 - Draw $y_{l,d} | a_{l,d}, \eta_l \sim \mathcal{N}(\mathbf{z}_{l,d}^T \eta_l, 1) \mathbb{I}(a_{l,d} > 0)$

Here $\mathbf{z}_d^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The $(\cdot)_-$ in the subscript means the count resulting from summing over the omitted subscript variable. Also \mathcal{L}_d is the set of labels which are observed for document d .

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model will be biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

3 Inference

In this section we provide the conditional distributions required to Gibbs sample the HSLDA posterior distribution. Note that like in collapsed Gibbs samplers for LDA [2], we have analytically marginalized out the parameters $\phi_{v,d}$ and $\theta_d|d$. In the following \mathbf{a} is the set of all auxiliary variables, \mathbf{w} is the set of all words, η is the set of all regression coefficients, and $z_{n,d}$ is the $z_{n,d}$ with element $z_{n,d}$ removed.

3.1 Data and Pre-Processing

5.1 Diagnosis Prediction

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient's complaints, findings and diagnoses, along with treatment and hospital course. For each admission trained medical coders review the information in the discharge summary and assign a series of diagnoses codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a child-parents relationship based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for "Pneumonia due to adenovirus" is a child of the code for "Viral pneumonia," where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [8], and sometimes make mistakes [10].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP competition challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [7, 11, 9]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naïve Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [14, 18, 16, 19, 15].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 average ICD-9 codes on average (std dev=5.01) and contain an average of 53.67 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² Vocabulary was determined as the top 10,000 tokens with highest frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Category Prediction

In this experiment, we look at product descriptions and their categorizations according to a product hierarchy. Product ID's and categories were obtained from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were crawled from the amazon.com website directly. We were able to deduce the structure of the hierarchy for the Amazon.com products directly since all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD. "The Time Machine."

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

5.2 Comparison Models

We applied HSLDA, along with three other closely related models, to the clinical data and the retail product data. Specifically, we evaluate models including SLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. We will refer to the three comparison models as the SLDA model, the separate HSLDA model, and the random HSLDA model, respectively. These models were chosen as to highlight performance in absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter for regression parameters.

5.3 Evaluation

For each dataset, a held out set of 1,000 documents and accompanying labels were used for evaluation. The two main methods of evaluation for the model are prediction and topic quality. To evaluate predictive performance for all comparison models equivalently, each model was evaluated with two methods. The first evaluation method augments the observed labels in the held out set with their ancestors and considers all other non-existent labels to be negative. The second method ignores the ancestors of the observed labels in the held out set and considers all other non-existent labels to be negative. This uniform treatment of ancestors allows for a fair comparison of the models.

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA
 {ajp9009@dbmi, bartlett@stat, noe@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal of this work, but we show how lower-dimensional representations of the bag-of-words data is also of interest and that using a size of size from hierarchical labels substantially improves out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include news feeds, web pages, and entries in hierarchical directories of the same [1], product descriptions and catalogs [3], and patient records and their associated billing codes. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [4] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [5] augmented with per document “supervision”; often taking the form of a single numerical or categorical “label.” More generally this supervision is just extra per document data: for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical (and often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$ is a parent node in \mathcal{L} and a parent node in \mathcal{L} also in the set of labels. We make no prior assumptions about whether this label set has “leaf” or “parent-child” constraints (explained later), although this assumption is not relevant at the cost of more complexity. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an $is-a$ label hierarchy are that if the l th label is applied to document d , i.e. $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{pa(l),d} = 1, y_{pa(pa(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

1

2

6 Results

7 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [4] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [5] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [6] Jonathan Chang and David M. Blei. Hierarchical models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [7] K Crammer, M Dredze, K Ganchev, PP Talukdar, and S Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [8] Yan Yan David S. Nilesen Martha J. Radford Brian F. Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [9] R Farkas and G Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [10] Mehrad Farzandipour, Abbas Sheikhtabar, and F. Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:78–84, 2010.
- [11] I Goldstein, A Arumugam, and O Uzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [12] T.L. Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [13] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [14] Leah Larkey and Bruce Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [15] LV Liita, S Yu, S Niculescu, and J. B. Large scale diagnostic code classification for medical patient records, 2008.
- [16] Serguei Pakhomov, James Buntrock, and Christopher Chute. Automating the assignment of diagnosis codes to patient encounters using example-based and machine learning techniques. *Journal of the American Medical Informatics Association (JAMIA)*, 13(5):516–525, 2006.
- [17] Daniel Ramage, David Hall, Ranesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL <http://portal.acm.org/citation.cfm?id=1699510.1699543>.
- [18] B RibeinNeto, AHF Laender, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and technology*, 52(3):391–401, 2001.
- [19] P Ravinder, J. Russell, I. Tsabirli, and A. Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic coding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [20] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476): 1566–1581, 2006.
- [21] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.
- [22] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

6

¹<http://www.cdc.gov/nchs/icd/icd9cm.htm>

²<http://www.nltk.org>

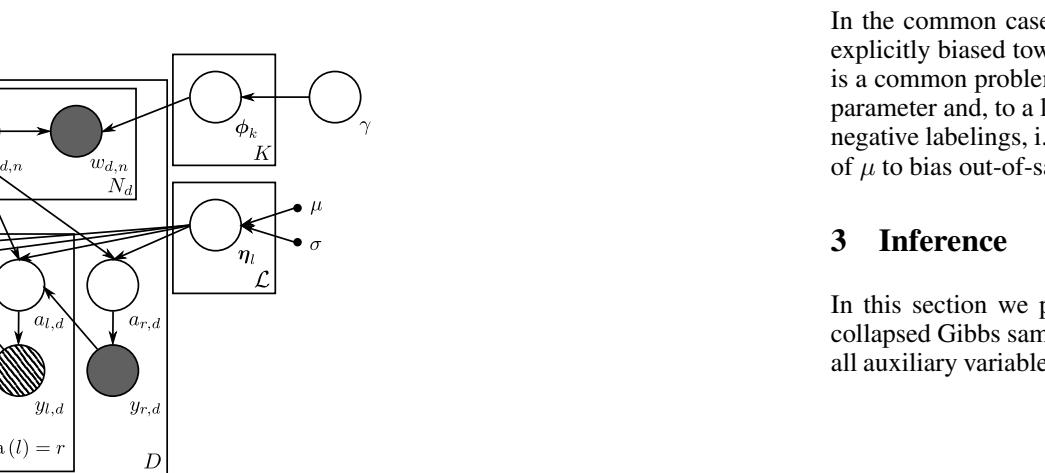


Figure 1: HSLDA graphical model

as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA “topics” (distributions over the elements of Σ). θ_d is a distribution over “words.” θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $\mathcal{I}_K(\cdot)$ is the K -dimensional Normal distribution, I_K is the K -dimensional identity matrix, $\mathbf{1}_d$ is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw $a_{n,d} | \theta_d, \eta_l \sim \text{Dir}(m_{(l),1} + \alpha', m_{(l),2} + \alpha', \dots, m_{(l),K} + \alpha')$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_{l,z_{n,d}} \sim \text{Multinomial}(\phi_{l,z_{n,d}})$
 - Set $y_{l,d} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_d, \eta_l, y_{pa(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}_d^T \eta_l, 1), & y_{pa(l),d} = 1 \\ \mathcal{N}(\mathbf{z}_d^T \eta_l, 1)(\mathbf{z}_d^T \eta_l) & y_{pa(l),d} = -1 \end{cases}$
 - Apply label l to document d according to $a_{l,d}$
 - $y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here $\mathbf{z}_d^T = [z_1, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k , $z_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generatively labeled using a hierarchy of conditionally dependent probit regressors [3, 7]. For every label l to be applied to document d and whether or not its parent label was applied (i.e. $\mathbb{I}(y_{pa(l),d} = 1)$) are used to decide whether or not label l is to be applied to document d and whether or not its parent label $pa(l)$ is to be applied to document d if its parent label $pa(l)$ is to be applied to document d . These regressions are specific to the constraints but can be modified to account for different constraint sets. The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model is biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5.1 Data and Pre-Processing

5.1.1 Diagnosis Prediction

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient's complaints, findings and diagnoses, along with treatment and hospital course. For each admission trained medical coders review the information in the discharge summary and assign a series of diagnoses codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a partial ordering of the diagnoses based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [8], and sometimes make mistakes [10].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP competition challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [7, 11, 9]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naïve Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [14, 18, 16, 19, 15].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=309.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² Vocabulary was determined as the top 10,000 tokens with highest frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Category Prediction

In this section we look at product descriptions and their categorizations according to a product hierarchy. Product ID's and categories were obtained from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were crawled from the amazon.com website directly. We have been shown to improve the quality and stability of inferred topics [21]. Sampling β is done using the “direct assignment” method of Teh et al. [20].

$$\beta | \mathbf{z}, \alpha', \alpha \sim \text{Dir}(m_{(1),1} + \alpha', m_{(1),2} + \alpha', \dots, m_{(1),K} + \alpha') \quad (4)$$

where $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β and are governed by the following function.

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + c_{d,k})} s(c_{d,k}^k) (\alpha \beta_k)^m \quad (5)$$

$s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

5.2 Comparison Models

We applied HSLDA, along with three other closely related

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
 Columbia University, New York, NY 10027, USA noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal, but we believe lower-dimensional representations of the bag-of-words data is also of interest. We show that using a size from hierarchical labels substantially improves our out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and entries in hierarchical directories of the same [2], product descriptions and catalogs [3], and patient records and their associated billing codes. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [4] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [5] augmented with per document “supervision”; often taking the form of a single numerical or categorical “label”. More generally this supervision is just extra per document data; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multi-staged) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathcal{W}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the states in the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$ is a parent of $l' \in \mathcal{L}$ if $l \in \text{pa}(l')$ $\in \mathcal{L}$ also in the set of labels. We note that for any parent label l that this label set has hard “leaf” parent-child constraints (explained later) although this assumption may be relaxed at the cost of more complexity. Such a label hierarchy forms a multiply-rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d , i.e. $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , i.e. $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation.

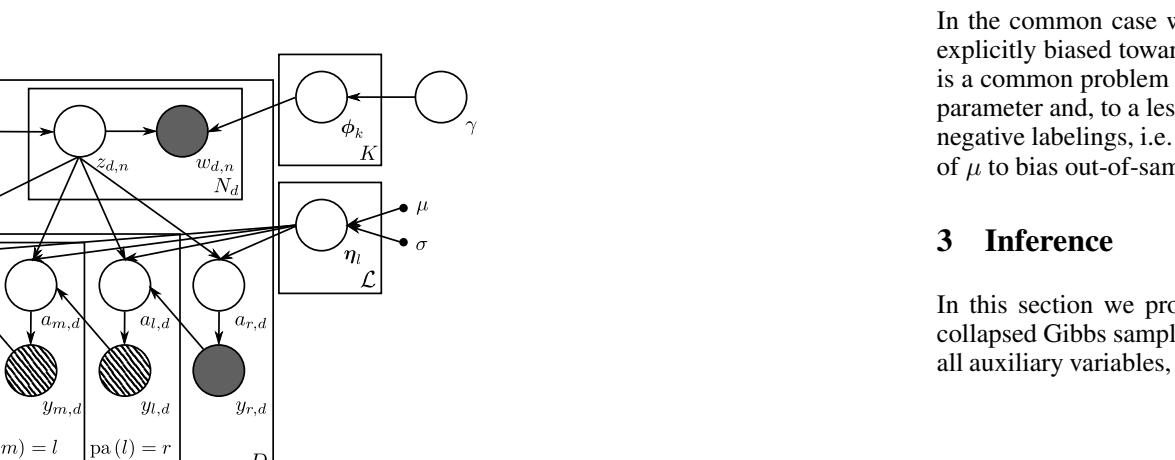


Figure 1: HSLDA graphical model

as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In the following, we provide the conditional distributions required to Gibbs sample the HSLDA posterior distribution. Note that, like in collapsed Gibbs samplers for LDA [12], we have analytically marginalized out the parameters $\phi_{1:D}$ and $\theta_{1:D}$. In the following \mathbf{a} is the set of all auxiliary variables, \mathbf{w} is the set of all words, $\boldsymbol{\eta}$ is the set of all regression coefficients, and \mathbf{z}_d is the set of $\mathbf{z}_{n,d}$ with element $z_{n,d}$ removed.

$$p(z_{n,d} = k | \mathbf{z}_d, \mathbf{z}_{n,d}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta) \propto \left(\frac{c_{k,n-d}^{a_{n,d}-a_{k,d}} + \alpha_k \beta_k}{c_{\cdot,n-d}^{a_{n,d}-a_{\cdot,d}} + \alpha \beta} \right) \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(z_{n,d}^T \boldsymbol{\eta} - a_{l,d})^2}{2} \right\} \quad (1)$$

Here, $c_{k,n-d}^{a_{n,d}-a_{k,d}}$ is the number of words of type k in document d assigned to topic k omitting the n th word of document d . The (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Also \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\boldsymbol{\eta}_l | \mathbf{z}, \mathbf{a}, \sigma) = \mathcal{N}(\boldsymbol{\mu}_l, \boldsymbol{\Sigma}) \quad (2)$$

where

$$\boldsymbol{\mu}_l = \Sigma \left(\frac{1}{\sigma^2} + \mathbf{Z}^T \mathbf{A}_l \right)^{-1} \Sigma^{-1} = \mathbf{I} \sigma^{-1} - \mathbf{Z}^T \mathbf{Z}.$$

Here \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{A}_l = (a_{l,1}, a_{l,2}, \dots, a_{l,D})^T$. The simplicity of this conditional distribution follows from the choice of probit regression [2]. The standard form of the update is a standard result from Bayesian normal linear regression [2]. That the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution

$$p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ - \frac{1}{2} (a_{l,d} - \eta_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (3)$$

is also a standard probit regression result [2].

HSLDA departs from stock LDA in that we estimate a hierarchical Dirichlet prior over topic assignments (i.e. β is a parameter in our model). This has been shown to improve the quality and stability of inferred topics [21]. Sampling β is done using the “direct assignment” method of Teh et al. [20]

$$\beta | \mathbf{z}, \alpha' \sim \text{Dir}(m_{(1),1} + \alpha', m_{(1),2} + \alpha', \dots, m_{(1),K} + \alpha') \quad (4)$$

where $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β and are governed by the following function.

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + c_{k,d}^m)} s(c_{k,d}^m, m) (\alpha \beta_k)^m \quad (5)$$

$s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

4 Related Work

Latent dirichlet allocation (LDA) is a generative probabilistic model of corpora that represents documents as a mixed membership bag-of-words. Also known as topic models, these models infer the latent structure, or topics, of documents in a corpus. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [5].

Supervised latent Dirichlet allocation (sLDA) builds on LDA by incorporating supervision in the form of an observed exponential family response variable per document. As a result, sLDA infers topics such that the model predicts the response variable while improving likelihood. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [4].

Other models that incorporate both and some form of supervision include LabeledLDA[17], DiscLDA[13], models applied to computer vision and document networks[22, 6].

The two measures for predictive performance used here include the true positive rate and the false positive rate evaluated based on $p(y_{l,d} | w_{1:N_d})$ for each label in each model.

[1] http://www.cdc.gov/nchs/icd/icd9cm.htm

[2] http://www.nitk.org

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , i.e. $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In this section we provide the conditional distributions required to Gibbs sample the HSLDA posterior distribution. Note that, like in collapsed Gibbs samplers for LDA [12], we have analytically marginalized out the parameters $\phi_{1:D}$ and $\theta_{1:D}$. In the following \mathbf{a} is the set of all auxiliary variables, \mathbf{w} is the set of all words, $\boldsymbol{\eta}$ is the set of all regression coefficients, and \mathbf{z}_d is the set of $\mathbf{z}_{n,d}$ with element $z_{n,d}$ removed.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be biased towards generating data that has negative labels in order to keep it from learning to assign all labels to the document. This is a common problem in modeling unbalanced data. To see this how this can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5.1 Data and Pre-Processing

5.1.1 Diagnosis Prediction

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient's complaints, findings and diagnosis, along with treatment and hospital course. For each admission trained medical coders review the information in the discharge summary and assign a series of diagnoses codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a labeling of a patient's diagnoses based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia”, where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [8], and sometimes make mistakes [10].

HSLDA	sLDA with Independent Regressors
children, animation, cartoon, kids, song, story, family, new, film, video, friends, fun, animals, family, adventure, minute, world, magic, animals, song, musical, action, voice, feature, along, baby, ages, magical, action, help, voice, toys, original, learn, characters, computer, sing, like, adults, parents, version, old, first, including, young, comedy, funny, comic, hilarious, series, gags, slapstick, satire, laugh, stars, films, comedians, first, median, funniest, show, involving, characters, cast, perfect, and, performance, play, write, together, scary, hilarious, ensemble, character, star, writer, hit, host, classic, romantic, mutants, boss, new, played, together, choice, like, karaoke, daughter, vehicle, news, splitting, actress, get, story, hilarious, wacky, featuring, could, side, comedies, comedian, horror, film, vampire, blood, one, dead, supernatural, killer, evil, thriller, gore, monster, young, night, terror, mysterious, director, movie, cult, creepiness, terrifying, original, dark, mysterious, death, revenge, violence, town, gang, turn, prison, noir, violent, deadly, scary, young, blood, body, vampire, night, murderer, kill, criminal, action, terrorism, crazy, gore, sex, victims, body, ghost, talk, death, predator, house, nightmare, budget, eerie, death, doctor, victim, later, effects	animation, songs, video, friends, fun, animals, family, adventure, little, music, like, along, baby, ages, magical, action, help, classic, world, magic, adventures, home, show, gang, animal, adults, toys, many, action, comedy, funny, comic, hilarious, series, gags, slapstick, satire, laugh, stars, films, comedians, first, median, funniest, show, involving, characters, cast, perfect, and, performance, play, write, together, scary, hilarious, ensemble, character, star, writer, hit, host, classic, romantic, mutants, boss, new, played, together, choice, like, karaoke, daughter, vehicle, news, splitting, actress, get, story, hilarious, wacky, featuring, could, side, comedies, comedian, horror, film, vampire, blood, one, dead, supernatural, killer, evil, thriller, gore, monster, young, night, terror, mysterious, director, movie, cult, creepiness, terrifying, original, dark, mysterious, death, revenge, violence, town, gang, turn, prison, noir, violent, deadly, scary, young, blood, body, vampire, night, murderer, kill, criminal, action, terrorism, crazy, gore, sex, victims, body, ghost, talk, death, doctor, victim, later, effects

The text of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP community challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a labeling of a patient's diagnoses based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia”, where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [8], and sometimes make mistakes [10].

HSLDA	sLDA with Independent Regressors
people, screwball, cast, best, including, sex, sketch, girlfriend, hilariously, charming, ensemble, comic, character, star, writer, hit, host, classic, romantic, mutants, boss, new, played, together, choice, like, karaoke, daughter, vehicle, news, splitting, actress, get, story, hilarious, wacky, featuring, could, side, comedies, comedian, horror, film, vampire, blood, one, dead, supernatural, killer, evil, thriller, gore, monster, young, night, terror, mysterious, director, movie, cult, creepiness, terrifying, original, dark, mysterious, death, revenge, violence, town, gang, turn, prison, noir, violent, deadly, scary, young, blood, body, vampire, night, murderer, kill, criminal, action, terrorism, crazy, gore, sex, victims, body, ghost, talk, death, doctor, victim, later, effects	workout, body, yoga, video, moves, dance, minute, fitness, easy, poses, muscles, routine, exercises, learn, strength, movements, practice, help, step, workouts, minutes, techniques, flexibility, get, exercise, one, instructor, first, series, time, instruction, two, fun, use, first, back, follow, steps, breathing, muscle, basic, using, II, beginners, three, warm, balance, technique, designed

workout, body, yoga, video, moves, dance, minute, fitness, easy, poses, muscles, routine, exercises, learn, strength, movements, practice, help, step, workouts, minutes, techniques, flexibility, get, exercise, one, instructor, first, series, time, instruction, two, fun, use, first, back, follow, steps, breathing, muscle, basic, using, II, beginners, three, warm, balance, technique, designed

program, fitness, easy, poses, muscles, exercises, learn, strength, movements, practice, help, step, workouts, minutes, techniques, flexibility, get, exercise, one, instructor, first, series, time, instruction, two, fun, use, follow, training, breathing, basic, back, II, muscle, fun, time, flexibility, exercise, instructor, instruction, one, two, fun, training, use, first, back, follow, steps, breathing, muscle, basic, using, II, beginners, three, warm, balance, technique, designed

workout, body, yoga, video, moves, dance, minute, fitness, easy, poses, muscles, routine, exercises, learn, strength, movements, practice, help, step, workouts, minutes, techniques, flexibility, get, exercise, one, instructor, first, series, time, instruction, two, fun, use, first, back, follow, steps, breathing, muscle, basic, using, II, beginners, three, warm, balance, technique, designed

program, fitness, easy, poses, muscles, exercises, learn, strength, movements, practice, help, step, workouts, minutes, techniques, flexibility, get, exercise, one, instructor, first, series, time, instruction, two, fun, use, first, back, follow, steps, breathing, muscle, basic, using, II, beginners, three, warm, balance, technique, designed

workout, body, yoga,

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include pages and their placement in Web directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show how leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to structure-agnostic models. Furthermore, we show evidence that the resulting HSLDA topics are more descriptive of the underlying data than sLDA topics, which ignore the label hierarchy.

1 Introduction

The task of multi-label classification, selecting the k-best labels for a given instance, has been a topic of research for several years. One simplistic way to carry out the classification is through a series of independent binary classifiers, but this ignores the many inherent dependencies among the labels. Thus, much work has been devoted on incorporating the co-occurrence patterns of the labels into the classification task. In this paper, we focus on multi-label classification, where the labels are organized in a hierarchical structure. Scenarios of use include, but are not limited to, placing webpages into manually curated Internet directories [2], categorizing images according to a taxonomy, tagging product descriptions with catalogue information [3], and assigning diagnosis codes to clinical records [1].

There are several challenges entailed in incorporating the hierarchical nature of labels into the classification task. One pertains to the labeling itself: in the datasets (especially real-world, noisy ones), for a given label, instances labeled with it contribute positive instance, but it is unclear how to determine the negative instances. In particular, how to treat the parent labels of the selected ones?

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. In particular, we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. We hypothesize that hierarchical label information provides more information about labeling than considering labels as a flat list.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Following the trend observed in supervised topic modeling, we note that the learned topic models are more representative of the underlying data in both of our datasets [5].

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d}$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. The set of labels is $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$ has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a response $y_{l,d} \in \{-1, 1\}$ to every label which indicates whether the label applies to document d or not.

The is-a hierarchical constraint is a hard constraint that document d has a positive response to label l_2 then it will also have a positive response to label l_1 . Conversely, if document d has a negative response to label l_1 then it will also have a negative response to l_2 . To capture this hierarchical structure we model the labeling of documents using a generative cascade of conditional probit regression models.

Each document is assigned a response of either -1 or 1 for at least one, but potentially many labels in \mathcal{L} . The label, l , for a document, d , will be used interchangeably to refer to the observed response of document d to label l .

The fixed parameters of the model are the number of topics K , the number of unique words in the vocabulary V , the number of documents D , as well as the mean, μ , and the standard deviation, σ , used in a normal prior distribution. The hyper-parameters α' , α , and γ are weight

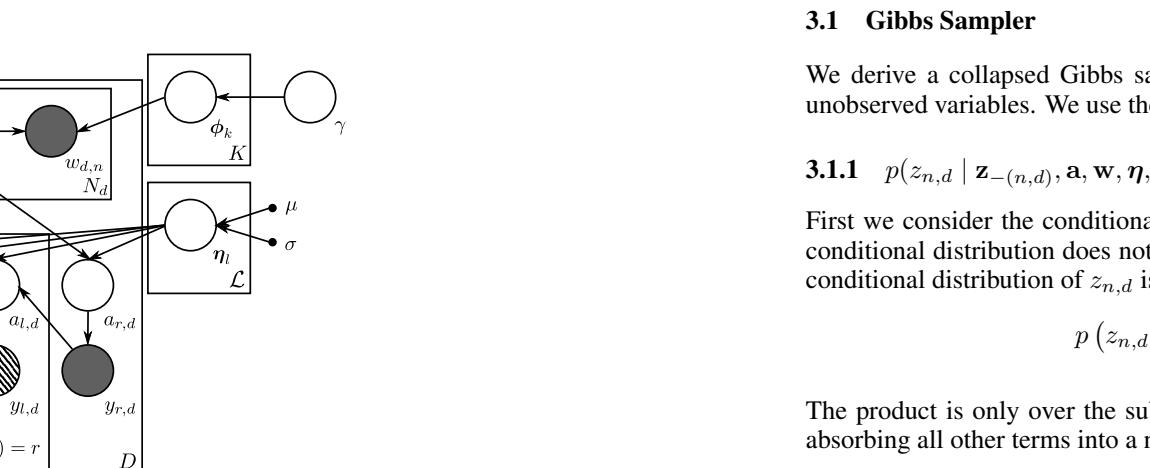


Figure 1: HSLDA graphical model

parameters for Dirichlet prior distributions. We will denote the K -dimensional Dirichlet distribution as $\text{Dir}_K(\cdot)$ and the K dimensional normal distribution as $\mathcal{N}_K(\cdot)$.

We will now describe the stochastic generative process which defines our model. The graphical model is show in Figure ??.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_{k,l} \sim \text{Dir}_{V_k}(\gamma \mathbf{I}_{V_k})$
2. For each label $l \in \mathcal{L}$
 - Draw a regression coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$, where \mathbf{I}_K is the K dimensional identity matrix
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $v = 1, \dots, V$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{1:D} \sim \text{Multinomial}(\phi_{1:D})$
 - For each label $l \in \mathcal{L}$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw $a_{l,d} \mid z_{n,d}, \eta_l \sim \begin{cases} \mathcal{N}(\bar{z}_{n,d} \eta_l, 1), & y_{p(l)} = 1 \\ \mathcal{N}(\bar{z}_{n,d} \eta_l, 1/(a_{l,d} < 0)), & y_{p(l)} = -1 \end{cases}$
where $\bar{z}_{n,d} = \frac{1}{N_d} \sum_{n=1}^{N_d} z_{n,d}$
 - Set the response variable
 $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \text{ and } y_{p(l),d} = 1 \\ -1 & \text{otherwise} \end{cases}$

This type of generative model is known as a probit regression model. Probit regression models are a type of discriminative probabilistic model similar to logistic regression. However, instead of using the logistic sigmoid as the link function, the probit regression model uses the CDF for a standard normal distribution – the inverse of which is known as the probit. In this case, the regression is conditional on the parents according to the constraints of the labeling hierarchy. The latent variables $a_{l,d}$ utilized here are also known as an auxiliary variables because they are introduced to make exact Gibbs sampling possible and are not of primary interest.

Given that negative labels are uncommon and that the absence of a label is not equivalent to a negative label, we apply an informative prior to the regression parameters, $\beta_{\mathcal{L}}$, in the form of a negative prior that encodes a bias towards being truly negative in the absence of a label.

In the Bayesian approach to statistical modeling, the primary task of inference is to find the posterior distribution over the unobserved parameters of the model. However, it is often possible and desirable to integrate over certain variables in the model, also known as collapsing. In this case, we collapse over the labels \mathcal{L} not fully observed for document d . We will define \mathcal{L}_d to be the subset of labels which have been observed for document d . It is straightforward to integrate out the variables $a_{l,d}$ for $l \in \mathcal{L} \setminus \mathcal{L}_d$ in the full generative model. We can also integrate out the parameters θ_d and $\theta_{l,D}$ in Griffiths and Steyvers [13]. Therefore, in our model the latent variables are $\mathbf{z} = \{z_{1:N_d,d}\}_{d=1:D}$, $\boldsymbol{\eta} = \{\eta_l\}_{l \in \mathcal{L}}$, $\theta = \{\theta_d\}_{d \in \mathcal{L}_d}$, β , α' , and γ .

The posterior distribution we seek cannot be solved in closed form. This is often the case in calculating posterior distributions of non-trivial probabilistic models. We will appeal to one of the common methods for approximating posterior distributions in the face of intractable normalization factors: Markov chain Monte Carlo (MCMC) sampling. Since in this model it is possible to sample from the conditional distributions for all variables we will use the Gibbs sampling algorithm to obtain an approximation to this posterior:

7 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...

- what about the nonparametric version of this?
- discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

References

- [1] The computational medicine center’s 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [4] James H Albert and Siddhartha Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669, 1993.
- [5] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet Allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [7] Jonathan Chang and David M. Blei. Hierarchical topic models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [8] K. Crammer, M. Dredze, K. Ganchev, PP. Talukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [9] Yan Yan David S. Silanes Martha J. Radford Brian F. Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [10] R. Farkas and G. Szarvas. Automatic construction of rule-based icd-9-cm coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [11] Mehrdad Farzandpour, Abbas Sheikhtabar, and F. Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:78–84, 2010.
- [12] I.Goldstein, A.Arztumyan, and O.Uzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [13] T.L. Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5233, 2004.
- [14] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [15] Leah Larkey and Bruce Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [16] LV. Lita, S. Yu, S. Niculescu, and J. Bi. Large scale diagnostic code classification for medical patient records, 2008.
- [17] Serguei Pakhomov, James Buntrock, and Christopher Chute. Automating the assignment of diagnosis codes to patient encounters using example-based and machine learning techniques. *Journal of the American Medical Informatics Association (JAMIA)*, 13(5):516–525, 2006.
- [18] Daniel Ramage, David Hall, Ramesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 – Volume 1*, pages 248–252. USA: Association for Computational Linguistics, ISBN 978-1-932432-59-6. URL <https://aclweb.org/anthology/E09-1030.pdf>; <https://aclweb.org/citation.cfm?id=1693943>.
- [19] B. RibeiroNeto, AHF Laender, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and Technology*, 52(5):391–401, 2001.
- [20] P. Ruch, J. Gobell, I. Thushri, and A. Grabischler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [21] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [22] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.
- [23] Chong Wang, David Blei, and Li Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

3.1 Gibbs Sampler

We derive a collapsed Gibbs sampler for this model by considering the individual conditional probability distributions for each of the unobserved variables. We use the notation $\mathbf{z}_{-(n,d)}$ to denote $\mathbf{z}_d \setminus z_{n,d}$.

3.1.1 $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$

First we consider the conditional distribution of the assignment variable for each word $n = 1, \dots, N_d$ in documents $d = 1, \dots, D$. The conditional distribution does not include $\theta_{l,D}$ and $\phi_{1:D}$ because they have been integrated out as in the collapsed Gibbs sampler [13]. The conditional distribution of $z_{n,d}$ is proportional to the joint distribution of its markov blanket.

$$p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \prod_{l \in \mathcal{L}_d} p(a_{l,d} \mid \mathbf{z}, \boldsymbol{\eta}_l) p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \alpha, \beta, \gamma). \quad (1)$$

The product is only over the subset of labels \mathcal{L}_d which have been observed for document d . By isolating terms that depend on $z_{n,d}$ and absorbing all other terms into a normalizing constant as in [13] we find

$$p(z_{n,d} = k \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \left(c_{(k-1)(n,d)}^{k-1} + \alpha_k \right) \frac{\left(\frac{a_{l,d}}{\mu_{l,d}} + \beta_k \right)^{-\frac{1}{2}}}{\left(\frac{a_{l,d}}{\mu_{l,d}} + \beta_k + V_k \right)} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(z_{n,d} - \mu_{l,d})^2}{2} \right\}. \quad (2)$$

Here, $c_{(k-1)(n,d)}$ represents the number of words of type v in document d assigned to topic k omitting the n th word of document d . The superscript (\cdot) in the subscript means the count resulting from summing over the omitted subscript variable. Given Equation 1, $p(z_{n,d} \mid \mathbf{z}_{-(n,d)}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma)$ can be sampled through enumeration.

3.1.2 $p(\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \mathbf{w}, \sigma)$

We now consider the conditional distribution of the regression coefficients $\boldsymbol{\eta}_l$ for $l \in \mathcal{L}$. Given that $\boldsymbol{\eta}_l$ and $a_{l,d}$ are distributed normally, the posterior distribution of $\boldsymbol{\eta}_l$ is normally distributed with mean $\hat{\boldsymbol{\eta}}_l$ and covariance Σ such that

$$\Sigma^{-1} = \mathbf{I} \sigma^{-2} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

$$\hat{\boldsymbol{\eta}}_l = \Sigma^{-1} \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right). \quad (4)$$

This is a standard result from normal Bayesian linear regression [2]. Here, \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is $\mathbf{a}_{l,d}$, and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$.

3.1.3 $p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta})$

The auxiliary variables $a_{l,d}$ must be sampled for documents $d = 1, \dots, D$ and $l \in \mathcal{L}_d$. The conditional posterior distribution of $a_{l,d}$ is the truncated normal distribution

$$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \boldsymbol{\eta}_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (5)$$

This conditional distribution can be sampled using an inverse CDF method.

###

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
 {ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply-labeled bag-of-word data. Examples of such data include web pages and their placement in link directories, product descriptions and placement(s) in product hierarchies, and free-text clinical records and diagnosis codes assigned to them. Out-of-sample label prediction is the primary goal we want to achieve. We show how improved lower-dimensional representations of the bag-of-words data is also achieved than using a size 1 vector from hierarchical labels substantially improves our out-of-sample label prediction in comparison to other models that don't utilize the structure of the labels. We demonstrate HSLDA on large-scale data from medical document labeling and retail product categorization tasks. We show improved label prediction performance and evidence that the learned topics also improve.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on hierarchical and multi-labeled data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and entries in medical record systems and their associated billing codes. In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new text documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

In this work we extend supervised latent Dirichlet allocation (sLDA) [4] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [5] augmented with per-document "supervision"; often taking the form of a single numerical or categorical "label." More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as having been conditionally drawn from some distribution that depends on the specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [4].

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browseable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unstructured labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is structured as follows. Section 2 introduces hierarchically supervised LDA (HSLDA). Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and Web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathcal{W}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label l except the root $\emptyset \in \mathcal{L}$ has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We say a label l is a child of another label if it has had a "seen" parent-child constraint (explained later) although the assumption is not required at the cost of more complex code. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{i,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d , i.e. $y_{r,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{p(l),d} = 1, y_{p(p(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

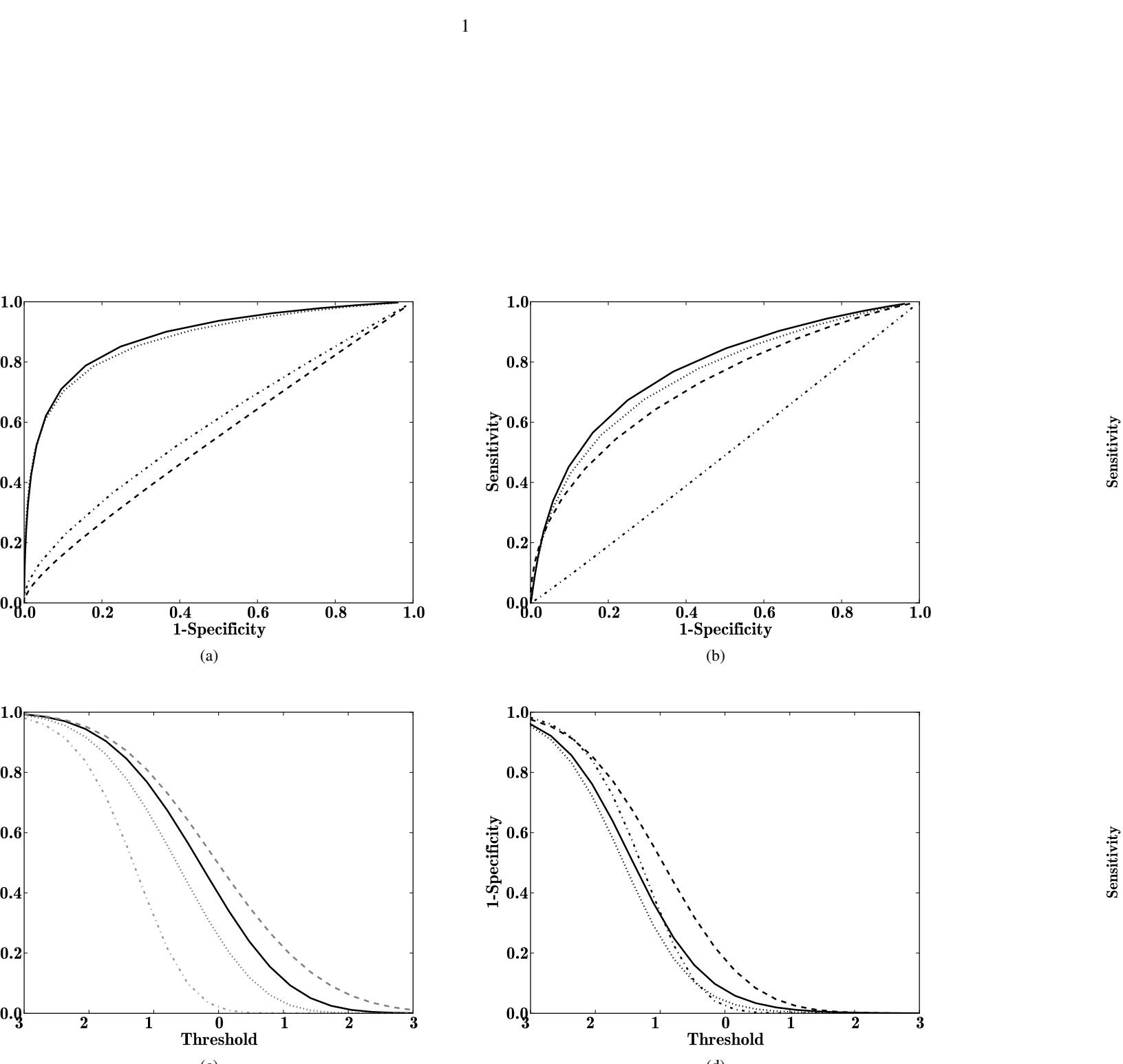


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

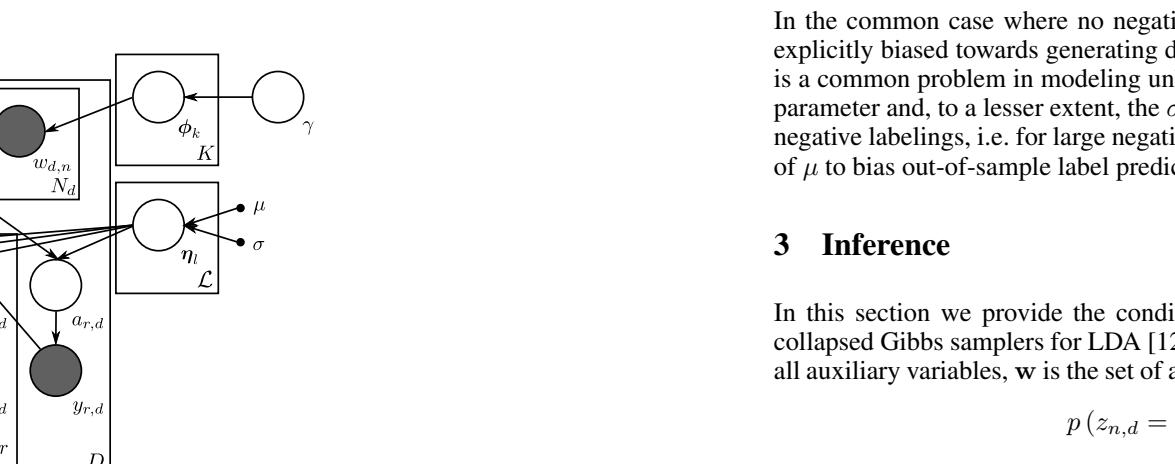


Figure 1: HSLDA graphical model

as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plotted variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA "topics" (distributions over the elements of Σ). θ_d is a distribution over "words," θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $\mathbf{I}_K(\cdot)$ is the K -dimensional identity matrix, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{I}_K, \sigma \mathbf{I})$
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{I})$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{l,z_{n,d}} \sim \text{Multinomial}(\phi_{l,z_{n,d}})$
 - Set $y_{r,d} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} \mid \mathbf{z}_d, \eta_l, y_{p(l),d} \sim \begin{cases} \mathcal{N}(\mathbf{z}_d^T \eta_l, 1), & y_{p(l),d} = 1 \\ \mathcal{N}(\mathbf{z}_d^T \eta_l, 1) \mathbb{I}(a_{l,d} < 0), & y_{p(l),d} = -1 \end{cases}$
 - Apply label l to document d according to $a_{l,d}$
 - $y_{l,d} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here $\mathbf{z}_d^T = [\bar{z}_1, \dots, \bar{z}_k, \dots, \bar{z}_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k , $\bar{z}_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generatively labeled using a hierarchy of conditionally dependent probit regressors [3?]. For every label $l \in \mathcal{L}$ both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $\mathbb{I}(y_{p(l),d} = 1)$) are used to decide whether or not label l is to be applied to document d as well. Note that $y_{p(l),d} = 1$ if its parent label $p(l)$ is applied (these expressions are specific to the constraints but can be modified to account for different constraints). The regression coefficients η_l are independent a priori; however, the hierarchical coupling in this model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

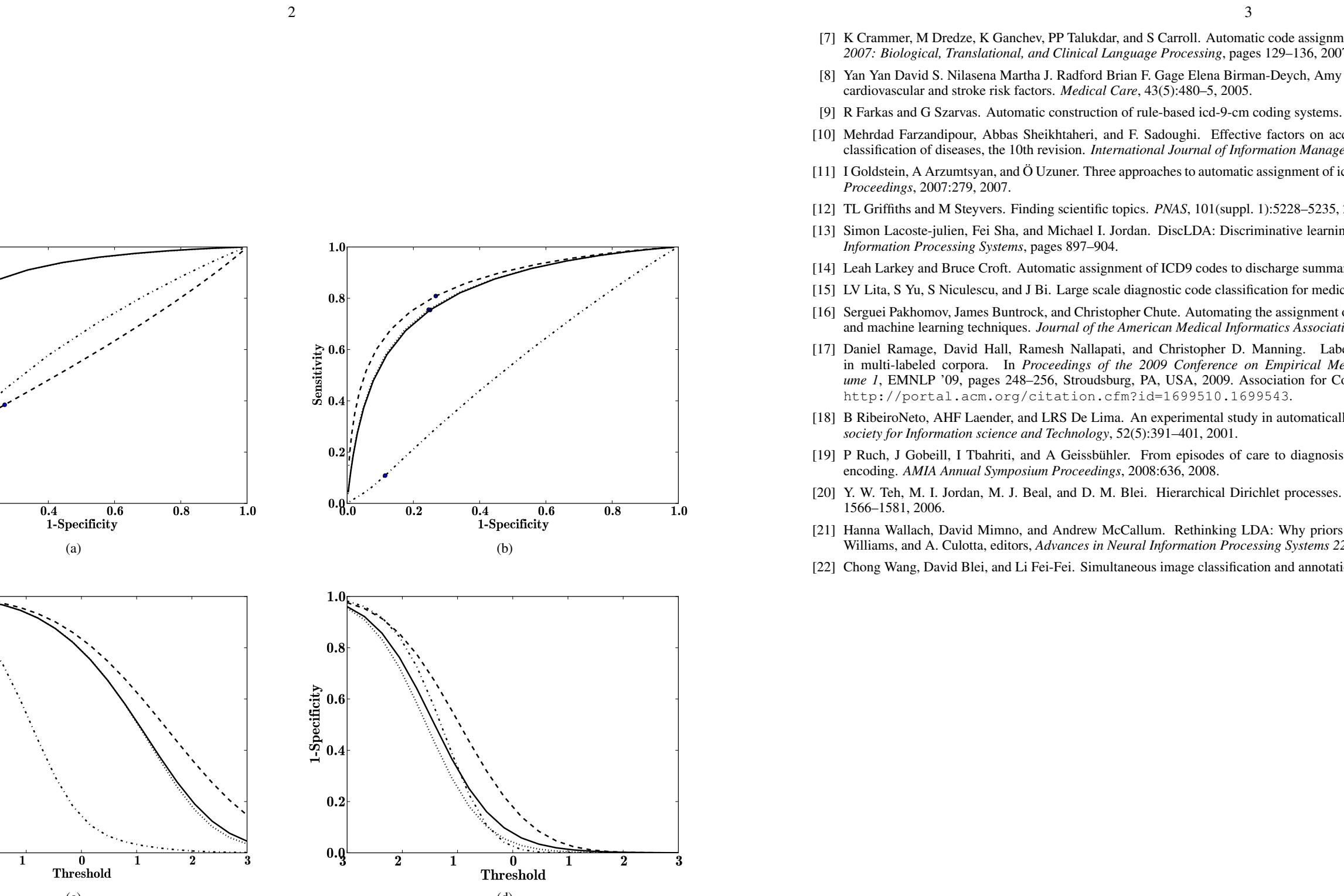


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be biased towards generating data that has negative labels in order to keep it from learning to assign all labels to the document. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because \mathbf{z}_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing

5.1.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient's complaints, findings and diagnoses, along with treatment and hospital course. For each admission trained medical coders review the information in the discharge summary and assign a series of diagnoses codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a labeling of a patient's diagnoses based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for "Pneumonia due to adenovirus" is a child of the code for "Viral pneumonia" where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [8], and sometimes make mistakes [10].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical community challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [7, 11, 19]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [14, 18, 16, 19, 15].

Our dataset was collected from clinical discharge summaries of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The text was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Description and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a hierarchy and provides product descriptions for most products in their catalog. Products can be discovered by users through query-based searching or through product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

We were able to deduce the structure of the hierarchy for the Amazon.com products because all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine."

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

7 Discussion

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...

• what about the nonparametric version of this?

• discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, is the first principled approach to doing so

References

[1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007/

[2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.

[3] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.

[4] James H. Albert and Sudhirtha Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669–693, 1993.

[4] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.

[5] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.

[6] Jonathan Chang and David M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.

5

4

3

2

1

0

5

4

3

2

1

0

5

4

3

2

1

0

5

4

3

2

1

0

5

4

3

2

1

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noemie@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include pages and their placement in Web directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to structure-agnostic models. Furthermore, we show evidence that the resulting HSLDA topics are more descriptive of the underlying data than sLDA topics, which ignore the label hierarchy.

1 Introduction

The task of multi-label classification, selecting the k -best labels for a given instance, has been a topic of research for several years. One simple way to carry out the classification is through a series of independent binary classifiers, but this ignores the many inherent dependencies among the labels. Thus, much work has been devoted to incorporating the co-occurrence patterns of the labels into the classification task. In this paper, we focus on multi-label classification, where the labels are organized in a hierarchical structure. Scenarios of use include, but are not limited to, placing webpages into manually curated Internet directories [2], categorizing images according to a taxonomy, tagging product descriptions with catalog information [3], and assigning diagnosis codes to clinical records [1].

There are several challenges entailed in incorporating the hierarchical nature of labels into the classification task. One pertains to the labeling itself: in the datasets (especially real-world, noisy ones), for a given label, instances labeled with it contribute positive instance, but it is unclear how to determine the negative instances. In particular, how to treat the parent labels of the selected ones?

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. In particular, we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. We hypothesize that hierarchical label information provides more information about labeling than considering labels as a flat list.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Following the trend observed in supervised topic modeling, we find that the learned topic models are more representative of the underlying data in both of our datasets [5].

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at a cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document d has a variable $y_{d,l} \in \{0, 1\}$ indicating whether it contains label l . Conversely, if a label l is marked as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

1

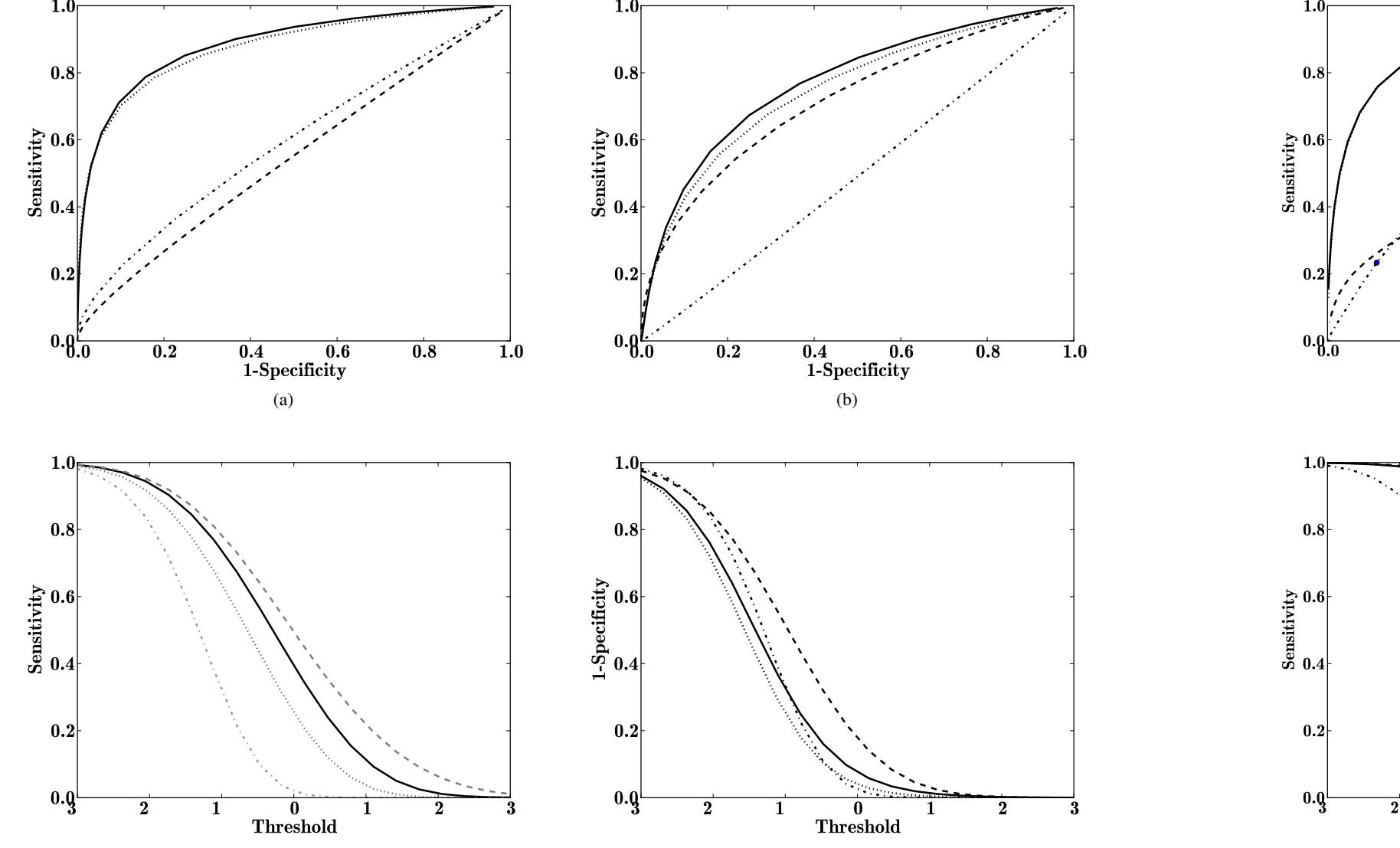


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

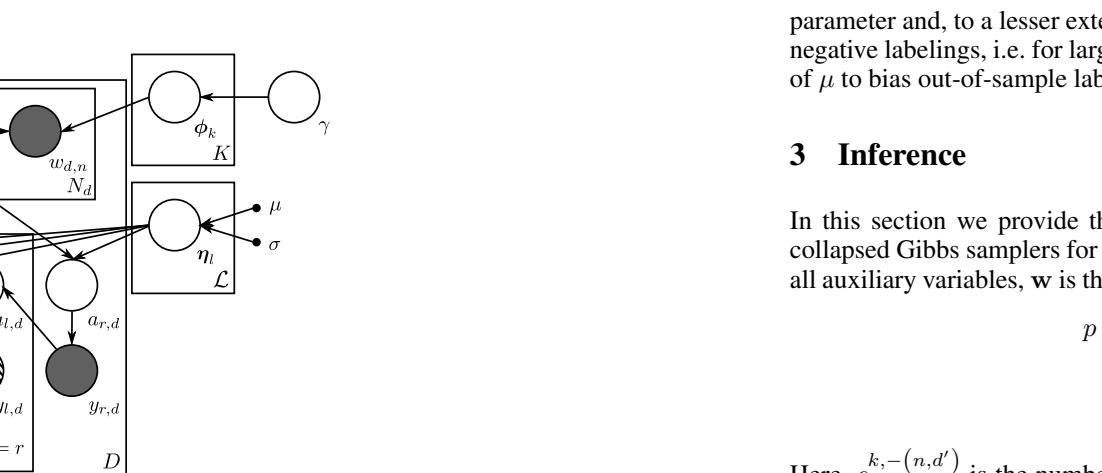


Figure 1: HSLDA graphical model

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA “topics” distributions over the elements of Σ , ϕ_k is a distribution over “words.” θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $\mathcal{N}_K(\cdot)$ is the K -dimensional Normal distribution, I_K is the K -dimensional identity matrix, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{I}_K, \sigma \mathbf{I}_K)$
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{I}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{I}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{1:K} \sim \text{Multinomial}(\phi_{z_{n,d}})$
 - Set $y_{d,l} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} \mid z_d, \eta_l, y_{p(l),d} \sim \begin{cases} \mathcal{N}(z_d^\top \eta_l, 1), & y_{p(l),d} = 1 \\ \mathcal{N}(z_d^\top \eta_l, 1)(1 - y_{p(l),d}), & y_{p(l),d} = 0 \end{cases}$
 - Apply label l to document d according to $a_{l,d}$

Here $z_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k ; $\tilde{z}_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generated and used as input to a tree-conditional regression problem [2, 7]. For every label $l \in \mathcal{L}$ we learn a conditional topic distribution for document d and whether or not it is applied (i.e. $\mathbb{I}(y_{d,l} = 1)$). If $pa(l) = r$ then $a_{l,d}$ is applied (these expressions are specific to is-a constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Other models that incorporate LDA and supervision include LabeledLDA[18], DiscLDA[14], and other models applied to computer vision and document networks[23, 7]. These models, however, do not implement constraints on the label space.

In other work, researchers have classified documents in a hierarchy with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces and has focused on single label classification without a model of documents such as LDA[2, 7, 22].

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from leaning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ

2

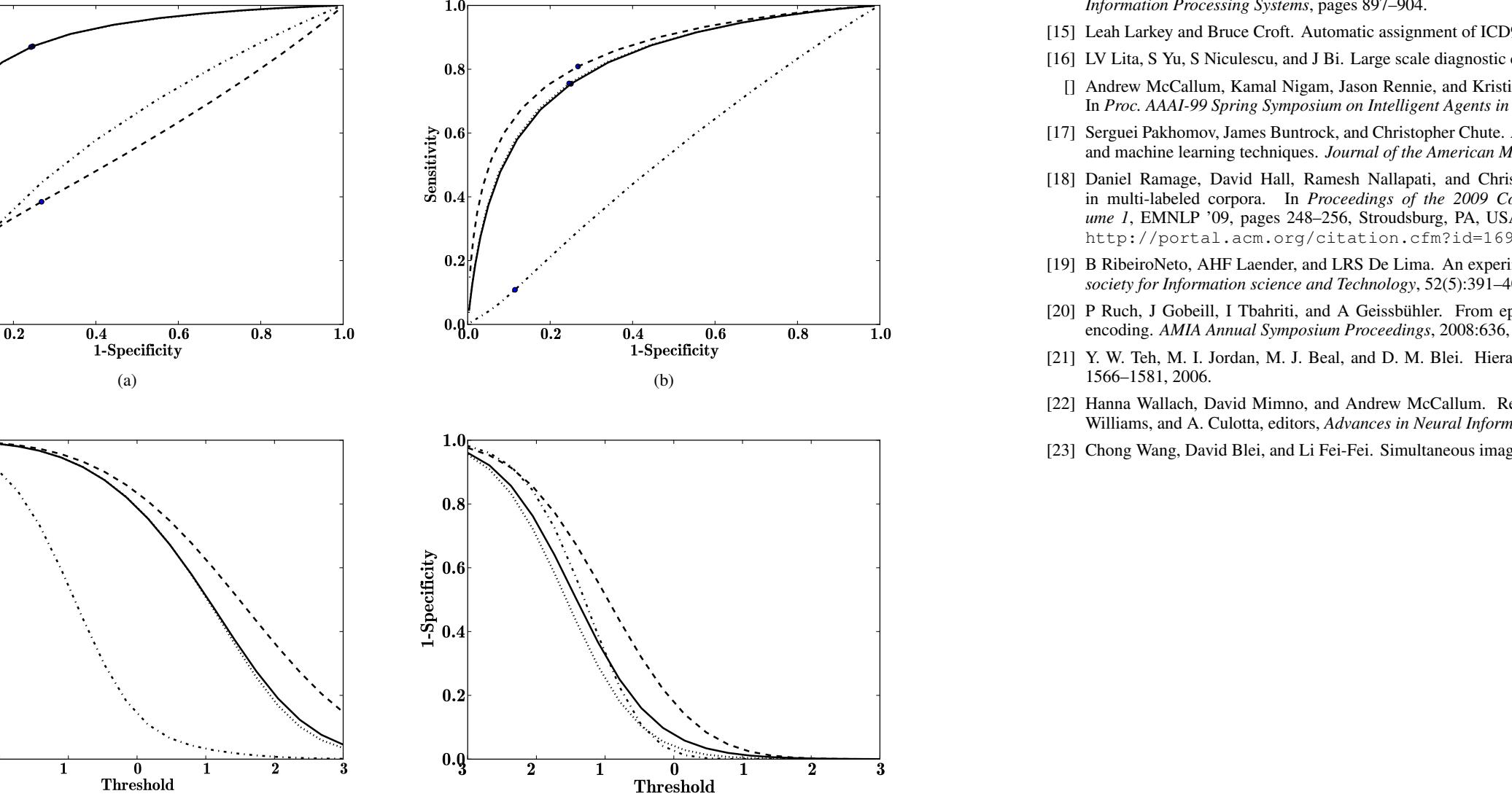


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

parameter and, to a lesser extent, the σ parameter above. Since z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

3 Data and Pre-Processing

3.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient complaints, findings and diagnoses, along with treatment and hospital course. For each admission trained medical coders review the information in the discharge summary and assign a series of diagnoses codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a labeling of a patient’s diagnoses based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia” where the former is the type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [9], and sometimes make mistakes [11].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP community challenge [1]. The data in the challenge, however, differs from ours in scope. The datasets were smaller (1,000 training and 1,000 testing) and focused on a restricted number of ICD-9 codes (45 of them, ranging from 3 to 7K in our dataset). Methods ranged from manual rules to online learning [8, 12, 10]. Other work has leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, as well as preexisting results [15, 19, 21, 16].

Our dataset was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall, representing all the discharges from the hospital in 2009. We have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of stop words and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The text was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

3.2 Product Descriptions and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a hierarchy and provides product descriptions for most products in their catalog. Product descriptions can be discovered by users through query-based searching or through product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy. This has been shown to improve the quality and stability of inferred topics [22]. Sampling β is done using the “direct assignment” method of Teh et al. [21]

$$\beta \mid \mathbf{z}, \alpha', \alpha \sim \text{Dir}(m_{(\cdot),1} + \alpha', m_{(\cdot),2} + \alpha', \dots, m_{(\cdot),K} + \alpha') \quad (4)$$

where $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β and are governed by the following function.

$$p(m_{d,k} \mid \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + c_{d,k})} s(c_{d,k}, m) (\alpha \beta_k)^m \quad (5)$$

$s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

4 Related Work

While there has been much work in multi-label classification of text and text modeling in general, we focus here on topic modeling approaches. Latent Dirichlet allocation (LDA) is a generative probabilistic model which represents documents as a mixed-membership bag of words. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include rating associated with an online review, grades for an essay, and number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

We evaluated model performance for all three models with a range of values for μ ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$), the mean prior parameter for regression performance.

Other models that incorporate LDA and supervision include LabeledLDA[18], DiscLDA[14], and other models applied to computer vision and document networks[23, 7]. These models, however, do not implement constraints on the label space.

In other work, researchers have classified documents in a hierarchy with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces and has focused on single label classification without a model of documents such as LDA[2, 7, 22].

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

3

- [9] Yan David S, Nilasena Martha J, Radford Brian F, Gage Elena Birman-Deych, Amy D. Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [10] R Farkas and G Szarvas. Automatic construction of rule-based icd-9 coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [11] Mehdi Farzandipour, Abbas Sheikhtaheri, and F. Soughaf. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:75–84, 2010.
- [12] I Goldstein, A Arzumanyan, and O Uzuner. Three approaches to automatic assignment of icd-9 codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279–207.
- [13] T.L Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl):5228–5235, 2004.
- [14] Daphne Koller and Mehran Sahami. Hierarchical classifying documents using very few words. *Technical Report 1997-75*, Stanford InfoLab, February 1997. Previous number is SIDL-WP-1997-009.
- [15] Simon Lacoste-Julien, Fei Sha, and Michael I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [16] Leah Larkey and Bruce Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [17] LV Lita, S Yu, S Niculescu, and JBL. Large scale diagnostic code classification for medical records. 2008.
- [18] Andrew McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Building domain-specific search engines with machine learning techniques. In *Proc. AAAI-99 Spring Symposium on Intelligent Agents in Cyberspace*, 1999.
- [19] Daniel Ramage, David Hall, Ramesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP ’09, pages 248–256, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL <http://portal.acm.org/citation.cfm?id=1699510.1699543>.
- [20] P. Ribeiro, S. Ribeiro, J. Tshirhi, and A. Geissbuhler. Prior episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [21] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [22] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.
- [23

Hierarchically Supervised Latent Dirichlet Allocation

Adler Perotte Nicholas Bartlett Noémie Elhadad Frank Wood
Columbia University, New York, NY 10027, USA
{ajp9009@dbmi, bartlett@stat, noe@dbmi, fwood@stat}.columbia.edu

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include pages and their placement in Web directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to structure-agnostic models. Furthermore, we show evidence that the resulting HSLDA topics are more descriptive of the underlying data than sLDA topics, which ignore the label hierarchy.

1 Introduction

The task of multi-label classification, selecting the k -best labels for a given instance, has been a topic of research for several years. One simple way to carry out the classification is through a series of independent binary classifiers, but this ignores the many inherent dependencies among the labels. Thus, much work has been devoted on incorporating the co-occurrence patterns of the labels into the classification task. In this paper, we focus on multi-label classification, where the labels are organized in a hierarchical structure. Scenarios of use include, but are not limited to, placing webpages into manually curated Internet directories [2], categorizing images according to a taxonomy, tagging product descriptions with category information [3], and assigning diagnosis codes to clinical records [1].

There are several challenges entailed in incorporating the hierarchical nature of labels into the classification task. One pertains to the labeling itself: in the datasets (especially real-world, noisy ones), for a given label, instances labeled with it contribute positive instance, but it is unclear how to determine the negative instances. In particular, how to treat the parent labels of the selected ones? It is also unclear how to determine the negative instances.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. In particular, we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. We hypothesize that hierarchical label information provides more information about labeling than considering labels as a flat list.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Following the trend observed in supervised topic modeling, we note that the learned topic models are more representative of the underlying data in both of our datasets [5].

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at a cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document d has a value $y_{d,l} \in \{0, 1\}$ indicating whether or not it contains label l . In most cases $y_{d,l}$ will be 0, but in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor cases it will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , i.e. $y_{d,l} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{d,pa(l)} = 1, y_{d,pa(pa(l))} = 1, \dots, y_{d,r} = 1$. Conversely, if a label l' is marked as not applying to a document d then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

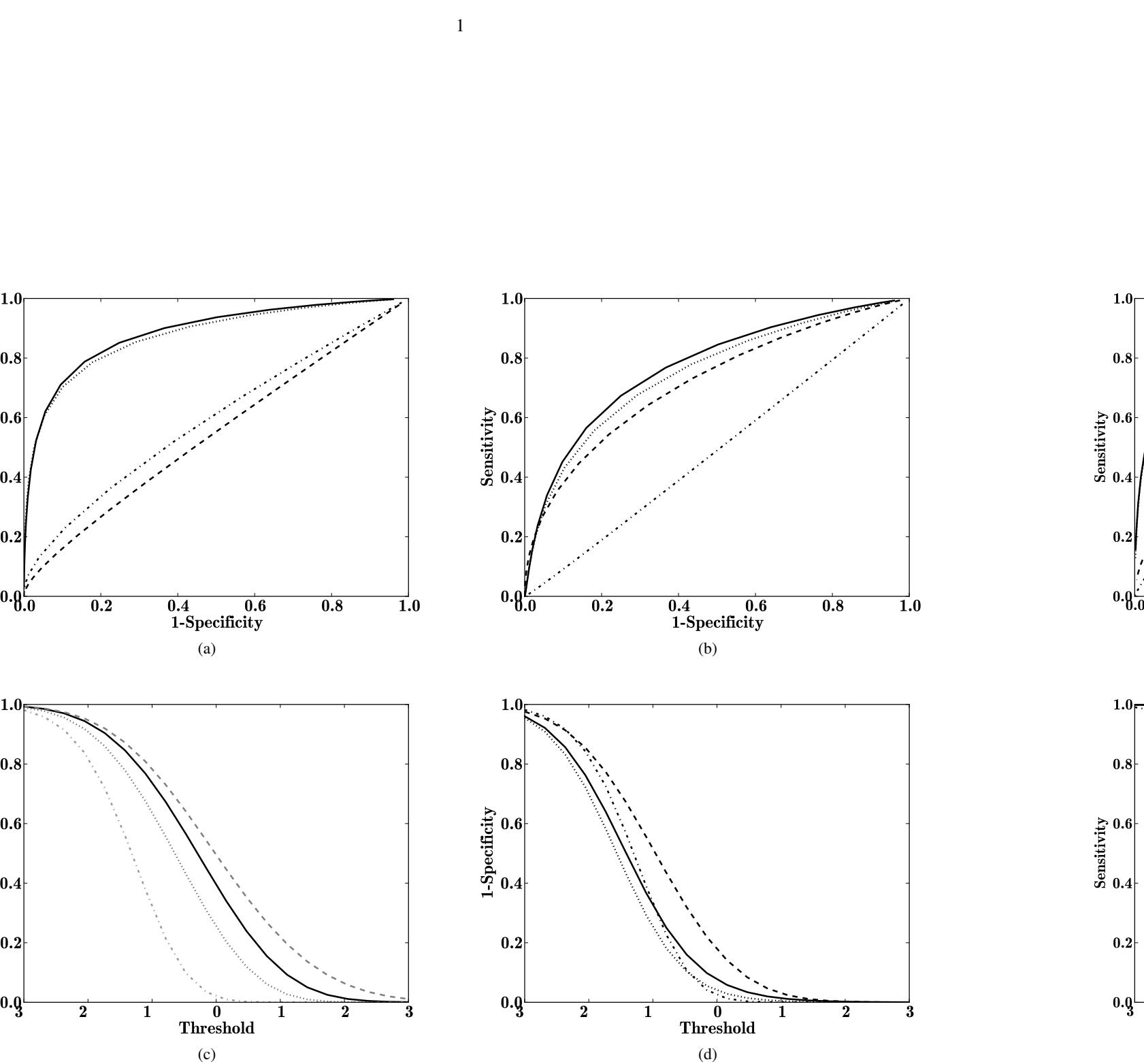


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

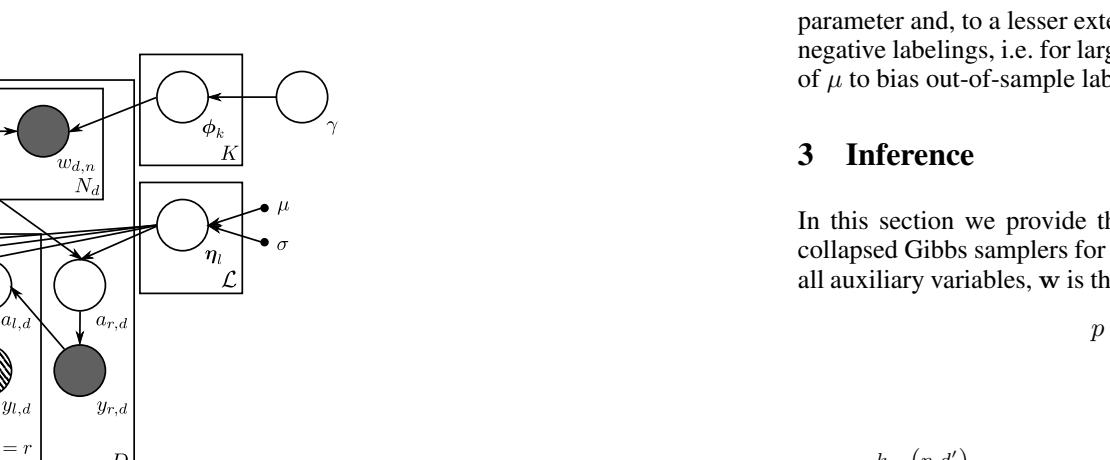


Figure 1: HSLDA graphical model

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA “topics” distributions over the elements of Σ , ϕ_k is a distribution over “words.” θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $I_d(\cdot)$ is the K -dimensional Normal distribution, $\text{Dir}_d(\cdot)$ is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_k \sim \text{Multinomial}(\phi_{z_{n,d}})$
 - Set $y_{d,l} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_d, \eta_l, y_{d,pa(l)} \sim \frac{\mathcal{N}(z_d^T \eta_l, 1)}{\mathcal{N}(z_d^T \eta_l, 1)(\mathbb{I}(a_{l,d} < 0))}, y_{d,pa(l),d} = 1$
 - Apply label l to document d according to $a_{l,d}$

$$y_{d,l} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$$

Here $z_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k ; $\tilde{z}_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is given an independent set of constraints of causally dependent probabilities [2, 7]. For every label l in the document, an explicit topic distribution for document d and whether or not it is positive is applied (i.e. $\mathbb{I}(y_{d,l}) = 1$) to decide whether or not label l is to be applied to document d as well. Note that label $y_{d,l}$ can only be applied to document d if its parent label $y_{d,pa(l)}$ is also applied (these expressions are specific to is-a constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from leaning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ

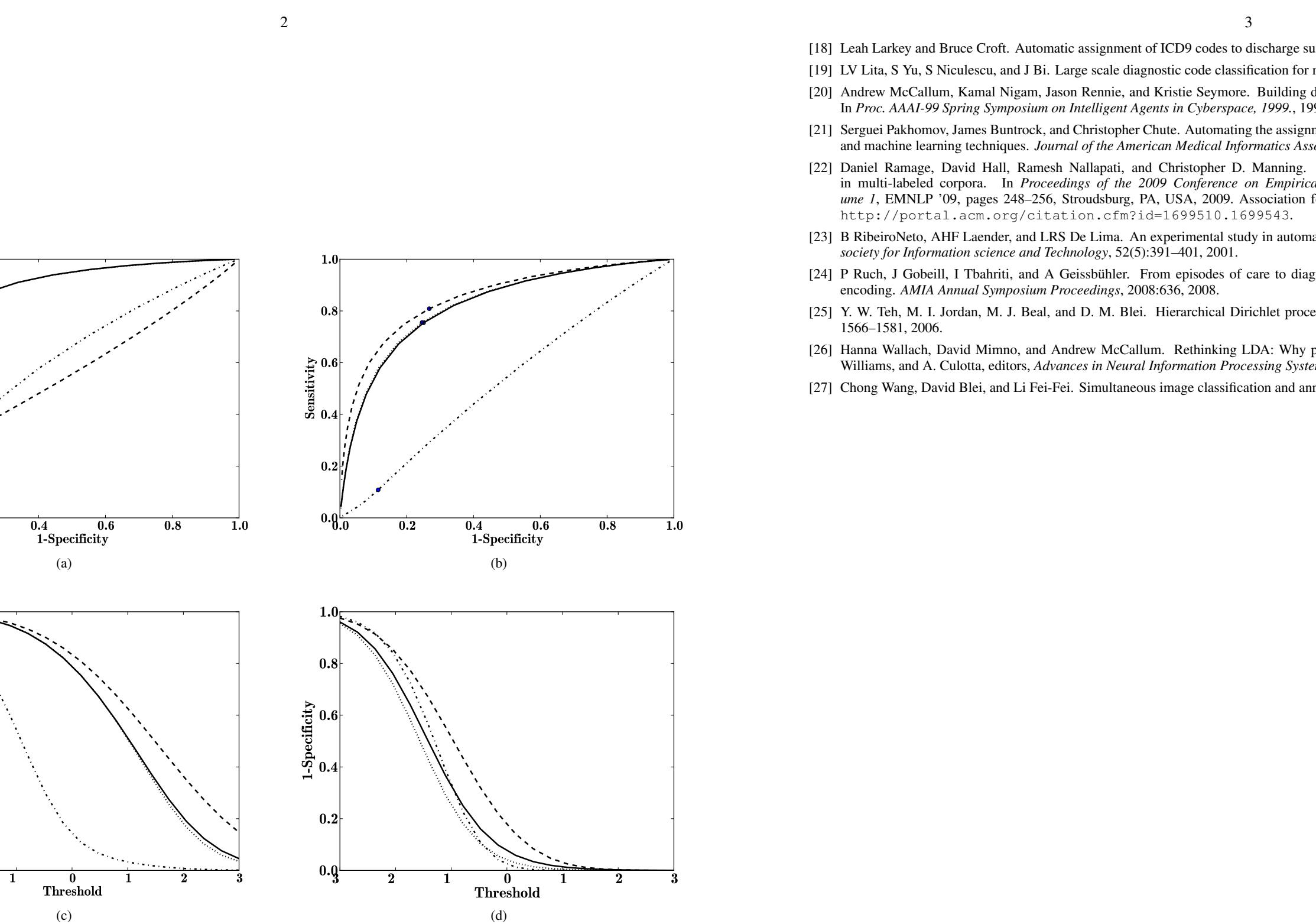


Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

parameter, and to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{d,l} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

3 Inference

In this section we provide the conditional distributions required to Gibbs sample the HSLDA posterior distribution. Note that, like in the discharge summaries and ICD-9 codes, the labels are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [11], and sometimes make mistakes [13].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge competition in the challenge, however different from ours in its scope. The datasets were very similar to ours (training and 1000 testing documents) and focused on radiology with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [9, 14, 12]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [18, 23, 21, 24, 19].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8,39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 53.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK 2. Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

4 Related Work

Amazon.com, an online retail store, organizes its catalog of products in a hierarchy and provides product descriptions for most products in their catalog. Products can be discovered by users through query-based searching or through product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

We were able to deduce the structure of the hierarchy for the Amazon.com products because all ancestors in the hierarchy were included with each category label. For example, “DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine.”

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

We evaluated HSLDA along with three other closely related models against these two datasets. The comparison models included sLDA with independent labels (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. These models were chosen to highlight several aspects of the model including performance in the absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition constraint imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance. Other models that incorporate LDA and supervision include LabeledLDA[22], DiscLDA[17], and other models applied to computer vision and document networks[27, 8]. These models, however, do not implement constraints on the label space. Other models that incorporate LDA and supervision include LabeledLDA[22], DiscLDA[17], and other models applied to computer vision and document networks[27, 8]. These models, however, do not implement constraints on the label space. In other work, researchers have classified documents in a hierarchy with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces and has focused on single label classification without a model of documents such as LDA[20, 10, 16, 7]. There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference procedures do not allow the response to reflect the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA[5]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nitk.org>

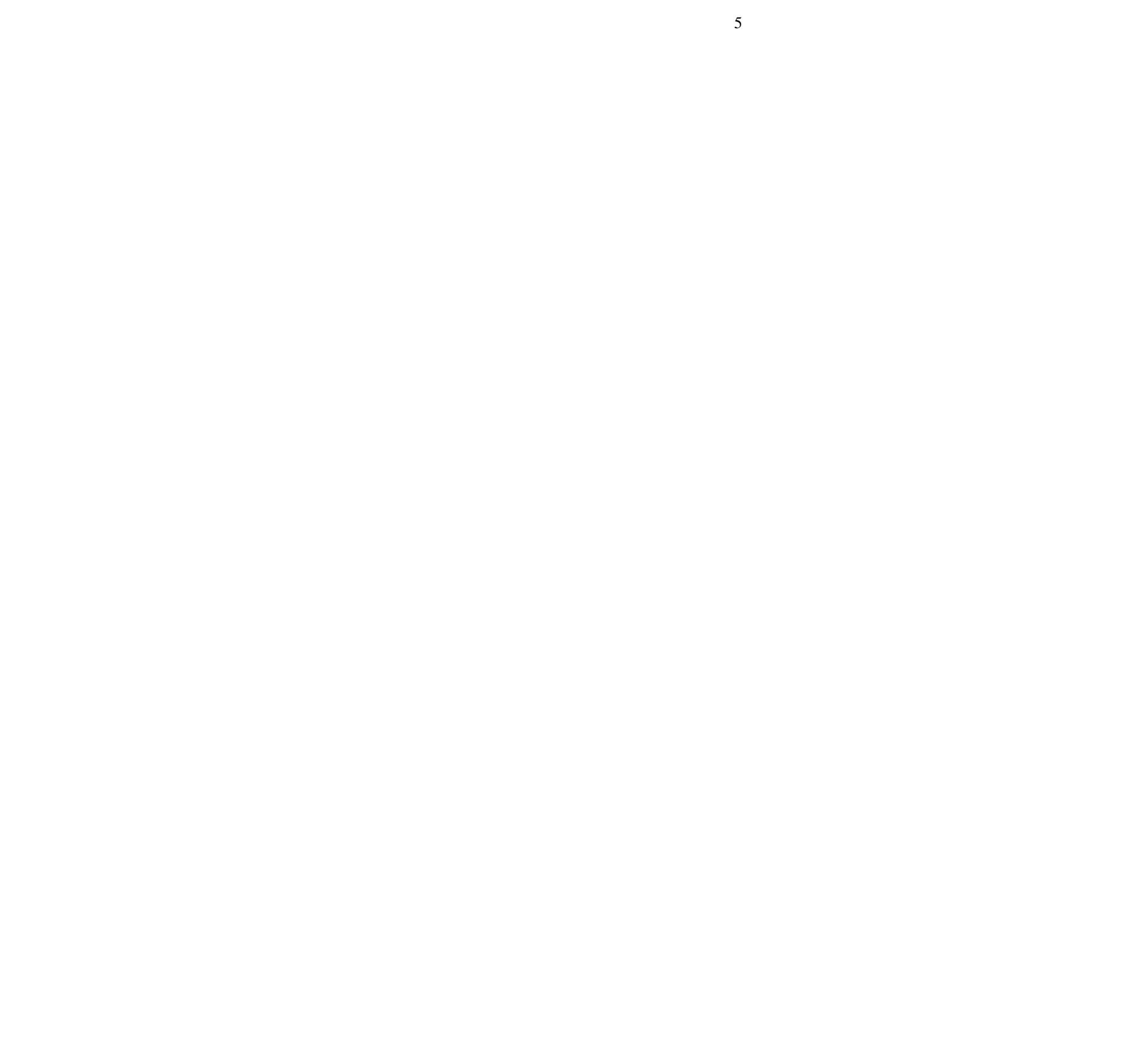


Figure 4: Comparison of HSLDA models across various datasets. The figure consists of a 4x2 grid of plots. Rows represent different datasets: (a) Medical Diagnosis Codes, (b) Amazon Product Categories, (c) Amazon Product Descriptions, and (d) Movie Reviews. Columns represent different models: (left) HSLDA fit by running LDA first then running tree-conditional regressions, and (right) HSLDA fit with fixed random regression parameters. Each plot shows Sensitivity vs 1-Specificity for both ancestor and leaf label prediction performance.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include pages and their placement in Web directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work. We propose a hierarchical topic model of bag-of-word data as a means to perform classification tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to structure-agnostic models. Furthermore, we show evidence that the resulting HSLDA topics are more descriptive of the underlying data than sLDA topics, which ignore the label hierarchy.

1 Introduction

The task of multi-label classification, selecting the k -best labels for a given instance, has been a topic of research for several years. One simplistic way to treat the classification is through a series of independent binary classifiers, but this ignores the many inherent dependencies among the labels. Thus, much work has been devoted on incorporating the co-occurrence patterns of the labels into the classification task. In this paper, we focus on multi-label classification, where the labels are organized in a hierarchical structure. Scenarios of use include, but are not limited to, placing webpages into manually curated Internet directories [2], categorizing images according to a taxonomy, tagging product descriptions with catalogue information [3], and assigning diagnosis codes to clinical records [1].

These are several challenges entailed in incorporating the hierarchical nature of labels into the classification task. One pertains to the labeling itself: in the datasets (especially real-world, noisy ones), for a given label, instances labeled with it contribute positive instance, but it is unclear how to determine the negative instances. In particular, how to treat the labels of the selected ones?

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. In particular, we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. We hypothesize that hierarchical label information provides more information about labeling than considering labels as a flat list.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Following the trend observed in supervised topic modeling, we note that the learned topic models are more representative of the underlying data in both of our datasets [5].

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each predictor is a variable $y_{l,d} \in \{0, 1\}$ indicating whether label l is present in document d . In most cases $y_{l,d}$ is 0. Note that the choice of variables $y_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $y_{l,d}$ ’s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from leaning to assign all labels to all documents. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include pages and their placement in Web directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, comparing the performance of HSLDA against the state-of-the-art bag-of-word data is also of interest. Our experiments show that HSLDA outperforms LDA in hospital document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to structure-agnostic models. Furthermore, we show evidence that the resulting HSLDA topics are more descriptive of the underlying data than sLDA topics, which ignore the label hierarchy.

1 Introduction

The task of multi-label classification, selecting the k -best labels for a given instance, has been a topic of research for several years. One simplistic way to carry out the classification is through a series of independent binary classifiers, but this ignores the many inherent dependencies among the labels. Thus, much work has been devoted on incorporating the co-occurrence patterns of the labels into the classification task. In this paper, we focus on multi-label classification, where the labels are organized in a hierarchical structure. Scenarios of use include, but are not limited to, placing webpages into manually curated Internet directories [2], categorizing images according to a taxonomy, tagging product descriptions with catalogue information [3], and assigning diagnosis codes to clinical records [1].

There are several challenges entailed in incorporating the hierarchical nature of labels into the classification task. One pertains to the labeling itself: in the datasets (especially real-world, noisy ones), for a given label, instances labeled with it contribute positive instance, but it is unclear how to determine the negative instances. In particular, how to treat the labels of the selected ones?

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. In particular, we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. We hypothesize that hierarchical label information provides more information about labeling than considering labels as a flat list.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Following the trend observed in supervised topic modeling, we note that the learned topic models are more representative of the underlying data in both of our datasets [5].

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

© 2009 The authors. ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2009, pp. 6–14. © 2009 ACM, Inc. This is the author's version of the work. It is available online at <http://www.acm.org>. Any changes made after acceptance by the editor will be noted in the journal version. DOI: <http://doi.acm.org/10.1145/1553376.1553387>

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each node l has a value $y_{l,d} \in \{0, 1\}$ corresponding to whether or not it is applied to document d . Note that label $y_{l,d}$ can only be applied to document d if its parent label $y_{p(l),d}$ is also applied (these constraints are specific to is-a constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence. The net effect of this is that label expressions are specific to is-a constraints but can be modified to accommodate different constraints.

Other models that incorporate LDA and supervision include LabeledLDA[22], DiscLDA[17], and other models applied to computer vision and document networks[27, 28]. These models, however, do not implement constraints on the label space.

In other work, researchers have classified documents into a hierarchy with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces and has focused on single label classification without a model of the remaining value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , i.e. $y_{l,d} = 1$, then all labels in the label hierarchy up to the root r are also applied to document d , i.e. $y_{p(l),d} = 1, y_{p(p(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l is marked as not applying to a document then no descendant of label l may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

Here $z_d^T = [z_{1,d}, \dots, z_{k,d}, \dots, z_{|\mathcal{L}|,d}]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k ; $\tilde{z}_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generated by first drawing a vector of conditionally dependent topic proportions $\gamma_d = [\gamma_{1,d}, \dots, \gamma_{k,d}]$, then drawing topic assignments for words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include ratings or associations with other documents, grades for essays, and the number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [2].

Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each node l has a value $y_{l,d} \in \{0, 1\}$ corresponding to whether or not it is applied to document d . Note that label $y_{l,d}$ can only be applied to document d if its parent label $y_{p(l),d}$ is also applied (these constraints are specific to is-a constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence. The net effect of this is that label expressions are specific to is-a constraints but can be modified to accommodate different constraints.

Other models that incorporate LDA and supervision include LabeledLDA[22], DiscLDA[17], and other models applied to computer vision and document networks[27, 28]. These models, however, do not implement constraints on the label space.

There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference procedures do not allow the response to reflect the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA [5]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from leaning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ

parameter, and to a lesser extent, the σ parameter above. Since z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

© 2009 The authors. ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2009, pp. 6–14. © 2009 ACM, Inc. This is the author's version of the work. It is available online at <http://www.acm.org>. Any changes made after acceptance by the editor will be noted in the journal version. DOI: <http://doi.acm.org/10.1145/1553376.1553387>

Figure 1: HSLDA graphical model

Figure 1 shows the HSLDA graphical model. The model consists of several components: a root node α' with an outgoing edge to a node β ; nodes θ_d and η_d with incoming edges from β ; a node γ_d with an outgoing edge to a node $\omega_{n,d}$; a node $\omega_{n,d}$ with an outgoing edge to a node $\omega_{k,d}$; a node $\omega_{k,d}$ with an outgoing edge to a node η_d ; a node η_d with an outgoing edge to a node μ ; a node μ with an outgoing edge to a node σ ; a node σ with an outgoing edge to a node γ ; and finally, a node γ with an outgoing edge to a node $z_{n,d}$. There are also several unlabeled nodes in the middle layer, represented by circles with diagonal hatching. The connections between these unlabeled nodes and the labeled nodes θ_d , η_d , and $\omega_{n,d}$ are indicated by dashed lines.

Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

Figure 2 consists of four plots arranged in a 2x2 grid. The top row contains plots (a) and (b), and the bottom row contains plots (c) and (d). All plots share the same axes: the x-axis for (a) and (c) is Sensitivity (0.0 to 1.0), and the x-axis for (b) and (d) is 1-Specificity (0.0 to 1.0). The y-axis for all plots is Threshold (0.3 to 3). Plot (a) shows Sensitivity vs. Threshold for ancestor prediction performance. Plot (b) shows 1-Specificity vs. Threshold for given leaf labels. Plot (c) shows Sensitivity vs. Threshold for given leaf labels, aligned on the threshold value from plot (b). Plot (d) shows 1-Specificity vs. Threshold for given leaf labels, aligned on the threshold value from plot (b). The plots show various curves representing different models: solid (HSLDA), dashed (independent regressors + sLDA), dotted (HSLDA fit by running LDA first then running tree-conditional regressions), and dot-dashed (HSLDA fit with fixed random regression parameters).

Figure 3: Out-of-sample Amazon product code predictions from product free-text descriptions. In all figures solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), dotted is HSLDA fit by running LDA first then running tree-conditional regressions, and dot-dashed is HSLDA fit with fixed random regression parameters. Top row: (a) includes ancestor prediction performance, (b) results are for given (leaf) labels alone. Bottom row: (c) are the sensitivity curves from (b) aligned on threshold value, (d) are the 1-specificity curves from (b) aligned on threshold value.

Figure 3 consists of four plots arranged in a 2x2 grid. The top row contains plots (a) and (b), and the bottom row contains plots (c) and (d). All plots share the same axes: the x-axis for (a) and (c) is Sensitivity (0.0 to 1.0), and the x-axis for (b) and (d) is 1-Specificity (0.0 to 1.0). The y-axis for all plots is Threshold (0.3 to 3). Plot (a) shows Sensitivity vs. Threshold for ancestor prediction performance. Plot (b) shows 1-Specificity vs. Threshold for given leaf labels. Plot (c) shows Sensitivity vs. Threshold for given leaf labels, aligned on the threshold value from plot (b). Plot (d) shows 1-Specificity vs. Threshold for given leaf labels, aligned on the threshold value from plot (b). The plots show various curves representing different models: solid (HSLDA), dashed (independent regressors + sLDA), dotted (HSLDA fit by running LDA first then running tree-conditional regressions), and dot-dashed (HSLDA fit with fixed random regression parameters).

Figure 4: Sensitivity and 1-specificity curves for HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 4 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 5: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 5 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 6: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 6 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 7: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 7 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 8: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 8 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 9: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 9 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 10: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 10 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 11: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 11 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 12: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 12 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 13: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 13 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 14: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Figure 14 consists of two rows of plots. The top row shows Sensitivity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid line) and sLDA (dashed line). Both rows have Threshold on the x-axis (0.3 to 3) and Sensitivity or 1-Specificity on the y-axis (0.0 to 1.0).

Figure 15: Comparison of HSLDA and sLDA on the Amazon dataset. The top row shows Sensitivity vs. Threshold for HSLDA (solid) and sLDA (dashed). The bottom row shows 1-Specificity vs. Threshold for HSLDA (solid) and sLDA (dashed).

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include pages and their placement in Web directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work. We propose a novel hierarchical Dirichlet model of bag-of-word data that is also of interest to domain-specific LDA on large-scale clinical document labeling and related product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to structure-agnostic models. Furthermore, we show evidence that the resulting HSLDA topics are more descriptive of the underlying data than sLDA topics, which ignore the label hierarchy.

1 Introduction

The task of multi-label classification, selecting the k-best labels for a given instance, has been a topic of research for several years. One simplistic way to approach the classification is through a series of independent binary classifiers, but this ignores the many inherent dependencies among the labels. Thus, much work has been devoted on incorporating the co-occurrence patterns of the labels into the classification task. In this paper, we focus on multi-label classification, where the labels are organized in a hierarchical structure. Scenarios of use include, but are not limited to, placing webpages into manually curated Internet directories [2], categorizing images according to a taxonomy, tagging product descriptions with catalogue information [3], and assigning diagnosis codes to clinical records [1].

There are several challenges entailed in incorporating the hierarchical nature of labels into the classification task. One pertains to the labeling itself: in the datasets (especially real-world, noisy ones), for a given label, instances labeled with it contribute positive instance, but it is unclear how to determine the negative instances. In particular, how to treat the parent labels of the selected ones?

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. In particular, we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. We hypothesize that hierarchical label information provides more information about labeling than considering labels as a flat list.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Following the trend observed in supervised topic modeling, we note that the learned topic models are more representative of the underlying data in both of our datasets [5].

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{d,l} \in \{0, 1\}$ indicating whether or not it contains label l . In most cases $y_{d,l}$ is 0. Note that the choice of variables $y_{d,l}$ and how they are distributed varies driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $y_{d,l}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from leaning to assign all labels to all documents. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

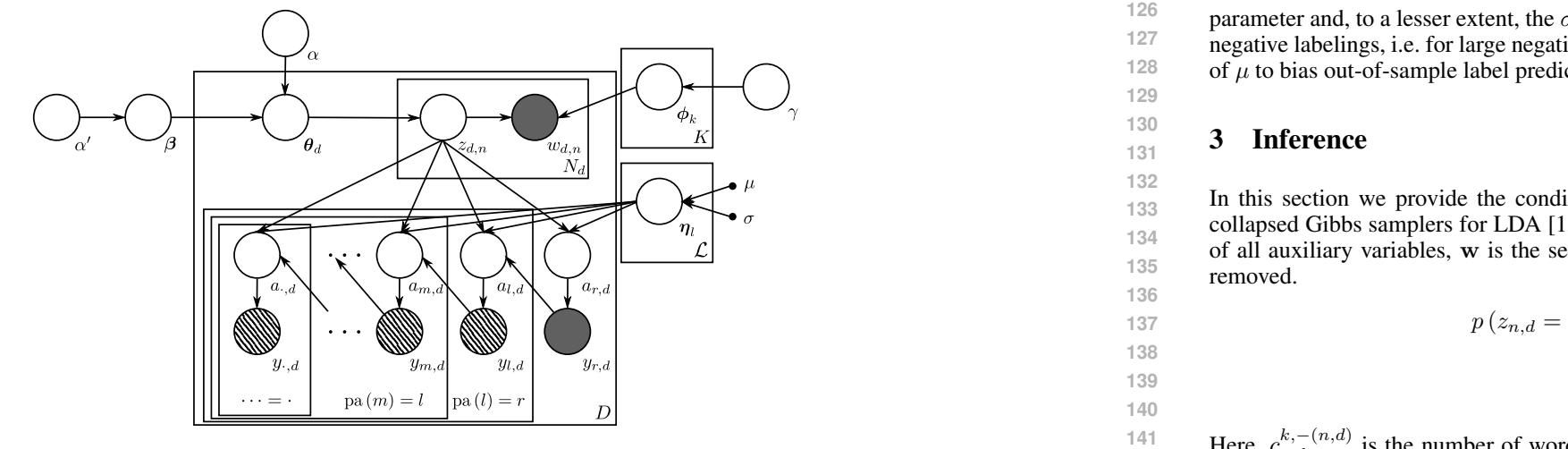


Figure 1: HSLDA graphical model

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA “topics” (distributions over the elements of Σ), ϕ_k is a distribution over “words,” θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $N_K(\cdot)$ is the K -dimensional Normal distribution, I_K is the K -dimensional identity matrix, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim N_K(\mu \mathbf{1}_K, \sigma I_K)$
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha')$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_k \sim \text{Multinomial}(\phi_{k,z_{n,d}})$
 - Set $y_{d,n} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{i,d} | \alpha, \mathbf{I}_d, \eta_l \sim \mathbb{I}(y_{d,l}) \mathcal{N}(\mathbf{0}, \mathbf{I}_{N_d})$
 - Draw $a_{i,d} | \alpha, \mathbf{I}_d, \eta_l \sim \mathbb{I}(y_{d,l}) \mathcal{N}(y_{d,l} \eta_l, \mathbf{I}_{N_d})$
 - Apply label l to document d according to $a_{i,d}$
 - 5. For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $y_{d,l} | \alpha, \mathbf{I}_d, \eta_l \sim \mathcal{N}(\eta_l, \mathbf{I}_{N_d})$
 - Draw $y_{d,l} | \alpha, \mathbf{I}_d, \eta_l \sim \mathcal{N}(\eta_l, \mathbf{I}_{N_d})$
 - Apply label l to document d according to $y_{d,l}$

Here $\mathbf{z}_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k ; $\tilde{z}_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generated by applying a generalization of conditionally dependent probit

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include pages and their placement in Web directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction of the primary goal of this work, we show that HSLDA improves upon the state-of-the-art bag-of-word data is also of interest in downstream NLP tasks like clinical document labeling and related product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to structure-agnostic models. Furthermore, we show evidence that the resulting HSLDA topics are more descriptive of the underlying data than sLDA topics, which ignore the label hierarchy.

1 Introduction

The task of multi-label classification, selecting the k-best labels for a given instance, has been a topic of research for several years. One simplistic way to approach the classification is through a series of independent binary classifiers, but this ignores the many inherent dependencies among the labels. Thus, much work has been devoted on incorporating the co-occurrence patterns of the labels into the classification task. In this paper, we focus on multi-label classification, where the labels are organized in a hierarchical structure. Scenarios of use include, but are not limited to, placing webpages into manually curated Internet directories [2], categorizing images according to a taxonomy, tagging product descriptions with catalogue information [3], and assigning diagnosis codes to clinical records [1].

There are several challenges entailed in incorporating the hierarchical nature of labels into the classification task. One pertains to the labeling itself: in the datasets (especially real-world, noisy ones), for a given label, instances labeled with it contribute positive instance, but it is unclear how to determine the negative instances. In particular, how to treat the parent labels of the selected ones?

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. In particular, we extend supervised latent Dirichlet allocation (sLDA) [5] to take advantage of hierarchical supervision. We hypothesize that hierarchical label information provides more information about labeling than considering labels as a flat list.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Following the trend observed in supervised topic modeling, we note that the learned topic models are more representative of the underlying data in both of our datasets [5].

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at a cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document d value $y_{d,l}$ is 1 if label l is applied to document d , i.e. $y_{d,l} = 1$, then all labels in the label hierarchy up to the root r are also applied to document d , i.e. $y_{d,p(l)} = 1, y_{d,p(p(l))} = 1, \dots, y_{d,r} = 1$. Conversely, if a label l is marked as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

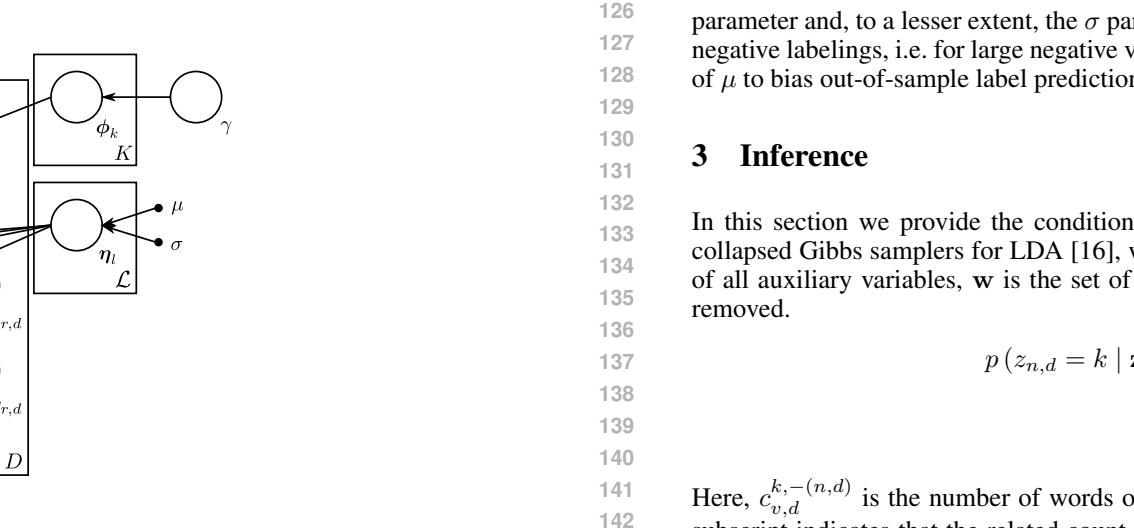


Figure 1: HSLDA graphical model

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA “topics” distributions over the elements of Σ , ϕ_k is a distribution over “words,” θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $N_K(\cdot)$ is the K -dimensional Normal distribution, I_K is the K -dimensional identity matrix, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l \mid \mu, \sigma \sim N_K(\mu_K, \sigma_K)$
 - 3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha')$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_k \sim \text{Multinomial}(\phi_{z_{n,d}})$
 - Set $y_{d,l} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} \mid z_{d,l}, \eta_l, y_{d,l} \sim \frac{\mathbb{I}(y_{d,l})}{\mathbb{I}(y_{d,l}) + 1} \mathcal{N}(\bar{z}_{d,l} \eta_l, 1)$, $y_{d,l} = 1$
 - Apply label l to document d according to $a_{l,d}$

$y_{d,l} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here $\bar{z}_{d,l}^T = [\bar{z}_{d,1}, \dots, \bar{z}_{d,1}, \dots, \bar{z}_{d,K}]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k ; $\bar{z}_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generated by first drawing a collection of conditionally dependent topic proportions [2, 7, 1]. For every label $l \in \mathcal{L}$, however, each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include ratings associated with reviews, grades for essays, and the number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5]. Other models that incorporate LDA and supervision include LabeledLDA[23], DiscLDA[18], and other models applied to computer vision and document networks[28, 8]. These models, however, do not implement constraints on the label space. There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant improvement over LDA, especially when the label space is large. The third comparison model examined not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$ ’s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from leaning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ

parameter, and to a lesser extent, the σ parameter above. Since z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{d,l} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5.1 Data and Pre-Processing

5.1.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient’s complaints, findings and diagnoses, along with treatment and hospital course. For each admission trained medical coders review the information in the discharge summary and assign a series of diagnoses codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a labeling of a patient’s diagnoses based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [11], and sometimes make mistakes [13].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP community challenge described in the challenge, however different from ours in its scope. The datasets were used for training and testing (400,000 training documents) and focused on radiology with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Models ranged from manual rules to online learning [9, 15, 12]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [19, 24, 22, 25, 20].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8,39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 53.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries and 1,000 for testing. Our dataset is a standard result from Bayesian normal linear regression [14]. It is a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution[4].

$$p(a_{l,d} \mid z_d \setminus z_{n,d}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \left(\frac{\mathbb{I}(a_{l,d} < 0)}{\mathbb{I}(a_{l,d} > 0)} + \alpha \beta_k \right) \frac{\mathbb{I}(a_{l,d} < 0)}{\binom{N_d}{a_{l,d}}} \prod_{n=1}^{N_d} \exp \left\{ - \frac{(a_{l,d} - \eta_l)^2}{2} \right\} \quad (1)$$

Here, $\binom{N_d}{a_{l,d}}$ is the number of words of type n in document d assigned to topic k omitting the n th word of document d . The (\cdot) in the subscript indicates that the related count is a sum over the omitted subscript variable. Also, \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\boldsymbol{\eta}_l \mid \mathbf{z}, \mathbf{a}, \sigma) = \mathcal{N}(\boldsymbol{\mu}_l, \Sigma) \quad (2)$$

where

$$\boldsymbol{\mu}_l = \Sigma \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right) \quad \Sigma^{-1} = \mathbf{I}_{|\mathcal{L}|} + \mathbf{Z}^T \mathbf{Z}$$

Here \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is z_d , and $\mathbf{a}_l = (a_{l,1}, a_{l,2}, \dots, a_{l,D})^T$. The simplicity of this conditional distribution follows from the choice of probit regression [4]; the specific form of the update is a standard result from Bayesian normal linear regression [14]. It is a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution[4].

$p(a_{l,d} \mid \mathbf{z}, \mathbf{Y}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ - \frac{1}{2} (a_{l,d} - \eta_l^T \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} > 0) \quad (3)$

where $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of $\boldsymbol{\beta}$ and are governed by the following function.

$$\boldsymbol{\beta} \mid \mathbf{z}, \alpha', \alpha \sim \text{Dir}(m_{(1),1} + \alpha', m_{(1),2} + \alpha', \dots, m_{(1),K} + \alpha') \quad (4)$$

Other inferences were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

We were able to deduce the structure of the hierarchy for the Amazon.com products because all ancestors in the hierarchy were included with each category label. For example, “DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine.”

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

We evaluated HSLDA along with three other closely related models against these two datasets. The comparison models included sLDA with independent regressors (no hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. These models were chosen to highlight several aspects of the model including performance in the absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance. Aside from this, all other features of sLDA are preserved between the two models.

Other models that incorporate LDA and supervision include LabeledLDA[23], DiscLDA[18], and other models applied to computer vision and document networks[28, 8]. These models, however, do not implement constraints on the label space.

There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference procedures do not allow the response to use the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA [5]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

¹<http://www.cdc.gov/nchs/icd/icd9cm.htm>

²<http://www.nitk.org>

Here the $\mathbf{z}_d^T = [z_{d,1}, \dots, z_{d,1}, \dots, z_{d,K}]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k ; $\bar{z}_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generated by first drawing a collection of conditionally dependent topic proportions [2, 7, 1]. For every label $l \in \mathcal{L}$, however, each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include ratings associated with reviews, grades for essays, and the number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5]. Other models that incorporate LDA and supervision include LabeledLDA[23], DiscLDA[18], and other models applied to computer vision and document networks[28, 8]. These models, however, do not implement constraints on the label space. There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant improvement over LDA, especially when the label space is large. The third comparison model examined not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$ ’s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from leaning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ

parameter, and to a lesser extent, the σ parameter above. Since z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{d,l} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

The constraints imposed by an is-a label hierarchy are that if the i th label is applied to document d , i.e. $y_{d,i} = 1$, then all labels in the label hierarchy up to the root r are also applied to document d , i.e. $y_{d,p(i)} = 1, y_{d,p(p(i))} = 1, \dots, y_{d,r} = 1$. Conversely, if a label i is marked as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Authors)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in product directories, product descriptions and associated categories, and product reviews. We introduce hierarchical free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on medical clinical data that has been at least in part manually categorized. Examples include but are not limited to webpages and treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [1]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let $z_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{z}_d = \{z_{1,d}, \dots, z_{N_d,d}\}$ be the set of N_d observations in document d . Here $\mathbf{z}_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k , $\bar{z}_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generatively labeled using a hierarchy of conditionally dependent probit regressors [2, 3]. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e., $\mathbb{I}(y_{\text{pa}(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{l,d}$ can only be applied to document d if its parent label $\text{pa}(l)$ is also applied (these expressions are specific to l constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in the model induces a posteriori dependence. The net effect of this is that label l , which is a label hierarchy, forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an l -a label hierarchy are that if the l th label is applied to document d , i.e., $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1$. Conversely, if a label l' is marked as not being applied to a document d then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In

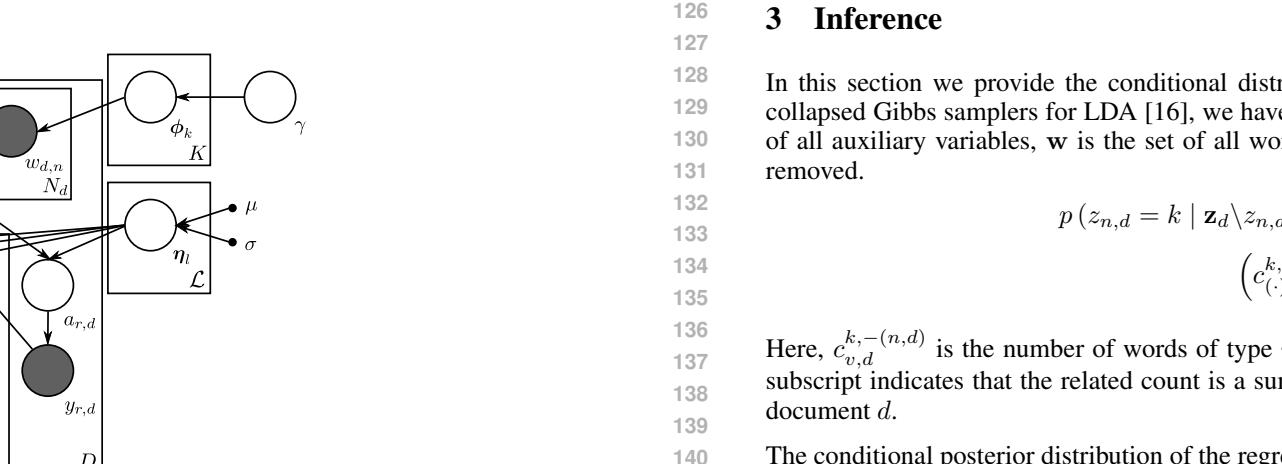


Figure 1: HSLDA graphical model

3 Inference

In this section we provide the conditional distributions required to Gibbs sample the HSLDA posterior distribution. Note that, like in collapsed Gibbs samplers for LDA [16], we have analytically marginalized out the parameters $\phi_{l,K}$ and $\theta_{l,D}$. In the following, \mathbf{a} is the set of all auxiliary variables, \mathbf{w} is the set of all words, η is the set of all regression coefficients, and $\mathbf{z}_d \setminus z_{n,d}$ is the set \mathbf{z}_d with element $z_{n,d}$ removed.

$$p(z_{n,d} = k | \mathbf{z}_d \setminus z_{n,d}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \left(c_{v,d}^{k-(n,d)} + \alpha \beta_v \right)^{\frac{c_{v,d}^{k-(n,d)} + \gamma}{(\beta_v + \gamma)}} \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(z_{l,d}^T \eta_l - a_{l,d})^2}{2} \right\} \quad (1)$$

Here, $c_{v,d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The (\cdot) in the subscript indicates that the related count is a sum over the omitted subscript variable. Also, \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\eta_l | \mathbf{z}_d, \mathbf{a}, \sigma) = \mathcal{N}(\mu_l | \mu_l, \Sigma_l) \quad (2)$$

where

$$\mu_l = \Sigma \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right) \quad \Sigma^{-1} = \mathbf{I} \sigma^{-1} + \mathbf{Z}^T \mathbf{Z}.$$

The text of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge, which focused on the challenge, however, did not include the clinical domain. The dataset was released after 14,000 training and 1,000 testing documents, and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [9, 15, 12]. Other work had leveraged large datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [19, 24, 22, 25, 20].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8,39

associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=309.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

(a)

Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (b) shows predictive performance as a function of the prior mean on regression parameters.

(a)

Figure 3: Out-of-sample Amazon product category predictions from product free-text descriptions. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (b) shows predictive performance as a function of the prior mean on regression parameters.

(b)

Figure 4: Related Work

Amazon.com, an online retail store, organizes its catalog of products in a hierarchy and provides product descriptions for most products in their catalog. Products can be discovered by users through query-based searching or through product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

We were able to deduce the structure of the hierarchy for the Amazon.com products because all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine."

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91-89 terms on average, std dev=533.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

5.2 Comparison Models

While there has been much work in multi-label classification of text and text modeling in general, we focus here on topic modeling approaches. Latent Dirichlet Allocation (LDA) is a generative probabilistic model which represents documents as a mixed-membership bag of words. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet Allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include ratings associated with online reviews, grades for essays, and the number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance. Aside from this, all other features of sLDA are preserved between the two models.

In other work, researchers have classified documents into a hierarchy with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces and has focused on single label classification without a model of the documents such as HSLDA [21, 10, 17, 7].

There are two components to HSLDA. LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the responses to influence the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA [5]. This comparison model examines not the structuring of the label space, but the benefit of predictors deeper in the label hierarchy are able to focus on finding specific, contextual labeling features. We believe this to be a significant leap in the combined inference procedure we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variables distributions for the $a_{l,d}$'s yield conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because \bar{z}_k is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where the auxiliary labels $a_{l,d}$ are observed, the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing

5.1.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient's complaints, findings and diagnoses, along with treatment and hospital course. For each admission trained medical coders review

the last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where the auxiliary labels $a_{l,d}$ are observed, the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

5.2 Comparison Models

While there has been much work in multi-label classification of text and text modeling in general, we focus here on topic modeling approaches. Latent Dirichlet Allocation (LDA) is a generative probabilistic model which represents documents as a mixed-membership bag of words. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet Allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include ratings associated with online reviews, grades for essays, and the number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance. Aside from this, all other features of sLDA are preserved between the two models.

In other work, researchers have classified documents into a hierarchy with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces and has focused on single label classification without a model of the documents such as HSLDA [21, 10, 17, 7].

There are two components to HSLDA. LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the responses to influence the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA [5]. This comparison model examines not the structuring of the label space, but the benefit of predictors deeper in the label hierarchy are able to focus on finding specific, contextual labeling features. We believe this to be a significant leap in the combined inference procedure we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variables distributions for the $a_{l,d}$'s yield conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because \bar{z}_k is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where the auxiliary labels $a_{l,d}$ are observed, the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

5.2 Comparison Models

While there has been much work in multi-label classification of text and text modeling in general, we focus here on topic modeling approaches. Latent Dirichlet Allocation (LDA) is a generative probabilistic model which represents documents as a mixed-membership bag of words. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet Allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include ratings associated with online reviews, grades for essays, and the number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance. Aside from this, all other features of sLDA are preserved between the two models.

In other work,

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in product directories, product descriptions and associated categories, and product reviews. We introduce hierarchical and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and we improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured data that has been at least in part manually categorized. Examples include but are not limited to weblogs and curated hierarchical structures of the same [2], product descriptions and catalogs (e.g. [2] as available from [3]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [2]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and Web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let Σ be the union of all words. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e., $\text{pa}(l,d) = 1$) are used to determine whether or not label l is applied to document d , i.e., $y_{l,d} = 1$ if $\text{pa}(l,d) = 1$. Conversely, if a label l' is marked as not applying to a document d no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the labeled variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In

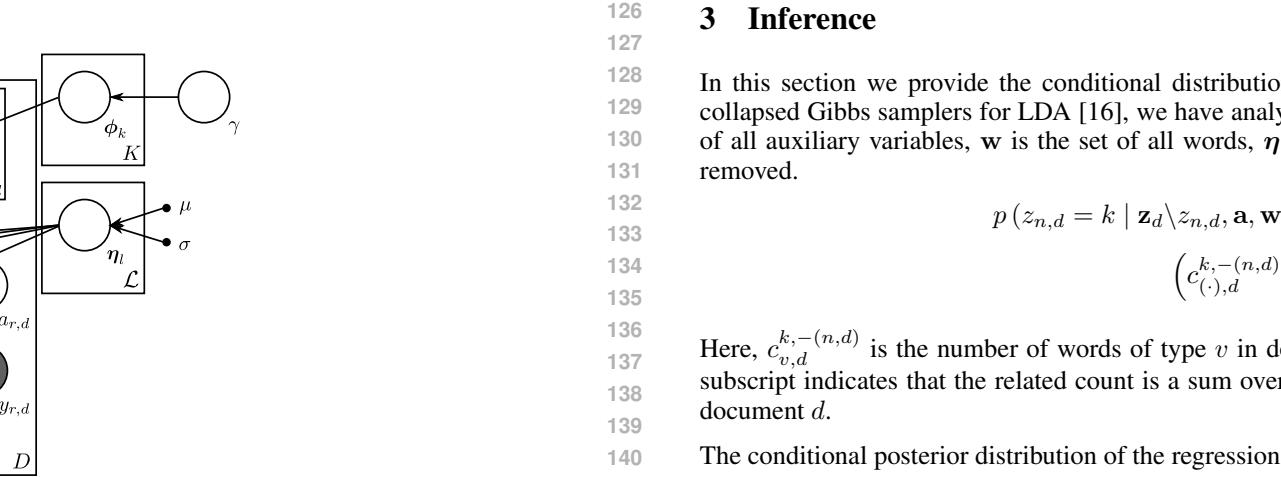


Figure 1: HSLDA graphical model

3 Inference

In this section we provide the conditional distributions required to Gibbs sample the HSLDA posterior distribution. Note that, like in collapsed Gibbs samplers for LDA [16], we have analytically marginalized out the parameters $\phi_{l,K}$ and $\theta_{l,D}$. In the following, a is the set of all auxiliary variables, w is the set of all words, η is the set of all regression coefficients, and $z_d \setminus z_{n,d}$ is the set z_d with element $z_{n,d}$ removed.

$$p(z_{n,d} = k | z_d \setminus z_{n,d}, a, w, \eta, \alpha, \beta, \gamma) \propto \left(c_{\gamma,d}^{k - (n,d)} + \alpha \beta_k \right) \frac{c_{\gamma,d}^{k - (n,d)} + \gamma}{\binom{N_d}{k} \cdot V \gamma} \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(z_{l,d}^T \eta_l - a_{l,d})^2}{2} \right\} \quad (1)$$

Here, $c_{\gamma,d}^{k - (n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The (\cdot) in the subscript indicates that the related count is a sum over the omitted subscript variable. Also, \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\eta_l | z_d, a, \sigma) = \mathcal{N}(\mu_l, \Sigma_l) \quad (2)$$

where

$$\mu_l = \sum \left(\frac{\mu}{\sigma} + Z^T a_l \right) \Sigma^{-1} = \mathbf{I} \sigma^{-1} + Z^T Z.$$

Here Z is a $D \times K$ matrix such that row d of Z is z_d , and $a_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$. The simplicity of this conditional distribution follows from the K -dimensional Normal distribution, $\mathcal{N}_K(\cdot)$ is the K -dimensional identity matrix, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

$$p(a_{l,d} | z_d, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ - \frac{1}{2} (a_{l,d} - \eta_l^T z_d) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0). \quad (3)$$

HSLDA implements an asymmetric prior as a hierarchical Dirichlet prior over topic assignments (i.e., β is a parameter in our model). This has been shown to improve the quality and stability of inferred topics [27]. Sampling β is done via a "direct assignment" method of Teh et al. [26]

$$\beta | z, \alpha', \alpha \sim \text{Dir} (m_{(1),1} + \alpha', m_{(1),2} + \alpha', \dots, m_{(1),K} + \alpha') \quad (4)$$

where $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β and are governed by the following function.

$$p(m_{d,k} = m | z_d, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha_k)}{\Gamma(\alpha_k + c_{d,k}^k)} \cdot \frac{(c_{d,k}^k)^m}{m!} \quad (5)$$

$s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

4 Related Work

While there has been much work in multi-label classification of text and text modeling in general, we focus here on topic modeling approaches. Latent Dirichlet allocation (LDA) is a generative probabilistic model which represents documents as a mixed-membership bag of word. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include ratings associated with online reviews, grades for essays, and the number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

Other models that incorporate LDA and supervision include LabelerLDA[23], DiscLDA[18], and other models applied to computer vision and document networks[28, 8]. These models, however, do not implement constraints on the label space.

In other work, researchers have classified documents into a hierarchy with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces and has focused on single label classification without a model of documents such as LDA[21, 10, 17, 7].

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing

5.1.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize the course of a hospitalized patient. The summaries typically contain a record of the patient's complaints, findings and diagnoses, along with treatment and hospital course. For each admission trained medical coders review

the information in the discharge summary and assign a series of diagnosis codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ As such, the ICD-9 codes constitute a labeling of a patient's diagnoses based on a discharge summary. The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for "Pneumonia due to adenovirus" is a child of the code for "Viral pneumonia," where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [11], and sometimes make mistakes [13].

The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics of performance - sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that model performance generalizes over two very disparate domains.

To evaluate the performance of these models, we establish a gold standard for comparison. In our evaluation of each dataset, a gold set of 1000 discharge summaries and 1000 test documents were used.

Specifically, ancestors of observed nodes were ignored, observed nodes were considered positive and unobserved nodes were considered to be negative. This method of defining positive and negative labels was chosen to be as fair as possible to all models being compared. In particular, since the SLDA model does not enforce the hierarchical constraints, the models can be compared on a more equal footing by focusing on the observed labels as positives. The gold standard defined in this way will likely lead to a slight overestimation of the number of false positives. It is known that ICD-9 codes lack sensitivity and their use as a gold standard could lead to correctly positive predictions being labeled as false positives. However, given that the label space is often large (as in our examples) it is a moderate assumption that erroneous false positives should not skew results significantly.

Predictive performance in HSLDA is evaluated by $p(y_{l,d} | w_1, N_d, d, w_1, N_d, l, d, y_{l,d}, l)$ where d represents the test document. For efficiency, the expectation of this probability distribution is estimated in the following way. Expectations of z_d and η_l were estimated from the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

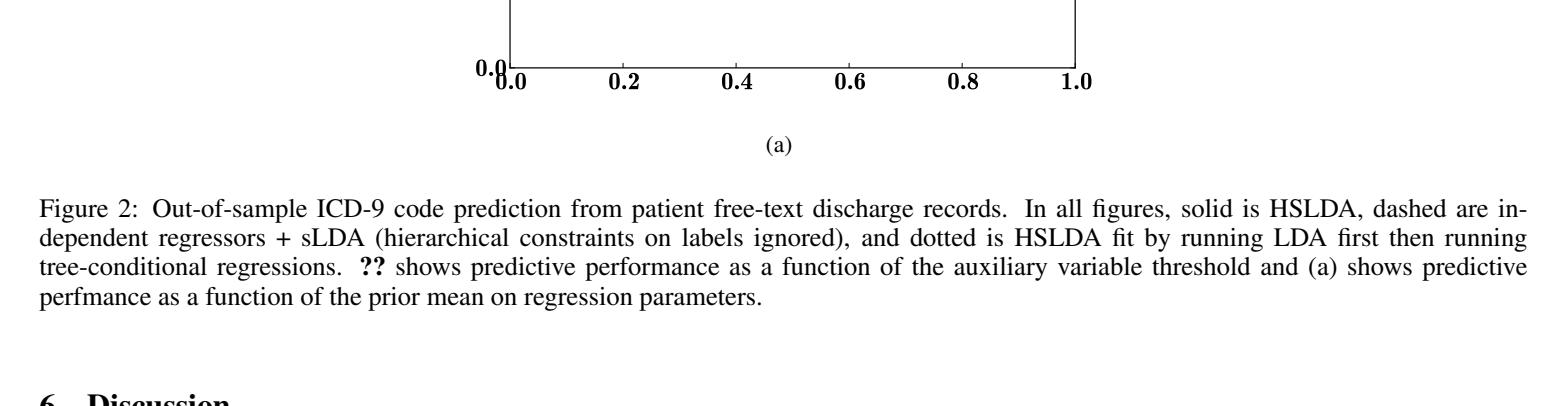
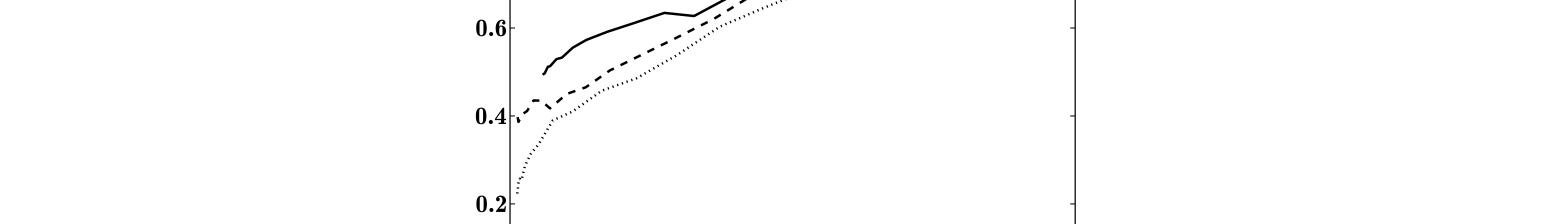


Figure 2: Out-of-sample Amazon product category predictions from product free-text descriptions. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

We have described a mixed membership model with hierarchical supervision. We have demonstrated this model in the context of document modeling with hierarchical multi-label supervision. Such a model is appropriate in domains where there are hierarchical constraints among the labels such as is the case in an IS-A hierarchy.

...
• what about the nonparametric version of this?
• discuss the broader goal, from the beginning of search engine time, to combine categorization and free text, this, to our knowledge, for all these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters are specific to is-a constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in the model induces a posterior correlation. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextual labeling features. We believe this is a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where the auxiliary variable is a categorical label, Gaussian priors over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

5 References
[1] The computational medicine center's 2007 medical natural language processing challenge. <http://www.computationalmedicine.org/challenge/previous-2007>
[2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.

5

6 Discussion

We evaluated HSLDA along with three other closely related models against these two datasets. The comparison models included sLDA with independent regressions (hierarchical constraints on labels ignored), HSLDA fit by performing LDA followed by tree-conditional regressions, and HSLDA fit with fixed random regression parameters. These models were chosen to highlight several aspects of the model including performance in the absence of hierarchical constraints, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance. Aside from this, all other features of sLDA are preserved between the two models.

Other models that incorporate LDA and supervision include LabelerLDA[23], DiscLDA[18], and other models applied to computer vision and document networks[28, 8]. These models, however, do not implement constraints on the label space.

In other work, researchers have classified documents into a hierarchy with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces and has focused on single label classification without a model of documents such as LDA[21, 10, 17, 7].

There are two components to HSLDA. LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the responses to influence the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA [5]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where the auxiliary variable is a categorical label, Gaussian priors over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

5 References
[1] The computational medicine center's 2007 medical natural language processing challenge. <http://www.computationalmedicine.org/challenge/previous-2007>
[2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.

5

5

5

5

5

5

5

5

5

5

5

5

5

5

5

5

5

5

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in directories, product descriptions and associated categories, and product hierarchy of free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and we improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on medical discharge summaries that have been at least in part manually categorized. Examples include but are not limited to Web pages and curated hierarchical directories of the same [2], product descriptions and catalogs (e.g., 12 [1] as available from [3]), and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [2]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that that label set is complete, i.e., every label is a child of some other label. This assumption is reasonable for most applications. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{i,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{i,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an i -level hierarchy are that if the j th label is applied to document d , i.e., $y_{j,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{i,d} = 1 \Rightarrow y_{j,d} = 1$. Conversely, if a label l' is marked as not applying to a document d then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In

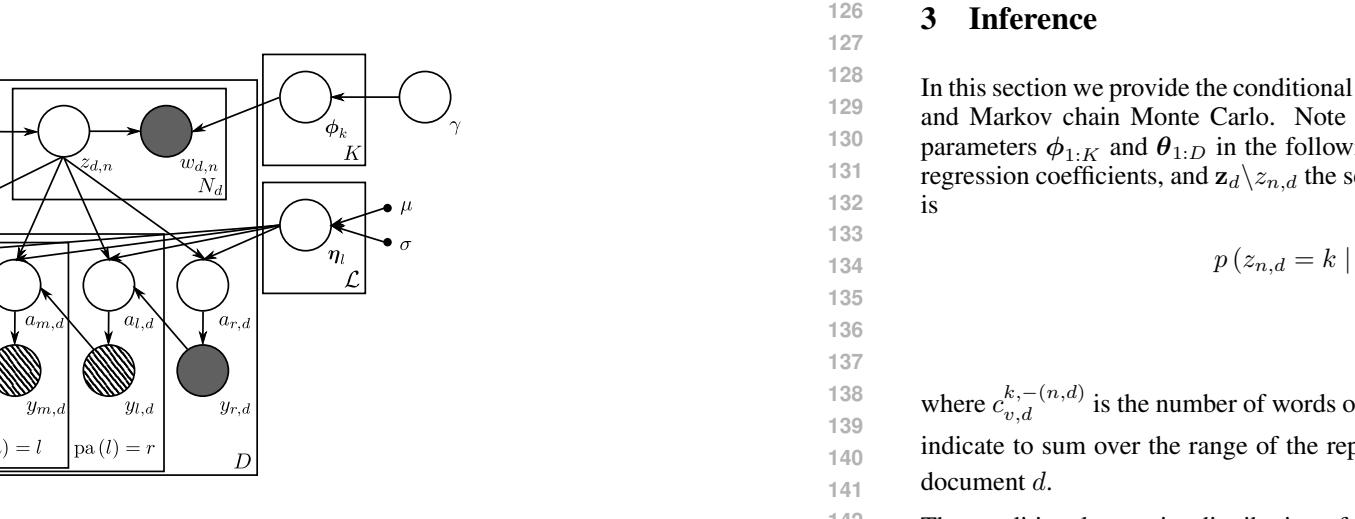


Figure 1: HSLDA graphical model

3 Inference

3.1 Data and Pre-Processing

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [16], we have analytically marginalized out parameters $\phi_{1,K}$ and $\theta_{1,D}$ in the following expressions. Let \mathbf{z} be the set of all auxiliary variables, \mathbf{w} the set of all words, η the set of all regression coefficients, and $\mathbf{z}_{n,d}$ the set \mathbf{z} with element $z_{n,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | \mathbf{z}_{\setminus n,d}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \left(c_{(1),d}^{k-(n,d)} + \alpha \beta_k \right) \frac{c_{(1),d}^{k-(n,d)+1}}{\binom{c_{(1),d}}{k-1} + V} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(\mathbf{z}_{l,d}^\top \eta_l - a_{l,d})^2}{2} \right\} \quad (1)$$

where $c_{(1),d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript (\cdot) 's indicate to sum over the range of the replaced variable, i.e. $c_{(w_n,d),d}^{k-(n,d)} = \sum_{w \in \mathbf{w}_{n,d}} c_{w,d}^{k-(n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\eta_l | \mathbf{z}, \mathbf{a}, \sigma) = \mathcal{N}(\hat{\mu}_l, \hat{\Sigma}) \quad (2)$$

where

$$\hat{\mu}_l = \Sigma \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right) \quad \Sigma^{-1} = \mathbf{I} \sigma^{-1} + \mathbf{Z}^T \mathbf{Z}$$

Here \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is $\mathbf{z}_{d,d}$ and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$. The simplicity of this conditional distribution follows from the choice of probit regression [4]; the specific form of the update standard from Bayesian normal linear regression [14]. It also is a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution [4]

$$p(a_{l,d} | \mathbf{z}, \mathbf{a}, \eta_l) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^\top \mathbf{z}_{d,d}) \right\} \mathbb{I}(a_{l,d} y_{l,d} > 0) \quad (3)$$

HSLDA employs a hierarchical Dirichlet prior over topic assignments (i.e., β is estimated from data rather than fixed a priori). This has been shown to improve the quality and stability of inferred topics [27]. Sampling β is done using the "direct assignment" method of Teh et al. [26]

$$\beta | \mathbf{z}, \alpha' \sim \text{Dir}(m_{(1),1} + \alpha', m_{(1),2} + \alpha', \dots, m_{(1),K} + \alpha') \quad (4)$$

Here $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β . Their conditional posterior distribution is sampled according to

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{\setminus(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + c_{(1),d}^k)} s(c_{(1),d}, m) (\alpha \beta_k)^m \quad (5)$$

where $s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

4 Related Work

While there has been much work in multi-label classification of text and text modeling in general, we focus here on topic modeling approaches. Latent Dirichlet allocation (LDA) is a generative probabilistic model which represents documents as a mixed-membership bag of words. Each document is represented as a collection of words, generated from a set of topic assignments (one for each word), where each topic assignment is drawn from a distribution over topics [6]. sLDA is latent Dirichlet allocation (LDA) [6] augmented with per-document labeling, often taking the form of a single numerical or categorical label. Examples of labels include ratings associated with online reviews, grades for essays, and the number of times a webpage is linked. This approach has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [5].

Other models that incorporate LDA and supervision include LabeledLDA[23], DiscLDA[18], and other models applied to computer vision and document networks[28, 8]. These models, however, do not implement constraints on the label space.

There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the response to see the low dimensional structure imposed by LDA. Combined inference has been shown to improve performance in sLDA [5]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nltk.org>

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$). The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics of performance - sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that model performance generalizes over two very disparate domains.

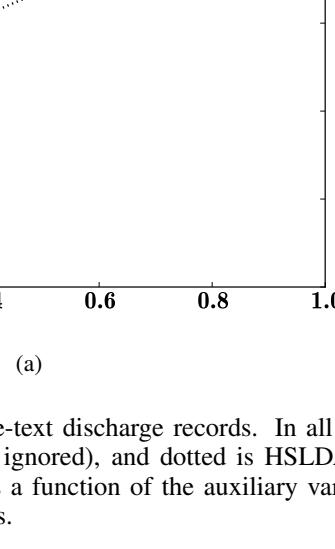
To evaluate the performance of these models, we establish a gold standard for comparison. In our evaluation of each dataset, a held out NLP domain classifier is used in the challenge, however different from ours in its scope. The datasets were used for training and testing (and 1000 testing documents) and focused on radiology with a restricted number of ICD-9 codes (45 of them, compared to 78+ in our dataset). Methods ranged from manual rules to online learning [9, 15, 12]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naïve Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [19, 24, 22, 25, 20].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 53.57 terms after preprocessing (std dev=30.29). We split our dataset into 5,000 discharge summaries and 1,000 for testing.

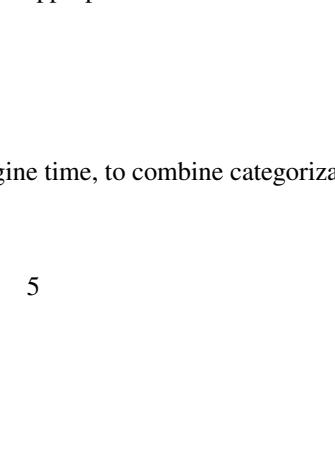
The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP domain challenge in the challenge, however different from ours in its scope. The datasets were used for training and testing by focusing on the observed labels as positive. The gold standard defined in this way will likely lead to a slight overestimation of the number of false positives. It is known that ICD-9 codes lack sensitivity and their use as a gold standard could lead to a correctly positive predictions labeled as false positives. However, given that the label space is often large (as in our example) it is a moderate assumption that erroneous false positives should not result significantly.

Predictive performance in HSLDA is evaluated by $p(y_{i,d} | w_{1:N_d,d}, w_{1:N_d,d} \in \mathcal{L}_{\setminus d}, D)$ where $\mathcal{L}_{\setminus d}$ represents the test document. For efficiency, the expectation of this probability distribution was estimated in the following way. Expectations of $\mathbf{z}_{i,d}$ and η_i were estimated with samples from the posterior. Using these expectations, we performed Gibbs sampling over the hierarchy to acquire predictive samples for the documents in the test set.

(a)



(b)



(c)



(d)



(e)

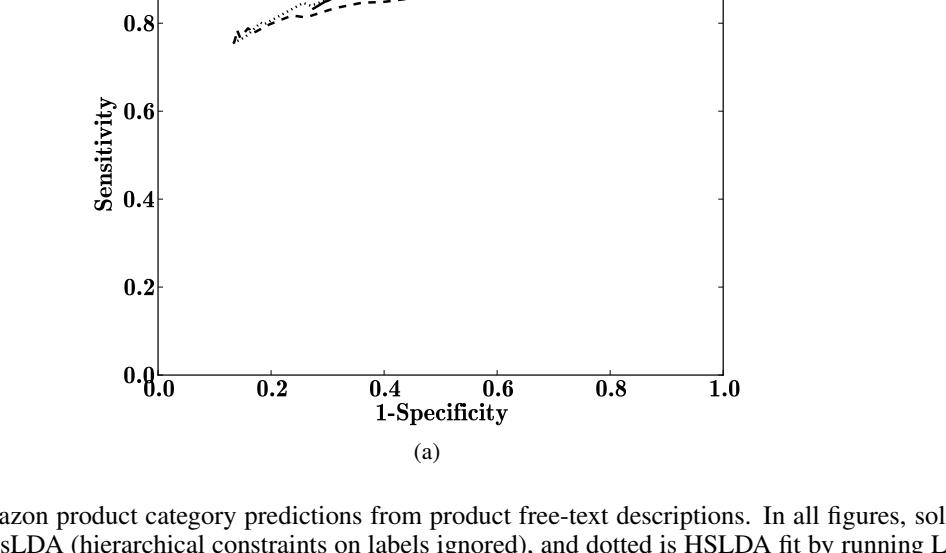


Figure 3: Out-of-sample Amazon product category predictions from product free-text descriptions. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.html, 2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] James H Albert and Siddhartha Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669, 1993.
- [6] David Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [7] Daniel Ramage, David Hall, Rameesh Nallapati, and Christopher D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, pages 248–256, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-59-6. URL <http://aclweb.org/EMNLP2009/pdf/corpora/corpora.pdf>.
- [8] B. RibeiroNeto, AHF Lacerda, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for Information science and Technology*, 52(3):391–401, 2001.
- [9] P. Rech, J. Gobello, I. Thashiti, and A. Grisoliahter. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [10] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(47

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in directories, product descriptions and associated categories, and product hierarchies of free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we consider unstructured textual data that have been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [3], product descriptions and catalogs [e.g. 11] and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [4]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will not for exposition purposes assume that that label set is complete; it is possible that some labels in the hierarchy are unlabeled. Such a label hierarchy forms a tree of unlabeled nodes. Each document has a variable $y_{i,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{i,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an i -level hierarchy are that if the j th label is applied to document d , i.e., $y_{j,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{i,d} = 1 \Rightarrow y_{pa(i),d} = 1$. Conversely, if a label l is not explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling imbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because \bar{y}_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{i,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

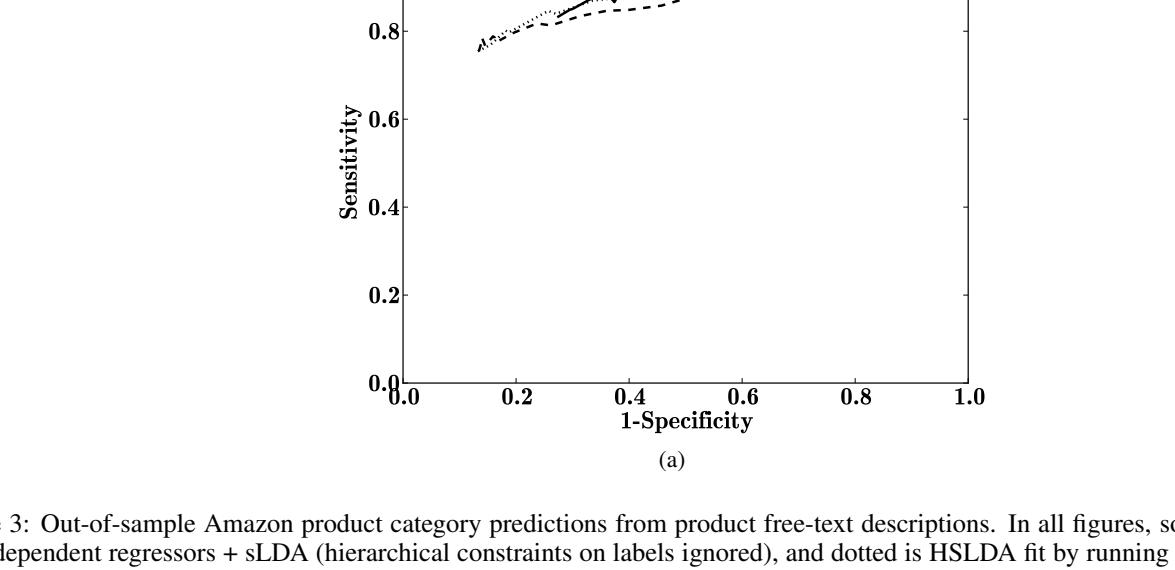


Figure 1: HSLDA graphical model

Here $z_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k . $\bar{y}_d = N_d^{-1} \sum_{n=1}^{N_d} I(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generatively labeled using a hierarchy of conditionally dependent probit regressors [? ?]. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e., $I(y_{pa(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{i,d}$ can only be applied to document d if its parent label $y_{pa(l),d}$ is also applied (these expressions are specific to i as constraints but can be modified to accommodate different constraints). The regression coefficients η_i are independent a priori, however, the hierarchical coupling in the model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually-labeling features. We believe this to be a significant improvement over previous work.

Note that the choice of variables $a_{i,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probability auxiliary variables for the $a_{i,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling imbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because \bar{y}_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{i,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In

1

(a)

Figure 2: Out-of-sample Amazon product category predictions from product free-text descriptions. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007. 2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>. 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA. 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>. 2004.
- [5] James H Albert and Siddhartha Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669, 1993.
- [6] K. Crotzer, M. Dror, & G. Ganchev. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Bioinformatics, Translation, and Clinical Language Processing*, pages 129–136, 2007.
- [7] S. Danai and Hao Chen. Hierarchical classification of web content. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR '00*, pages 258–263, New York, NY, USA, 2000. ACM.
- [8] Yan David S, Nilasena Mantha J, Bradford Brian F, Gage Elena Birman-Deych, Amy D Waterman. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–51, 2005.
- [9] R Farkas and G Srivastava. Automatic construction of rule-based icd-9-cm coding systems. *BMC Bioinformatics*, 9(Suppl 3):S10, 2008.
- [10] Mahdad Farzandifar, Abbas Ghahremani, and F. Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the tenth revision. *International Journal of Information Management*, 30:78–84, 2010.
- [11] Andrew Gelman, John B. Carlin, Hal S. Stern, and Donald B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2nd ed. edition, 2004.
- [12] I Goldstein, A Arzamustov, and O Uzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [13] TL Griffiths and M Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [14] Daphne Koller and Mehran Sahami. Hierarchically classifying documents using very few words. *Technical Report 1997-75*, Stanford InfoLab, February 1997. Previous number = SIDL-WP-1997-0059.
- [15] Simon Lacoste-julien, Fei Sha, and Michael I Jordan. DisLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [16] Leah Larkey and Bruce Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [17] LV Lita, S Yu, S Niculescu, and J Bi. Large scale diagnostic code classification for medical patient records. 2008.
- [18] Andrew McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Building domain-specific search engines with machine learning techniques. In *Proc. AAAI-99 Spring Symposium on Intelligent Agents in Cyberspace*, 1999, 1999.

2

(b)

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in directories, product descriptions and associated categories, and product hierarchies of free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we used medical discharge summaries and clinical records, which include but are not limited to webpages and curated hierarchical directories of the same [3], product descriptions and catalogs [4] and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [4]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified e.g. image catalogs with bag-of-feature image representations).

Our main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have both a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{z}_{n,d} = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$.

Each label except root labels $l \in \mathcal{L}$ has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will call for exposition purposes assume that the label set is complete, i.e., every label has a parent label. This assumption is reasonable since most labels are composed of several words.

Such a label hierarchy forms a tree with a single root $r \in \mathcal{L}$. Each

document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$

will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor

remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an is-a label hierarchy are that if the i th label is applied to document d , i.e. $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{p(l),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l ' is not applied to document d then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation.

Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1.

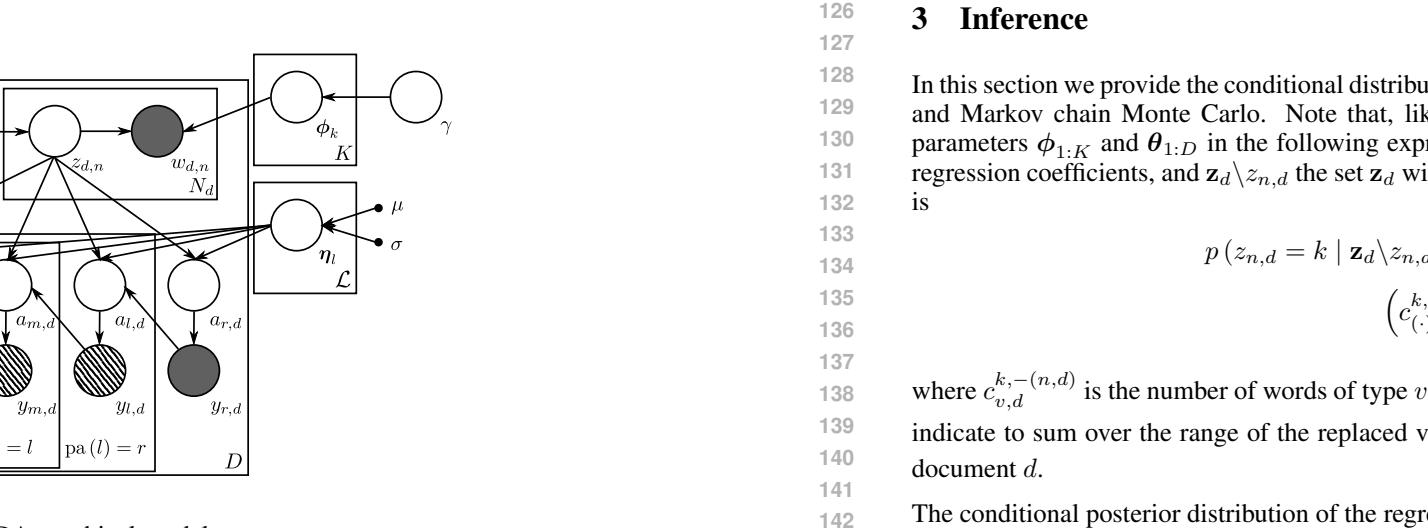


Figure 1: HSLDA graphical model

3 Inference

3.1 Data and Pre-Processing
In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [18], we have analytically marginalized out the parameters $\phi_{l,d}$ and $\theta_{l,D}$ in the following expressions. Let \mathbf{z} be the set of all auxiliary variables, \mathbf{w} the set of all words, η the set of all regression coefficients, and $\mathbf{z}_{n,d}$ the set \mathbf{z}_d with element $z_{n,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | \mathbf{z}_d \setminus z_{n,d}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \\ \left(\frac{c_{k,(n,d)}^{k-(n,d)}}{c_{(n,d)}^{k-(n,d)} + V} \right)^{\eta_{k,(n,d)}} \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(y_{l,d}^k - \eta_{l,(n,d)})^2}{2} \right\} \quad (1)$$

where $c_{v,(n,d)}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript (\cdot) 's indicate to sum over the range of the replaced variable, i.e. $c_{w_n,d}^{k-(n,d)} = \sum_{w \in \mathcal{W}} c_{w_n,d}^{k-(n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\eta_l | \mathbf{z}_d, \mathbf{a}, \sigma) = \mathcal{N}(\mu_l, \Sigma_l) \quad (2)$$

where

$$\mu_l = \Sigma \left(\frac{1}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right) \quad \Sigma^{-1} = \mathbf{I}\sigma^{-1} + \mathbf{Z}^T \mathbf{Z}$$

Here \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is $\mathbf{z}_{d,d}$ and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$. The simplicity of this conditional distribution follows from the choice of probit regression [6]; the specific form of the update is a standard result from Bayesian normal linear regression [16]. It is also a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution [6]

$$p(a_{l,d} | \mathbf{z}_d, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ - \frac{1}{2} (a_{l,d} - \eta_l^T \mathbf{z}_{d,d}) \right\} \mathbb{I}(a_{l,d} > 0) \quad (3)$$

HSLDA employs a hierarchical Dirichlet prior over topic assignments (β is estimated from data rather than fixed a priori). This has been shown to improve the quality and stability of inferred topics [29]. Sampling β is done using the “direct assignment” method of Teh et al. [28]

$$\beta | \mathbf{z}, \alpha, \alpha' \sim \text{Dir}(m_{1,1} + \alpha, m_{1,2} + \alpha', \dots, m_{K,1} + \alpha') \quad (4)$$

Here $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β . Their conditional posterior distribution is sampled according to

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + c_{(d,k)}^k)} s(c_{(d,k),m}^k) (\alpha \beta_k)^m \quad (5)$$

where $s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

4 Related Work

In this work, we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document “supervision”, often taking the form of a single numbered or categorical label. More generally this supervision is just extra per-document data; for instance, in quality or relevance (e.g., online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as a conditional regression, and are often fit with tree-based random regression parameters. These models were chosen to highlight several aspects of the model including performance from some of the topics; topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [7]. It has also been demonstrated that sLDA has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [7]. sLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we contrast HSLDA with sLDA applied in this way.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the responsibility for the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA [7]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing
For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$). The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics of performance – sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that model performance generalizes over two very disparate domains.

To evaluate the performance of the models, we generate a gold-standard in each dataset, and then compare the predicted against the observed labels. Specifically, ancestors of observed nodes were ignored; observed nodes were considered positive and unobserved nodes were considered to be negative. This method of defining positive and negative labels was chosen to be as fair as possible to all models being compared. In particular, since the sLDA model does not enforce the hierarchical constraints, the models can be compared on a more equal footing by promising results [21, 26, 24, 27, 22].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional identity matrix, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 • Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma \mathbf{I}_V)$

2. For each label $l \in \mathcal{L}$
 • Draw a label application coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{I}_K, \sigma \mathbf{I}_K)$

3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}(\alpha' \mathbf{I}_K)$

4. For each document $d = 1, \dots, D$
 • Draw topic proportions $\theta_{d,l} | \beta, \alpha \sim \text{Dir}_K(\alpha \beta_l)$

• For $n = 1, \dots, N_d$
 • Draw topic assignment $z_{n,d} | \theta_{d,l} \sim \text{Multinomial}(\theta_{d,l})$

• Draw word $w_{n,d} | z_{n,d}, \phi_k \sim \text{Multinomial}(\phi_{k,z_{n,d}})$

• Set $y_{p(l),d} = 1$
 • For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r

• Draw $a_{l,d} | \eta_l, y_{p(l),d} \sim \mathcal{N}(y_{p(l),d} \eta_l, 1)$, $y_{p(l),d} = 1$
 • Draw word $w_{n,d} | z_{n,d}, \phi_k \sim \text{Multinomial}(\phi_{k,z_{n,d}})$
 • Apply label l to document d according to $a_{l,d}$

$y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here $z_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k , $z_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generatively labeled using a hierarchy of conditionally dependent probit regressors [$\mathbf{?}$?]. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $\mathbb{I}(y_{p(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{p(l),d}$ can only be applied if its parent label $y_{p(p(l),d)}$ is also applied (these expressions are specific to \mathcal{L} as constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in the model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually labeling features. We believe this to be a significant predictor of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit auxiliary variables for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because \bar{z}_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability

of μ to bias out-of-sample label prediction performance in Section 5.

...
 http://www.cdc.gov/nchs/icd/icd9.htm
 http://www.nitk.org

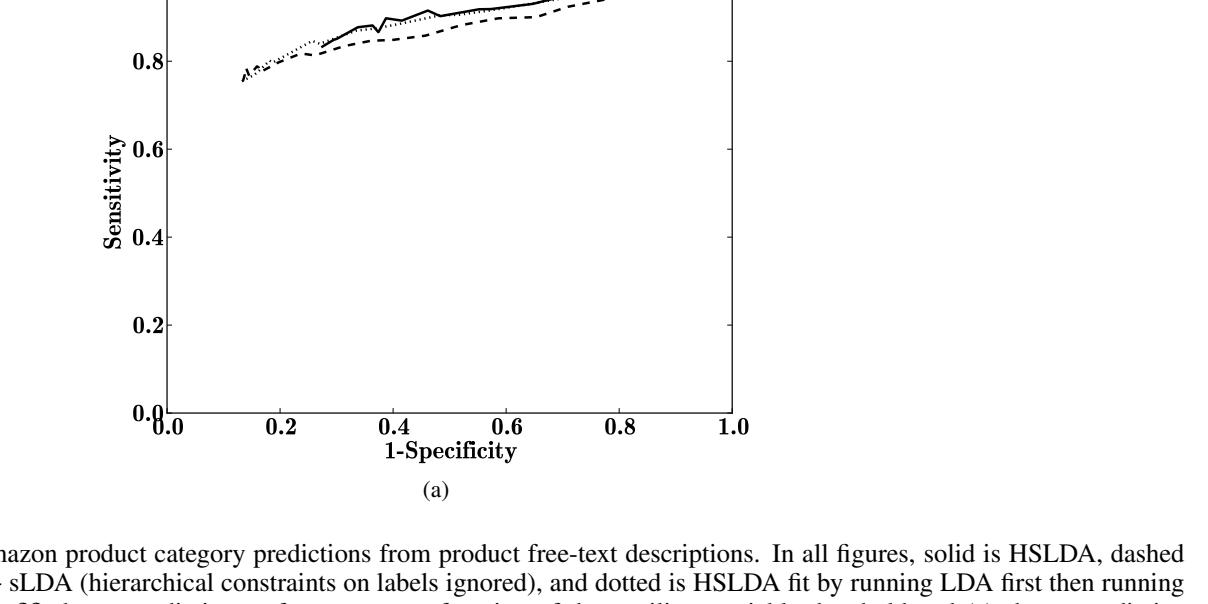


Figure 3: Out-of-sample Amazon product category predictions from product free-text descriptions. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

References

- [1] K. Craske, M. Dredze, & Grice, PP. Latent Dirichlet Allocation for Medical

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in directories, product descriptions and associated categories, and product hierarchies of free-text medical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we use unstructured textual data that have been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [3], product descriptions and catalogs (e.g. [1]) as available from [5]; and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., hospital discharge records with International Classification of Disease, Ninth Revision, Clinical Modification (ICD-9-CM) codes assigned [4]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

The main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data.

We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\omega_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents, and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$.

Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will not for exposition purposes assume that this label set is complete. This means that some labels may have no parent label, and some may be unlabeled. We will consider three types of constraints imposed by a label hierarchy: (1) a label l is a child of a label l' if $l' = pa(l)$; (2) a label l is a sibling of a label l' if $pa(l) = pa(l')$; and (3) a label l is a parent of a label l' if $pa(pa(l)) = l$. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{i,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{i,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an i -level hierarchy are that if the l th label is applied to document d , i.e., $y_{i,d} = 1$, then all labels in the label hierarchy up to the i th root are also applied to document d , i.e., $y_{j,d} = 1$ for all $j < i$. Conversely, if a label l' is marked as having been applied (diagonal hashing) it indicates that potentially some of the plated variables are observed.

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1.

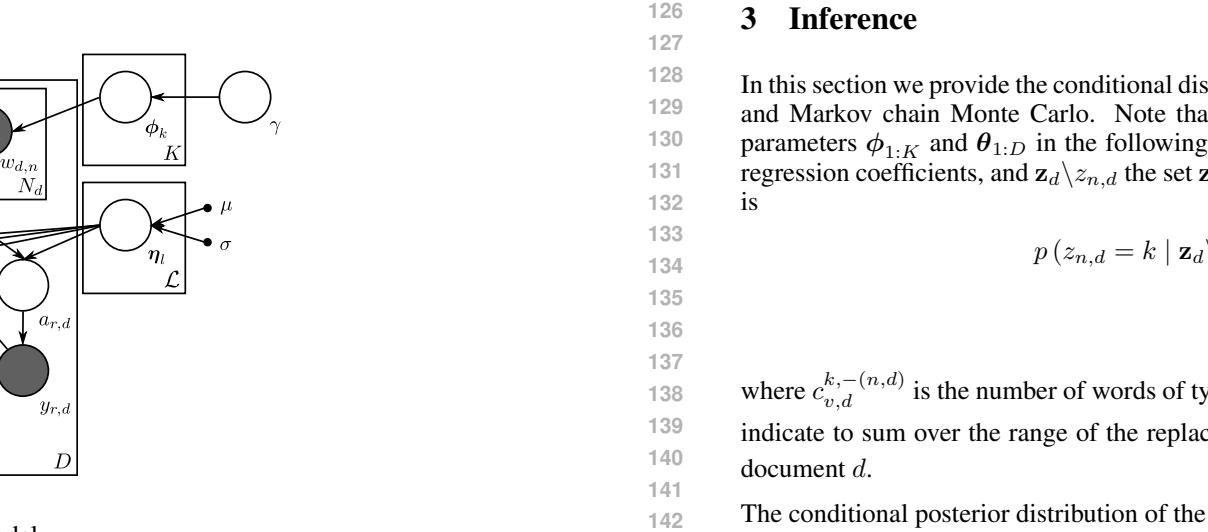


Figure 1: HSLDA graphical model

3 Inference

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [18], we have analytically marginalized out parameters $\phi_{1,K}$ and $\theta_{1,D}$ in the following expressions. Let \mathbf{z} be the set of all auxiliary variables, \mathbf{w} the set of all words, $\boldsymbol{\eta}$ the set of all regression coefficients, and $\mathbf{z}_{n,d}$ the set \mathbf{z}_d with element $z_{n,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | \mathbf{z}_d \setminus z_{n,d}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \\ \left(\frac{c_{(k),d}^{k-(n,d)} + \eta}{c_{(k),d}^{k-(n,d)} + V^{-1}} \right)^{\frac{n_d}{V}} \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(y_{l,d} \eta_{l,d} - n_{l,d})^2}{2} \right\} \quad (1)$$

where $c_{v,d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript ‘(\cdot)’ indicate to sum over the range of the replaced variable, i.e. $c_{w_n,d}^{k-(n,d)} = \sum_v c_{w_n,v,d}^{k-(n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\boldsymbol{\eta}_l | \mathbf{z}, \mathbf{a}, \sigma) = \mathcal{N}(\mu_l, \Sigma) \quad (2)$$

where

$$\hat{\mu}_l = \Sigma \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right) \quad \Sigma^{-1} = \mathbf{I}\sigma^{-1} + \mathbf{Z}^T \mathbf{Z}$$

Here \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$. The simplicity of this conditional distribution follows from the choice of probit regression [6]; the specific form of the update is a standard result from Bayesian normal linear regression [16]. It is also a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution [6].

The text of the discharge summaries was tokenized with NLTK 2 Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Descriptions and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a hierarchy and provides product descriptions for most products in their catalog. Products can be discovered by users through query-based searching or through product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [5]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

These results indicate that the full HSLDA model predicts more of the correct labels at the cost of an increase in the number of false positives should one trust the comparison models.

4 Related Work

In this work, we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document “supervision”, often taking the form of a single numerical or categorical label. More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g., online review scores), marks given to written work (e.g., essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as a conditional distribution for document d and whether or not its parent label was applied (i.e., $I(y_{pa(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{pa(l),d}$ may only be applied to document d if its parent label $y_{pa(l),d}$ is also applied (these expressions are specific to l is a constraint but can be modified to accommodate different constraints). The regression coefficients $\boldsymbol{\eta}_l$ are independent a priori, however, the hierarchical coupling in the model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually labeling features. We believe this to be a significant improvement in label prediction we observe experimentally. We test this hypothesis in Section 3.

Here $z_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k , $\bar{z}_k = N_d^{-1} \sum_{n=1}^{N_d} I(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generatively labeled using a hierarchy of conditionally dependent probit regressors [7, ?]. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e., $I(y_{pa(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{pa(l),d}$ may only be applied to document d if its parent label $y_{pa(l),d}$ is also applied (these expressions are specific to l is a constraint but can be modified to accommodate different constraints). The regression coefficients $\boldsymbol{\eta}_l$ are independent a priori, however, the hierarchical coupling in the model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually labeling features. We believe this to be a significant improvement in label prediction we observe experimentally. We test this hypothesis in Section 3.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in performance.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DisseLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nltk.org>

5.2 Comparison Models

We evaluated HSLDA along with three other closely related models against these two datasets. The comparison models included sLDA with independent regressors (this model constrains on labels ignored), LDA followed by tree-conditional regressions, and sLDA with hierarchical constraints on labels ignored.

More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g., online review scores), marks given to written work (e.g., essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as a conditional distribution for document d and whether or not its parent label was applied (i.e., $I(y_{pa(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{pa(l),d}$ may only be applied to document d if its parent label $y_{pa(l),d}$ is also applied (these expressions are specific to l is a constraint but can be modified to accommodate different constraints). The regression coefficients $\boldsymbol{\eta}_l$ are independent a priori, however, the hierarchical coupling in the model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually labeling features. We believe this to be a significant improvement in label prediction we observe experimentally. We test this hypothesis in Section 3.

More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g., online review scores), marks given to written work (e.g., essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as a conditional distribution for document d and whether or not its parent label was applied (i.e., $I(y_{pa(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{pa(l),d}$ may only be applied to document d if its parent label $y_{pa(l),d}$ is also applied (these expressions are specific to l is a constraint but can be modified to accommodate different constraints). The regression coefficients $\boldsymbol{\eta}_l$ are independent a priori, however, the hierarchical coupling in the model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually labeling features. We believe this to be a significant improvement in label prediction we observe experimentally. We test this hypothesis in Section 3.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics of performance - sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that model performance generalizes over two very disparate domains.

To evaluate the performance of these models, we establish a gold standard for comparison. In our evaluation of each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against the observed labeling. Specifically, ancestors of observed nodes were ignored, observed nodes were considered positive and unobserved nodes were considered to be negative. The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

5.4 Discussion

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms more straightforward approaches should be of interest to practitioners.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

...

5

Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

(a)

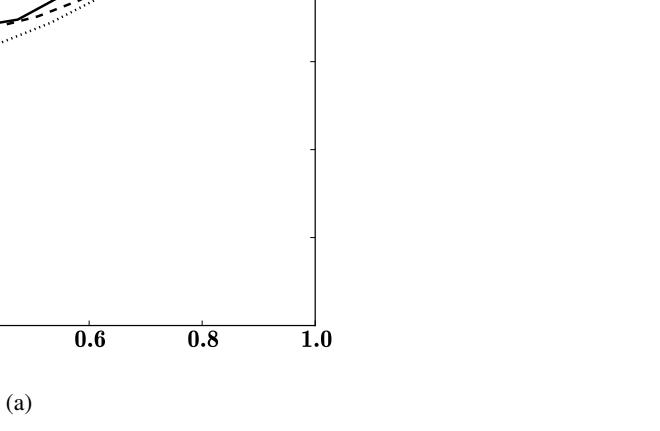


Figure 2(b): A line graph showing Sensitivity on the y-axis (0.0 to 1.0) versus 1-Specificity on the x-axis (0.0 to 1.0). Three curves are plotted: a solid line for HSLDA, a dashed line for independent regressors + sLDA, and a dotted line for HSLDA fit by running LDA first then running tree-conditional regressions. The HSLDA curve is the highest, followed by the solid line, then the dotted line.

(b)

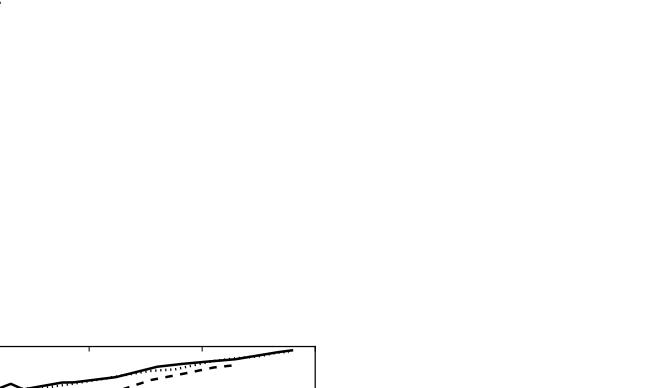


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

(a)

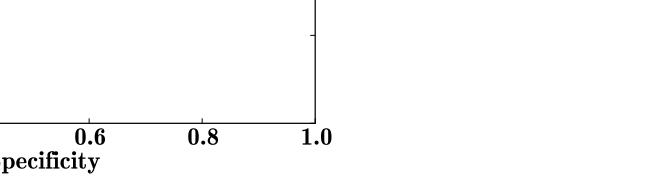


Figure 3(b): A line graph showing Sensitivity on the y-axis (0.0 to 1.0) versus 1-Specificity on the x-axis (0.0 to 1.0). Three curves are plotted: a solid line for HSLDA, a dashed line for independent regressors + sLDA, and a dotted line for HSLDA fit by running LDA first then running tree-conditional regressions. The HSLDA curve is the highest, followed by the solid line, then the dotted line.

(b)



Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in directories, product descriptions and associated categories, and product hierarchies of free-text medical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we use an unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [3], product descriptions and catalogs (e.g. [1]) as available from [5]; and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [4]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

The task of supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data.

We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$.

Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that that label set is complete, i.e., every label has a parent label. This assumption is not required for the proposed model.

Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{i,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{i,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an i -level hierarchy are that if the i th label is applied to document d , i.e. $y_{i,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{j,d} = 1 \Rightarrow y_{pa(j),d} = 1$. Conversely, if a label l' is marked as having been applied (diagonal hashing) indicates that potentially some of the plated variables are observed.

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1.

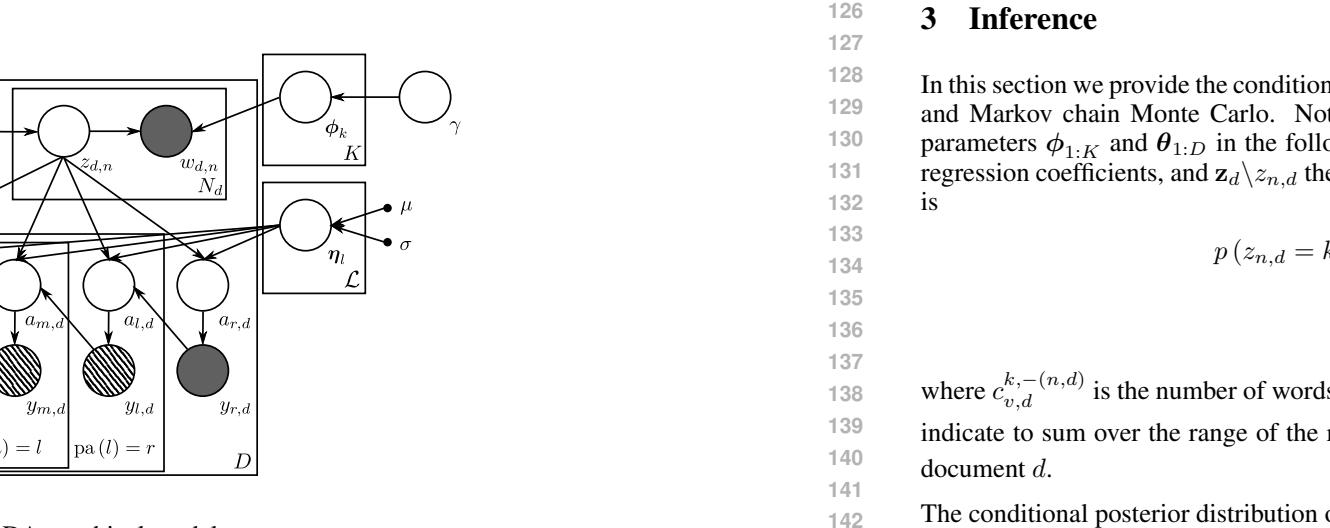


Figure 1: HSLDA graphical model

3 Inference

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [18], we have analytically marginalized out parameters $\phi_{L,K}$ and $\theta_{1,D}$ in the following expressions. Let \mathbf{z} be the set of all auxiliary variables, \mathbf{w} the set of all words, $\boldsymbol{\eta}$ the set of all regression coefficients, and $\mathbf{c}_{w_n,d}$ the set $z_{n,d}$ with element $z_{n,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | \mathbf{z} \setminus z_{n,d}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \\ \left(\frac{c_{w_n,d}^{k-(n,d)} + \gamma}{c_{w_n,d}^{k-(n,d)} + V \gamma} \right)^{\frac{1}{V}} \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(x_{l,d}^k \eta_{l,d} - a_{l,D})^2}{2} \right\} \quad (1)$$

where $x_{l,d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript ‘ (\cdot) ’ indicate to sum over the range of the replaced variable, i.e. $c_{w_n,d}^{k-(n,d)} = \sum_l c_{w_n,d}^{k-(n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\boldsymbol{\eta}_l | \mathbf{z}, \mathbf{a}, \sigma) = \mathcal{N}(\boldsymbol{\mu}_l, \Sigma) \quad (2)$$

where

$$\boldsymbol{\mu}_l = \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right) \quad \Sigma^{-1} = \mathbf{I} \sigma^{-1} + \mathbf{Z}^T \mathbf{Z}$$

Here \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is $\mathbf{z}_{d,:}$ and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$. The simplicity of this conditional distribution follows from the choice of probit regression [6]; the specific form of the update is a standard result from Bayesian normal linear regression [16]. It is also a standard probit regression result that the conditional posterior distribution of $a_{l,D}$ is a truncated normal distribution [6].

The text of the discharge summaries was tokenized with NLTK 2 Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Descriptions and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a hierarchy and provides product descriptions for most products in their catalog. Products can be discovered by users through query-based searching or through product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [5].

Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

These results indicate that the full HSLDA model predicts more of the correct labels at the cost of an increase in the number of false positives should one predict the wrong labels.

Predictive performance in HSLDA is evaluated by $p(y_{i,d} | w_{1:N_d,d}, w_{1:N_d,d}, y_{l \in \mathcal{L}:l \neq d})$ where d represents the test document. Efficiency, the expectation of this probability distribution is estimated in the following way. Expectations of \bar{z}_d and η_d were estimated with samples from the posterior. Using these expectations, we performed Gibbs sampling over the hierarchy to acquire predictive samples for the documents in the test set.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

These results indicate that the full HSLDA model predicts more of the correct labels at the cost of an increase in the number of false positives relative to the comparison models.

The following are the predictive performance results for the Amazon.com products because all ancestors in the hierarchy were included with each category label. For example, “DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine.”

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 words on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

4 Related Work

In this work, we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document “supervision”, often taking the form of a single numbered or categorical label. More generally this supervision is just extra per-document data, for instance its variable or relevance (e.g. online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as a conditional distribution for document d and whether or not its parent label was applied (i.e. $I(y_{pa(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{pa(l),d}$ may only be applied to document d if its parent label $y_{pa(l),d}$ is also applied (these expressions are specific to l is constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in the model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually labeling features. We believe this to be a significant strength of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

SLDA with independent regressors and LDA and supervision is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

These two components are important relational features in the model.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in directories, product descriptions and associated categories, and product hierarchies of free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we use an unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [3], product descriptions and catalogs (e.g. [1]) as available from [5]; and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g., hospital discharge records with International Classification of Disease, Ninth Revision, Clinical Modification (ICD-9-CM) codes assigned [4]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified e.g. image catalogs with bag-of-feature image representations).

The task of supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data.

We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathcal{L}_d = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$ be the set of labels for document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$.

Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that that label set is complete, i.e., every label in \mathcal{L} is a child of some other label in \mathcal{L} . Each label is marked with a single root $r \in \mathcal{L}$. Each

document has a variable $y_{i,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{i,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an i -level hierarchy are that if the j th label is applied to document d , i.e., $y_{j,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{i,d} = 1 \Rightarrow y_{pa(i),d} = 1$. Conversely, if a label l' is marked as having been applied (diagonal hashing) indicates that potentially some of the plated variables are observed.

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1.



Figure 1: HSLDA graphical model

3 Inference

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [18], we have analytically marginalized over the parameters $\phi_{1,K}$ and $\theta_{1,D}$ in the following expressions. Let \mathbf{z} be the set of all auxiliary variables, \mathbf{w} the set of all words, η the set of all regression coefficients, and $\mathbf{z}_{n,d}$ the set z_d with element $z_{n,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | \mathbf{z} \setminus z_{n,d}, \mathbf{a}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \\ \left(\frac{c_{(k),d}^{k-(n,d)} + \gamma}{c_{(k),d}^{k-(n,d)} + V \gamma} \right) \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(x_{l,d}^k \eta_l - a_{l,d})^2}{2} \right\} \quad (1)$$

where $c_{(k),d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript ‘(·)’ indicate to sum over the range of the replaced variable, i.e. $c_{(k),d}^{k-(n,d)} = \sum_{v \in w_{n,d}} c_{w_{n,d}}^{k-(n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\boldsymbol{\eta}_l | \mathbf{z}, \mathbf{a}, \sigma) = \mathcal{N}(\boldsymbol{\mu}_l, \Sigma_l) \quad (2)$$

where

$$\boldsymbol{\mu}_l = \left(\frac{\mu}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right) \quad \Sigma_l^{-1} = \mathbf{I} \sigma^{-1} + \mathbf{Z}^T \mathbf{Z}$$

Here \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is \mathbf{z}_d , and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$. The simplicity of this conditional distribution follows from the choice of probit regression [6]; the specific form of the update is a standard result from Bayesian normal linear regression [16]. It is also a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution [6].

The text of the discharge summaries was tokenized with NLTK 2 Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Descriptions and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a hierarchy and provides product descriptions for most products in their catalog. Products can be discovered by users through query-based searching or through product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [5].

Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

These results indicate that the full HSLDA model predicts more of the correct labels at a cost of an increase in the number of false positives should one trust the comparison models.

Predictive performance in HSLDA is evaluated by $p(y_{i,d} | w_{1:N_d}, y_{1:N_d}, \mathbf{a}, \mathbf{y}_{\mathcal{L}}, \mathbf{z}, d)$ where d represents the test document. Efficiency, the expectation of this probability distribution is evaluated in the following way. Expectations of \bar{z}_d and η_d were estimated with samples from the posterior. Using these expectations, we performed Gibbs sampling over the hierarchy to acquire predictive samples for the documents in the test set.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

These results indicate that the full HSLDA model predicts more of the correct labels at a cost of an increase in the number of false positives relative to the comparison models.

The following are the predictive performance results for the Amazon.com products because all ancestors in the hierarchy were included with each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine".

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

4 Related Work

In this work, we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document "supervision", often taking the form of a single numbered or categorical label. More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g., online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as a conditional distribution for document d and whether or not its parent label was applied (i.e., $I(y_{pa(d)}=1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{pa(d)}$ can only be applied if its parent label $y_{pa(d)}$ is also applied (these expressions are specific to l is a constraint but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in the model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually labeling features. We believe this to be a significant improvement in label prediction we observe experimentally. We test this hypothesis in Section 3.

SLDA is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is that sLDA with independent regressors and LDA and a hierarchically constrained response. The second comparison model is sLDA fit by performing LDA followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the response to incorporate the low dimensional structure inferred by LDA. Combined inference has been demonstrated to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [7]. sLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we constrain HSLDA with sLDA applied in this way.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiselLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the response to incorporate the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA [7]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

...

5 Discussion

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms more straightforward approaches should be of interest to practitioners.

Figures 2 and 3 show the predictive performance of HSLDA relative to the three comparison models.

6 Conclusion

The HSLDA model is a member of the SLDA model family, which can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms more straightforward approaches should be of interest to practitioners.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

...

(a)

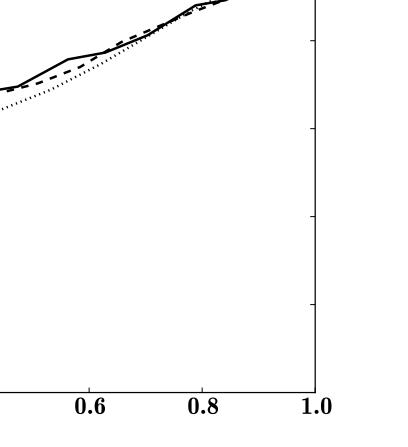


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is sLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

(a)

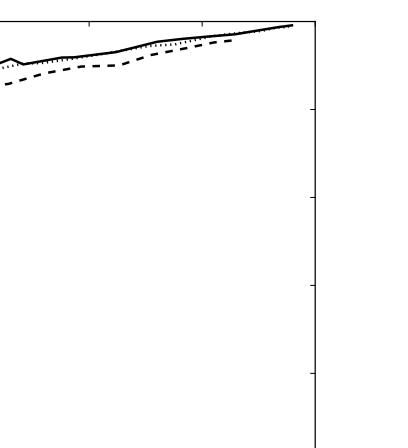


Figure 3: Out-of-sample Amazon product category predictions from product free-text descriptions. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is sLDA fit by running LDA first then running tree-conditional regressions. ?? shows predictive performance as a function of the auxiliary variable threshold and (a) shows predictive performance as a function of the prior mean on regression parameters.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multi-labeled bag-of-word data. Examples of such data include Web pages and their placement in directories, product descriptions and associated categories, and product hierarchies of free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, and improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we use unstructured textual data that have been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [3], product descriptions and catalogs (e.g. [1]) available from [5]; and patient hospital treatment transcripts and codes applied to them for bookkeeping and insurance purposes (e.g. hospital discharge records with International Classification of Disease 9th Revision, Clinical Modification (ICD-9-CM) codes assigned [4]). In this work we show how to combine these two sources of information using a single model that allows one to automatically categorize new documents, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g. image catalogs with bag-of-feature image representations).

The main contribution is to show how to utilize supervision in the form of hierarchical and (often) multiple labelings in a similar manner. Consider Web retail data. Web retailers often have a browsable product hierarchy and free-text descriptions for all products they sell. The situation of each product in a product hierarchy (often multiply situated) constitutes a multiple, hierarchical labeling of the free-text product descriptions. We hypothesize that such hierarchical labels should, at least in theory, provide better supervision than the simpler unlabeled labels previously considered. Results from applying our model to both medical record and Web retail data suggests that this is likely the case. In particular, we observe gains in our primary goal of out-of-sample label prediction that result specifically from leveraging hierarchical supervision.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and Web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data.

We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{y}_{d,l} = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d .

Each label except root labels $l \in \mathcal{L}$ has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will let $\mathbf{y}_{d,l}$ denote the set of labels that are descendants of label l . This is denoted by $\mathbf{y}_{d,l}$ if l is a simple root label, and by $\mathbf{y}_{d,l}$ if l is a multiple rooted node. Without loss of generality we will consider a document d with a single root $r \in \mathcal{L}$. Each

label l in \mathcal{L} has a variable $y_{d,l} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{d,l}$ will be 0, in some cases it will be 1, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an l -level hierarchy are that if the l th label is applied to document d , i.e. $y_{d,l} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e. $y_{d,p(l)} = 1$. Conversely, if a label l' is not applied to document d then none of its descendants are applied to document d . This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1.

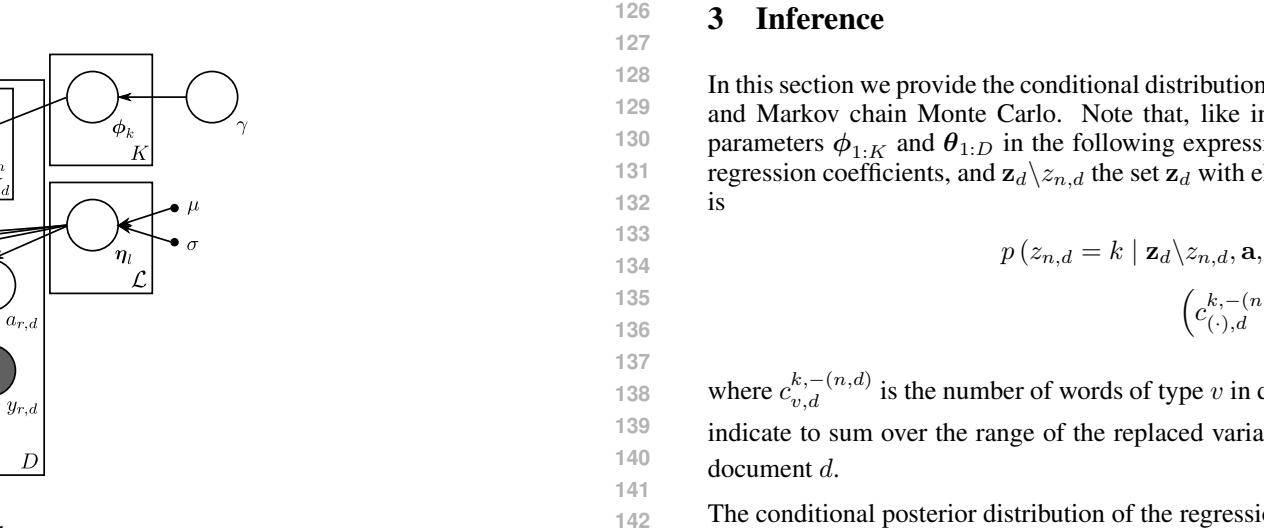


Figure 1: HSLDA graphical model

3 Inference

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [18], we have analytically marginalized out the parameters $\phi_{l,K}$ and $\theta_{l,D}$ in the following expressions. Let \mathbf{z} be the set of all auxiliary variables, \mathbf{w} the set of all words, $\boldsymbol{\eta}$ the set of all regression coefficients, and \mathbf{y}_d the set $y_{d,l}$ with element $z_{n,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | \mathbf{z}_{-n,d}, \mathbf{w}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto$$

$$\left(\frac{c_{(k),d}^{k-(n,d)} + \gamma}{c_{(k),d}^{k-(n,d)} + V_\gamma} \right) \prod_{l \in \mathcal{L}} c_{(k),l}^{y_{d,l}(n,d)} \exp \left\{ - \frac{(x_{d,l}^T \boldsymbol{\eta}_l - n_{l,d})^2}{2} \right\} \quad (1)$$

where $c_{v,d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript ‘ (\cdot) ’ is used to indicate to sum over the range of the replaced variable, i.e. $c_{w_n,d}^{k-(n,d)} = \sum_l c_{w_n,d,l}^{k-(n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\boldsymbol{\eta}_l | \mathbf{z}, \mathbf{a}, \sigma) = \mathcal{N}(\boldsymbol{\mu}_l, \Sigma_l) \quad (2)$$

where

$$\boldsymbol{\mu}_l = \left(\frac{1}{\sigma} + \mathbf{Z}^T \mathbf{a}_l \right) \Sigma^{-1} = \mathbf{a}_l^{-1} + \mathbf{Z}^T \mathbf{Z} \quad (3)$$

Here \mathbf{Z} is a $D \times K$ matrix such that row d of \mathbf{Z} is $\mathbf{z}_{d,l}$ and $\mathbf{a}_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$. The simplicity of this conditional distribution follows from the choice of probit regression [6]; the specific form of the update is a standard result from Bayesian normal linear regression [16]. It is also a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution [6].

The text of the discharge summaries was tokenized with NLTK 2 Vocabulary was determined as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

Predictive performance in HSLDA is evaluated by $p(y_{d,l} | w_{1:N_d,d}, \mathbf{y}_{d,l}, \mathbf{z}_{d,l}, \mathbf{a}_l)$ where d represents the test document. For efficiency, the expectation of this probability distribution was estimated in the following way. Expectations of \bar{z}_d and $\boldsymbol{\eta}_d$ were estimated with samples from the posterior. Using these expectations, we performed Gibbs sampling over the hierarchy to acquire predictive samples for the documents in the test set.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

These results indicate that the full HSLDA model predicts more of the correct labels at a cost of an increase in the number of false positives relative to the comparison models.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -2.2. The full HSLDA model had a true positive rate of 0.85 and a false positive rate of 0.30, the sLDA model had a true positive rate of 0.78 and a false positive rate of 0.14, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.77 and a false positive rate of 0.16. These results follow a similar pattern to the clinical data.

To further explore the tradeoffs between the true positive and false positive rates, we evaluated predictive performance for a range of values for the different parameters – the prior mean for the regression coefficients and the threshold for the auxiliary variables. The goal of this analysis was to evaluate the performance of these models subject to more or less stringent requirements for predicting positive labels.

The hyperparameters α , σ , and γ are sampled using Metropolis-Hastings.

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

We were able to deduce the structure of the hierarchy for the Amazon.com products because all ancestors in the hierarchy were included with each category label. For example, “DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi” is a single product category for the DVD, “The Time Machine.”

The hyperparameters α , σ , and γ are sampled using Metropolis-Hastings.

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [5]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

Here $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β . Their conditional posterior distribution is sampled according to

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha\beta_m)}{\Gamma(\alpha\beta_k + c_{(k),d}^k)} s(c_{(k),d}^k, m) (\alpha\beta_k)^m \quad (5)$$

where $s(n, m)$ represents stirling numbers of the first kind.

Draw $a_{l,d} | \mathbf{z}_d, \eta_l, y_{p(l),d} \sim \mathcal{N}(\bar{z}_d \eta_l, 1)$ if $y_{p(l),d} < 0$, $y_{p(l),d} = 1$ otherwise

Apply label l to document d according to $a_{l,d}$

$$y_{d,l} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$$

Here $\mathbf{z}_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k , $\bar{z}_d = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here each document is generatively labeled using a hierarchy of conditionally dependent probit regressors [16]. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $\mathbb{I}(y_{p(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{d,l}$ can only be applied to document d if its parent label $y_{p(l),d}$ is also applied (these expressions are specific to l as constraints but can be modified to accommodate different constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in the model induces a posteriori dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, contextually labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

SLDA is an extension of supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. SLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document “supervision”, often taking the form of a single numbered or categorical label. More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g., online review scores), marks given to written work (e.g. essay grades), or the number of times a web page is linked. These labels are usually generatively modeled as a conditional distribution from some distribution that depends on each document-specific topic mixture. It has been demonstrated that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [7]. It has also been shown that SLDA has been to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [7]. SLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we contrast HSLDA with SLDA applied in this way.

SLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and SLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiscLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the response to use the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in SLDA [7]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

Other models that incorporate LDA and supervision include LabeledLDA[25] and DiscLDA[20]. Various applications of these models to computer vision and document networks have been explored [30, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[23, 12, 19, 9].

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nltk.org>

...

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

5 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics - sensitivity (true positive rate) and specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that our model generalizes to two different domains.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard was derived from the observed data.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents were drawn from the challenge. The gold standard

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of file [3] and medical records and diagnosis codes assigned to them for hospital admissions. In this work we focus on ICD-9-CM codes assigned to hospital discharge summaries and Amazon.com product descriptions. Our approach leverages the underlying structure of the data to incorporate hierarchical information into the model. In the case of simplicity, we focus on ICD-9 categories, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (LDA) [8] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical information into the model. We hypothesize that the context of labels in a joint model, while leveraging the hierarchical structure of the labels. For the sake of simplicity, we focus on ICD-9 categories, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (LDA) [8] to take advantage of hierarchical supervision. The task of word-document membership is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA "topics" (distributions over the elements of Σ). ϕ_k is a distribution over "words." θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $\mathcal{N}_K(\cdot)$ is the K -dimensional Normal distribution, \mathbf{I}_K is the K -dimensional identity matrix, $\mathbf{1}_d$ is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l \mid \mu, \sigma \sim \mathcal{N}_K(\mu\mathbf{1}_K, \sigma\mathbf{I}_K)$
3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha'\mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha\beta)$
 - For $n = 1, \dots, N_d$
 - Draw a topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} \mid z_{n,d}, \phi_{z_{n,d}} \sim \text{Multinomial}(\phi_{z_{n,d}})$
 - Set $y_{d,l} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} \mid z_d, \eta_l; y_{pa(l),d} \sim \begin{cases} \mathcal{N}(\bar{\eta}_l, 1), & y_{pa(l),d} = 1 \\ \mathcal{N}(\bar{\eta}_l, 1)\mathbb{I}(y_{pa(l),d} < 0), & y_{pa(l),d} = -1 \end{cases}$
 - Apply label l to document d according to $a_{l,d}$

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has "hard is-a" parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{d,l} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{d,l}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

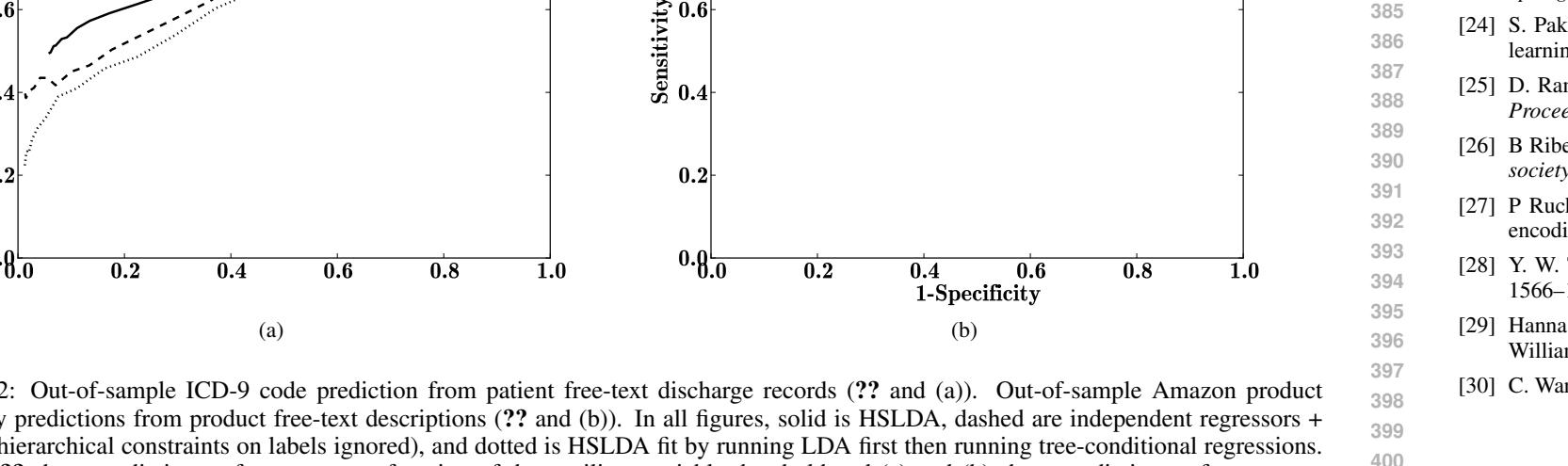


Figure 1: HSLDA graphical model

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [2] The computational medicine center's 2007 medical natural language processing challenge. <http://www.computationalmedicine.org/challenge/previous.2007>.
- [3] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [4] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [5] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [6] J.H Albert and S. Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669, 1993.
- [7] E. Birman-Deych, A. D. Waterman, Y. Yan, D. S. Nilusena, M. J. Bradford, and B. F. Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [8] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [9] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [10] J. Chakrabarti, B. Dom, R. Agrawal, and P. Raghavan. Scalable feature selection, classification and signature generation for organizing large text databases into hierarchical topic taxonomies. *The VLDB Journal*, 7:163–178, August 1998. ISSN 1066-8888.
- [11] J. Chang and D. M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [12] K. Crammer, M. Dredze, K. Ganchev, P. Tkalčík, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [13] Susan Dumais and Hao Chen. Hierarchical classification of web content. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR '00*, pages 256–263, New York, NY, USA, 2000. ACM.
- [14] R. Farkas and G. Szarvas. Automatic construction of rule-based ICD-9-CM coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [15] M. Farzandipour, A. Sheikhtabar, and F. Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:78–84, 2010.
- [16] A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2nd ed. edition, 2004.
- [17] I. Goldstein, A. Arzumanyan, and Ö. Üzuner. Three approaches to automatic assignment of icd-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [18] T. L. Griffiths and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [19] D. Koller and M. Sahami. Hierarchical classifying documents using very few words. Technical Report 1997-75, Stanford InfoLab, February 1997. Previous number = SIDL-WP-1997-059.

6

7

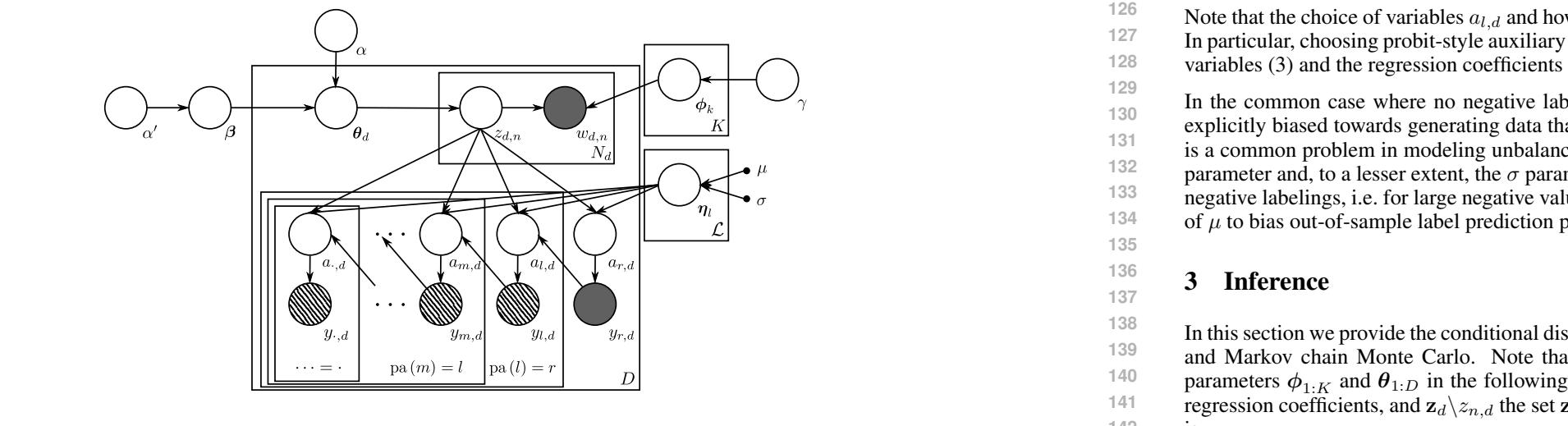
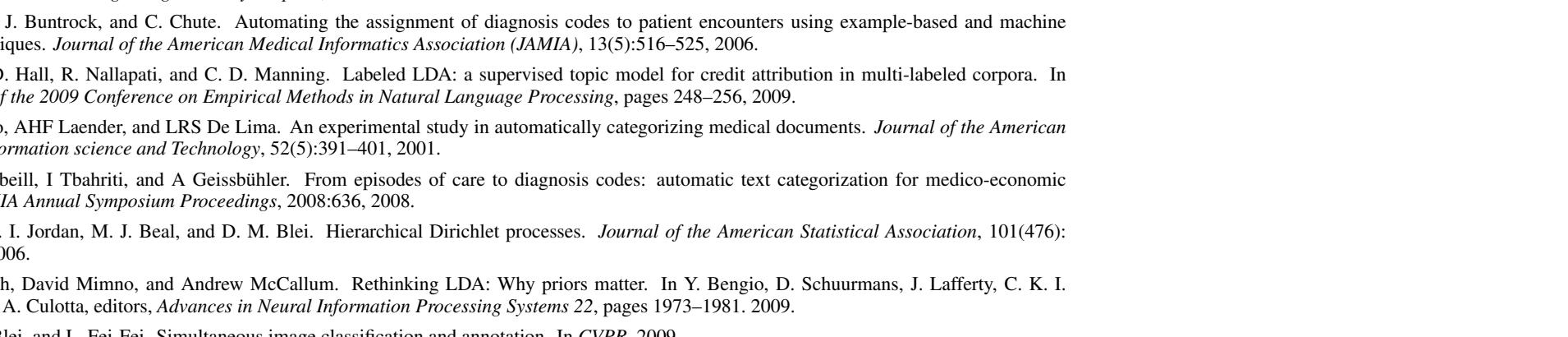


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (?? and (a)). Out-of-sample Amazon product category predictions from predicted free-text descriptions (?? and (b)). In all figures, solid is HSLDA, dashed are independent regressors + SLDA hierarchical constraints on labels ignored, and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? and ?? show predictive performance as a function of the auxiliary variable threshold and (a) and (b) show predictive performance as a function of the prior mean on regression coefficients.



(a)



(b)

Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (?? and (a)). Out-of-sample Amazon product category predictions from predicted free-text descriptions (?? and (b)). In all figures, solid is HSLDA, dashed are independent regressors + SLDA hierarchical constraints on labels ignored, and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? and ?? show predictive performance as a function of the auxiliary variable threshold and (a) and (b) show predictive performance as a function of the prior mean on regression coefficients.

- [20] S. Lacoste-Julien, F. Sha, and M. I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [21] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [22] LV Lita, S. Yu, S. Niculescu, and J. Bi. Large scale diagnostic code classification for medical patient records. In *Proceedings of the 3rd International Joint Conference on Natural Language Processing (IJCNLP'08)*, 2008.
- [23] A. McCallum, K. Nigam, J. Rennie, and K. Seymore. Building domain-specific search engines with machine learning techniques. In *Proc. AAAI-99 Spring Symposium on Intelligent Agents in Cyberspace*, 1999.
- [24] S. Pahikkala, J. Buntrock, and C. Chuote. Automating the assignment of diagnosis codes to patient encounters using example-based and machine learning techniques. *Journal of the American Medical Informatics Association (JAMIA)*, 13(5):516–525, 2006.
- [25] D. Ramage, D. Hall, R. Nallapati, and C. D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 248–256, 2009.
- [26] B. RibeiroNeto, AHF Laender, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and technology*, 52(5):391–401, 2001.
- [27] P. Ruch, J. Gobeli, I. Tshirkin, and A. Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [28] Y. W. Teh, M. I. Jordan, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [29] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981. 2009.
- [30] C. Wang, D. Blei, and L. Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

4

5

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labeling, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{d,l} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

3 Inference

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note, like in collapsed Gibbs samplers for LDA [18], we have analytically marginalized out the parameters $\phi_{k,l}$ and $\theta_{l,D}$ in the following expressions. Let \mathbf{a} be the set of all auxiliary variables, \mathbf{w} the set of all words, \mathbf{y} the set of all regression coefficients, and $\mathbf{z}_{d,\mathcal{L}}$ the set z_d with element $z_{d,l}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{d,l} = k \mid \mathbf{z}_{d,\mathcal{L}}, \mathbf{a}, \mathbf{w}, \mathbf{y}, \boldsymbol{\eta}, \alpha, \beta) \propto \left(\frac{c_{w,n,(n,d)}^{k-(n,d)} + \alpha \beta_k}{c_{w,n,(n,d)}^{k-(n,d)} + V \mathbf{1}^T} \right) \prod_{l \in \mathcal{L}, l \neq k} \exp \left\{ - \frac{(\mathbf{z}_{d,l}^T \boldsymbol{\eta} - a_{l,d})^2}{2} \right\} \quad (1)$$

where $c_{w,n,(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript (\cdot) indicates to sum over the range of the replaced variable, i.e. $c_{w,n,(n,d)} = \sum_{l \in \mathcal{L}} c_{w,n,(n,d,l)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [21]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

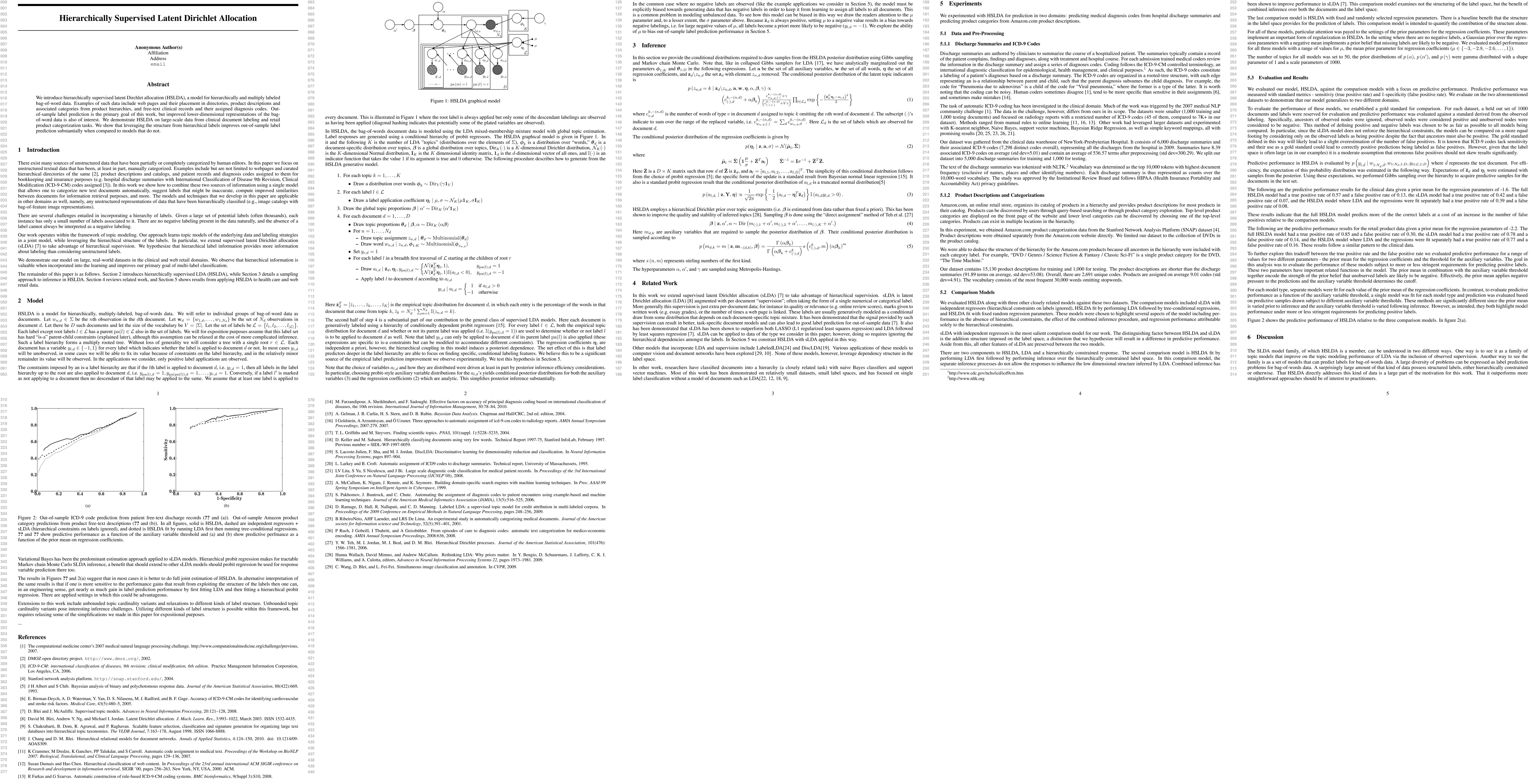
To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special care was taken to ensure that the labels were representative of the true clinical domain.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from



Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of file (see [2]), product descriptions and catalogs, and patient records and discharge codes assigned to them for hospital admissions. In this work we focus on the hospital discharge summaries from the 2007 medical natural language processing challenge (ICD-9-CM) codes assigned ([3]). In this work we show how to combine these two sources of information using a single model that allows one to categorize new text documents automatically; suggest labels that might be inaccurate; compute improved similarities between documents for information retrieval purposes; and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, and the absence of a label cannot always be interpreted as a negative labeling.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. For the sake of simplicity, we focus on ICD-9 codes, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (LDA) [7] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical information into the model. We hypothesize that the context of labels in a joint model outperforms a classification with unstructured labels as well as a disjoint model, where the topic modeling and the hierarchical classification are carried out independently of each other.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label except root labels $l \in \mathcal{L}$ has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a single root $r \in \mathcal{L}$. Each document has a variable $y_{d,l} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document or not. In most cases $y_{d,l}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

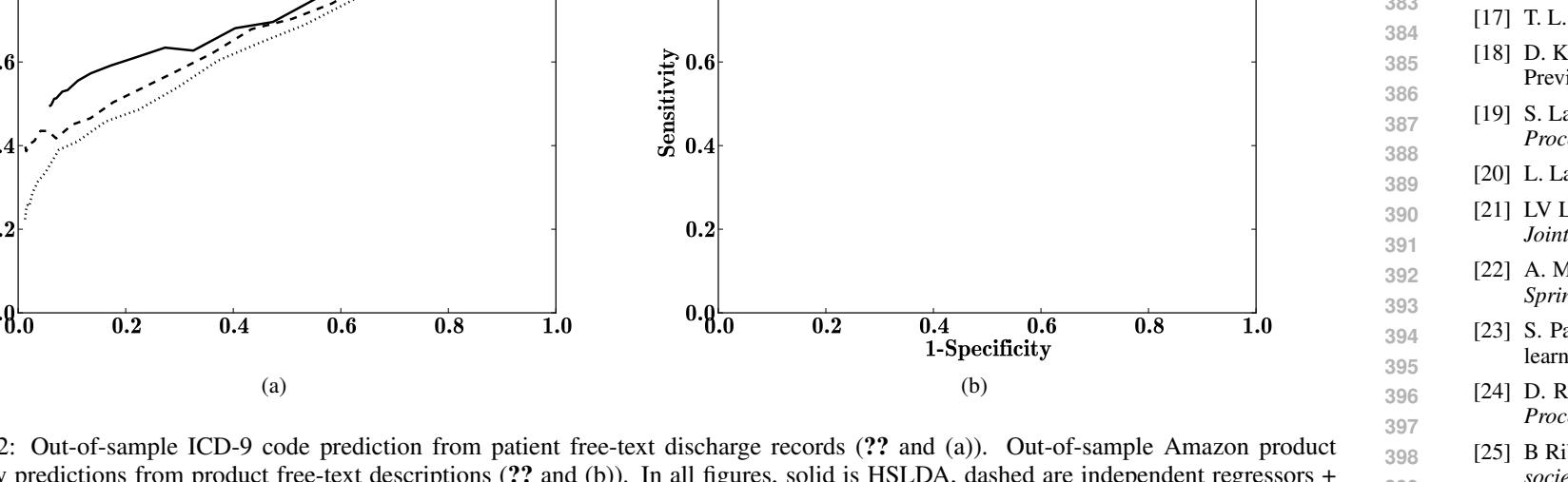


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (?? and (a)). Out-of-sample Amazon product category predictions from product free-text descriptions (?? and (b)). In all figures, solid is HSLDA, dashed are independent regressors + SLDA hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions. ?? and ?? show predictive performance as a function of the auxiliary variable threshold and (a) and (b) show predictive performance as a function of the prior mean on regression coefficients.

Variational Bayes has been the predominant estimation approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo sLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

The results in Figures ?? and ?? suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that one is more sensitive to the performance gains that result from exploiting the structure of the labels than one, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. There are get settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [1] The computational medicine center’s 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] J H Albert and S Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669–683, 1993.
- [6] E. Birman-Deych, A. D. Waterman, Y. Yan, D. S. Nilasena, M. J. Bradford, and B. F. Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [7] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [8] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [9] S. Chakrabarti, B. Dom, R. Agrawal, and P. Raghavan. Scalable feature selection, classification and signature generation for organizing large text databases into hierarchical topic taxonomies. *The VLDB Journal*, 7:163–178, August 1998. ISSN 1066-8888.
- [10] J. Chang and D. M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [11] K. Cranner, M. Dredze, K. Ganchev, PP Talukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [12] Susan Dumais and Hao Chen. Hierarchical classification of web content. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR ’00*, pages 256–263, New York, NY, USA, 2000. ACM.
- [13] R Farkas and G Szarus. Automatic construction of rule-based ICD-9-CM coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.

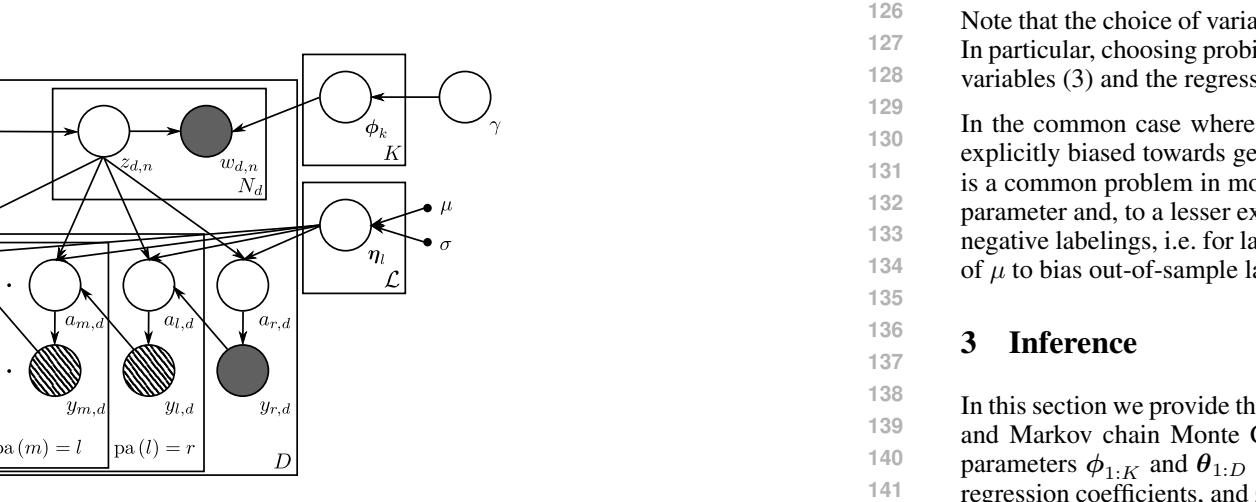


Figure 1: HSLDA graphical model

Note that the choice of variables $a_{n,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{n,d}$ ’s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because $z_{d,l}$ is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{d,l} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. This is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

3 Inference

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [17], we have analytically marginalized out the parameters $\theta_{1,K}$ and $\theta_{1,D}$ in the following expressions. Let a be the set of all auxiliary variables, w the set of all words, η the set of all regression coefficients, and z_d the $z_{d,l}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | z_d \setminus z_{n,d}, a, w, \eta, \alpha, \beta, \gamma) \propto \left(\frac{c_{w,n,d}^{k-(n,d)} + \alpha \beta_k}{c_{w,n,d}^{k-(n,d)} + V \gamma} \right)^{\frac{1}{V}} \prod_{l \in \mathcal{L}_d} \exp \left\{ - \frac{(\mathbf{z}_{d,l}^T \eta - a_{l,d})^2}{2} \right\} \quad (1)$$

where $c_{w,n,d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript ‘(·)’ indicate to sum over the range of the replaced variable, i.e. $c_{w,n,d}^{k-(n,d)} = \sum_l c_{w,n,d,l}^{k-(n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

Our dataset was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes on average (std dev=5.01) and contain a total of 1536.57 terms after preprocessing (std dev=309.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The task of automatic ICD-9 coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 test documents) and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special attention was given to the evaluation of negative labels, which are often the most difficult to predict. Our work leveraged larger datasets and experimented with labeling. Special attention was given to the evaluation of negative labels, which are often the most difficult to predict. The models can be compared on a more equal basis. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as being positive despite the fact that the ancestor must also be positive. The gold standard has been applied (diagonal hashing) indicates that potentially some of the predicted variables are observed.

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA “topics” (distributions over the elements of Σ), ϕ_k is a distribution over “words,” θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $N_K(\cdot)$ is the K -dimensional Normal distribution, I_K is the K -dimensional identity matrix, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 • Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma_V)$

2. For each label $l \in \mathcal{L}$
 • Draw a label application coefficient $\eta_l \mid \mu, \sigma \sim N_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$

3. Draw the global topic proportions $\beta \mid \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$

4. For each document $d = 1, \dots, D$
 • Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$

• Draw topic proportions $\theta_d \mid \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$

• For each label $l = 1, \dots, N_d$
 • Draw topic assignment $z_{n,d} \mid \theta_d \sim \text{Multinomial}(\theta_d)$

• Draw word $w_{n,d} \mid z_{n,d}, \phi_l \sim \text{Multinomial}(\phi_{z_{n,d}, l})$

• Set $y_{d,l} = 1$
 • For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r

– Draw $a_{l,d} \mid z_d, \eta_l; y_{pa(l),d} \sim \begin{cases} N(\mathbf{z}_d^T \eta_l, 1), & y_{pa(l),d} = 1 \\ N(\mathbf{z}_d^T \eta_l, 1), & y_{pa(l),d} < 0, \end{cases}$

– Apply label l to document d according to $a_{l,d}$

$y_{d,l} \mid a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here \mathbf{z}_d^T is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k , $\mathbf{z}_d = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a special part of our contribution to the general class of supervised LDA models. Here each document is generally labeled using a hierarchy of conditionally dependent probit regressors [15]. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $\mathbb{I}(y_{pa(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that $y_{d,l}$ can only be applied to document d if its parent label $y_{pa(l),d}$ is also applied (these express “is-a” constraints that can be moved to accommodate different constraints). The regression coefficients η_l are independent of $\eta_{pa(l),d}$, however, the hierarchical coupling in this model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

In this work we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document “supervision”; often taking the form of a single numerical or categorical label. More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g. online review score), marks given to writing, or popularity of the post. In this work we use the term “supervision” to denote general external knowledge that is used to guide the inference process. We were able to demonstrate that the signal provided by such supervision can result in better, task-specific document models and can also lead to good label prediction for out-of-sample data [7]. It has also been demonstrated that sLDA has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [7]. sLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we constrain HSLDA with LDA applied in this way.

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[22, 12, 18, 9].

SLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and SLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance. Aside from this, all other features of SLDA are preserved between the two models.

There are two components to HSLDA. LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the prior mean on regression coefficients and separate inference processes do not allow the responses to influence the low dimensional structure inferred by LDA. Combined inference has

been demonstrated that HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms more straightforward approaches should be of interest to practitioners.

5 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics – sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the aforementioned datasets to demonstrate that our model generalizes to two different domains.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the labels. Special attention was given to the evaluation of negative labels, which are often the most difficult to predict. Our work leveraged larger datasets and experimented with labeling. Special attention was given to the evaluation of negative labels, which are often the most difficult to predict. The models can be compared on a more equal basis. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as being positive despite the fact that the ancestor must also be positive. The gold standard has been applied (diagonal hashing) indicates that potentially some of the predicted variables are observed.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. This is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placements in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of file [2], product descriptions and catalogs, patient records and diagnosis codes assigned to them for hospital admissions, hospital discharge summaries, and clinical documents. The National Clinical Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine these two sources of information using a single model that allows one to categorize new text documents automatically; suggest labels that might be inaccurate; compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable in other domains as well, namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, and the absence of a label cannot always be interpreted as a negative labeling. The text of the discharge summaries was tokenized with NLTK [2]. Vocabulary was estimated as the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). Each discharge summary is thus represented as counts over the 10,000-word vocabulary. The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

In HSLDA, the bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA "topics" (distributions over the elements of Σ). ϕ_k is a distribution over "words." θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $N_K(\cdot)$ is the K -dimensional Normal distribution, I_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA graphical model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma_1 \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mathbf{w}, \sigma \sim N_K(\mu_l \mathbf{1}_K, \sigma \mathbf{I}_K)$
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_{k,d} \sim \text{Multinomial}(\phi_{k,d})$
 - Set $y_{d,l} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | \mathbf{z}_d, \eta_l; y_{\text{pa}(l),d} \sim \begin{cases} N(\mathbf{z}_d^\top \eta_l, 1), & y_{\text{pa}(l),d} = 1 \\ N(\mathbf{z}_d^\top \eta_l, 1)(a_{l,d} < 0), & y_{\text{pa}(l),d} = -1 \end{cases}$
 - Apply label l to document d according to $a_{l,d}$

For each model type, separate models were fit for each value of the prior mean of the regression coefficients. In contrast, to evaluate predictive performance as a function of the auxiliary variable threshold, a single model was fit for each model type and prediction was evaluated based on predictive samples drawn subject to different auxiliary variable thresholds. These methods are significantly different since the prior mean is varied prior to inference and the auxiliary variable threshold is varied following inference. However, as intended, they both highlight model performance under more or less stringent requirements for predicting positive labels.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-words data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$.

Each label except root labels $l \in \mathcal{L}$ has a parent $\text{pa}(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has "hard is-a" parent-child constraints (explained later), although this assumption can be relaxed at the cost of more complicated inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{d,l} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{d,l}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

Figure 2 shows the results of running HSLDA on medical discharge summaries and Amazon product category predictions. In both figures, solid is HSLDA, dashed are independent sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.

The results in Figures ?? and ?? suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that one is more sensitive to the performance gains that result from exploiting the structure of the labels than one, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. There are applied settings in which this could be advantageous.

Extensions to this work include unbouded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

...

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision, clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] J H Albert and S Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669–693.
- [6] E Birman-Deych, A D Waterman, Y Yan, D Nilasena, M J Radford, and B F Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [7] D Blei and J McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [8] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [9] S Chakrabarti, B Dom, R Agrawal, and P Raghavan. Scalable feature selection, classification and signature generation for organizing large text databases via hierarchical topic taxonomies. *The VLDB Journal*, 7:163–178, August 1998. ISSN 1066-8888.
- [10] J Chang and D. M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [11] K Crammer, M Dredze, K Ganchev, PP Talukdar, and S Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biologic, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [12] Susan Dumais and Hao Chen. Hierarchical classification of web content. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR '00*, pages 256–263, New York, NY, USA, 2000. ACM.
- [13] R Farkas and G Sarvaras. Automatic construction of rule-based ICD-9-CM coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [14] M Farzanpour, A Sheikholeslami, and F Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:78–84, 2010.
- [15] A Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2nd ed. edition, 2004.

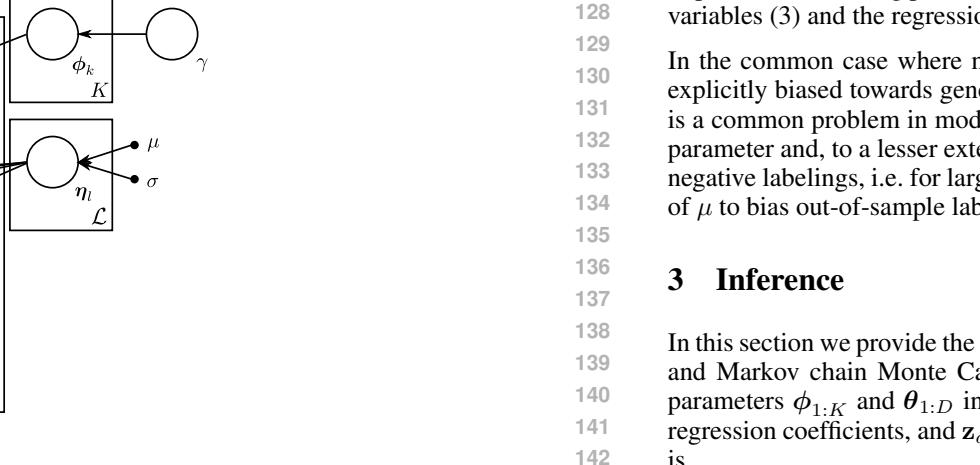


Figure 1: HSLDA graphical model

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because $z_{d,l}$ is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{d,l} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics - sensitivity (true positive rate) and specificity (false positive rate). We evaluated on the two aforementioned datasets to demonstrate that our model generalizes to two different domains.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the observed hierarchy up to the root are also applied to document d , i.e. $y_{d,\text{pa}(l),d} = 1, \dots, y_{d,l,d} = 1$. Conversely, if a label l' is marked as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied. Diagonal hashing indicates that potentially some of the plotted variables are observed.

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with Naïve Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [20, 25, 23, 26, 21].

Our dataset was gathered from the clinical data warehouse of New York-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=309.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

Predictive performance in HSLDA is evaluated by $p(y_{d,l} | \mathbf{z}_d, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto$ $(\mathbf{z}_d^{\mathbf{k}_{(l)} - (n,d)} + \alpha \mathbf{z}_d) \frac{c_{w,n,d}^{k_{(l)} - (n,d) + \gamma}}{(c_{w,n,d}^{k_{(l)} - (n,d) + \gamma} + V)^{\gamma}} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(\mathbf{z}_d^\top \eta_l - a_{l,d})^2}{2} \right\}$ (1)

where $c_{w,n,d}^{k_{(l)} - (n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript (\cdot) indicate to sum over the range of the replaced variable, i.e. $c_{w,n,d}^{k_{(l)} - (n,d)} = \sum_{k \in \mathcal{L}_d} c_{w,n,d}^{k_{(l)} - (n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the latent topic indicators is

$$p(z_{d,l} = k | \mathbf{z}_d \setminus z_{d,l}, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto p(\eta_l | \mathbf{z}_d \setminus z_{d,l}, \mathbf{a}, \sigma) = \mathcal{N}(\mu_l, \Sigma) \quad (2)$$

where

$$\mu_l = \Sigma \left(\frac{1}{\mu_l + Z^T \mathbf{a}_l} \right) \Sigma^{-1} \mathbf{1} \sigma^{-1} + Z^T \mathbf{z}_d.$$

Here Z is a $D \times K$ matrix such that row d of Z is \mathbf{z}_d , and $\mathbf{a}_l = (a_{l,1}, a_{l,2}, \dots, a_{l,D})^T$. The simplicity of this conditional distribution follows from the choice of probit regression [5]; the specific form of the update is a standard result from Bayesian normal linear regression [15]. It also is a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution [5].

$$p(a_{l,d} | \mathbf{z}, \mathbf{Y}, \eta) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (a_{l,d} - \eta_l^\top \mathbf{z}_d) \right\} \mathbb{I}(a_{l,d} > 0). \quad (3)$$

HSLDA employs a hierarchical Dirichlet prior over topic assignments (i.e. β is estimated from data rather than fixed a priori). This has been shown to improve the quality and stability of inferred topics [28]. Sampling β is done using the "direct assignment" method of Teh et al. [27].

$$\beta | \mathbf{z}, \alpha', \alpha \sim \text{Dir}((m_{(1,1)} + \alpha', m_{(1,2)} + \alpha', \dots, m_{(1,K)} + \alpha')^T) \quad (4)$$

Here $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β . Their conditional posterior distribution is sampled according to

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + c_{(k),d}^m)} s((c_{(k),d}^m) / (\alpha \beta_k)^m) \quad (5)$$

where $s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

For each model type, separate models were fit for each value of the prior mean of the regression coefficients. In contrast, to evaluate predictive performance as a function of the auxiliary variable threshold, a single model was fit for each model type and prediction was evaluated based on predictive samples drawn subject to different auxiliary variable thresholds. These methods are significantly different since the prior mean is varied prior to inference and the auxiliary variable threshold is varied following inference. However, as intended, they both highlight model performance under more or less stringent requirements for predicting positive labels.

Figure ?? shows the predictive performance of HSLDA relative to the three comparison models as a function of the prior mean on regression coefficient. Figure ?? demo.

6 Discussion

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as prediction problems for bag-of-words data. A surprisingly large amount of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms more straightforward approaches should be of interest to practitioners.

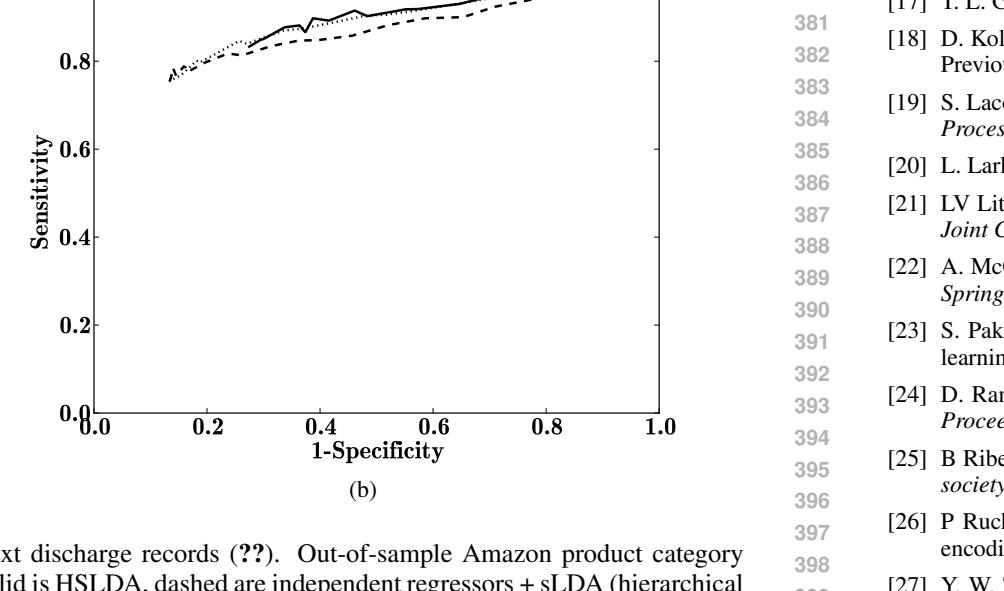


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records [??]. Out-of-sample Amazon product category predictions from product free-text descriptions [??]. In both figures, solid is HSLDA, dashed are independent sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.

V

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placements in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data is also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured textual data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs, and patient records and diagnosis codes assigned to them for bookings and charges. The challenge of learning from this data is that it is often unlabeled. For example, the ICD-9-CM codes assigned to a hospital discharge record are not labeled with the corresponding clinical diagnosis. Coding follows the ICD-9-CM controlled terminology, and the ICD-9-CM codes are organized in a tree-structured hierarchy. A diagnosis code such as ICD-9-CM code 78.4 corresponds to a labeling of a patient's diagnosis, based on a discharge summary. The ICD-9-CM codes are organized in a tree-structured hierarchy, with each node representing an *is-a* relationship between the parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for "Pneumonia due to adenovirus" is a child of the code for "Viral pneumonia," where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

2

The task of automatic ICD-9-Coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them), compared to 7K+ in our dataset. Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with different models [10]. In particular, since the sLDA model does not respect the hierarchical constraints, the models can be compared on a more equal footing by considering only the observed labels as being positive despite the fact that ancestors must also be positive. The gold standard promising results [20, 23, 25, 26, 21].

3

Our dataset was gathered from the clinical data warehouse of NewYork-Presbyterian Hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 average words per document and 1,000 tokens per sentence. The datasets were split into training (80%) and testing (20%) sets. The bag-of-words document data is modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In it and the following K is the number of LDA "topics" (distributions over the elements of Σ). ϕ_k is a distribution over "words." θ_k is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $N_K(\cdot)$ is the K -dimensional Normal distribution, I_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}_V(\gamma_1 \mathbf{1}_V)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim N_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $n = 1, \dots, N_d$
 - Draw a topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,d} | z_{n,d}, \phi_{z_n} \sim \text{Multinomial}(\phi_{z_n})$
 - Set $y_{d,l} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_d, \eta_l; y_{\text{pa}(l),d} \sim \begin{cases} N(\bar{z}_d^T \eta_l, 1), & y_{\text{pa}(l),d} = 1 \\ N(\bar{z}_d^T \eta_l, 1)(a_{l,d} < 0), & y_{\text{pa}(l),d} = -1 \end{cases}$
 - Apply label l to document d according to $a_{l,d}$

Here $\bar{z}_d^T = [\bar{z}_1, \dots, \bar{z}_k, \dots, \bar{z}_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k . $\bar{z}_k = \frac{1}{N_d} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a subtask of our part of contribution to the general class of supervised LDA models. Here each document is generally labeled using a hierarchy of conditionally dependent probit regressors [15]. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $\mathbb{I}(y_{\text{pa}(l),d} = 1)$) are used to determine whether or not label l is to be applied to document d as well. Note that label $y_{d,l}$ can only be applied to document d if its parent label $y_{\text{pa}(l)}$ is also applied (these same label hierarchy forms a multiply rooted tree). Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{d,l} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{d,l}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

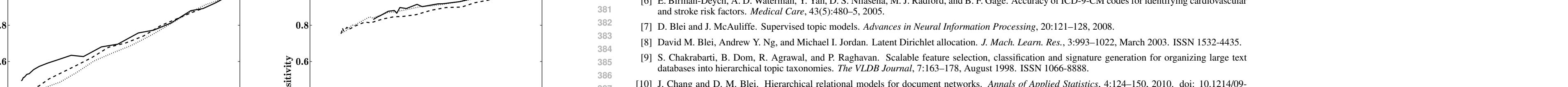


Figure 1: HSLDA graphical model

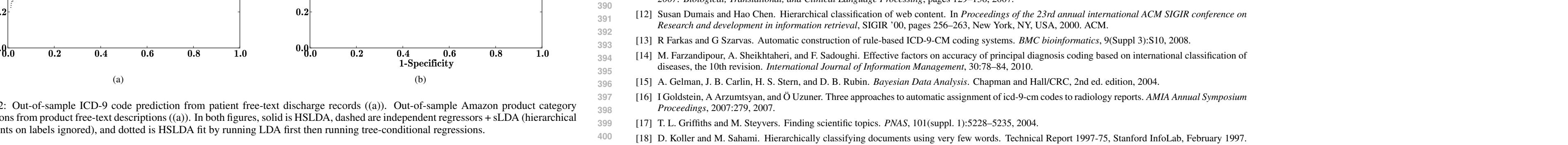


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records ((a)). Out-of-sample Amazon product category predictions from product free-text descriptions ((a)). In both figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.

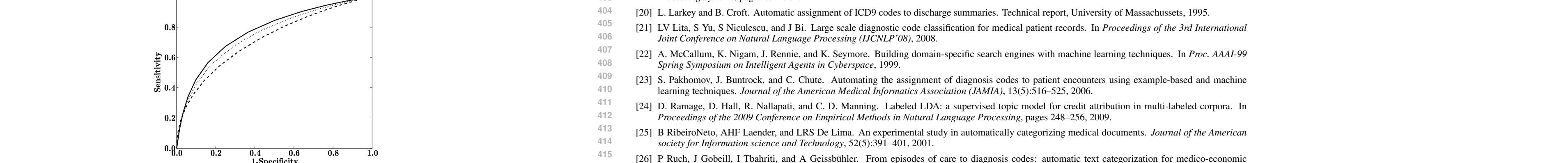


Figure 3: ROC Curve for clinical data

Variational Bayes has been the predominant estimation approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo sLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

The results in Figures ?? and 2(a) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that one is more sensitive to the performance gains that result from exploiting the structure of the labels then one, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. There are applied settings in this which could be advantageous.

Extensions to this work include unbounded topic cardinality constraints and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics - sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluated on the two aforementioned datasets to demonstrate that our model generalizes to two different domains.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the observed labeling. Specifically, ancestors of observed nodes were considered positive and unobserved nodes were considered negative.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the readers attention to the μ parameter and, to a lesser extent, the σ parameter above. Because $z_{d,l}$ is always positive, setting μ to a negative value results in a bias towards negative labeling, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{d,l} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

3 Inference

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note, like in collapsed Gibbs samplers for LDA [17], we have analytically marginalized out the parameters $\phi_{k,l}$ and $\theta_{l,D}$ in the following expressions. Let \mathcal{A} be the set of all auxiliary variables, \mathcal{W} the set of all words, \mathcal{Y} the set of all regression coefficients, and $\mathcal{Z}_{d,l}$ the set $z_{d,l}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{d,l} = k | \mathcal{Z}_{d,l}, \mathcal{A}, \mathcal{W}, \mathbf{y}, \boldsymbol{\eta}, \alpha, \beta, \gamma) \propto \left(\frac{c_{w_{n,d},l}^{k-(n,d)} + \eta_l}{c_{w_{n,d},l}^{k-(n,d)} + V} \right)^{\gamma} \prod_{i \in \mathcal{L}_{d,l}} \exp \left\{ - \frac{(\bar{z}_{d,l}^T \boldsymbol{\eta}_i - a_{i,d})^2}{2} \right\} \quad (1)$$

where $c_{w_{n,d},l}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript (\cdot) indicate to sum over the range of the replaced variable, i.e. $c_{w_{n,d},l}^{k-(n,d)} = \sum_{v \in \mathcal{V}} c_{w_{n,d},l}^{k-(n,d)}$. Here $\mathcal{L}_{d,l}$ is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\boldsymbol{\eta}_i | \mathcal{Z}_{d,l}, \mathcal{A}, \sigma) = \mathcal{N}(\boldsymbol{\mu}_i, \Sigma_i) \quad (2)$$

where

$$\boldsymbol{\mu}_i = \left(\frac{1}{\mu} + \mathbf{Z}' \mathbf{Z}_i \right)^{-1} \mathbf{Z}' \mathbf{Z}_i \mathbf{y}_i \quad \Sigma_i = \mathbf{I} \sigma^{-2} + \mathbf{Z}' \mathbf{Z}$$

Here \mathbf{Z}_i is a $D \times K$ matrix such that row d of \mathbf{Z}_i is $z_{d,l}$ and $\mathbf{a}_i = (a_{1,d}, a_{2,d}, \dots, a_{|L|,d})^T$. The simplicity of this conditional distribution follows from the choice of probit regression [5]; the specific form of the update is a standard result from Bayesian normal linear regression [15]. It is also a standard probit regression result that the conditional posterior distribution of $a_{i,d}$ is a truncated normal distribution [5].

$$p(a_{i,d} | \mathcal{Z}_{d,l}, \mathbf{Y}, \boldsymbol{\eta}) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ - \frac{1}{2} (a_{i,d} - \eta_i^T \mathbf{z}_{d,l}) \right\} \mathbb{I}(a_{i,d} y_{d,l} > 0) \quad (3)$$

HSLDA employs a hierarchical Dirichlet prior over topic assignments (i.e. β is estimated from data rather than fixed a priori). This has been shown to improve the quality and stability of inferred topics [28]. Sampling β is done using the "direct assignment" method of Teh et al. [27].

$$\beta | \mathbf{z}, \alpha' \sim \text{Dir}((m_{(1,1)} + \alpha', m_{(1,2)} + \alpha', \dots, m_{(1,K)} + \alpha')) \quad (4)$$

Here $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β . Their conditional posterior distribution is sampled according to

$$p(m_{d,k} = m | \mathcal{Z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha \beta)}{\Gamma(\alpha \beta + c_{(1),d}^k)} s(c_{(1),d}^k, m) (\alpha \beta)_m^m \quad (5)$$

where $s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α' , α , and γ are sampled using Metropolis-Hastings.

4 Related Work

In this work we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document "supervision"; often taking the form of a single numerical or categorical label. More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g. online review score), marks given to writing, or popularity of the document. In general, this supervision is often generated by a human editor and is often added to the training data. These models were chosen to highlight several aspects of the model including performance relative to the comparison models.

For each model type, separate models were fit for each category label. For example, "DVD / Genres / Science Fiction & Fantasy / Classic Sci-Fi" is a single product category for the DVD, "The Time Machine."

To further explore this tradeoff between the true positive rate we evaluated predictive performance for a range of values of α for two different parameter sets - the prior mean for the regression coefficients and the threshold for the auxiliary variables. The goal in this analysis was to evaluate the performance of the model as the strength of the prior increased.

These two parameters have implicit related functions in the model. It has been demonstrated that the signal provided by such supervision can result in better task-specific document models and can also lead to good label prediction for out-of-sample data [7]. It has also been demonstrated that sLDA has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [7]. sLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we constrain HSLDA with sLDA applied in this way.

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[22, 12, 18, 9].

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nikit.org>

5 Discussion

6.1 Discussion

6.1.1 Model Comparison

6.1.1.1 HSLDA vs. sLDA

6.1.1.1.1 Performance

6.1.1.1.1.1 Sensitivity

6.1.1.1.1.1.1 Sensitivity

6.1.1.1.1.1.1.1 Sensitivity

6.1.1.1.1.1.1.1.1 Sensitivity

6.1.1.1.1.1.1.1.1.1 Sensitivity

6.1.1.1.1.1.1.1.1.1.1 Sens

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placements in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories, lists of products and categories, patient records and discharge summaries from hospital discharge summaries.

In this work we show how to combine these two sources of information using a single model that allows one to categorize new text documents automatically; suggest labels that might be inaccurate; compute improved similarities in other data as well; namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, and the absence of a label cannot always be interpreted as a negative labeling.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. For the sake of simplicity, we focus on *a*-is-hierarchies, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical information into the model. We hypothesize that the context of labels within the hierarchy provides valuable information about labeling.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Our results show that a joint, hierarchical model outperforms a classification with unstructured labels as well as a disjoint model, where the topic modeling and the hierarchical classification are carried out independently of each other.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d .

Let \mathcal{L} be the set of documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except root, has a parent $pa(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has an hard “ $is-a$ ” parent-child constraint (explained later), although this assumption can be relaxed at the cost of more computationally complex inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{d,l} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{d,l}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remaining its value will be observed. In the applications we consider, only positive label applications are observed.

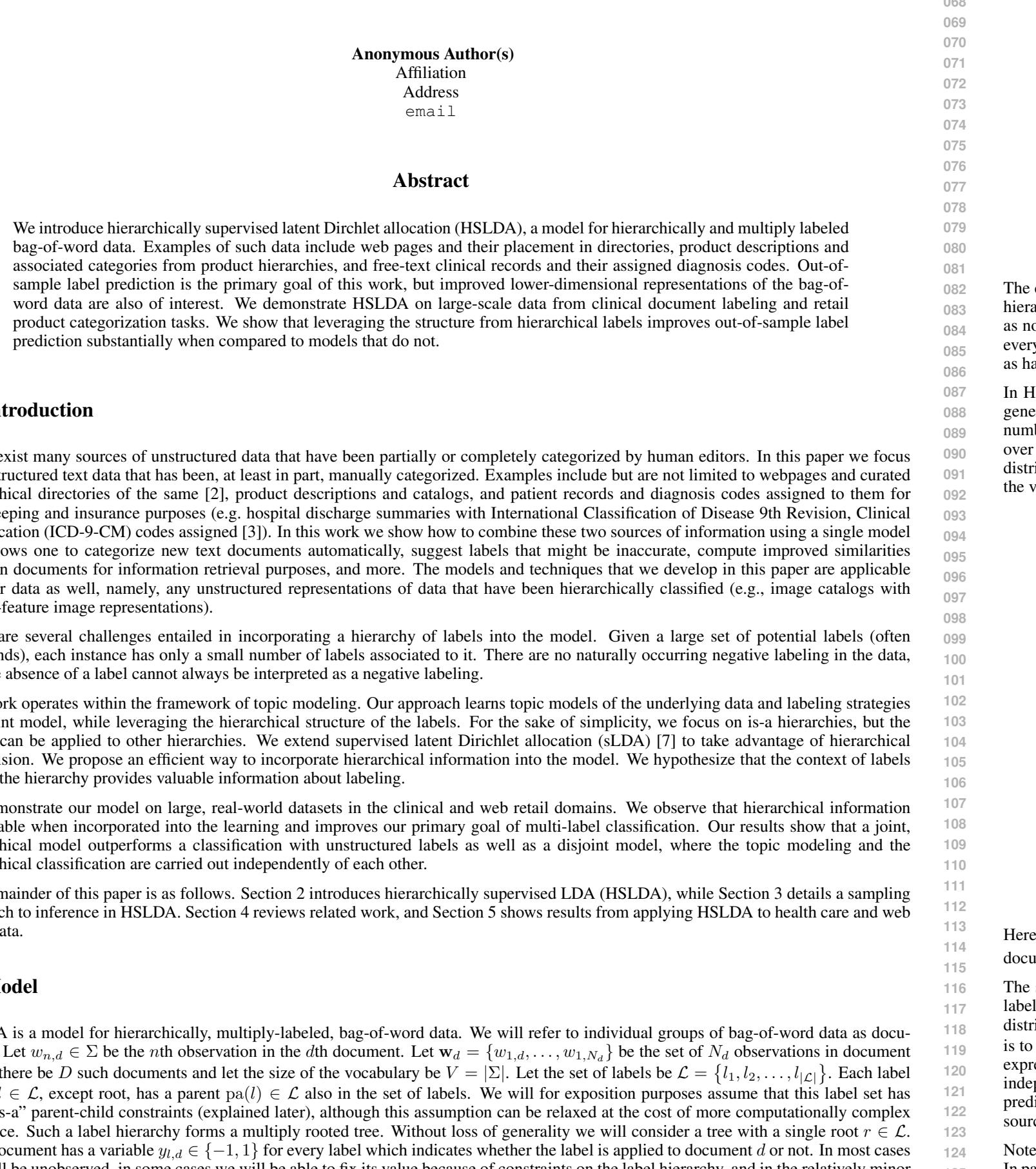


Figure 1: HSLDA graphical model

Here $z_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k . $z_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(a_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here, each document is labeled generatively using a conditionally dependent probit regression (15). For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $(y_{pa(l),d} = 1)$) are used to determine whether or not label l is applied to document d and whether or not its parent label was applied to document d if its parent label $pa(l)$ is not applied (these are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence). The regression coefficients η_l are specific to few constants but can be modified to accommodate different constants. The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{i,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{i,d}$ yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

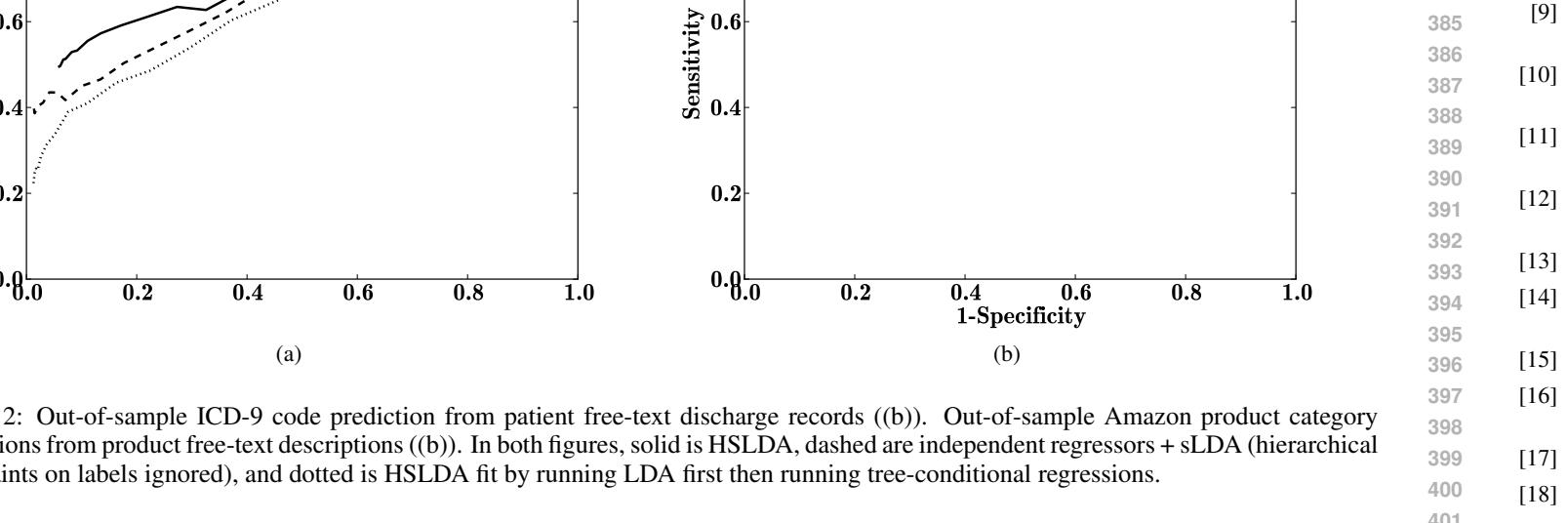


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (2(a)). Out-of-sample Amazon product category predictions from product free-text descriptions (2(b)). In both figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running treo-conditional regressions.

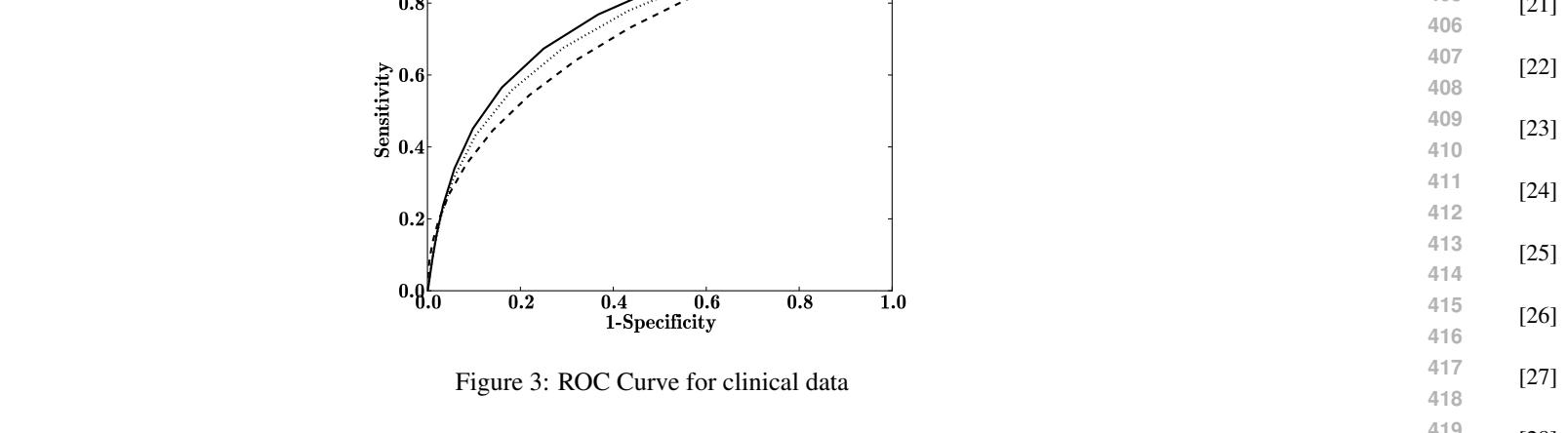


Figure 3: ROC Curve for clinical data

Bayesian has been the predominant estimation approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo sLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that if one is more sensitive to the performance gains that result from exploiting the structure of the labels then one can, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. There are applied settings in this which could be advantageous.

Extensions to this work include unbounded topic cardinality constraints and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.

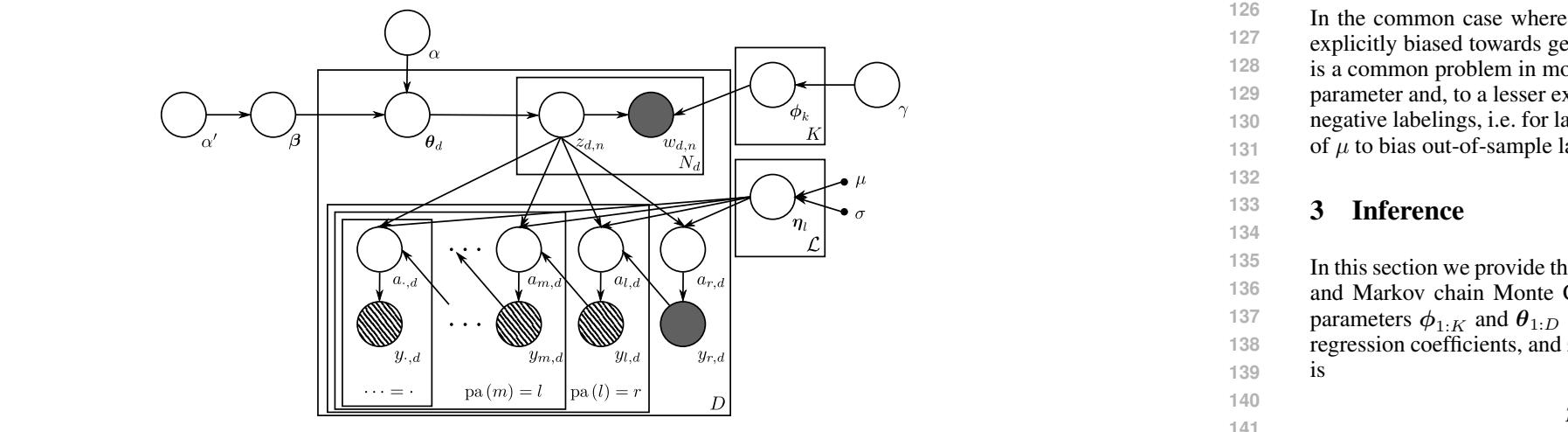


Figure 4: HSLDA graphical model

The constraints imposed by an a -is-a label hierarchy are that if the l th label is applied to document d , i.e., $y_{d,l} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{d,pa(l)} = 1, \dots, y_{d,r} = 1$. Conversely, if a label l' is marked as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labels are observed being applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA “topics” (distributions over the elements of Σ), ϕ_k is a distribution over “words”. θ_d is a document-specific distribution over topics, β is a global distribution over topics, $Dir_K(\cdot)$ is a K -dimensional Dirichlet distribution, $N_K(\cdot)$ is the K -dimensional Normal distribution, I_K is the K dimensional identity matrix, a_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 • Draw a distribution over words $\phi_k \sim Dir(\gamma \mathbf{1}_V)$

2. For each label $l \in \mathcal{L}$
 • Draw a label application coefficient $\eta_l | \mu, \sigma \sim N_K(\mu \mathbf{1}_K, \sigma I_K)$

3. Draw the global topic proportions $\beta | \alpha' \sim Dir_K(\alpha' \mathbf{1}_K)$

4. For each document $d = 1, \dots, D$
 • Draw topic proportions $\theta_d | \beta, \alpha \sim Dir_K(\alpha \beta)$

• For $a = 1, \dots, N_d$
 • Draw topic assignment $z_{a,d} | \theta_d \sim Multinomial(\theta_d)$

• Draw word $w_{a,d} | z_{a,d}, \phi_{k,a} \sim Multinomial(\phi_{k,a})$

• Set $y_{d,l} = 1$

• For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r

– Draw $a_{l,d} | z_{a,d}, \eta_l \sim N(\eta_l z_{a,d}, 1)$, $y_{pa(l),d} = 1$

– Draw $a_{l,d} | z_{a,d}, \eta_l \sim N(\eta_l z_{a,d}, 1) \mathbb{I}(a_{l,d} < 0)$, $y_{pa(l),d} = -1$

– Apply label l to document d according to $a_{l,d}$

$y_{d,l} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here $z_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k . $z_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(a_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here, each document is labeled generatively using a conditionally dependent probit regression (15). For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $(y_{pa(l),d} = 1)$) are used to determine whether or not label l is applied to document d and whether or not its parent label was applied to document d if its parent label $pa(l)$ is not applied (these are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence). The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{i,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{i,d}$ yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

In this work we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document “supervision”; often taking the form of a single numerical or categorical label. More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g. online review score), marks given to writing samples, or other document properties. The sLDA generative model is identical to the standard LDA model except that the prior $p(\theta_d | \beta)$ is replaced by $p(\theta_d | \beta, y_d)$ where y_d is the supervisory label. The sLDA posterior distribution is $p(\theta_d | \beta, y_d, \theta_r)$ where θ_r is the global topic distribution. These models were chosen to highlight several aspects of the model including performance under different types of supervision, the effect of the combined inference procedure, and regression performance attributable solely to the hierarchical constraints.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesize will result in a difference in predictive performance.

Aside from this, all other features of sLDA are preserved between the two models.

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA[22, 12, 18, 9].

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

The number of topics for all models was set to 30, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics - sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that our model generalizes to two different domains.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the observed labeling. Specifically, ancestors of observed nodes were ignored, observed nodes were considered positive and unobserved nodes were considered to be negative. This method of defining positive and negative labels was chosen to be as fair as possible to all models being compared. In particular, since the sLDA model does not respect the hierarchical constraints, the models can be compared on a more equal footing by considering only the observed labels as being positive despite the fact that ancestors must also be positive. The gold standard defined in this way will likely lead to a slight overestimation of the number of false positives. It is known that the ICD-9 codes lack sensitivity and their associated ICD-9 codes on average (std dev=0.01) and contain a range of values (std dev=0.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

Predictive performance in HSLDA is evaluated by $p(y_d | w_{1:N_d}, \beta, \theta_d, \alpha, \alpha')$ where β represents the test document. For efficiency, the expectation of this probability distribution was estimated in the following way. Expectations of z_d and η_d were estimated from the choice of probit regression [5]; the specific form of the update is a standard result from Bayesian normal regression [15]. It also is a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution [5].

$$p(a_{l,d} | z_d, \mathbf{Y}, \eta_d) \propto \frac{1}{\sqrt{2\pi}} \exp \left\{ \frac{1}{2} (a_{l,d} - \eta_d)^2 \right\} \mathbb{I}(a_{l,d} \geq 0). \quad (3)$$

HSLDA employs a hierarchical Dirichlet prior over topic assignments (i.e., β is estimated from data rather than fixed a priori). This has been shown to improve the quality and stability of inferred topics [28]. Sampling β is done using the “direct assignment” method of Teh et al. [27].

$$\beta | \mathbf{z}, \alpha, \alpha' \sim Dir(m_{(1),1} + \alpha', m_{(1),2} + \alpha', \dots, m_{(1),K} + \alpha'). \quad (4)$$

Here $m_{(1),k}$ are auxiliary variables that are required to sample the posterior distribution of β . Their conditional posterior distribution is sampled according to

$$p(m_{(1),k} = m | \mathbf{z}, \mathbf{m}_{-(1,k)}, \beta) = \frac{\Gamma(\alpha_k)}{\Gamma(\alpha_k + c_{(1),k}^T \mathbf{z}_d)} s(c_{(1),k}^T \mathbf{z}_d) (\alpha_k \beta_k)^m \quad (5)$$

where $s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

4 Related Work

In this work we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document “supervision”; often taking the form of a single numerical or categorical label. More generally this supervision is just extra per-document data; for instance its quality or relevance (e.g. online review score), marks given to writing samples, or other document properties.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of file (e.g., product descriptions and categories, medical patient records and discharge summaries) used for bookkeeping and administrative purposes. In this work we focus on ICD-9-CM codes, which are a standard way of classifying medical (ICD-9-CM) codes assigned [3].

(3)

In this work we show how to combine these two sources of information using a single model. In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA “topics” (distributions over the elements of Σ). ϕ_k is a distribution over “words.” θ_d is a document-specific distribution over topics. β is a global distribution over topics. D_{ik} is a K -dimensional Dirichlet distribution. $N_K(\cdot)$ is the K -dimensional Normal distribution. I_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(1_{\mathcal{V}})$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim N_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$
 - 3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
 - 4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$
 - For $a = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,a} | z_{n,d}, \phi_{k,n} \sim \text{Multinomial}(\phi_{k,n})$
 - Set $y_{t,d} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{t,d} | z_{d,l}, \eta_l, y_{p(a_l),d} \sim N(\bar{z}_{d,l} \eta_l, 1)$, $y_{p(a_l),d} = 1$
 - Draw $a_{t,d} | z_{d,l}, \eta_l, y_{p(a_l),d} \sim N(\bar{z}_{d,l} \eta_l, 1) \mathbb{I}(a_{t,d} < 0)$, $y_{p(a_l),d} = -1$
 - Apply label l to document d according to $a_{t,d}$

$$y_{t,d} | a_{t,d} = \begin{cases} 1 & \text{if } a_{t,d} > 0 \\ -1 & \text{otherwise} \end{cases}$$

Here $\bar{z}_{d,l}^T = [\bar{z}_{d,1}, \dots, \bar{z}_{d,L}, \bar{z}_{d,K}]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic l .

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $(y_{p(l),d} = 1)$) are used to determine whether or not label l is to be applied. Note that only if a parent label $p(l)$ is not applied (these are specific to less constrained models) can it lead to good label prediction for out-of-sample data [7]. It has been demonstrated that SLDA has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [7]. SLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical coupling in this model to avoid a posterior dependence. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

3

Note that the choice of variables $a_{t,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{t,d}$ ’s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

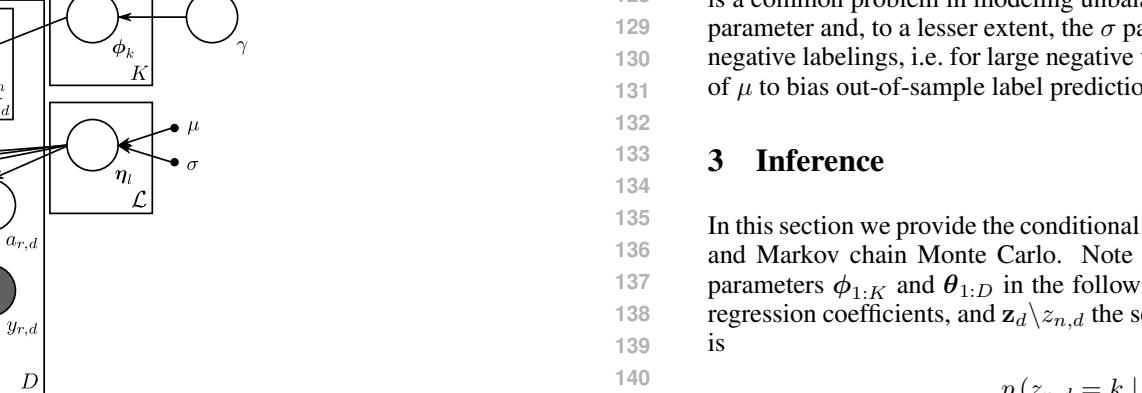


Figure 1: HSLDA graphical model

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it learning to assign all labels to all documents. This is a common problem in modeling unlabeled data. To see how this model can be biased in this way we draw the reader’s attention to the μ parameter and, to a lesser extent, the σ parameter above. Because ζ_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{t,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics – sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that our model generalizes to two different domains.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics – sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that our model generalizes to two different domains.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the observed labeling. To make the comparison as fair as possible, ancestors of observed nodes were ignored, observed nodes were considered positive and descendants of observed nodes were considered negative. This method is similar to the one used in the HSLDA graphical model to be as simple as possible while maintaining consistency. In particular, this model does not enforce any hierarchical constraints, the ancestors must also be positive. The gold standard defined in this way will likely lead to a slight overestimation of the number of false positives. It is known that ICD-9-CM codes lack sensitivity and their use as a gold standard could lead to correctly positive predictions being labeled as false positives. However, given that the label space is often large (as in our examples) it is a moderate assumption that erroneous false positives should not skew results significantly.

6

Predictive performance in HSLDA is evaluated by $p(y_{t,d} | w_{1:N_d}, d, y_{1:L-1:D})$ where \hat{d} represents the test document. For efficiency, the expectation of this probability distribution was estimated in the following way. Expectations of ζ_d and η_l were estimated from samples from the posterior. Using these expectations, we performed Gibbs sampling over the hierarchy to acquire predictive samples for the documents in the test set.

7

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

8

These results indicate that the full HSLDA model predicts more of the correct labels at a cost of an increase in the number of false positives relative to the comparison models.

9

The following are the predictive performance results for the retail product data given a prior mean for the regression parameters of -2.2. The full HSLDA model had a true positive rate of 0.85 and a false positive rate of 0.30, the sLDA model had a true positive rate of 0.78 and a false positive rate of 0.14, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.77 and a false positive rate of 0.16. These results follow a similar pattern to the clinical data.

10

To further explore this tradeoff between the true positive rate and the false positive rate we evaluated predictive performance for a range of values for two different parameters - the prior mean for the regression coefficients and the threshold for the auxiliary variables. The goal in this analysis was to evaluate the performance of these models subject to more or less stringent requirements for predicting positive labels. Products can exist in multiple locations in the hierarchy.

11

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [4]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

12

Product descriptions were fit for each model type, separate models were fit for each category. The prior mean in combination with the auxiliary variable thresholds for predicting positive labels. These two parameters have important related relations in the model. The prior mean applies negative pressure to the predictions and the auxiliary variable threshold determines the cutoff.

13

For each model type, separate models were fit for each category. The prior mean in combination with the auxiliary variable thresholds for predicting positive labels. These two parameters have important related relations in the model. The prior mean applies negative pressure to the predictions and the auxiliary variable threshold determines the cutoff.

14

Figure 2 shows the predictive performance of HSLDA relative to the three comparison models as a function of the prior mean on regression coefficients as a receiver operating characteristic (ROC) curve. For low values of the auxiliary variable threshold, the models predict labels in independent regresses (hierarchical constraints on labels ignored). HSLDA fit by first performing LDA then fitting tree-conditional regressions. These models were chosen to highlight several aspects of HSLDA including performance in the absence of hierarchical constraints, the effect of the combined inference, and regression performance attributable solely to the hierarchical constraints.

15

HSLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesized would result in a difference in predictive performance. There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the responses to influence the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in SLDA [7]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

16

For other models, particularly LDA, the prior mean was paid to settings of the prior parameters for the regression coefficients. These parameters are a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms most straightforward approaches should be of interest to practitioners.

17

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

18

Variational Bayes has been the predominant approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo SLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

19

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that if one is more sensitive to the performance gains that result from exploiting the structure of the labels then one can.

20

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms most straightforward approaches should be of interest to practitioners.

21

For other models, particularly LDA, the prior mean was paid to settings of the prior parameters for the regression coefficients. These parameters are a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms most straightforward approaches should be of interest to practitioners.

22

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

23

Variational Bayes has been the predominant approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo SLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

24

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that if one is more sensitive to the performance gains that result from exploiting the structure of the labels then one can.

25

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms most straightforward approaches should be of interest to practitioners.

26

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

27

Variational Bayes has been the predominant approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo SLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

28

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that if one is more sensitive to the performance gains that result from exploiting the structure of the labels then one can.

29

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms most straightforward approaches should be of interest to practitioners.

30

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

31

Variational Bayes has been the predominant approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo SLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

32

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that if one is more sensitive to the performance gains that result from exploiting the structure of the labels then one can.

33

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms most straightforward approaches should be of interest to practitioners.

34

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placements in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories, lists of names [2], product descriptions and categories, and patient records and diagnosis codes assigned to them for hospital discharge summaries. The ICD-9-CM controlled terminology [1] is a well-known example of this kind of data. Modification (ICD-9-CM) codes assigned [3]. In this work we show how to combine these two sources of information using a single model.

In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA “topics” (distributions over the elements of Σ). ϕ_k is a distribution over “words”. θ_d is a document-specific distribution over topics. β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $\mathcal{N}_K(\cdot)$ is the K -dimensional Normal distribution. \mathbf{I}_K is the K dimensional identity matrix, $\mathbf{1}_d$ is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 • Draw a distribution over words $\phi_k \sim \text{Dir}(\gamma \mathbf{1}_V)$

2. For each label $l \in \mathcal{L}$
 • Draw a label application coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}_K(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$

3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$

4. For each document $d = 1, \dots, D$
 • Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$

• For $a = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$

- Draw word $w_{n,d} | z_{n,d}, \phi_{k,n} \sim \text{Multinomial}(\phi_{k,n})$

• Set $y_{l,d} = 1$
 • For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r

- Draw $a_{l,d} | z_{d,l}, \eta_l \sim \mathcal{N}(z_{d,l} \eta_l, 1)$, $y_{p(l),d} = 1$
 - Draw $a_{l,d} | z_{d,l}, \eta_l \sim \mathcal{N}(z_{d,l} \eta_l, 1) \mathbb{I}(a_{l,d} < 0)$, $y_{p(l),d} = -1$

- Apply label l to document d according to $a_{l,d}$

$y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here $\mathbf{z}_d^T = [\bar{z}_1, \dots, \bar{z}_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k . $z_{d,l} = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here, each document is labeled generatively using a conditionally dependent probit regression. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $(y_{p(l),d} = 1)$) are used to determine whether or not label l is applied. This is done by applying label l to document d if its parent label $p(l)$ is not applied (i.e. $y_{p(l),d} = 0$). The regression coefficients η_l are independent a priori; however, the hierarchical coupling in this model induces a posterior dependence. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$ yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

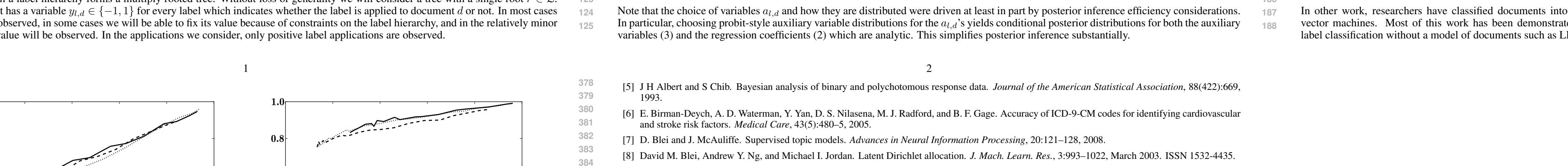


Figure 1: HSLDA graphical model

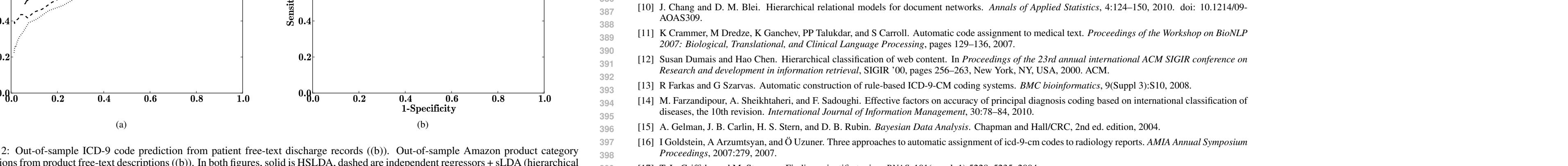
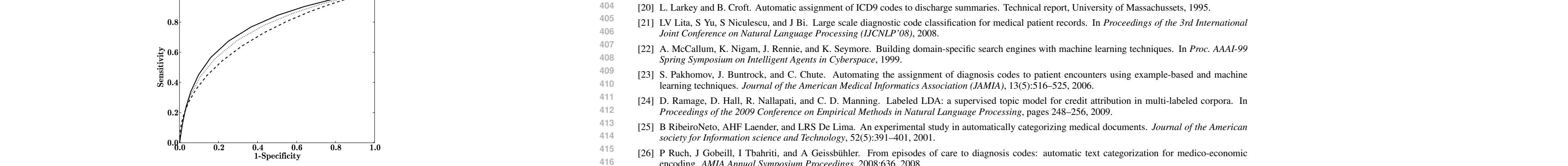


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (b)). Out-of-sample Amazon product category predictions from product free-text descriptions (b)). In both figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.



Variational Bayes has been the predominant estimation approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo sLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that if one is more sensitive to the performance gains that result from exploiting the structure of the labels then one can, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. There are applied settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality constraints and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unlabeled data. To see how this model can be biased in this way we draw the reader's attention to the μ parameter and, to a lesser extent, the σ parameter above. Because $z_{n,d}$ is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We experimented with HSLDA for prediction in two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

be shown to improve performance in sLDA [7]. This comparison model examines not the structuring of the label space, but the benefit of combined inference over both the documents and the label space.

The last comparison model is HSLDA with fixed and randomly selected regression parameters. There is a baseline benefit that the structure in the label space provides for the prediction of labels. This comparison model is intended to quantify the contribution of the structure alone.

For all of these models, particular attention was paid to the settings of the prior parameters for the regression coefficients. These parameters implement an important form of regularization in HSLDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in \{-3, -2.8, -2.6, \dots, 1\}$).

The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics - sensitivity (true positive rate) and specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that our model generalizes to two different domains.

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the observed labeling. Specifically, ancestors of observed nodes were ignored, observed nodes were considered positive and unobserved nodes were considered to be negative. This method of defining positive and negative labels was chosen to be as fair as possible to all models being compared. In particular, since the sLDA model does not respect the hierarchical constraints, its models can be compared on a more equal footing by considering only the observed labels as being positive despite the fact that ancestors must also be positive. The gold standard defined in this way will likely lead to a slight overestimation of the number of false positives. It is known that the gold standard is often large (as in our examples) and it is a moderation assumption that erroneous false positives should not skew results significantly.

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9-CM coding has been investigated in the clinical domain. Much of the work was triggered by the 2007 medical NLP challenge [1]. The data in the challenge, however, differs from ours in its scope. The datasets were smaller (1,000 training and 1,000 testing documents) and focused on radiology reports with a restricted number of ICD-9 codes (45 of them, compared to 7K+ in our dataset). Methods ranged from manual rules to online learning [11, 16, 13]. Other work had leveraged larger datasets and experimented with labeling of patient diagnoses based on a discharge summary. The ICD-9 codes are organized in a tree-structured hierarchy, with each node representing an is-a relationship between a parent and child, such that the parent diagnosis subserves the child diagnosis. For example, the code for “Pneumonia due to adenovirus” is a child of the code for “Viral pneumonia,” where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of file (see [2]), product descriptions and catalogues, and patient records and discharge codes assigned to them for hospital admissions. In this work we leverage the structure of these data to build a model for label prediction.

In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA “topics” (distributions over the elements of Σ). ϕ_k is a distribution over “words.” θ_d is a document-specific distribution over topics. β is a global distribution over topics. $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution. $N_K(\cdot)$ is the K -dimensional Normal distribution. $\text{I}_{\{ \cdot \}}$ is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 • Draw a distribution over words $\phi_k \sim \text{Dirv}(\gamma_{1k})$

2. For each label $l \in \mathcal{L}$
 • Draw a label application coefficient $\eta_l | \mu, \sigma \sim N_K(\mu \mathbf{I}_K, \sigma \mathbf{I}_K)$

3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$

4. For each document $d = 1, \dots, D$
 • Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$

• For $a = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$

- Draw word $w_{n,a} | z_{n,d}, \phi_{z_{n,d}} \sim \text{Multinomial}(\phi_{z_{n,d}})$

• Set $y_{l,d} = 1$

• For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_{d,l}, \eta_l, y_{p(l),d} \sim \begin{cases} N(\bar{z}_{d,l} \eta_l, 1), & y_{p(l),d} = 1 \\ N(\bar{z}_{d,l} \eta_l, 1) \mathbb{I}(a_{l,d} < 0), & y_{p(l),d} = -1 \end{cases}$

- Apply label l to document d according to $a_{l,d}$
 $y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except root, has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has no hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at the cost of more computationally complex inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remaining its value will be observed. In the applications we consider, only positive label applications are observed.

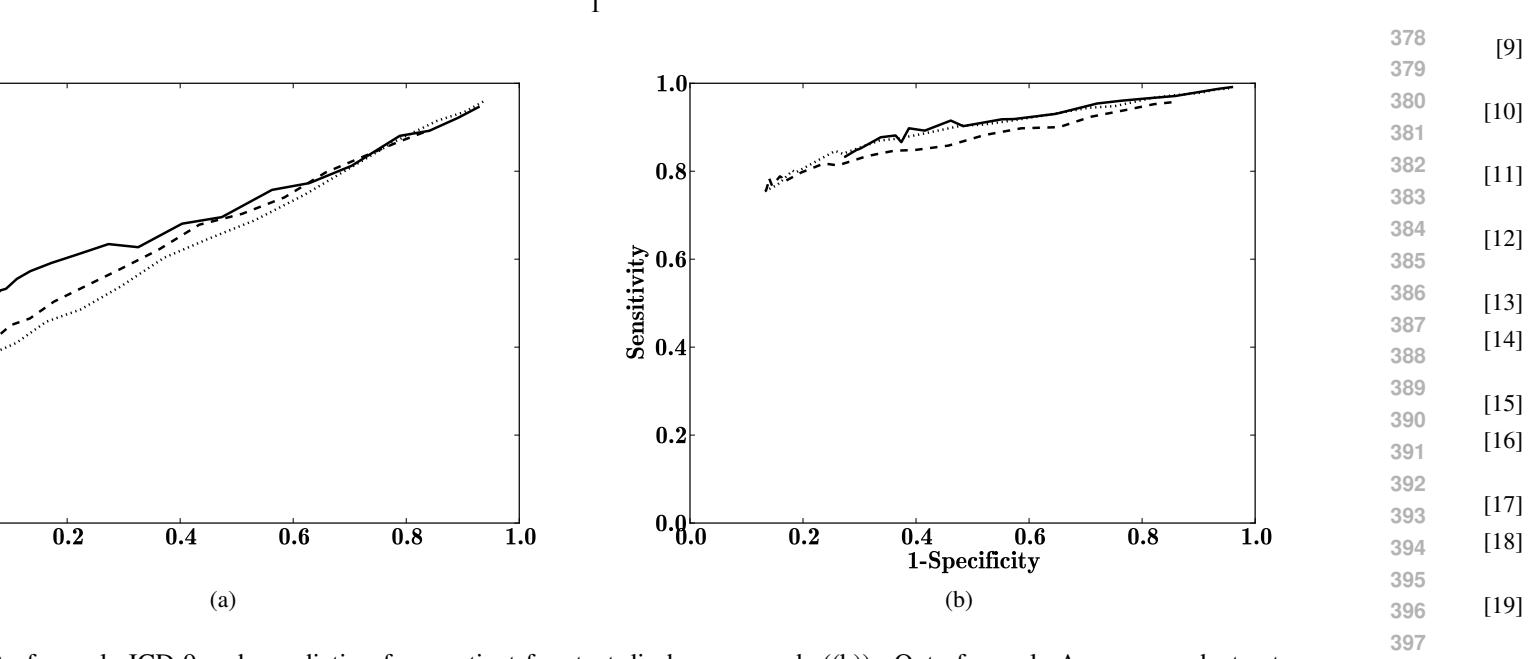


Figure 1: HSLDA graphical model

Here $\mathbf{z}_d^T = [\bar{z}_1, \dots, \bar{z}_k, \dots, \bar{z}_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k .

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here, each document is labeled generatively using a hierarchy of conditionally dependent probit regressors. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d and whether or not its parent label was applied (i.e. $(y_{p(l),d} = 1)$) are used to determine whether or not label l is to be applied. Specifically, if $y_{p(l),d} = 1$ then $a_{l,d}$ is applied to document d if its parent label $p(l)$ is applied (these regressions are specific to less constrained constraints). The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependence. The net effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$ ’s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

2.1 Related Work

In this work we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document “supervision,” often taking the form of a single numerical or categorical label.

It has been demonstrated that the signal provided by such supervision can result in better task-specific document models and can also lead to good label prediction for out-of-sample data [7]. It also has been demonstrated that sLDA has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [7], sLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we contrast HSLDA with sLDA applied in this way.

For other models, particularly those that were paying attention to settings of the prior parameters for the regression coefficients, these parameters implement an important role of regularization in LDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters is a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in [-3, -2.8, -2.6, \dots, 1]$).

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA [22, 12, 18, 9].

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nitk.org>

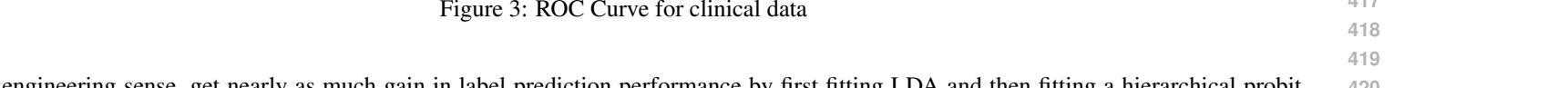


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (a) and Amazon product category predictions from product free-text descriptions (b). In both figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.

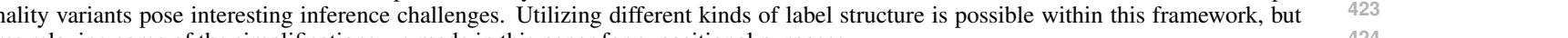


Figure 3: ROC Curve for clinical data

in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. There are applied settings in which this can be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

...
References

- [1] The international medicine center’s 2007 medical natural language processing challenge. <http://www.computationalmedicine.org/challenge/previous/2007>.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] J.H Albert and S Chiba. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669–693.
- [6] E. Birman-Deych, A. D. Waterman, Y. Yan, D. S. Nilusena, M. J. Radford, and B. F. Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [7] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing Systems*, 20:121–128, 2008.
- [8] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.

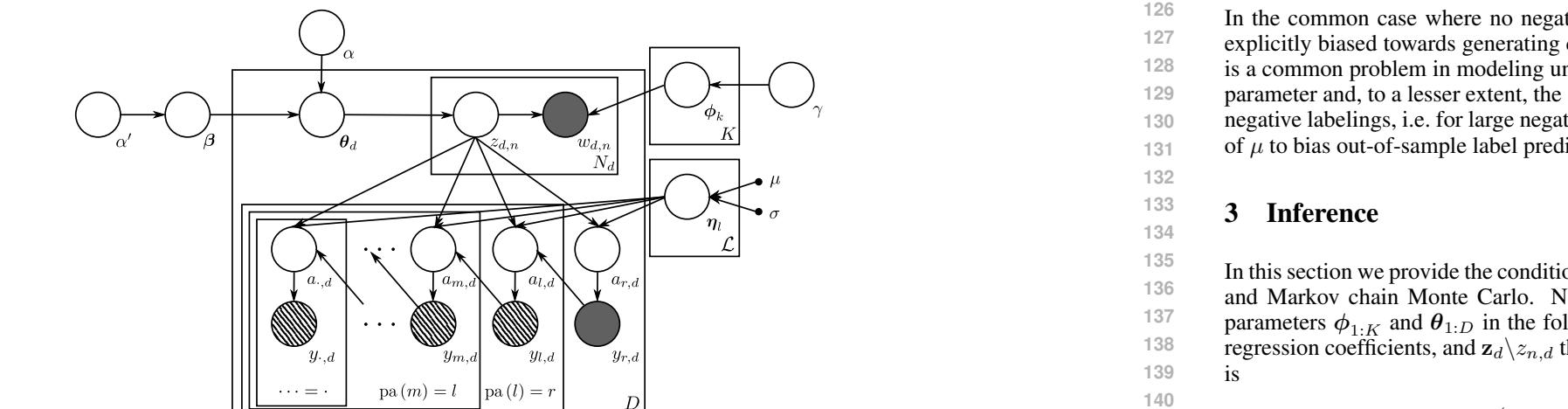


Figure 1: HSLDA graphical model

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unlabeled data. To see how this model can be biased in this way we draw the reader’s attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We applied HSLDA to data from two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the observed labeling. To make the comparison as fair as possible, ancestors of observed nodes were ignored, observed nodes were considered positive and descendants of observed nodes were considered negative. This method of evaluation is appropriate for the clinical domain, as it is common for medical models being compared. In particular, this model does not enforce a hierarchical constraint, the ancestors must also be positive. The gold standard defined in this way will likely lead to a slight overestimation of the number of false positives. It is known that ICD-9 codes lack sensitivity and their use as a gold standard could lead to correctly positive predictions being labeled as false positives. However, given that the label space is often large (as in our examples) it is a moderate assumption that erroneous false positives should not skew results significantly.

Predictive performance in HSLDA is evaluated by $p(y_{l,d} | w_{1:N_d,d}, w_{1:N_d,d}, \theta_d, \eta_l, \alpha, \beta, \gamma)$ where \hat{d} represents the test document. For efficiency, the expectation of this probability distribution was estimated in the following way. Expectations of z_d and η_l were estimated from the samples from the posterior. Using these expectations, we performed Gibbs sampling over the hierarchy to acquire predictive samples for the documents in the test set.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

These results indicate that the full HSLDA model predicts more of the correct labels at a cost of an increase in the number of false positives relative to the comparison models.

The following are the predictive performance results for the retail product data given a prior mean for the regression parameters of -2.2. The full HSLDA model had a true positive rate of 0.85 and a false positive rate of 0.30, the sLDA model had a true positive rate of 0.78 and a false positive rate of 0.14, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.77 and a false positive rate of 0.16. These results follow a similar pattern to the clinical data.

To further explore this tradeoff between the true positive rate and the false positive rate we evaluated predictive performance for a range of values for two different parameters - the prior mean for the regression coefficients and the threshold for the auxiliary variables. The goal in this analysis was to evaluate the performance of these models subject to more or less stringent requirements for predicting positive labels. Products can exist in multiple locations in the hierarchy. In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [4]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 words on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

For each model type, separate models were fit for each value of the prior mean of the regression coefficients. In contrast, to evaluate predictive performance as a function of the auxiliary variable threshold, the models predict labels in two different ways.

Figure 2 shows the predictive performance of HSLDA relative to the three comparison models as a function of the prior mean on regression coefficients as a receiver operating characteristic (ROC) curve. For low values of the auxiliary variable threshold, the models predict labels in an more sensitive and less specific manner, creating the points in the upper right corner of the ROC curve. As the auxiliary variable threshold is increased, the models predict in a less sensitive and more specific manner, creating the points in the lower left hand corner of the ROC curve.

For all values of the prior mean in both datasets HSLDA outperforms sLDA with independent regressors. In the case of HSLDA with separate trained regression, HSLDA outperforms in the clinical dataset but performs equally well across the board with the retail product dataset.

6 Discussion

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the topic modeling performance of LDA via the inclusion of observed supervision. Another way to see the family is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or not. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms sLDA applied in this way.

For other models, particularly those that were paying attention to settings of the prior parameters for the regression coefficients, these parameters implement an important role of regularization in LDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in [-3, -2.8, -2.6, \dots, 1]$).

Variational Bayes has been the predominant approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo SLDA inference, a benefit that should extend to other SLDA models. Hierarchical probit regression is used for response variable prediction there too.

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that if one is more sensitive to the performance gains that result from exploiting the structure of the labels then one can.

5 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance.

Predictive performance was evaluated by $p(y_{l,d} | w_{1:N_d,d}, w_{1:N_d,d}, \theta_d, \eta_l, \alpha, \beta, \gamma)$ where \hat{d} represents the test document.

Expectation of this probability distribution was estimated in the following way. Expectations of z_d and η_l were estimated from the samples from the posterior. Using these expectations, we performed Gibbs sampling over the hierarchy to acquire predictive samples for the documents in the test set.

The following are the predictive performance results for the clinical data given a prior mean for the regression parameters of -1.6. The full HSLDA model had a true positive rate of 0.57 and a false positive rate of 0.13, the sLDA model had a true positive rate of 0.42 and a false positive rate of 0.07, and the HSLDA model where LDA and the regressions were fit separately had a true positive rate of 0.39 and a false positive rate of 0.08.

These results indicate that the full HSLDA model predicts more of the correct labels at a cost of an increase in the number of false positives relative to the comparison models.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of file (see [2]). Product descriptions and catalogues, medical patient records and discharge codes assigned to them for hospital admissions, and hospital discharge summaries are examples of this type of data. In this work we focus on ICD-9-CM codes assigned to hospital discharge summaries. In this work we show how to combine these two sources of information using a single model.

In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA “topics” (distributions over the elements of Σ). ϕ_k is a distribution over “words.” θ_d is a document-specific distribution over topics. β is a global distribution over topics. $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution. $N_K(\cdot)$ is the K -dimensional Normal distribution. $\text{I}_{\{ \cdot \}}$ is the D -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dirv}(\mathbf{1}_{\mathcal{V}})$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim N_K(\mu \mathbf{1}_K, \sigma \text{I}_K)$
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $a = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,a} | z_{n,d}, \phi_{k,n} \sim \text{Multinomial}(\phi_{k,n})$
 - Set $y_{l,d} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_{d,l}, \eta_l, y_{p(l),d} \sim N(\eta_l z_{d,l}, 1)$, $y_{p(l),d} = 1$
 - Draw $a_{l,d} | z_{d,l}, \eta_l, y_{p(l),d} \sim N(\eta_l z_{d,l}, 1) \mathbb{I}(a_{l,d} < 0)$, $y_{p(l),d} = -1$
 - Apply label l to document d according to $a_{l,d}$

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, and the absence of a label cannot always be interpreted as a negative labeling.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. For the sake of simplicity, we focus on α -hierarchies, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical information into the model. We hypothesize that the context of labels within the hierarchy provides valuable information about labeling.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Our results show that a joint, hierarchical model outperforms a classification with unstructured labels as well as a disjoint model, where the topic modeling and the hierarchical classification are carried out independently of each other.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except root, has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard “is-a” parent-child constraints (explained later), although this assumption can be relaxed at the cost of more computationally complex inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$ will be known, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remaining its value will be observed. In the applications we consider, only positive label applications are observed.

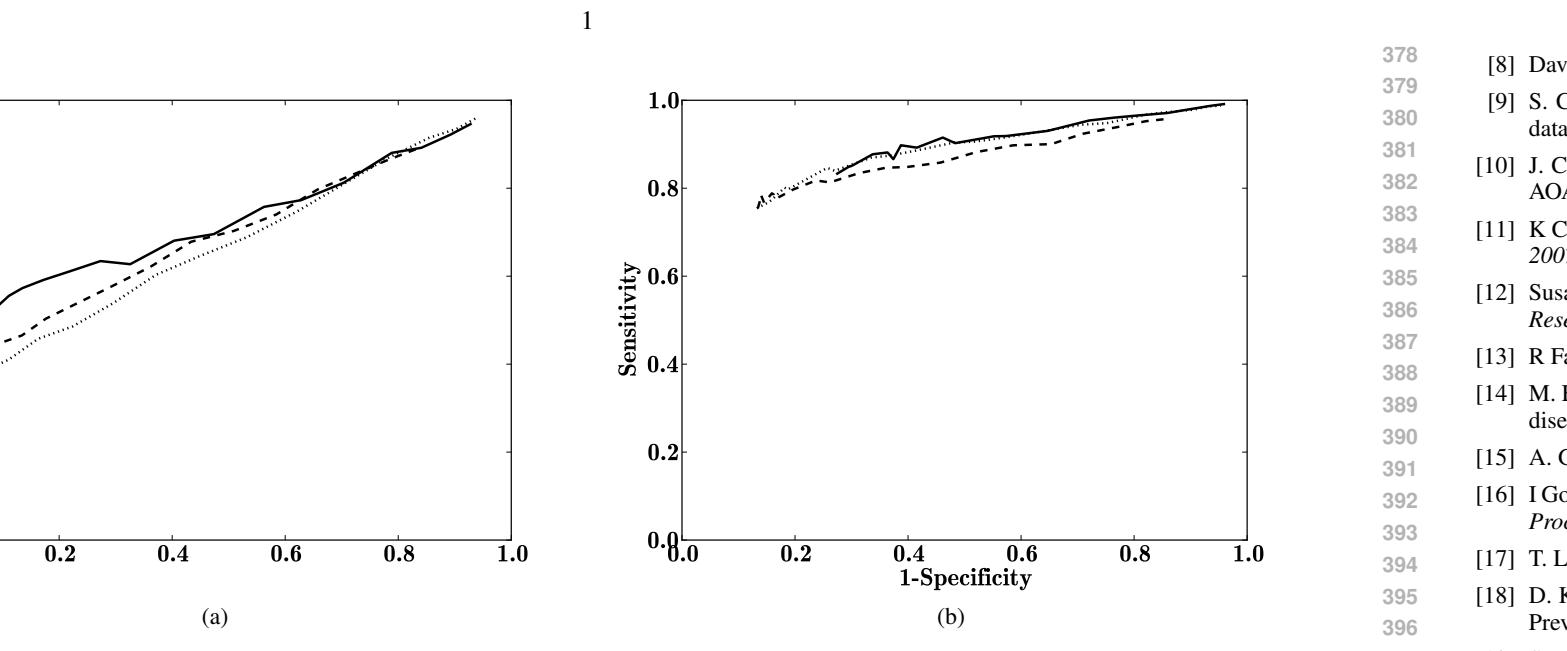
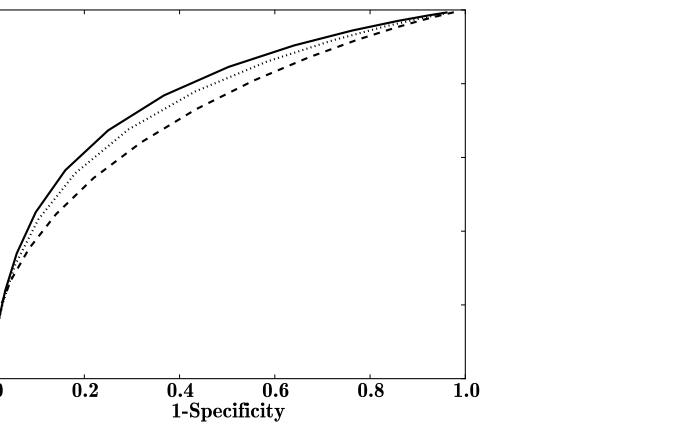


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (2(a)). Out-of-sample Amazon product category predictions from product free-text descriptions (2(b)). In both figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.



The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. In alternative interpretation of the same results is that if one is more sensitive to the performance gains that result from exploiting the structure of the labels then one, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. These are applied settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [1] The computational medicine center’s 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] J H Albert and S Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669, 1993.
- [6] E Birman-Deych, A D Waterman, Y Yan, D S Nilasena, M J Radford, and B F Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [7] D Blei and J McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.

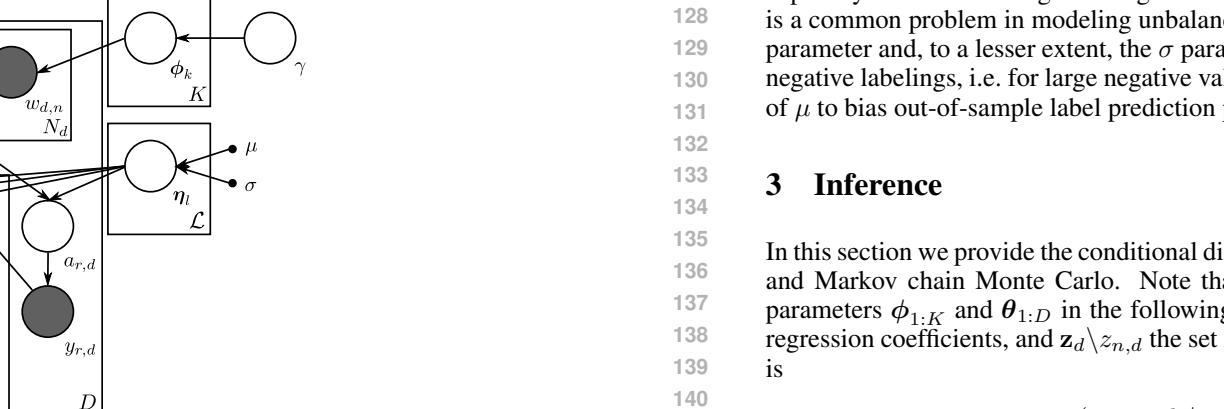


Figure 1: HSLDA graphical model

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unbalanced data. To see how this model can be biased in this way we draw the reader’s attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We applied HSLDA to data from two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing

5.1.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize patient hospitalization course. The summaries typically contain a record of patient complaints, findings and diagnoses, along with treatment and hospital course. For each hospitalization, trained medical coders review and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [17], we have analytically marginalized out the parameters $\phi_{k,K}$ and $\theta_{1,D}$ in the following expressions. Let \mathbf{a} be the set of auxiliary variables, \mathbf{w} the set of all words, η the set of all regression coefficients, and $\mathbf{z}_{\setminus d}$ the set \mathbf{z}_d with element $z_{n,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | \mathbf{z}_{\setminus d}, \mathbf{z}_{n,d}, \mathbf{w}, \mathbf{\eta}, \alpha, \beta, \gamma) \propto \\ (c_{(i),d}^{k-(n,d)} + \alpha \beta_i) \frac{e^{\frac{\alpha}{\sigma^2} n_{i,d} + \gamma}}{(c_{(i),d}^{k-(n,d)} + V + \gamma)} \prod_{l \in \mathcal{L}} \exp \left\{ -\frac{(z_{l,d} - \mu_{l,d})^2}{2 \sigma^2} \right\} \quad (1)$$

where $c_{(i),d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The false positive rate $\text{I}_{(i),d}$ is the probability that a word v is a positive label for document d .

The task of automatic ICD-9 coding has been investigated in the clinical domain. Methods range from manual rules to online learning [11, 16, 17]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Ridge Regression, as well as simple keyword mappings, all with promising results [20, 23, 26, 21].

Our dataset was gathered from the clinical data warehouse of a large metropolitan hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² A fixed vocabulary was formed by taking the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Descriptions and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a multiply-rooted hierarchy and provides textual product descriptions for most products. Products can be discovered by users through free-text search and product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

As sensitivity and specificity can always be traded off we examined sensitivity for a range of values for two different parameters - the prior means for the regression coefficients and the threshold for the auxiliary variables. The prior in this analysis was to evaluate the performance of these models subject to more or less stringent requirements for predicting positive labels. These two parameters have important relatedness in the model. The prior mean in combination with the auxiliary variable threshold together encode the strength of the prior belief that unobserved labels are likely to be negative. Effectively, the prior mean applies negative pressure to the predictions and the auxiliary variable threshold determines the cutoff.

For each model type, separate models were fit for each value of the prior mean of the regression coefficients. This is a proper Bayesian sensitivity analysis. In contrast, to evaluate predictive performance as a function of the auxiliary variable threshold, a single model was fit for each model type and prediction was evaluated based on predictive samples drawn subject to different auxiliary variable thresholds. These methods are significantly different since the prior mean is varied prior to inference and the auxiliary variable threshold is varied following inference. However, as intended, they both highlight model performance under more or less stringent requirements for predicting positive labels.

Figure 2 shows the predictive performance of HSLDA relative to the three comparison models as a function of the prior mean on regression coefficients as a receiver operating characteristic (ROC) curve. For low values of the auxiliary variable threshold, the models predict label in independent regressors (hierarchical constraints on labels ignored). These models were chosen to highlight several aspects of HSLDA including performance in the absence of hierarchical constraints, the effect of the combined inference, and regression performance attributable solely to the hierarchical constraints.

HSLDA employs a hierarchical Dirichlet prior over topic assignments (i.e., β is estimated from data rather than fixed a priori). This has been shown to improve the quality and stability of inferred topics [28]. Sampling β is done using the “direct assignment” method of Teh et al. [27].

$\beta | \mathbf{z}, \mathbf{Y}, \eta, \alpha \sim \text{Dir}(m_{(i),1} + \alpha, m_{(i),2} + \alpha, \dots, m_{(i),K} + \alpha)$ (4)

Here $m_{(i),k}$ are auxiliary variables that are required to sample the posterior distribution of β . Their conditional posterior distribution is sampled according to

$$p(m_{(i),k} = m | \mathbf{z}, \mathbf{y}_{\setminus d}, \beta) \propto \frac{\Gamma(\alpha \beta_k)}{\Gamma(\alpha \beta_k + c_{(i),d}^k)} s(c_{(i),d}^k, m) (\alpha \beta_k)^m \quad (5)$$

where $s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

4 Related Work

We evaluated HSLDA along with three other closely related models against these two datasets. The comparison models included sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by first performing LDA then fitting tree-conditional regressions. These models were chosen to highlight several aspects of HSLDA including performance in the absence of hierarchical constraints, the effect of the combined inference, and regression performance attributable solely to the hierarchical constraints.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesized would result in a difference in predictive performance.

There are two components to HSLDA, LDA and a hierarchically constrained response. In this comparison model, the separate inference processes do not allow the responses to influence the low dimensional structure imposed by LDA. Combined inference has been shown to improve performance in sLDA [7]. This comparison model examines not the structure of the label space, but the benefit of combined inference over both the documents and the label space.

For other models, particularly those that were payed attention to the prior parameters for the regression coefficients. These parameters implement an important role of regularization in LDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters is specific to less constraints but can be applied to accommodate different constraints. The regression coefficients η_l are independent a priori; however, the hierarchical coupling in this model induces a posterior dependence. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$ yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

<http://www.cdc.gov/nchs/icd/icd9cm.htm>

<http://www.nitk.org>

6 Discussion

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of models that improve on the topic modeling performance of LDA via the inclusion of observed supervisory. Another way to see the model is as a set of models that can predict labels for bag-of-words data. A large diversity of problems can be expressed as label prediction problems for bag-of-words data. A surprisingly large amount of that kind of data possess structured labels; either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms more straightforward approaches should be of interest to practitioners.

Variational Bayes has been the predominant estimation approach applied to sLDA models. Hierarchical probit regression makes for tractable Markov chain Monte Carlo sLDA inference, a benefit that should extend to other sLDA models should probit regression be used for response variable prediction there too.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)
Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of file (see [2]), product descriptions and catalogues, patient records and discharge codes assigned to them for hospital admissions, and hospital discharge summaries. In addition, there are many other types of structured data, such as ICD-9-CM codes assigned (3)).

In

HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA "topics" (distributions over the elements of Σ). ϕ_k is a distribution over "words." θ_d is a document-specific distribution over topics. β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $N_K(\cdot)$ is the K -dimensional Normal distribution. \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value one if its argument is true and zero otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$

• Draw a distribution over words $\phi_k \sim \text{Dir}(1_{\mathcal{V}})$

2. For each label $l \in \mathcal{L}$

• Draw a label application coefficient $\eta_l | \mu, \sigma \sim N_K(\mu\mathbf{I}_K, \sigma\mathbf{I}_K)$

3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$

4. For each document $d = 1, \dots, D$

• Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{1}_K)$

• For $a = 1, \dots, N_d$

– Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$

– Draw word $w_{n,d} | z_{n,d}, \phi_{1,K} \sim \text{Multinomial}(\phi_{1,K})$

• Set $y_{l,d} = 1$

• For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r

– Draw $a_{l,d} | z_{d,l}, \eta_l, y_{l,d} \sim N(\bar{z}_{d,l} \eta_l, 1)$, $y_{l,d} = 1$

– $y_{l,d} < 0$, $y_{l,d} = -1$

– Apply label l to document d according to $a_{l,d}$

$y_{l,d} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here $\bar{z}_{d,l} = [\bar{z}_1, \dots, \bar{z}_k, \dots, \bar{z}_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k . $z_{d,l} = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{1,N_d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except root, has a parent $\text{pa}(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has cardinality $|l|$ and that the label l is marked as being positive if it is applied to document d if its parent label $\text{pa}(l)$ is not applied to document d . The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependency. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

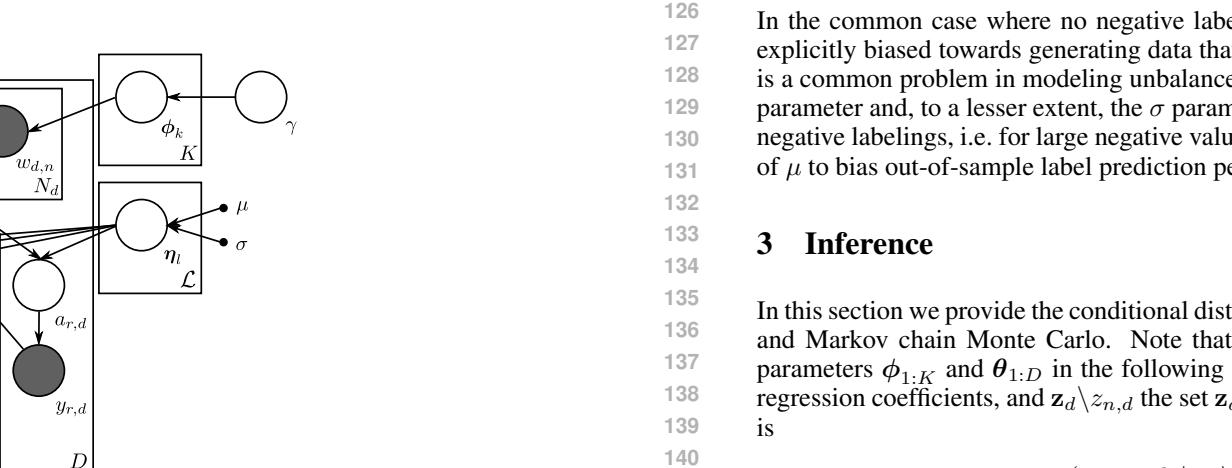


Figure 1: HSLDA graphical model

In the common case where no negative labels are observed (like the example application we consider in Section 5), the model must be slightly biased towards generating data that has negative labels in order to keep it learning to assign all labels to all documents. This is a common problem in modeling unlabeled data. To see how this model can be biased in this way we draw the reader's attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We applied HSLDA to data from two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance. Predictive performance was measured with standard metrics – sensitivity (true positive rate) and 1-specificity (false positive rate). We evaluate on the two aforementioned datasets to demonstrate that our model generalizes to two different domains.

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the observed labeling. To make the comparison as fair as possible, ancestors of observed nodes were ignored, observed nodes were considered positive and descendants of observed nodes were considered negative. This method is appropriate for the medical domain, but may not be as appropriate for the Amazon product domain. In particular, the HSLDA model does not enforce a hierarchical constraint, whereas the comparison models do.

The information in the discharge summary and assign a series of diagnoses codes. Following the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for "Pneumonia due to adenovirus" is a child of the code for "Viral pneumonia," where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [6], and sometimes make mistakes [14].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Methods range from manual rules to online learning [11, 16, 13]. Other work have leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Tree Regression, as well as simple keyword mappings, all with promising results [20, 23, 26, 21].

Our dataset was gathered from the clinical data warehouse of a large metropolitan hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² A fixed vocabulary was formed by taking the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Descriptions and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a multiply-rooted hierarchy and provides textual product descriptions for most products. Products can be discovered by users through query-based searching or through product category exploration. Top-level product categories are displayed on the front page of the website and lower-level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

To further explore this tradeoff between the true positive rate and the false positive rate we evaluated predictive performance for a range of values for two different parameters – the prior mean for the regression coefficients and the threshold for the auxiliary variables. The goal in this analysis was to evaluate the performance of these models subject to more or less stringent requirements for predicting positive labels.

These two parameters have important related relations in the model. The prior mean in combination with the auxiliary variable threshold pressure to the predictions and the auxiliary variable threshold determines the cutoff.

For each model type, separate models were fit for each value of the prior mean of the regression coefficients. In contrast, to evaluate predictive performance as a function of the auxiliary variable threshold, the prior mean in combination with the auxiliary variable threshold pressure to the predictions and the auxiliary variable threshold determines the cutoff.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [4]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

5.2 Comparison Models

We evaluated HSLDA along with three other closely related models against these two datasets. The comparison models included sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by first performing LDA then fitting tree-conditional regressions. These models were chosen to highlight several aspects of HSLDA including performance in the absence of hierarchical constraints, the effect of the combined inference, and regression performance attributable solely to the hierarchical constraints.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesized would result in a difference in predictive performance.

There are two components to HSLDA, LDA and a hierarchically constrained response. The second comparison model is HSLDA fit by performing LDA first followed by performing inference over the hierarchically constrained label space. In this comparison model, the separate inference processes do not allow the responses to influence the low dimensional structure inferred by LDA. Combined inference has been shown to improve performance in sLDA [7]. This comparison model examines not the structure of the label space, but the benefit of combining inference over both the documents and the label space.

For all of these models, particularly HSLDA, was paid to settings of the prior parameters for the regression coefficients. These parameters implement an important role of regularization in LDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters is specific to less constrained but can be modified to accommodate current constraints. The regression coefficients η_l are independent a priori, however, the hierarchical coupling in this model induces a posterior dependency. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$ yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

Figure 2 shows the predictive performance of HSLDA relative to the three comparison models as a function of the prior mean on regression coefficients as a receiver operating characteristic (ROC) curve. For low values of the auxiliary variable threshold, the models predict labels in an more sensitive and less specific manner, creating the points in the upper right corner of the ROC curve. As the auxiliary variable threshold is increased, the models predict in a less sensitive and more specific manner, creating the points in the lower left hand corner of the ROC curve. For all values of the prior mean in both datasets HSLDA outperforms sLDA with independent regressors. In the case of HSLDA with separately trained regression, HSLDA outperforms in the clinical dataset but performs equally well across the board with the retail product dataset.

6 Discussion

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that impose on the topic modeling performance of LDA via the inclusion of observed supervision. An alternative, complementary way is to see it as a set of models that can predict labels for bag-of-word data. A large diversity of products can be expressed as label prediction problems for bag-of-word data. A surprisingly large amount of that kind of data possess structured labels, either hierarchically constrained or otherwise. That HSLDA directly addresses this kind of data is a large part of the motivation for this work. That it outperforms sLDA applied in this way.

Other models that incorporate LDA and supervision include LabelledLDA[24] and DiscLDA[19]. Various applications of these models to computer vision and document networks have been explored [29, 10]. None of these models, however, leverage dependency structure in the label space.

Variational Bayes has been the predominant approach applied to sLDA models. Hierarchical probit regression makes for tractable inference. Such a label hierarchy forms a multiply rooted tree. The number of topics for all models was set to 50, the prior distributions of $p(\alpha)$, $p(\alpha')$, and $p(\gamma)$ were gamma distributed with a shape parameter of 1 and a scale parameter of 1000.

¹<http://www.cdc.gov/nchs/icd/icd9cm.htm>

²<http://www.nltk.org>

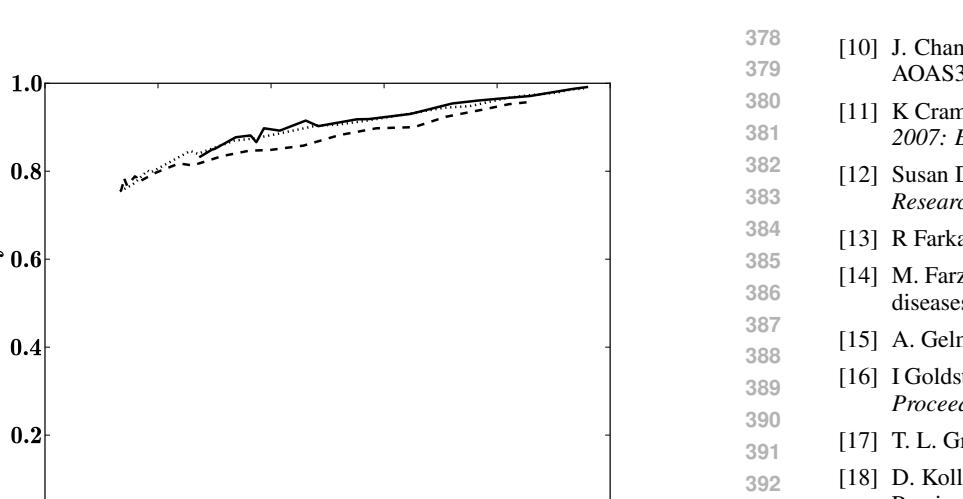


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (b)). Out-of-sample Amazon product category predictions from product free-text descriptions (b)). In both figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is sLDA fit by running LDA first then running tree-conditional regressions.

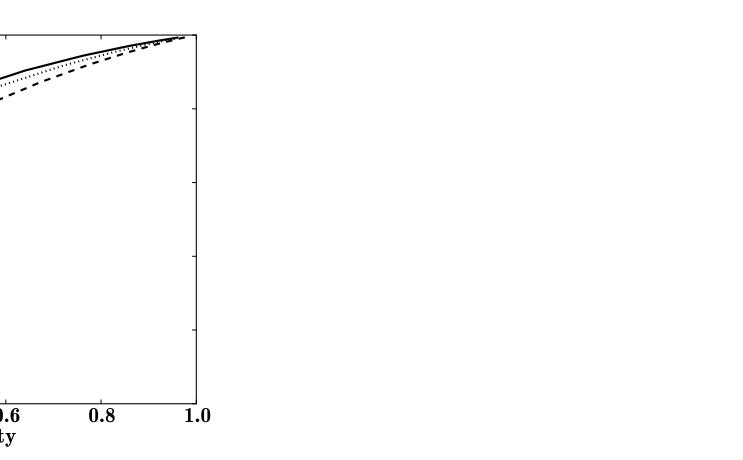


Figure 3: ROC Curve for clinical data

can, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. There are applied settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

- References
- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
 - [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
 - [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
 - [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
 - [5] J H Albert and S Chib. Bayesian analysis of binomial and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669-693.
 - [6] E Birman-Deych, A D Waterman, Y Yan, D S Nilsson, M J Ralford, and F Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480-5, 2005.
 - [7] D Blei and J McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121-128, 2008.
 - [8] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993-1022, March 2003. ISSN 1532-4435.
 - [9] S Chakrabarti, B Dom, R Agrawal, and P Raghavan. Scalable feature selection, classification and signature generation for organizing large text databases into hierarchical topic taxonomies. *The VLDB Journal*, 7:163-178, August 1998. ISSN 1066-8888.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Authors)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of file (see [2]). Product descriptions and catalogues, medical patient records and discharge documents to them for hospital admissions, and hospital discharge summaries are examples of this type of data. In this work we focus on ICD-9-CM codes assigned to hospital admissions, and medical patient records. We also focus on product descriptions and retail product categorization tasks.

In this work we show how to combine these two sources of information using a single model.

In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA "topics" (distributions over the elements of Σ), ϕ_k is a distribution over "words," θ_d is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $N_K(\cdot)$ is the K -dimensional Normal distribution, I_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(1_{\text{V}})$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim N(\mu \mathbf{1}_K, \sigma \text{I}_K)$
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \beta)$
 - For $a = 1, \dots, N_d$
 - Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,a} | z_{n,d}, \phi_{k,n} \sim \text{Multinomial}(\phi_{k,n})$
 - Set $y_{l,d} = 1$
 - For each label l in a breadth first traversal of L starting at the children of root r
 - Draw $a_{l,d} | z_d, \eta_l, y_{p(l),d} \sim N(z_d^\top \eta_l, 1)$, $y_{p(l),d} = 1$
 - Draw $a_{l,d} | z_d, \eta_l, y_{p(l),d} \sim N(z_d^\top \eta_l, 1) \mathbb{I}(a_{l,d} < 0)$, $y_{p(l),d} = -1$
 - Apply label l to document d according to $a_{l,d}$

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, and the absence of a label cannot always be interpreted as a negative labeling.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. For the sake of simplicity, we focus on α -hierarchies, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical information into the model. We hypothesize that the context of labels within the hierarchy provides valuable information about labeling.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Our results show that a joint, hierarchical model outperforms a classification with unstructured labels as well as a disjoint model, where the topic modeling and the hierarchical classification are carried out independently of each other.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except root, has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has hard "is-a" parent-child constraints (explained later), although this assumption can be relaxed at the cost of more computationally complex inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$.

Each document has a variable $a_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$ will be known, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remaining its value will be observed. In the applications we consider, only positive label applications are observed.

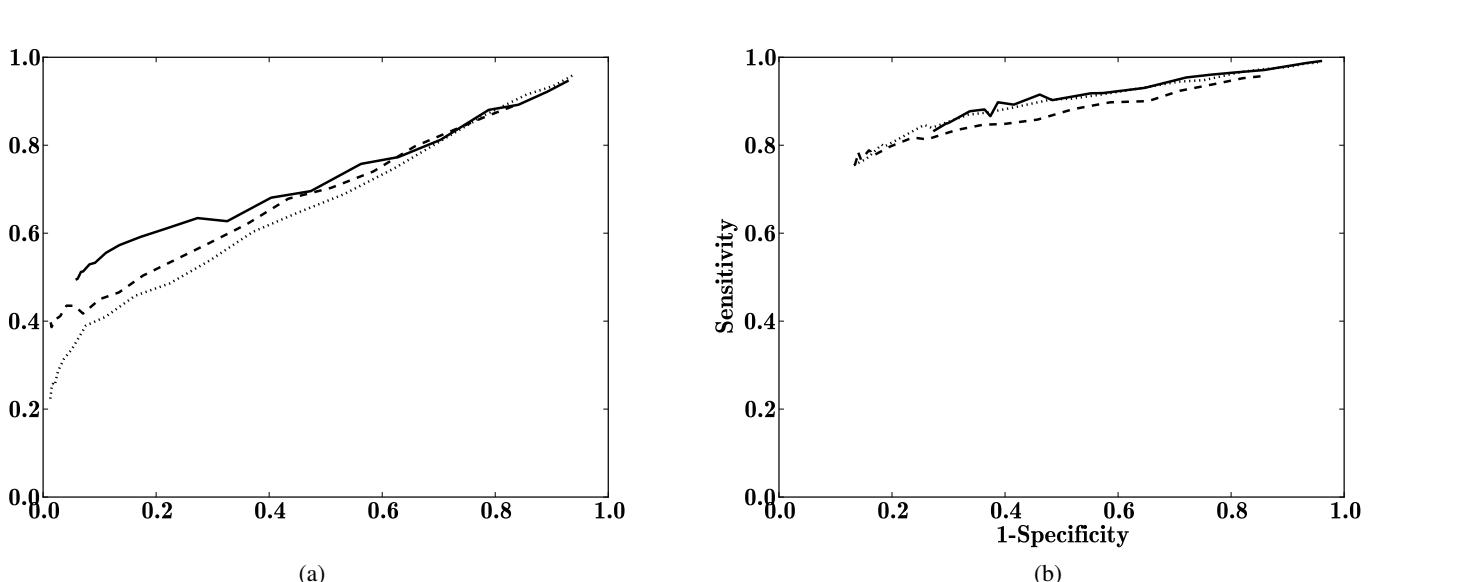
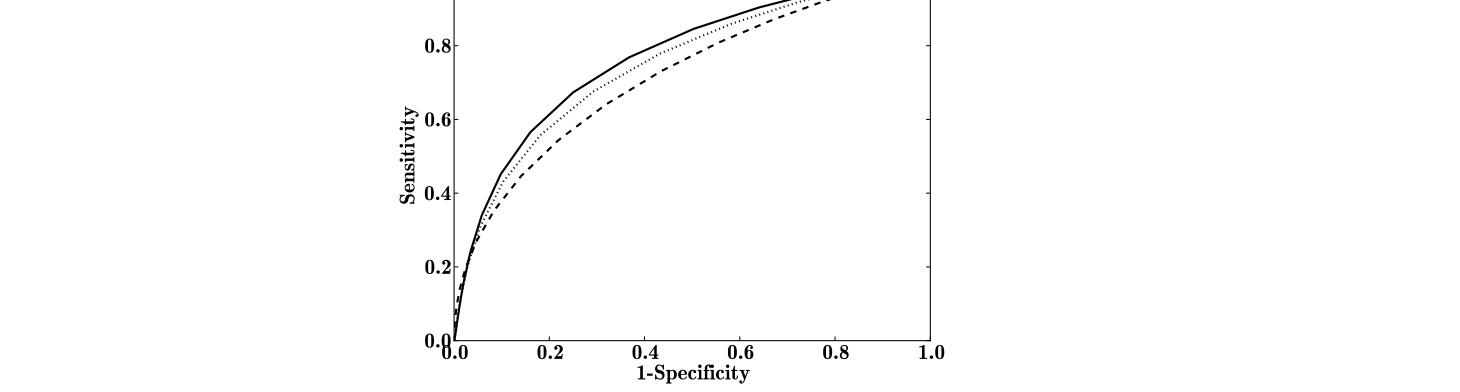


Figure 2: Out-of-sample ICD-9 code prediction from patient free-text discharge records (b)). Out-of-sample Amazon product category predictions from product free-text descriptions (b)). In both figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.



The results in Figures 2(b) and 3(b) suggest that in most cases it is better to do full joint estimation of HSLDA. An alternative interpretation of the same results is that, if one is more sensitive to the performance gains that result from exploiting the structure of the labels, then one can, in engineering terms, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. These are applied settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

- ## References
- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007
 - [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
 - [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
 - [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
 - [5] J H Albert and S Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669, 1993.
 - [6] E Birman-Deych, A D Waterman, Y Yan, D S Nilusena, M J Radford, and B F Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480-5, 2005.
 - [7] D Blei and J McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.

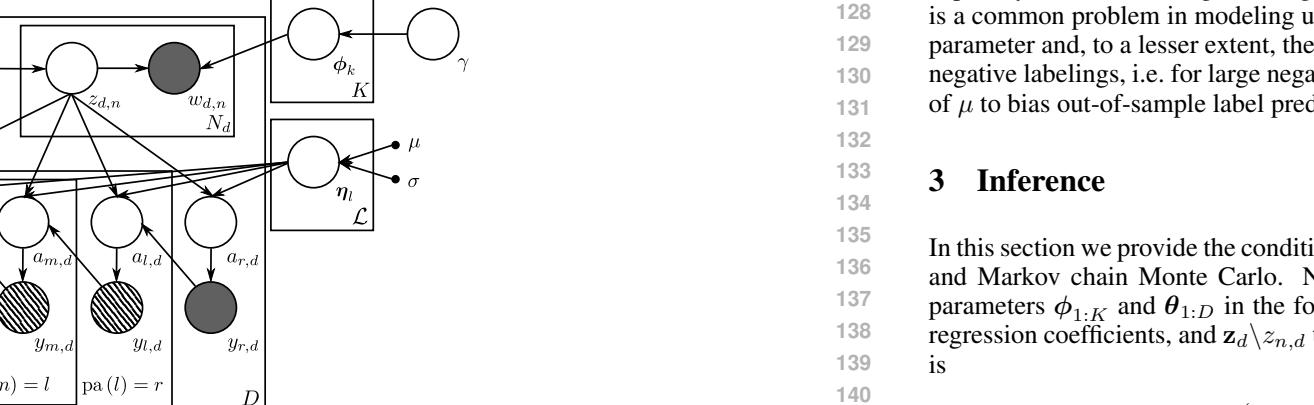


Figure 1: HSLDA graphical model

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unlabeled data. To see how this model can be biased in this way we draw the reader's attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We applied HSLDA to data from two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.3 Evaluation and Results

We evaluated our model, HSLDA, against the comparison models with a focus on predictive performance on held-out data. Prediction performance was measured with standard metrics - sensitivity (true positive rate) and 1-specificity (false positive rate).

To evaluate the performance of these models, we established a gold standard for comparison. For each dataset, a held out set of 1000 documents and labels were reserved for evaluation and predictive performance was evaluated against a standard derived from the observed labeling. To make the comparison as fair as possible, ancestors of observed nodes were ignored, observed nodes were considered positive and descendants of observed nodes were considered to be negative. This method of defining positive and negative labels was chosen to be as fair as possible to all models being compared. In particular, since the sLDA model does not enforce the hierarchical constraints, its performance was found to be so poor that it predicted that it was included. The models can be compared on a more equal footing by considering only the labels in the tree that are not descendants of the root node. The gold standard model may well have a false positive rate that will likely inflate the number of false positives because the labels applied to any particular document are usually not as complete as they could be. ICD-9 codes lack sensitivity and their use as a gold standard could lead to correctly positive predictions being labeled as false positives [1]. However, given that the label space is often large (as in our examples) it is a moderate assumption that erroneous false positives should not skew results significantly.

Predictive performance in HSLDA is evaluated by $p(y_{l,d} | w_1, \dots, w_{N_d}, d, y_{e,1}, \dots, y_{e,L_e})$ for each test document, d . For efficiency, the expectation of this probability distribution was estimated in the following way. Expectations of z_d and η_l were estimated with samples from the posterior. Using these expectations, we performed Gibbs sampling over the hierarchy to acquire predictive samples for the documents in the test set. The true positive rate was calculated as the average expected labeling for gold standard positive labels. The false positive rate is $1 - \text{true negative rate}$.

Our dataset was gathered from the clinical data warehouse of a large metropolitan hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² A fixed vocabulary was formed by taking the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize patient hospitalization course. The summaries typically contain a record of patient complaints, findings and diagnoses, along with treatment and hospital course. For each hospitalization, trained medical coders review and Markov chain Monte Carlo. Note that, like in collapsed Gibbs samplers for LDA [17], we have analytically marginalized out the parameters $\phi_{k,K}$ and $\theta_{1,D}$ in the following expressions. Let a be a set of auxiliary variables, w the set of all words, η the set of all regression coefficients, and $z_d \setminus z_{n,d}$ the set z_d with element $z_{n,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | z_d \setminus z_{n,d}, w, \eta, \alpha, \beta, \gamma) \propto \left(c_{(k),d}^{k-(n,d)} + \alpha \beta_k \right) \frac{c_{(k),d}^{k-(n,d)+1}}{\left(c_{(k),d}^{k-(n,d)} + V \right)^{\gamma+1}} \prod_{l \in \mathcal{L}_d} \exp \left\{ -\frac{(z_{l,d}^T \eta - a_{l,d})^2}{2} \right\} \quad (1)$$

where $c_{(k),d}^{k-(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The subscript (\cdot) indicates to sum over the range of the replaced variable, i.e. $c_{w,n,d}^{k-(n,d)} = \sum_d c_{w,n,d}^{k-(n,d)}$. Here \mathcal{L}_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\eta_l | z, a, \sigma) = \mathcal{N}(\hat{\mu}_l, \hat{\Sigma}) \quad (2)$$

where

$$\mu_l = \Sigma \left(\frac{\mu}{\sigma} + Z^T a_l \right) \quad \Sigma^{-1} = \mathbf{I}_{|\mathcal{L}|} + Z^T Z$$

Here Z is a $D \times K$ matrix such that row d of Z is z_d , and $a_l = [a_{l,1}, a_{l,2}, \dots, a_{l,D}]^T$. The simplicity of this conditional distribution follows from the choice of probit regression [5]; the specific form of the update is a standard result from Bayesian normal linear regression [15]. It is also a standard probit regression result that the conditional posterior distribution of $a_{l,d}$ is a truncated normal distribution [5].

As sensitivity and specificity can always be traded off we examined sensitivity for a range of values for two different parameters - the prior means for the regression coefficients and the threshold for the auxiliary variables. The goal in this analysis is to evaluate the performance of these models subject to more or less stringent requirements for predicting positive labels. These two parameters have important related

functions in the model. The prior mean in combination with the auxiliary variable threshold together encode the strength of the prior belief that unobserved labels are likely to be negative. Effectively, the prior mean applies negative pressure to the predictions and the auxiliary variable threshold determines the cutoff.

For each model type, separate models were fit for each value of the prior mean of the regression coefficients. This is a proper Bayesian sensitivity analysis. In contrast, to evaluate predictive performance as a function of the auxiliary variable threshold, a single model was fit for each model type and prediction was evaluated based on predictive samples drawn subject to different auxiliary variable thresholds. These methods are significantly different since the prior mean is varied prior to inference and the auxiliary variable threshold is varied following inference. However, as intended, both high model performance under more or less stringent requirements for predicting positive labels.

Figure 2 shows the predictive performance of HSLDA relative to the three comparison models as a function of the prior mean on regression coefficients as a receiver operating characteristic (ROC) curve. For low values of the auxiliary variable threshold, the models perform similarly. These models were chosen to highlight several aspects of HSLDA including performance in the absence of hierarchical constraints, the effect of the combined inference, and regression performance attributable solely to the hierarchical constraints.

HSLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesized would result in a difference in predictive performance. For each model type, separate models were fit for each value of the prior mean of the regression coefficients. These methods are significantly different since the prior mean is varied prior to inference and the auxiliary variable threshold is varied following inference. However, as intended, both high model performance under more or less stringent requirements for predicting positive labels.

For each of these models, particular attention was paid to settings of the prior parameters for the regression coefficients. These parameters implement an important role of regularization in LDA. In the setting where there are no negative labels, a Gaussian prior over the regression parameters are specific to less than constraints but can be applied to accommodate different constraints. The regression coefficients η_l are independent a priori; however, the hierarchical coupling in this model induces a posterior dependence. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant hard "is-a" parent-child constraint (explained later), although this assumption can be relaxed at the cost of more computationally complex inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$.

In this work we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document "supervision", often taking the form of a single numerical or categorical label. It has been demonstrated that the signal provided by such supervision can result in better task-specific document models and can also lead to good label prediction for out-of-sample data [7]. It has also been demonstrated that sLDA has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [7]. sLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we constrain HSLDA with hierarchical constraints on labels to be as similar to sLDA as possible. The second comparison model is HSLDA with independent regressors. These models are specifically designed to handle structured labels, either hierarchically constrained or otherwise. That sLDA directly addresses this kind of data is a large part of the motivation for this work. It outperforms sLDA with independent regressors in this way.

Other models that incorporate LDA and supervision include LabeledLDA[24] and DiscLDA[19]. Various applications of these models to different fields with a range of values for μ , the mean prior parameter for regression coefficients ($\mu \in [-3, -2.8, -2.6, \dots, 1]$). None of these models, however, leverage dependency structure in the label space.

In other work, researchers have classified documents into a hierarchy (a closely related task) with naive Bayes classifiers and support vector machines. Most of this work has been demonstrated on relatively small datasets, small label spaces, and has focused on single label classification without a model of documents such as LDA [12, 18, 9].

²<http://www.cdc.gov/nchs/icd/icd9cm.htm>

³<http://www.nitk.org>

4 Related Work

In this work we extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [8] augmented with per-document "supervision", often taking the form of a single numerical or categorical label. It has been demonstrated that the signal provided by such supervision can result in better task-specific document models and can also lead to good label prediction for out-of-sample data [7]. As the auxiliary variable threshold is increased, the models predict in a less sensitive and more specific manner, moving the points in the lower left hand corner of the ROC curve. For all values of the prior mean in both datasets HSLDA outperforms sLDA with independent regressors. In the case of HSLDA with separately trained regression, HSLDA outperforms in the clinical dataset but performs equally well across the board with the retail product dataset.

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Authors)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs, and patient records and diagnosis codes assigned to them for bookkeeping and insurance purposes. In this work we show how to combine these two sources of information using a single model that can automatically categorize new text documents automatically, suggest labels that might be inaccurate, compare improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable to other data as well; namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hatching indicates that potentially some of the plated variables are observed).

In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressions. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA "topics" (distributions over the elements of Σ), ϕ_k is a distribution over "words," θ_k is a document-specific distribution over topics, β is a global distribution over topics, $\text{Dir}_K(\cdot)$ is a K -dimensional Dirichlet distribution, $\mathcal{N}_K(\cdot)$ is the K -dimensional Normal distribution, \mathbf{I}_K is the K -dimensional identity matrix, \mathbf{I}_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 • Draw a distribution over words $\phi_k \sim \text{Dir}(\gamma \mathbf{1}_K)$

2. For each label $l \in \mathcal{L}$
 • Draw a label assignment coefficient $\eta_l | \mu, \sigma \sim \mathcal{N}(\mu \mathbf{1}_K, \sigma \mathbf{I}_K)$

3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}_K(\alpha' \mathbf{1}_K)$

4. For each document $d = 1, \dots, D$
 • Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha \mathbf{B})$

• For $i = 1, \dots, N_d$
 – Draw topic assignment $z_{n,d} | \theta_d \sim \text{Multinomial}(\theta_d)$

– Draw word $w_{n,d} | z_{n,d}, \phi_{l,K} \sim \text{Multinomial}(\phi_{l,z_{n,d}})$

• Set $y_{l,d} = 1$
 • For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r

– Draw $a_{l,d} | z_d, \eta_l, y_{\text{pa}(l),d} \sim \begin{cases} \mathcal{N}(z_d^T \eta_l, 1), & y_{\text{pa}(l),d} = 1 \\ \mathcal{N}(z_d^T \eta_l, 1) \mathbb{I}(a_{l,d} < 0), & y_{\text{pa}(l),d} = -1 \end{cases}$

– Apply label l to document d according to $a_{l,d}$

$y_{l,d} = a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

Here $z_d^T = [z_1, \dots, z_k, \dots, z_K]$ is the empirical topic distribution for document d , in which each entry is the percentage of the words in that document that come from topic k . $z_k = N_d^{-1} \sum_{n=1}^{N_d} \mathbb{I}(z_{n,d} = k)$.

The second half of step 4 is a substantial part of our contribution to the general class of supervised LDA models. Here, each document is labeled generatively using a hierarchy of conditionally dependent probit regressions [15]. For every label $l \in \mathcal{L}$, both the empirical topic distribution for document d whether or not its parent label was applied (i.e. $\mathbb{I}(y_{\text{pa}(l),d} = 1)$) are used to determine whether or not label l is to be applied. Specifically, if a parent label l' is applied, then label l is only applied to document d if its parent label $y_{\text{pa}(l),d}$ is applied. These parameters are specific to less complex constraints but can be applied to more complex constraints. The regression coefficients η_l are implemented in an important way of regularizing the LDA. In the setting where there are no negative labels, a Gaussian prior over these regression parameters with a negative mean implements a prior belief that missing labels are likely to be negative. We evaluated model performance for all three models with a range of values for μ ; the mean prior parameter for regression coefficients ($\mu \in [-3, -2.8, -2.6, \dots, 1]$). Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases labeling will consider a tree with a single root $r \in \mathcal{L}$. The effect of this is that label predictors deeper in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{\text{pa}(l),d} = 1, y_{\text{pa}(\text{pa}(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

as not being applied to document d , then none of its descendants are applied to document d .

The constraints imposed by an is-a label hierarchy are that if the l th label is

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical structures of the same [2]. Medical patient records and disease codes assigned to them for hospital admissions and hospital discharge summaries are another example. Both of these are examples of medical codebooks (ICD-9-CM) codes assigned [3].

In this work we show how to combine these two sources of information using a single model.

In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA “topics” (distributions over the elements of Σ). ϕ_k is a distribution over “words.” θ_k is a document-specific distribution over topics. β is a global distribution over topics. $D_{k,d}$ is a K -dimensional Dirichlet distribution. $N_K(\cdot)$ is the K -dimensional Normal distribution. I_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value of 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$
 - Draw a distribution over words $\phi_k \sim \text{Dir}(\gamma_1)$
2. For each label $l \in \mathcal{L}$
 - Draw a label application coefficient $\eta_l | \mu, \sigma \sim N(\mu I_K, \sigma I_K)$
3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}(\alpha' I_K)$
4. For each document $d = 1, \dots, D$
 - Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha I_K)$
 - For $a = 1, \dots, N_d$
 - Draw topic assignment $z_{n,a,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
 - Draw word $w_{n,a,d} | z_{n,a,d}, \phi_{k,z_{n,a,d}} \sim \text{Multinomial}(\phi_{k,z_{n,a,d}})$
 - Set $y_{d,l} = 1$
 - For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r
 - Draw $a_{l,d} | z_{d,l}, \eta_l, y_{p(a_l),d} \sim N(z_{d,l}^\top \eta_l, 1)$, $y_{p(a_l),d} = 1$
 - Draw $a_{l,d} | z_{d,l}, \eta_l, y_{p(a_l),d} \sim N(z_{d,l}^\top \eta_l, 1) \mathbb{I}(a_{l,d} < 0)$, $y_{p(a_l),d} = -1$
 - Apply label l to document d according to $a_{l,d}$
 - $y_{d,l} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, and the absence of a label cannot always be interpreted as a negative labeling.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a joint model, while leveraging the hierarchical structure of the labels. For the sake of simplicity, we focus on α -is hierarchies, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical information into the model. We hypothesize that the context of labels within the hierarchy provides valuable information about labeling.

We demonstrate our model on large, real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Our results show that a joint, hierarchical model outperforms a classification with unstructured labels as well as a disjoint model, where the topic modeling and the hierarchical classification are carried out independently of each other.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $w_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of the vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except root, has a parent $p(l) \in \mathcal{L}$ also in the set of labels. We will for exposition purposes assume that this label set has no “is-a” parent-child constraints (explained later), although this assumption can be relaxed at the cost of more computationally complex inference. Such a label hierarchy forms a multiply rooted tree. Without loss of generality we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{d,l} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{d,l}$ will be known, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

3 Inference

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. An alternative interpretation of the same results is that, if one is more sensitive to the performance gains that result from exploiting the structure of the labels, then one can, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. These are applied settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

4 Related Work

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. An alternative interpretation of the same results is that, if one is more sensitive to the performance gains that result from exploiting the structure of the labels, then one can, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. These are applied settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inference challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [1] The computational medicine center’s 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007/
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [4] H. Albert and S. Chiba. Bayesian analysis of binary and polytomous response data. *Journal of the American Statistical Association*, 88(422):669–693, 1993.
- [5] E. Birman-Deych, A. D. Wissner-Gross, Y. Yu, D. S. Nilsson, M. J. Radford, and B. F. Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [6] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [7] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [8] S. Chakrabarti, B. Dom, R. Agrawal, and P. Raghavan. Scalable feature selection, classification and signature generation for organizing large text databases into hierarchical topic taxonomies. *The VLDB Journal*, 7:163–178, August 1998. ISSN 1066-8888.
- [9] C. Chang and D. M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS309.
- [10] M. Crasmer, M. Dredze, K. Gauché, P.P. Talukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2009: Biomedical Text Mining and Computational Linguistics*, pages 129–136, 2009.
- [11] Susan Dumais and Hao Chen. Hierarchical classification of web content. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR ’00*, pages 256–263, New York, NY, USA, 2000. ACM.
- [12] R. Farkas and G. Székely. Automatic construction of rule-based ICD-9-CM coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [13] M. Farzanehpour, A. Sheikholeslami, and F. Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:78–84, 2010.
- [14] A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2nd ed. edition, 2004.
- [15] J. Goldstein, A. Arzumanyan, and O. Uzuner. Three approaches to automatic assignment of ICD-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [16] T. L. Griffith and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):S228–S235, 2004.
- [17] D. Koller and M. Sahami. Hierarchical classifying documents using very few words. Technical Report 1997-57, Stanford InfoLab, February 1997. Previous number = SIDL-WP-1997-0059.
- [18] S. Lacoste-Julien, F. Sha, and M. I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904, 2010.
- [19] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [20] L.Y. Lita, S. Yu, S. Niculescu, and J. Bi. Large scale diagnostic code classification for medical patient records. In *Proceedings of the 3rd International Joint Conference on Natural Language Processing (IJCNLP’08)*, 2008.
- [21] A. McCallum, N. Freitag, and K. Seymore. Building domain-specific search engines with machine learning techniques. In *Proc. AAAI-99 Spring Symposium: Intelligent Agents in Cyberspace*, 1999.
- [22] S. Pabstow, J. Barroso, and C. Chute. Automating the assignment of diagnosis codes to patient encounters using example-based and machine learning techniques. *Journal of the American Medical Informatics Association (JAMIA)*, 13(5):516–525, 2006.
- [23] D. Ramage, D. Hall, R. Nallapati, and C. D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 248–256, 2009.
- [24] B. RibeiroNeto, AHF Lacerda, and LRS Da Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and technology*, 52(3):391–401, 2001.
- [25] P. Ruch, J. Gobell, I. Thushri, and A. Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [26] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [27] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1973–1981, 2009.
- [28] C. Wang, D. Blei, and L. Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

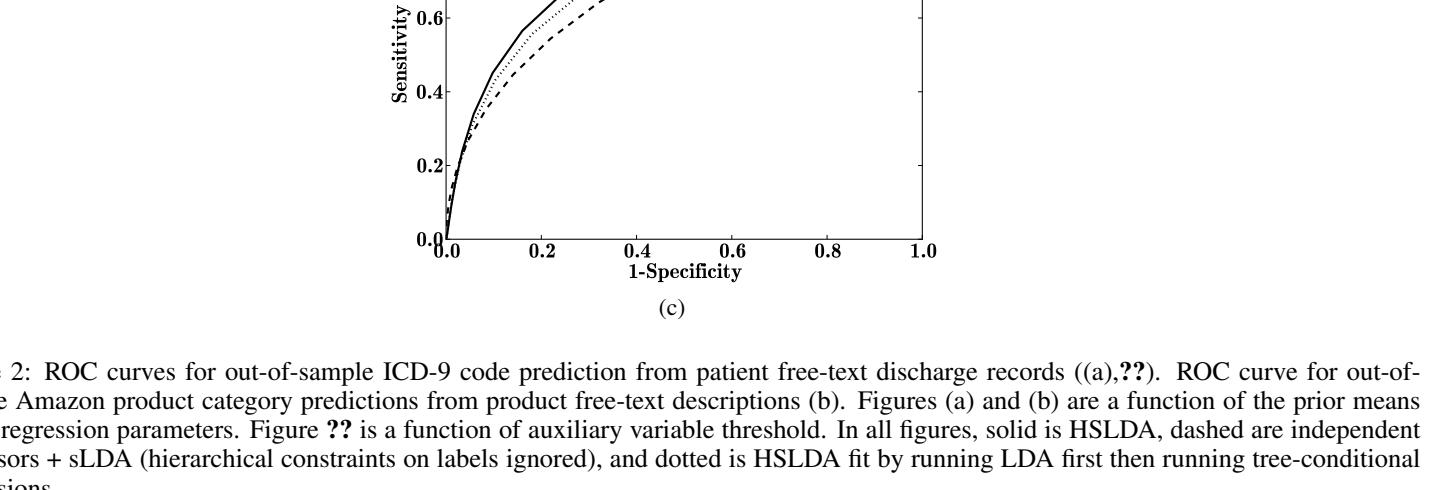


Figure 2: ROC curves for out-of-sample ICD-9 code prediction from patient free-text discharge records ((a),(b)). ROC curve for out-of-sample Amazon product category predictions from product free-text descriptions (b). Figures (a) and (b) are a function of the prior means of the regression parameters. Figure ?? is a function of auxiliary variable threshold. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.

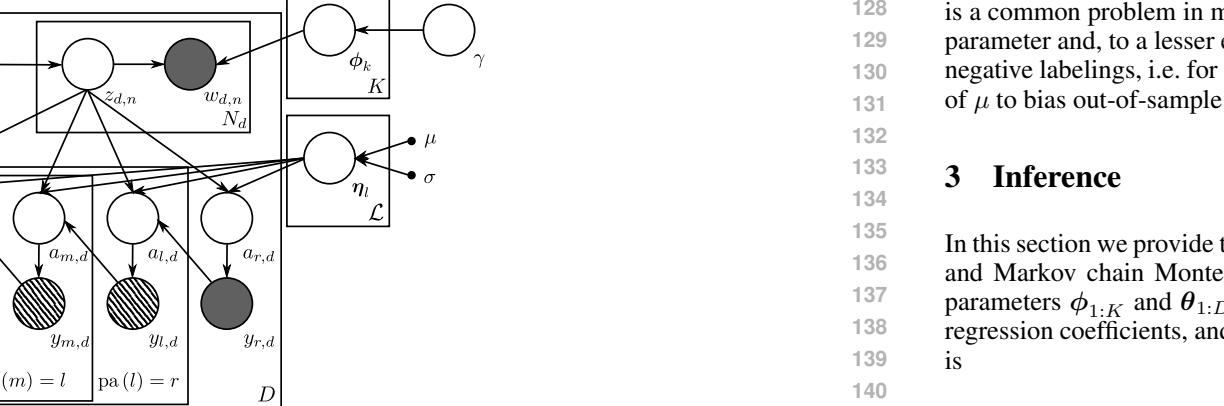


Figure 1: HSLDA graphical model

The constraints imposed by an is-a label hierarchy are that if the l th label is applied to document d , i.e., $y_{d,l} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{p(a_l),d} = 1, \dots, y_{p(a_1),d} = 1$. Conversely, if a label l' is marked as not applying to a document then no descendant of that label may be applied to the same. We assume that at least one label is applied to every document. This is illustrated in Figure 1 where the root label is always applied but only some of the descendant labelings are observed as having been applied (diagonal hashing indicates that potentially some of the plated variables are observed).

In HSLDA, documents are modeled using the LDA mixed-membership mixture model with global topic estimation. Label responses are generated using a conditional hierarchy of probit regressors. The HSLDA graphical model is given in Figure 1. In the model, K is the number of LDA “topics” (distributions over the elements of Σ). ϕ_k is a distribution over “words.” θ_k is a document-specific distribution over topics. β is a global distribution over topics. $D_{k,d}$ is a K -dimensional Dirichlet distribution. $N_K(\cdot)$ is the K -dimensional Normal distribution. I_d is the d -dimensional vector of all ones, and $\mathbb{I}(\cdot)$ is an indicator function that takes the value of 1 if its argument is true and 0 otherwise. The following procedure describes how to generate from the HSLDA generative model.

1. For each topic $k = 1, \dots, K$

- Draw a distribution over words $\phi_k \sim \text{Dir}(\gamma_1)$

2. For each label $l \in \mathcal{L}$

- Draw a label application coefficient $\eta_l | \mu, \sigma \sim N(\mu I_K, \sigma I_K)$

3. Draw the global topic proportions $\beta | \alpha' \sim \text{Dir}(\alpha' I_K)$

4. For each document $d = 1, \dots, D$

- Draw topic proportions $\theta_d | \beta, \alpha \sim \text{Dir}_K(\alpha I_K)$

• For $a = 1, \dots, N_d$

- Draw topic assignment $z_{n,a,d} | \theta_d \sim \text{Multinomial}(\theta_d)$
- Draw word $w_{n,a,d} | z_{n,a,d}, \phi_{k,z_{n,a,d}} \sim \text{Multinomial}(\phi_{k,z_{n,a,d}})$

• Set $y_{d,l} = 1$

• For each label l in a breadth first traversal of \mathcal{L} starting at the children of root r

- Draw $a_{l,d} | z_{d,l}, \eta_l, y_{p(a_l),d} \sim N(z_{d,l}^\top \eta_l, 1)$, $y_{p(a_l),d} = 1$
- Draw $a_{l,d} | z_{d,l}, \eta_l, y_{p(a_l),d} \sim N(z_{d,l}^\top \eta_l, 1) \mathbb{I}(a_{l,d} < 0)$, $y_{p(a_l),d} = -1$
- Apply label l to document d according to $a_{l,d}$

$y_{d,l} | a_{l,d} = \begin{cases} 1 & \text{if } a_{l,d} > 0 \\ -1 & \text{otherwise} \end{cases}$

5. Inference

5.1 Data and Pre-Processing

5.1.1 Discharge Summaries and ICD-9 Codes

In this section we provide the conditional distributions required to draw samples from the HSLDA posterior distribution using Gibbs sampling and Markov chain Monte Carlo. Note that like in collapsed Gibbs samplers for LDA [17], we have analytically marginalized out the parameters ϕ_k and $\theta_{l,D}$ in the following expressions. Let a be a set of auxiliary variables, w the set of all words, η the set of all regression coefficients, and $z_{d,\setminus d}$ the set z_d with element $z_{d,d}$ removed. The conditional posterior distribution of the latent topic indicators is

$$p(z_{n,d} = k | z_{d,\setminus d}, w, \mathbf{a}, \mathbf{w}, \eta, \alpha, \beta, \gamma) \propto \frac{c_{k,n-(a,d)}^{k-(n,a,d)+\gamma}}{\binom{c_{(k-1),n-(a,d)}}{c_{(k-1),n-(a,d)+V}}} \prod_{l \in \mathcal{L}, d} \exp \left\{ -\frac{(z_{d,l}^\top \eta_l)^2}{2} \right\} \quad (1)$$

where $c_{v,(n,d)}$ is the number of words of type v in document d assigned to topic k omitting the n th word of document d . The false positive rate α indicate to sum over the range of the replaced variable, i.e. $c_{w,n-(a,d)}^{k-(n,d)}$. Here Z_d is the set of labels which are observed for document d .

The conditional posterior distribution of the regression coefficients is given by

$$p(\eta_l | \mathbf{z}, \mathbf{a}, \sigma) = \mathcal{N}(\hat{\eta}_l, \hat{\Sigma}) \quad (2)$$

where

$$\mu_l = \frac{1}{\sigma} \left(\frac{\$$

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs, and patient records and diagnosis codes assigned to them for bookkeeping and insurance purposes. In this work we show how to combine these two sources of information using a single model that can automatically categorize new text documents automatically, suggest labels that might be inaccurate, compare improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable to other data as well; namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, and a label of a label cannot always be interpreted as a negative labeling.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a hierarchical manner. For the sake of simplicity, we focus on α -hierarchies, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (sLDA) [6] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical structure into the model. We hypothesize that the context of labels within the hierarchy provides valuable information about labeling constraints.

We demonstrate our model on large real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Our results show that a joint, hierarchical model outperforms a classification with unstructured labels as well as a disjoint model, where the topic modeling and the hierarchical classification are carried out independently of each other.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except for a particular label L also has a set of labels \mathcal{L}' that are its children. A label l is a parent of all its child labels. If l is a parent of l' , then l' is a child of l . This forms a multiply-rooted tree. Without loss of generality, we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an α -label hierarchy are that if the l th label is applied to document d , i.e., $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{l',\text{pa}(l),d} = 1, y_{l'',\text{pa}(l'),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

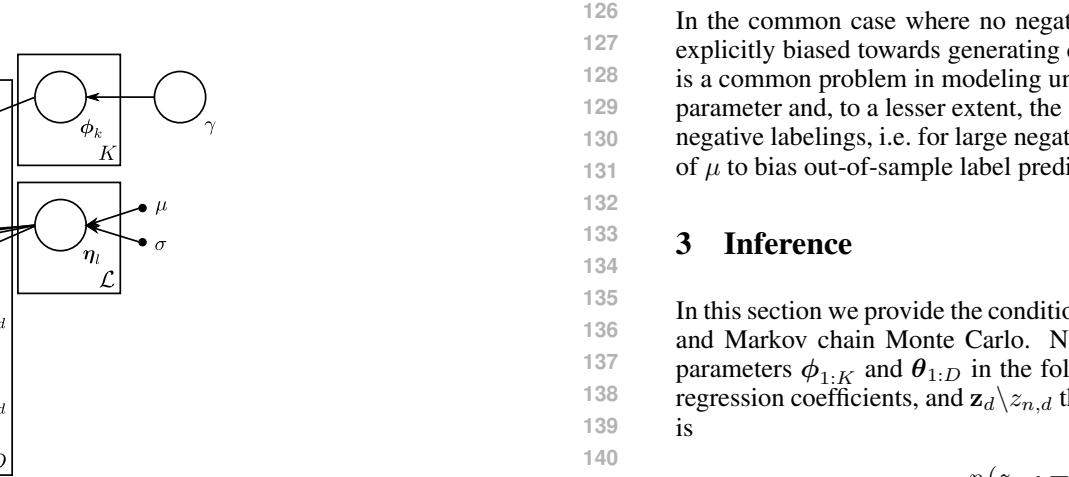


Figure 1: HSLDA graphical model

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unlabeled data. To see how this model can be biased in this way we draw the reader's attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We applied HSLDA to data from two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing

5.1.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize patient hospitalization course. The summaries typically contain a record of patient complaints, findings and diagnoses along with treatment and hospital course. For each hospitalization, trained medical coders review the information in the discharge summary and assign a series of diagnosis codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for "Pneumonia due to adenovirus" is a child of the code for "Viral pneumonia," where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [5], and sometimes make mistakes [13].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Methods range from manual rules to online learning [10, 15, 12]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Linear Regression, as well as simple keyword mappings, all with promising results [19, 24, 22, 20].

Our dataset was gathered from the clinical data warehouse of a large metropolitan hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² A fixed vocabulary was formed by taking the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Descriptions and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a multiply-rooted hierarchy and provides textual product descriptions for most products. Products can be discovered by users through free-text search and product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

As sensitivity and specificity can always be traded off we examined sensitivity for a range of values for two different parameters - the prior mean for the regression coefficients and the threshold for the auxiliary variables. The goal in this analysis was to evaluate the performance of these models subject to more or less stringent requirements for predicting positive labels. These two parameters have important related effects on the quality of the predictions.

In this experiment, we obtained Amazon.com product categorization data from the Stanford Network Analysis Platform (SNAP) dataset [3]. Product descriptions were obtained separately from the Amazon.com website directly. We limited our dataset to the collection of DVDs in the product catalog.

HSLDA employs a hierarchical Dirichlet prior over topic assignments (i.e., β is estimated from data rather than fixed a priori). This has been shown to improve the quality and stability of inferred topics [27]. Sampling β is done using the "direct assignment" method of Teh et al. [26]

$\beta | \mathbf{z}, \mathbf{Y}, \eta \sim \text{Dir}(m_{1,\cdot,1} + \alpha', m_{1,\cdot,2} + \alpha', \dots, m_{1,\cdot,K} + \alpha')$

Here $m_{d,k}$ are auxiliary variables that are required to sample the posterior distribution of β . Their conditional posterior distribution is sampled according to

$$p(m_{d,k} = m | \mathbf{z}, \mathbf{m}_{-(d,k)}, \beta) = \frac{\Gamma(\alpha\beta_k)}{\Gamma(\alpha\beta_k + c_{(k),d}^k)} s(c_{(k),d}, m)^{\alpha\beta_k} m^{c_{(k),d}}$$

where $s(n, m)$ represents stirling numbers of the first kind.

The hyperparameters α , α' , and γ are sampled using Metropolis-Hastings.

4 Related Work

In this work we extend supervised latent Dirichlet allocation (sLDA) [6] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [17] augmented with per-document "supervision", often taking the form of a single numerical or categorical label.

It has been demonstrated that sLDA can result in better task-specific document models and can also lead to good label prediction for out-of-sample data [6]. It has also been demonstrated that sLDA has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [6]. sLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we constrain HSLDA with hierarchical coupling in this model to avoid this posterior decoupling. The effect of this is that label predictors in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed were driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$'s yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

¹<http://www.cdc.gov/nchs/icd/icd9cm.htm>

²<http://www.nltk.org>

Figure 2: ROC curves for out-of-sample ICD-9 code prediction from patient free-text discharge records ((a),(c)). ROC curve for out-of-sample Amazon product category predictions from product free-text descriptions (b). Figures (a) and (b) are a function of the prior means of the regression parameters. Figure (c) is a function of auxiliary variable threshold. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLD fit by running LDA first then running tree-conditional regressions.

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. An alternative interpretation of the same results is that, if one is more sensitive to the performance gains that result from exploiting the structure of the labels, then one can, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. These are applied settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inferences challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] ICD-9-CM: international classification of diseases, 9th revision; clinical modification, 6th edition. Practice Management Information Corporation, Los Angeles, CA, 2006.
- [4] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [5] J. H. Albert and S. Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669–693, 1993.
- [6] E. Birman-Deych, A. D. Waterman, Y. Yan, D. S. Nilasena, M. J. Radford, and B. F. Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [7] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [8] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [9] S. Chakrabarti, B. Dom, T. Agrawal, and P. Raghavan. Scalable selection, classification and signature generation for organizing large text documents into hierarchical topic taxonomies. *The VLDB Journal*, 7:163–178, August 1998. ISSN 0949-8888.
- [10] J. Chu and D. M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi:10.1214/09-AOS780.
- [11] K. Crammer, M. Dredze, K. Gauché, P. Talukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2007: Biomedical, Translational, and Clinical Language Processing*, pages 129–136, 2007.
- [12] Susan Dumais and Hao Chen. Hierarchical classification of web content. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR '00*, pages 256–263, New York, NY, USA, 2000. ACM.
- [13] R. Farkas and G. Starović. Automatic construction of rule-based ICD-9-CM coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [14] M. Farzandipour, A. Sheikholeslami, and F. Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:78–84, 2010.
- [15] A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2nd ed. edition, 2004.
- [16] I. Goldstein, A. Arzamustyan, and O. Uzuner. Their approaches to automatic assignment of ICD-9-cm codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [17] T. L. Griffith and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [18] T. L. Griffith and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [19] T. L. Griffith and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [20] T. L. Griffith and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [21] T. L. Griffith and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [22] S. Lacoste-Julien, F. Sha, and M. I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904.
- [23] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [24] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [25] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [26] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [27] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [28] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [29] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [30] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [31] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [32] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [33] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [34] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [35] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [36] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [37] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [38] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [39] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [40] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [41] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [42] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [43] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [44] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [45] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [46] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [47] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [48] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [49] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [50] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [51] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [52] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [53] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [54] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [55] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [56] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [57] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Author(s)

Affiliation
Address
email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs, and patient records and diagnosis codes assigned to them for bookkeeping and insurance purposes. In this work we show how to combine these two sources of information using a single model that can automatically learn new text documents automatically, suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable to other data as well; namely, any unstructured representations of data that have been hierarchically categorized (e.g., image catalogs with bag-of-feature image representations).

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often times, each instance has only a small number of labels associated to it). There are no naturally occurring negative labeling in the data, and the absence of a label cannot always be interpreted as a negative labeling.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in parallel. For the sake of simplicity, we focus on *a-s* hierarchies, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (sLDA) [7] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical information into the model. We hypothesize that the context of labels within the hierarchy provides valuable information about labeling constraints.

We demonstrate our model on large real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Our results show that a joint, hierarchical model outperforms a classification with unstructured labels as a disjoint model, where the topic modeling and the hierarchical classification are carried out independently of each other.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d \in \Sigma^D$ be the set of N_d observations in document d . Let there be D documents and let the size of vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except for a particular label \mathcal{L} itself, has a parent label l' . A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label l is said to be a root if it has no parents. A label l is said to be a child of l' if $l \in \mathcal{L}_{l'}$. A label l is said to be a parent of l' if $l' \in \mathcal{L}_l$. A label l is said to be a sibling of l' if they share a common parent. A label l is said to be a leaf if it has no children. A label

Hierarchically Supervised Latent Dirichlet Allocation

Anonymous Authors)

Affiliation

Address

email

Abstract

We introduce hierarchically supervised latent Dirichlet allocation (HSLDA), a model for hierarchically and multiply labeled bag-of-word data. Examples of such data include web pages and their placement in directories, product descriptions and associated categories from product hierarchies, and free-text clinical records and their assigned diagnosis codes. Out-of-sample label prediction is the primary goal of this work, but improved lower-dimensional representations of the bag-of-word data are also of interest. We demonstrate HSLDA on large-scale data from clinical document labeling and retail product categorization tasks. We show that leveraging the structure from hierarchical labels improves out-of-sample label prediction substantially when compared to models that do not.

1 Introduction

There exist many sources of unstructured data that have been partially or completely categorized by human editors. In this paper we focus on unstructured text data that has been, at least in part, manually categorized. Examples include but are not limited to webpages and curated hierarchical directories of the same [2], product descriptions and catalogs, and patient records and diagnosis codes assigned to them for bookkeeping and insurance purposes. In this work we show how to combine these two sources of information using a single model that can automatically learn new text documents, automatically suggest labels that might be inaccurate, compute improved similarities between documents for information retrieval purposes, and more. The models and techniques that we develop in this paper are applicable to other data as well; namely, any unstructured representations of data that have been hierarchically classified (e.g., image catalogs with bag-of-feature image representations).

There are several challenges entailed in incorporating a hierarchy of labels into the model. Given a large set of potential labels (often thousands), each instance has only a small number of labels associated to it. There are no naturally occurring negative labeling in the data, and labeling of a label cannot always be interpreted as a negative labeling.

Our work operates within the framework of topic modeling. Our approach learns topic models of the underlying data and labeling strategies in a hierarchical manner. For the sake of simplicity, we focus on α -hierarchies, but the model can be applied to other hierarchies. We extend supervised latent Dirichlet allocation (sLDA) [6] to take advantage of hierarchical supervision. We propose an efficient way to incorporate hierarchical structure into the model. We hypothesize that the context of labels within the hierarchy provides valuable information about labeling constraints.

We demonstrate our model on large real-world datasets in the clinical and web retail domains. We observe that hierarchical information is valuable when incorporated into the learning and improves our primary goal of multi-label classification. Our results show that a joint, hierarchical model outperforms a classification with unstructured labels as well as a disjoint model, where the topic modeling and the hierarchical classification are carried out independently of each other.

The remainder of this paper is as follows. Section 2 introduces hierarchically supervised LDA (HSLDA), while Section 3 details a sampling approach to inference in HSLDA. Section 4 reviews related work, and Section 5 shows results from applying HSLDA to health care and web retail data.

2 Model

HSLDA is a model for hierarchically, multiply-labeled, bag-of-word data. We will refer to individual groups of bag-of-word data as documents. Let $w_{n,d} \in \Sigma$ be the n th observation in the d th document. Let $\mathbf{w}_d = \{w_{1,d}, \dots, w_{N_d,d}\}$ be the set of N_d observations in document d . Let there be D such documents and let the size of vocabulary be $V = |\Sigma|$. Let the set of labels be $\mathcal{L} = \{l_1, l_2, \dots, l_{|\mathcal{L}|}\}$. Each label $l \in \mathcal{L}$, except for a particular root label \mathcal{L} itself, has a set of children labels. A label l is a parent of all its child labels. If a label l is a child of a parent p , then l is a child of p if and only if l is a child of p . This assumption can be relaxed at the cost of more computationally complex inference. Such a label hierarchy forms a multiply-rooted tree. Without loss of generality, we will consider a tree with a single root $r \in \mathcal{L}$. Each document has a variable $y_{l,d} \in \{-1, 1\}$ for every label which indicates whether the label is applied to document d or not. In most cases $y_{l,d}$ will be unobserved, in some cases we will be able to fix its value because of constraints on the label hierarchy, and in the relatively minor remainder its value will be observed. In the applications we consider, only positive label applications are observed.

The constraints imposed by an α -label hierarchy are that if the l th label is applied to document d , i.e., $y_{l,d} = 1$, then all labels in the label hierarchy up to the root are also applied to document d , i.e., $y_{pa(l),d} = 1, y_{pa(pa(l)),d} = 1, \dots, y_{r,d} = 1$. Conversely, if a label l' is marked

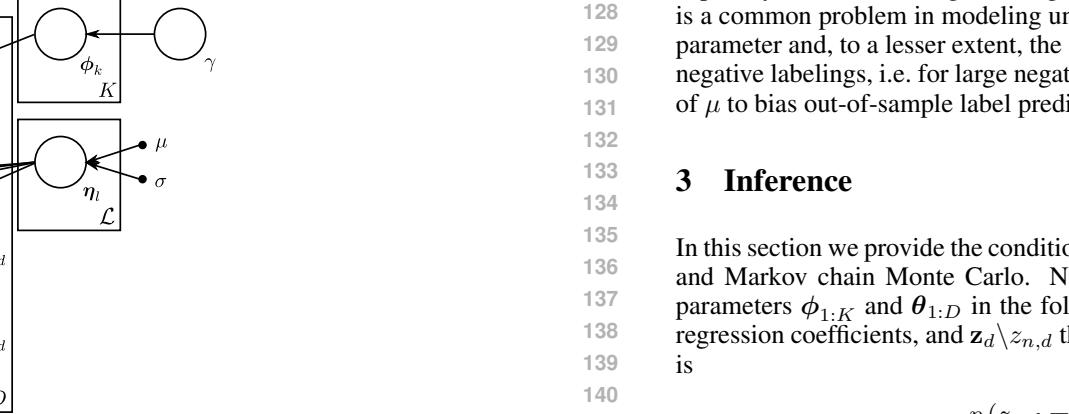


Figure 1: HSLDA graphical model

In the common case where no negative labels are observed (like the example applications we consider in Section 5), the model must be explicitly biased towards generating data that has negative labels in order to keep it from learning to assign all labels to all documents. This is a common problem in modeling unlabeled data. To see how this model can be biased in this way we draw the reader's attention to the μ parameter and, to a lesser extent, the σ parameter above. Because z_d is always positive, setting μ to a negative value results in a bias towards negative labelings, i.e. for large negative values of μ , all labels become a priori more likely to be negative ($y_{l,d} = -1$). We explore the ability of μ to bias out-of-sample label prediction performance in Section 5.

5 Experiments

We applied HSLDA to data from two domains: predicting medical diagnosis codes from hospital discharge summaries and predicting product categories from Amazon.com product descriptions.

5.1 Data and Pre-Processing

5.1.1 Discharge Summaries and ICD-9 Codes

Discharge summaries are authored by clinicians to summarize patient hospitalization course. The summaries typically contain a record of patient complaints, findings and diagnoses, along with treatment and hospital course. For each hospitalization, trained medical coders review the information in the discharge summary and assign a series of diagnosis codes. Coding follows the ICD-9-CM controlled terminology, an international diagnostic classification for epidemiological, health management, and clinical purposes.¹ The ICD-9 codes are organized in a rooted-tree structure, with each edge representing an is-a relationship between parent and child, such that the parent diagnosis subsumes the child diagnosis. For example, the code for "Pneumonia due to adenovirus" is a child of the code for "Viral pneumonia," where the former is a type of the latter. It is worth noting that the coding can be noisy. Human coders sometimes disagree [1], tend to be more specific than sensitive in their assignments [5], and sometimes make mistakes [13].

The task of automatic ICD-9 coding has been investigated in the clinical domain. Methods range from manual rules to online learning [10, 15, 12]. Other work had leveraged larger datasets and experimented with K-nearest neighbor, Naive Bayes, support vector machines, Bayesian Linear Regression, as well as simple keyword mappings, all with promising results [19, 24, 22, 20].

Our dataset was gathered from the clinical data warehouse of a large metropolitan hospital. It consists of 6,000 discharge summaries and their associated ICD-9 codes (7,298 distinct codes overall), representing all the discharges from the hospital in 2009. Summaries have 8.39 associated ICD-9 codes on average (std dev=5.01) and contain an average of 536.57 terms after preprocessing (std dev=300.29). We split our dataset into 5,000 discharge summaries for training and 1,000 for testing.

The text of the discharge summaries was tokenized with NLTK.² A fixed vocabulary was formed by taking the top 10,000 tokens with highest document frequency (exclusive of names, places and other identifying numbers). The study was approved by the Institutional Review Board and follows HIPAA (Health Insurance Portability and Accountability Act) privacy guidelines.

5.1.2 Product Descriptions and Categorizations

Amazon.com, an online retail store, organizes its catalog of products in a multiply-rooted hierarchy and provides textual product descriptions for most products. Products can be discovered by users through free-text search and product category exploration. Top-level product categories are displayed on the front page of the website and lower level categories can be discovered by choosing one of the top-level categories. Products can exist in multiple locations in the hierarchy.

As sensitivity and specificity can always be traded off we examined sensitivity for a range of values for two different parameters - the prior mean for the regression coefficients and the threshold for the auxiliary variables. The goal in this analysis was to evaluate the performance of these models subject to more or less stringent requirements for predicting positive labels. These two parameters have important related effects.

Our dataset contains 15,130 product descriptions for training and 1,000 for testing. The product descriptions are shorter than the discharge summaries (91.89 terms on average, std dev=53.08). Overall, there are 2,691 unique codes. Products are assigned on average 9.01 codes (std dev=4.91). The vocabulary consists of the most frequent 30,000 words omitting stopwords.

5.2 Comparison Models

We evaluated HSLDA along with three closely related models against these two datasets. The comparison models included sLDA with independent regressors (hierarchical constraints on labels ignored), HSLDA fit by first performing LDA then fitting tree-conditional regressions. These models were chosen to highlight several aspects of HSLDA including performance in the absence of hierarchical constraints, the effect of the combined inference, and regression performance attributable solely to the hierarchical constraints.

sLDA with independent regressors is the most salient comparison model for our work. The distinguishing factor between HSLDA and sLDA is the addition structure imposed on the label space, a distinction that we hypothesized would result in a difference in predictive performance.

For each model type, separate models were fit for each value of the prior mean of the regression coefficients. This is a proper Bayesian sensitivity analysis. In contrast, to evaluate predictive performance as a function of the auxiliary variable threshold, a single model was fit for each model type and prediction was evaluated based on predictive samples drawn subject to different auxiliary variable thresholds. These methods are significantly different since the prior mean is varied prior to inference and the auxiliary variable threshold is varied following inference. However, as intended, they both highlight model performance under more or less stringent requirements for predicting positive labels.

Figure 2 shows the predictive performance of HSLDA relative to the three comparison models as a function of the prior mean on regression coefficients as a receiver operating characteristic (ROC) curve. For low values of the auxiliary variable threshold, a single model is fit for each model type and prediction was evaluated based on predictive samples drawn subject to different auxiliary variable thresholds. These methods are significantly different since the prior mean is varied prior to inference and the auxiliary variable threshold is varied following inference. However, as intended, they both highlight model performance under more or less stringent requirements for predicting positive labels.

In this work we extend supervised latent Dirichlet allocation (sLDA) [6] to take advantage of hierarchical supervision. sLDA is latent Dirichlet allocation (LDA) [7] augmented with per-document "supervision", often taking the form of a single numerical or categorical label. It has been demonstrated that LDA with such supervision can result in better task-specific document models and can also lead to good label prediction for out-of-sample data [6]. It has also been demonstrated that sLDA has been shown to outperform both LASSO (L1 regularized least squares regression) and LDA followed by least squares regression [6]. sLDA can be applied to data of the type we consider in this paper; however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we constrain HSLDA with hierarchical constraints on labels to ensure that the signal provided by such supervision can result in better task-specific document models and can also lead to good label prediction for out-of-sample data [6]. This comparison model examines not the structure of the label space, but the benefit of combined inference over both the documents and the label space.

For most models, particularly sLDA, the prior was paid little attention to the prior parameters for the regression coefficients. These parameters are to be applied uniformly to all documents. In the setting where there are no negative labels, a Gaussian prior over the regression coefficients can be used to regularize the model. The prior is applied to the entire set of labels. The regression coefficients η_l are assumed to be independent a priori; however, the hierarchical coupling in this model makes it difficult to implement a posterior distribution.

The effect of this is that label predictors in the label hierarchy are able to focus on finding specific, conditional labeling features. We believe this to be a significant source of the empirical label prediction improvement we observe experimentally. We test this hypothesis in Section 5.

Note that the choice of variables $a_{l,d}$ and how they are distributed are driven at least in part by posterior inference efficiency considerations. In particular, choosing probit-style auxiliary variable distributions for the $a_{l,d}$ yields conditional posterior distributions for both the auxiliary variables (3) and the regression coefficients (2) which are analytic. This simplifies posterior inference substantially.

¹<http://www.cdc.gov/nchs/icd/icd9cm.htm>

²<http://www.nltk.org>

1

The results in Figures 2(a) and 2(b) suggest that in most cases it is better to do full joint estimation of HSLDA. An alternative interpretation of the same results is that, if one is more sensitive to the performance gains that result from exploiting the structure of the labels, then one can, in an engineering sense, get nearly as much gain in label prediction performance by first fitting LDA and then fitting a hierarchical probit regression. These are applied settings in which this could be advantageous.

Extensions to this work include unbounded topic cardinality variants and relaxations to different kinds of label structure. Unbounded topic cardinality variants pose interesting inferences challenges. Utilizing different kinds of label structure is possible within this framework, but requires relaxing some of the simplifications we made in this paper for expositional purposes.

References

- [1] The computational medicine center's 2007 medical natural language processing challenge. http://www.computationalmedicine.org/challenge/previous_2007.html
- [2] DMOZ open directory project. <http://www.dmoz.org/>, 2002.
- [3] Stanford network analysis platform. <http://snap.stanford.edu/>, 2004.
- [4] H. Albert and S. Chiba. Bayesian analysis of binary and polytomous response data. *Journal of the American Statistical Association*, 88(422):669–693, 1993.
- [5] E. Birman-Deych, A. D. Wimmer, Y. Yan, D. S. Nilsson, M. J. Radford, and B. F. Gage. Accuracy of ICD-9-CM codes for identifying cardiovascular and stroke risk factors. *Medical Care*, 43(5):480–5, 2005.
- [6] D. Blei and J. McAuliffe. Supervised topic models. *Advances in Neural Information Processing*, 20:121–128, 2008.
- [7] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003. ISSN 1532-4435.
- [8] S. Chakrabarti, B. Dom, R. Agrawal, and P. Raghavan. Scalable feature selection, classification and signature generation for organizing large text databases into hierarchical topic taxonomies. *The VLDB Journal*, 7:163–178, August 1998. ISSN 1066-8888.
- [9] J. Chang and D. M. Blei. Hierarchical relational models for document networks. *Annals of Applied Statistics*, 4:124–150, 2010. doi: 10.1214/09-AOAS399.
- [10] M. Crasmer, M. Dredze, K. Gauchard, P.P. Talukdar, and S. Carroll. Automatic code assignment to medical text. *Proceedings of the Workshop on BioNLP 2009: BioText Mining and Computational Linguistics*, pages 129–136, 2009.
- [11] Susan Dumais and Hao Chen. Hierarchical classification of web content. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR '00*, pages 256–263, New York, NY, USA, 2000. ACM.
- [12] R. Farkas and G. Székely. Automatic construction of rule-based ICD-9-CM coding systems. *BMC bioinformatics*, 9(Suppl 3):S10, 2008.
- [13] M. Farzanehpour, A. Sheikholeslami, and F. Sadoughi. Effective factors on accuracy of principal diagnosis coding based on international classification of diseases, the 10th revision. *International Journal of Information Management*, 30:78–94, 2010.
- [14] A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2nd ed. edition, 2004.
- [15] J. Goldstein, A. Arzumanyan, and O. Ozmen. Three approaches to automatic assignment of ICD-9-CM codes to radiology reports. *AMIA Annual Symposium Proceedings*, 2007:279, 2007.
- [16] T. L. Griffith and M. Steyvers. Finding scientific topics. *PNAS*, 101(suppl. 1):5228–5235, 2004.
- [17] D. Koller and M. Sahami. Hierarchical classifying documents using very few words. Technical Report 1997-57, Stanford InfoLab, February 1997. Previous number = SIDL-WP-1997-0059.
- [18] S. Lacoste-Julien, F. Sha, and M. I. Jordan. DiscLDA: Discriminative learning for dimensionality reduction and classification. In *Neural Information Processing Systems*, pages 897–904, 2011.
- [19] L. Larkey and B. Croft. Automatic assignment of ICD9 codes to discharge summaries. Technical report, University of Massachusetts, 1995.
- [20] L.Y. Lita, S. Yu, S. Niculescu, and J. Bi. Large scale diagnostic code classification for medical patient records. In *Proceedings of the 3rd International Joint Conference on Natural Language Processing (IJCNLP'08)*, 2008.
- [21] J. McNamee, N. Nigam, and R. Reilly. Building domain-specific search engines with machine learning techniques. In *Proc. AAAI-99*.
- [22] S. Pankonen, J. Barroso, and C. Chute. Automating the assignment of diagnosis codes to patient encounters using example-based and machine learning techniques. *Journal of the American Medical Informatics Association (JAMIA)*, 13(5):516–525, 2006.
- [23] D. Ramage, D. Hall, R. Nallapati, and C. D. Manning. Labeled LDA: a supervised topic model for credit attribution in multi-label corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 248–256, 2009.
- [24] B. RibeiroNeto, AHF Lacerda, and LRS De Lima. An experimental study in automatically categorizing medical documents. *Journal of the American society for information science and technology*, 52(2):391–401, 2001.
- [25] P. Ruch, J. Gobell, I. Thushri, and A. Geissbuhler. From episodes of care to diagnosis codes: automatic text categorization for medico-economic encoding. *AMIA Annual Symposium Proceedings*, 2008:636, 2008.
- [26] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [27] Hanna Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why priors matter. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. Williams, and A. Culotta, editors, *Advances in Information Processing*, 22, pages 1973–1981, 2009.
- [28] C. Wang, D. Blei, and L. Fei-Fei. Simultaneous image classification and annotation. In *CVPR*, 2009.

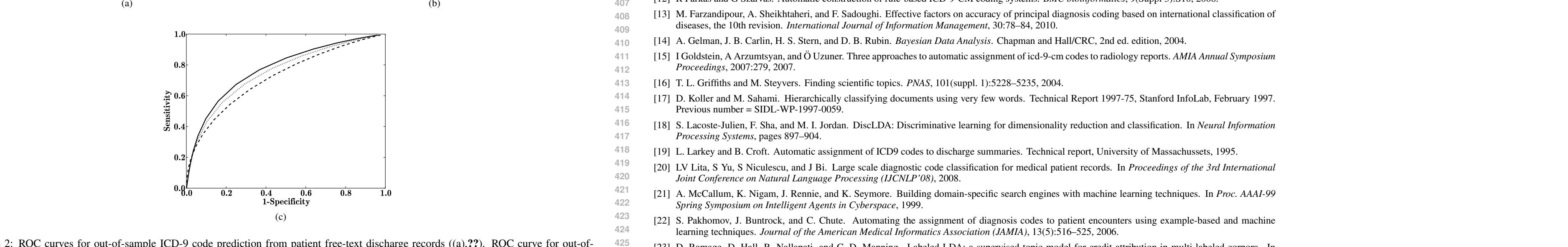


Figure 2: ROC curves for out-of-sample ICD-9 code prediction from patient free-text discharge records ((a)-(c)). ROC curve for out-of-sample Amazon product category predictions from product free-text descriptions (b). Figures (a) and (b) are a function of the prior means of the regression parameters. Figure (c) is a function of auxiliary variable threshold. In all figures, solid is HSLDA, dashed are independent regressors + sLDA (hierarchical constraints on labels ignored), and dotted is HSLDA fit by running LDA first then running tree-conditional regressions.

6 Discussion

The SLDA model family, of which HSLDA is a member, can be understood in two different ways. One way is to see it as a family of topic models that improve on the pure model by including the hierarchical structure imposed by LDA. Combined inference is the second way.

In this paper, however, doing so requires ignoring the hierarchical dependencies amongst the labels. In Section 5 we constrain HSLDA with hierarchical constraints on labels to ensure that the signal provided by such supervision can result in better task-specific document models and can also lead to good label prediction for out-of-sample data [6]. This comparison model examines not the structure of the label space, but the benefit of combined inference over both the documents and the label space.

For most models, particularly sLDA, the prior was paid little attention to the prior parameters for the regression coefficients. These parameters are to be applied uniformly to all documents. In the setting where there are no negative labels, a Gaussian prior over the regression coefficients can be used to regularize the model. The prior is applied to the entire set of labels. The regression coefficients η_l are assumed to be independent a priori; however, the hierarchical coupling in this model makes it difficult to implement a posterior distribution.

The effect of this is that label predictors in the label hierarchy

