

Analysing Narratives: Automatic Descriptive Feature Extraction Through Latent Entity modelling

Adam Slack

Declaration

I declare that this work is my own and not someone elses. May god beat down upon me with hellfire should this declaration be falsly claimed.

Abstract

Latent Entities are approximations of entities found within a corpus allowing for similar entities to be grouped into classes. Through the use of a document segmentation method processing entities and associated terms, LDA can be applied with the intention of extracting topics from each entity segmentation.

Contents

Declaration	1
Abstract	2
1 Introduction	6
1.1 Introduction	6
1.2 Aims and Objectives	7
1.3 Motivation	8
2 Literature Review	9
2.1 Introduction	9
2.2 Part-Of-Speech Tagging	10
2.2.1 POS Tagging Techniques and Tools	10
2.3 Named Entity Recognition	11
2.3.1 NER Techniques	12
2.3.2 NER Tools	12
2.4 Topic Models	13
2.4.1 Latent Semantic Analysis	14
2.4.2 Probabilistic Latent Semantic Indexing	14
2.4.3 Latent Dirichlet Allocation	15
2.5 Evaluating Topic Models	17
2.5.1 Perplexity	17
2.5.2 Interpretability and Coherence	18
2.5.3 Comparing Models	20
2.6 Entity Modeling	21
2.6.1 Entity Linking	21
2.6.2 Entity Topic Models	22
3 New Ideas	23
3.1 Introduction	23
3.2 Gaussian Entity-Term Matrix	23
3.3 Latent Entities	24
3.4 Entity Topic Model	25
3.4.1 Evaluating ETMs	26
3.5 Analysis Pipeline	26
3.6 Topic Model Explorer	26
4 Implementation	28

4.1	Introduction	28
4.2	Architecture Overview	28
4.3	NLP Toolkit	29
4.4	Corpus Pre-Processing	29
4.4.1	HMM POS Tagger	30
4.5	GETM Derivation	30
4.6	Entity-Topic Derivation	30
4.7	Latent Entity Clustering	30
4.8	Latent Entity Explorer	30
5	Results and Evaluation	31
5.1	Introduction	31
5.2	Gaussian Entity-Term Matrix	31
5.2.1	Experiment Setup	31
5.2.2	Results	31
5.3	Entity Topic-Models	31
5.3.1	Experiment Setup	31
5.3.2	Results	31
5.4	Latent Entities	31
6	Conclusion	32
6.1	Introduction	32

List of Figures

1	LDA Plate Notation	15
2	(Left) Bipartite topic conceptual similarity graph. (Middle, Right) Document relationships formed by LSA and LDA topic modelling methods	20
3	Example entity-term associations between three entities	25
4	Example entity-Topic distributions for entities in Figure 3 and the associated Latent Entity.	25
5	Analysis Pipeline for the Extraction of Latent Entity Information.	26
6	Architecture overview of the Latent Entity Model Analysis and Exploration System.	29

Chapter 1

Introduction

1.1 Introduction

The Web; A seemingly ever-expanding resource, with data being generated and information published at an accelerating rate [1]. As more gain access to the internet, the rate at which new information is made available will only increase. Whether the origins of data be, social media, news outlets, or e-commerce reviews, much of the resources on the web exist in a natural language format. From this data there exists underlying information that can be extracted and utilised in decision making process. Given the amount of processing required to consume data on this scale, it is necessary to ensure that computational methods exist that can accommodate for individual or business needs to understand data. Many methods for Methods that allow the processing of data to extract surface level information exist, however there is room for further understanding. By relating distinct pieces of information or subsets of data, it is possible to frame or contextualise a solution to a problem within a wider setting. This paper aims to provide a novel method capable of providing succinct summaries of documents in terms of entity topics.

Information Retrieval (IR), Machine Learning (ML), and Natural Language Processing (NLP) when considered in conjunction with each other, concern themselves with the extracting of information from data in a natural language texts. Relating textual data yields information valuable to many different entities, including businesses when marketing products, individuals when choosing what to read, and even researchers considering the current state of a research area. For a business, relating entities extracted from texts, can help identify target groups to aim products at. For an individual, relating entities can assist in decisions on what to read, buy,

and watch based on similarities between things that they do and don't like. For research, the identification of common themes, or prominent authors can be achieved through the relating of entities.

A Latent Entity in this investigation is defined as an abstract representation of some entity that can be used to describe one or more concrete entities. Thus the term Latent Entity Modelling is defined as a transformation from a corpus of text to a collection of Latent Entities describing a corpus. Latent Entities represent a layer of abstraction from a corpus, wherein it is possible to describe the corpus as a whole in terms of the Latent Entities.

Entity models build upon the concept of topic modelling, particularly the kinds of topics derived through Latent Dirichlet Allocation. A Topic in this sense is a probability distribution of words, such that the degree of membership for each word in a topic indicates the probability of that word being an indicator of that topic. By performing clustering on entities expressed in terms of entity topic models, this paper investigates if it is possible to

1.2 Aims and Objectives

The aims of this project are two-fold. The first aim is to utilise and extend natural language processing methods to model and extract latent information describing entities within narratives. The second aim is to develop a system within which corpora can be easily and effectively analysed and explored when methods developed for aim one are applied.

Natural Language Processing is a large field with many analytical methods available for modelling purposes, and as such the scope of the first aim is required to be well defined. Objectives for the first aim include:

- To survey the fields of Natural Language Processing and Topic Modelling to identify methods valuable for the task of Entity-Topic Modelling.
- To Define a Entity-Topic Modelling method capable of deriving interpret-able topics which describe entities within narratives.

The second aim involves the development of a software environment, the objectives of which include:

- To develop a front-end user interface for the visual exploration of entity-topic models.

- To implement back-end infrastructure allowing for dynamic visualisations in front-end user interfaces.
- To create a software pipeline that extracts entity-topic data from a provided corpus.

1.3 Motivation

Motivations for this project include the need for more effective methods and tools capable of analysing data in natural language form. Additionally the development of tools utilising these methods can assist in providing others with means of understanding more about themselves, their habits, and their preferences.

Chapter 2

Literature Review

2.1 Introduction

Natural Language Processing (NLP) is a vast field of study within Computer Science and Computational Intelligence, the focus of which is the development of intelligent systems capable of handling data in the form of natural language. The task of extracting Latent Entity Models is one that will touch upon many areas of NLP, including Part-of-Speech (POS) tagging, Named Entity Recognition (NER), and Latent Topic Modelling. Given that the focus of this paper is the extraction of descriptive information of narratives, it is necessary to consider the field of Information Retrieval and existing methods of descriptive summarisation as well.

It is possible to divide NLP into two large schools of thought; One involves the processing of natural language such that machine usable resources are available, the other is concerned with the application of information resulting from processed natural language. Whilst many NLP tasks require methods for parsing, understanding, or synthesising spoken words, this paper is only considering written texts. As such, the spoken language aspects of NLP will be overlooked.

Parsing written texts begins with tokenization. This is the division of a single string of characters in to strings representing sentences or words. Often text is tokenized into sentences, and each sentence is then tokenized into words. Once a text is expressed with basic sentence and word structure it is possible to apply additional processing steps. POS tagging is the application of tags to sequences of words, each tag and resulting sequence represents the grammatical structure of a sequence of words. NER is the extraction of entities from natural language, entities are tagged depending on their type. NER is performed after tokenization, however it doesn't necessarily rely on text being POS tagged first.

There are a range of NLP tasks involving the applications of natural language parsing. From Customer relation Chatbots, to Document Summarisation and Relation. Whilst many applications involve the use of parsed text features, many also orient themselves around additional features. For example, tools for relating and summarising distinct documents may utilise topic models and topic modelling techniques.

2.2 Part-Of-Speech Tagging

POS Tagging, is a useful task for many NLP problems. It forms a solid platform from which many investigations can be launched. The quality of a POS tagger can make or break a study. There is a range of POS tagging methods to choose from, many the highest performing taggers employ Maximum Entropy (MaxEnt) models or Hidden Markov Models (HMM). However Rule and transformation based taggers often suffice [2–5]. There even exists hybrid models that use probabilistic or stochastic methods in conjunction with rule sets. UCREL CLAWS tagger is an example hybrid model. [6] As long as the POS tagging tool utilised performs comparably to those in literature, labouring over the type of POS tagger that is used will not overly affect the results of this study. It might be more useful to consider the subjective merits of any POS tagging libraries that already exist. The role POS tagging plays in this study is to produce a corpus of words that can be filtered by type. When building a model of words associated with entities and deriving topics that describe entities, words such as ‘The’, ‘To’ and ‘where’ won’t provide much information about entities within a corpus. Removing them may see an improvement in the quality of any derived models.

2.2.1 POS Tagging Techniques and Tools

The Natural Language Toolkit (NLTK) by default utilises a Maximum Entropy POS Tagger using the Penn-Treebank tagset. The main benefit to this POS tagger is its ease of use and accessibility. The performance of the MaxEnt tagger used in NLTK was reported score an accuracy of 96.64% on all words, and 85.56% on unknown words. [7] However, in similar implementations, certain tags had error rates of 100% , this was likely a result of the absence of certain features necessary to correctly tag specific parts of speech like ‘TO’ which occurred

14,748 times in a study and was never correctly classified [8].

By comparison, the CLAWS tagger achieved a reported 96% accuracy, though accuracy fell to 82% when text was not preprocessed to filter spelling variants and shakespearean english words. The challenges with POS tagging is the variance of sentence structure, as well as ambiguity that can occur within the english language. Fictional narratives like that of shakespeare frequently contribute to the open classes of words (Adjectives, Nouns, etc), meaning that there is an increased likelihood of unknown or spelling variant words. Whilst modern texts are not likely to be as creative with the english language as shakespeare, it would be naive to assume the CLAWS tagger would achieve the upper bounds for its accuracy. Estimating the performance of the CLAWS tagger on the corpora being used in this investigation makes the performance of both the CLAWS tagger and NLTK’s MaxEnt tagger roughly comparable.

In much of the literature, POS tagging techniques were evaluated using corpora oriented around a specific subject or domain, like news, or personal tweets. Using corpora limited to specific domains means that any given POS tagging method might only perform as reported in literature if the same or similar corpus is used. This project will benefit from a POS tagger that performs well on general purpose text, meaning that considering many of the tools used in literature opens a potential point of failure. As POS tagging is one of the first steps performed when processing a corpus, errors could be introduced into the system early on should the POS taggers used in literature proved to be inadequate on the chosen corpus for this project. Additionally, the effect of literature tending to focus on domain limited corpora means that studies on the quality of general purpose cross-domain POS taggers are somewhat unreported.

2.3 Named Entity Recognition

The task of extracting a list of entities from text is a similar task to that of POS tagging. It involves parsing natural language and applying entity labels to entity words. A recent survey of the field summarised that Named Entity Recognition (NER) - the task of identifying entities and labelling them as ‘Person’, ‘Organisation’, or ‘Location’ - as being essential to many tasks of computational linguistics [9]. Similarly to POS tagging, ambiguity prevents NER from having

a simple solution. It is not always clear which entity some text may be referring to, especially in situations where an entity is not referred to directly, or when entities are named in unusual ways. [10] Its application in this study is to provide a set of entities from which topic models can be derived from and applied too.

2.3.1 NER Techniques

NER can be carried out using a range of different statistical methods. HHMs can be used to classify entities as either a name (person, organisation, or location), time, or numerical quantity. HMM-based Chunk Taggers have a reported accuracy ranging from 87% to 94% depending on the size of the training set [11].

Similarly to POS Tagging, statistical methods based on Maximum Entropy can be applied with similar levels of accuracy as other methods [12, 13]. Hybrid approach to NER have been investigated, by utilising HMM, MaxEnt and transformation-based learning, error rates were reduced by as much as 15% when used with the english language [14].

CRFs are a form of statistical model that is particularly useful for applying labels to sequence data, when applied to POS tagging or NER, CRF systems can attain error rates as low as 5.55% and 15.96% respectively [15, 16].

It is worth noting that the drawbacks that applied to publications regarding POS tagging, also apply to studies on NER. Notably, NER methods typically only concern themselves with labelling words in text, and provide no means of extracting additional information about a recognised named entity. Drawbacks common though many of the techniques revolve around the resolution of which entities are actually the same. Entities within books can be referred to directly, or indirectly, and even be addressed with different names.

2.3.2 NER Tools

Many free NER tools exist on the web, including Stanford NER, Illinois NER, OpenCalais NER, and Alias-i LingPipe. In a comparison between the relative performances of these tools to classify entity types (Person, Location, Organisation) the Stanford NER achieved the second

highest precision and the highest recall rates. [17]. In separate comparison of NER tagging tools, more mixed results were received for the Stanford NER tool. It performed comparably to generalised NER tools found in the NLTK and Apache OpenNLP toolkits. On specialised datasets, specialised NER taggers tuned to the task at hand predictably outperformed an untuned implementation of the Stanford NER.

Focusing on the Stanford NER [18] an NER tool that utilises conditional random fields (CRF). A Java implementation, the tool exists as part of a larger CoreNLP toolkit created at Stanford University. The self-reported performance of the tool ranges from 92.15% to 85.6% and 92.39% to 85.53% for precision and recall respectively. The upper bounds for possible precision and recall relied on additional processing for handling specific features of text. [19]. Additionally, these results are from tests occurring in 2006. Recent advances in CRFs have been utilised in subsequent versions of the Stanford NER.

2.4 Topic Models

Topic models are a way of representing a corpus of text in terms of latent or abstract topics. A topic is defined as the degrees of membership each term in a collection has to a topic. In methods like Latent Dirichlet Allocation, a corpus is expressed as a probability distribution of topics, which in turn are expressed as a probability distribution of words. Describing a specific element of a corpus in terms of topic models requires matching terms in the text to terms in the various topics, after the derivation of a set of models.

Deriving topics from a corpus is commonly done through statistical techniques such as Latent Dirichlet Allocation (LDA), Hierarchical LDA, Latent Semantic Indexing (LSI), or Probabilistic LSI (pLSI). However, many existing methods are in fact built upon the ideas used in LSA, for example, the valuable LDA method arose out of the shortcomings of the comparatively simple LSI [20].

2.4.1 Latent Semantic Analysis

Sometimes referred to as LSI, Latent Semantic Analysis (LSA) is a technique that provides vector representations of documents in a corpus. These vector representations allow for quantitative comparisons of different documents. The implementation of the technique is oriented around Single Value Decomposition (SVD), and attempts to model terms in a document as being averages of the document passages which in which it occurs [21]. Additionally, for NLP tasks It is possible to view LSA as an extension of the term frequency-inverse document frequency (tf-idf) method, this is due to the method producing a subset of the reduction carried out by tf-idf such that term occurrences with the most variation between documents are retained [22]. The obvious shortcoming of LSA is focused purely on documents at a word level. It bears no notion of topics or themes through which documents can be related. It is possible to apply additional post processing on the output of LDA. Additionally, for this study the method is totally inadequate for the analysis of entities within the documents. The possible use of the method with regards to the automatic extraction of Latent Entity Models could relate to the extraction of entity-term associations. It has been reported that the method performs well with highly dimensional data, and comparably better than standard vector space models like tf-idf. [23] As such, LSA would likely be useful as nothing more Principal Component Analysis style reduction.

2.4.2 Probabilistic Latent Semantic Indexing

Probabilistic LSA (pLSA) is derived as an effort to mathematically formalise the the ideas patented by Deerwester et. al. The formal statistical nature of the method means that the method can be used in conjunction with other models in a more predictable manner. The probabilistic approach changes the meaning of output produced by pLSA, resulting in models which numerically relate terms to some latent variable. When used in tasks of prediction, pLSA models outperformed models produced by LSA. pLSA, reducing the perplexity of any produced models. Notably, there is a greater difference in performance on more general purpose information retrieval tasks, than on restricted corpora. [24]. This is likely due to the effect that sparsity has on the ability of each method. Whilst still worse than pLSA, LSA actually performed better

as the sparsity of training data reduced whereas pLSA's performance decreased. This suggests that the mixture models used in pLSA perhaps overfit as the sparsity of data decreases. As this study is focused on written narratives, any corpora used will be diverse and intrinsically sparse. As such the overfitting problem potentially reducing the effectiveness of pLSA might not be so much of an issue. However, whilst it does begin to numerically relate variables, there is still no explicit relation between documents expressed within the models meaning that adding additional documents to the corpus would mean the algorithm needs to be re-run. Additionally, there remains no direct means of relating entities or expressing entities in terms of elements within a corpus.

2.4.3 Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) uses mixture models to generate a collection of topics which can be used to describe individual documents within a corpus. By representing individual topics as probability distributions of words found within the corpus, additional documents can be expressed as distributions of the topics already derived. This approach overcomes some of the drawbacks of the pLSA method, whilst providing a somewhat human-interpretable representation of the corpus. [20]

Needs reworking...

For each document w in a corpus D , LDA produces a mixture model representing k topics for w . Where each word w_n of length N , has a set denoting the probability of w_n belonging to each of k topics. Collections of models could be compared conceptually to weighted groups produced by some clustering method.

The simplified process for generating an LDA model across m documents is as follows:

1. Choose $N \sim \text{Poisson}(\xi)$
2. Choose $\theta_m \sim \text{Dirichlet}(\alpha)$
3. For each of the N_m words $w_{m,n}$:
 - a) Choose a topic $z_{m,n} \sim \text{Multinomial}(\theta_m)$
 - b) Choose a word $w_{m,n} \sim \text{Multinomial}(\phi_{z_{m,n}})$

Figure 1: LDA Plate Notation

The total probability of an LDA model is defined as:

$$P(\underline{W}, \underline{Z}, \underline{\theta}, \phi | \alpha, \beta) = \prod_{i=1}^K P(\phi_i | \beta) \prod_{j=1}^M P(\theta_j | \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) P(W_{j,t} | \phi_{z_{j,t}})$$

In LDA, the lengths of documents are said to be Poisson distributed, and that the set of probabilities generated for each word w is a random multinomial dirichlet distribution. There is little reasoning for these assumptions. The use of Poisson distributions in literature seems to be a standard distribution for estimating document lengths. However, given that the lengths of documents are known parameters, the first step in the process is redundant leaving only steps 2 and 3.

The versatility of LDA is one of its strongest characteristics. The method has seen use in a range of different settings, from filtering web spam, to the semantic annotation of satellite images. As a bayesian statistical process, it is possible to utilise other methods with relative ease given its modularity.

In web spam filtering, LDA has been combined with tf-idf and expanded to take into consideration links that exist between documents. For spam classification, the use of LDA as a bayesian network performed worse than baseline machine learning methods, however improvements were seen when the method was expanded and combined with other statistical techniques. [25]

For automatic image annotation, additional machine learning and computer vision methods can be applied to images in order transform images into collections of image features. Topics derived from a corpus of these images can be used to annotate new images introduced. Attempts at using LDA in image annotation have seen some success, whilst not all annotations from topics were correct, they were often at least conceptually related. [26,27]

LDA doesn't relate words semantically. The meanings of words in a topic are not necessarily relate in any way, which is demonstrated in studies using LDA to annotate images. This means that derived topics might not provide any meaningful grouping of words. The intuition is that words in books and documents about a specific subject will somewhat relate to other words in the document. LDA might not produce thematically related topics if the documents used to derive topic distributions are varied in content. For this investigation, the potential for

thematically unrelated topics is potentially a major drawback. Entities within narratives might be present throughout a whole series of books. Meaning that they may be associated with a wide range of topics. It is unclear at this stage what effect this may have on the quality of derived topics. It might be necessary to regard the same entity across two books to actually be distinct.

As a bag of words method, each unigram is treated as being unrelated to other words. Which is not true in the case of written narratives. Much of the meaning of words and relations between words are lost when used in ‘bag of word’ models. The drawback for this study is that relations between terms and entities are required. If topics were directly derived from the documents which entities were extracted from, then there would be no meaningful way of identifying how these topics apply to entities.

2.5 Evaluating Topic Models

Topic models are machine representations of word groups that have been deemed to be somehow related. The challenge for evaluating topic models is that it depends on the task which they are being applied to. Directly evaluating the accuracy of topic models is tough, as an unsupervised bayesian method, the probabilities are not interpretable, they are what they are and no ground truth exist. Where produced models can be effectively evaluated is on their efficacy at some secondary task, this could be to predict or estimate topics that might be present in documents not present in the sample used to train models.

2.5.1 Perplexity

Perplexity is a commonly used metric accross topic modelling literature, and is an estimate of how well a probability distribution or probabilistic model predicts a sample. In the literature for pLSA and LDA, perplexity is one of the main metrics used to evaluate the quality of the methods. [20, 24]

Perplexity can be calculated by subsampling a corpus into training and testing sets. By using the test set as unlablled data the perplexity of a model is measured as its ability to estimate

the probabilities density of the M documents in the test set.

For LDA, the perplexity is defined as:

$$perplexity(D_{test}) = exp\{-\frac{\sum_{d=1}^M \log(p(w_d))}{\sum_{d=1}^M N_d}\}$$

The probability of a word w_d is seemingly glossed over in many papers, however some methods for estimating the value have been proposed. The harmonic mean is often used as an estimator due to its relative simplicity and computational efficiency, however it has been subject to criticism regarding its suitability for these. Other methods including 'left-to-right' and a chib-style estimator have been shown to be more effective estimators. [28–31]

The use of perplexity as a metric for evaluating topic models does not necessarily mean that the model is correct or accurate, but rather that the model is capable of predicting a sample, not that the prediction is correct. If a topic model assigned equal probabilities across all words in a vocabulary, then the perplexity would be high, indicating that the topic model provides no insight into the relationships between documents.

2.5.2 Interpretability and Coherence

For human interpretable topic models, it is necessary that the terms within each topic be semantically related, or that it is clear why words in a topic co-exist. Whilst attempts have been made to utilise machine learning methods in the production of coherent topic models, the evaluation of them is somewhat lax. Whilst on the surface, models appear to be more coherent with the inclusion of methods like Markov Random Fields, the coherence is often done subjectively using the opinions of human judges. [32] Even though the use of human judges may be an imprecise method of evaluation, it does provide insight into the interpretability of derived models. For tools that aim to provide users with insight into large corpora of data, interpretability of the results is a key element defining the success of the tool. With low coherence of topics and thus low levels of interpretability, communicating any extracted information becomes increasingly difficult.

Assessing the interpretability and coherence of topic models has been somewhat formalised

through the proposal of a ‘word intrusion’ and ‘topic intrusion’ tests. The word intrusion test assesses the coherence of a topic by introducing words that don’t belong in a topic and asking humans to identify the misplaced word. A Topic Intrusion test is performed by introducing an incorrect topic and asking humans to identify the misplaced topic. It was found that in cases where topics are incoherent or difficult to interpret, that humans would tend to choose a word seemingly at random. [33]. This method has been extended by asking humans to select two intruder words, where only one of them is actually an intrusion. The idea is that in well defined topics not only will the correct intruding word be selected, but also there will be no means for individuals to distinguish between the strength of belonging for the remaining words. If participants can’t decide on what is an intruding word, then there should be a seemingly randomised selection. [34] Both approaches result in different assessments of any derived topics, there is no indication however as to which metric is a better estimate of the coherence of topics.

The precision metrics defined for word and topic intrusion tests may be of value when evaluating the quality of derived Entity Models, though they rely heavily on human evaluation of topics. Use of judges is time consuming and potentially unreliable. A more formal method of evaluating coherence was proposed by comparing word vectors formed in a semantic space built from wikipedia articles. By considering the distributional similarity of word vectors formed by pairs of words in topics, the semantic cohesion of topics can be estimated using one of a few different metrics. When considering the inter-rater agreement between the calculations and human raters, the spearmann rank correlation values had an average of $\bar{x} = 0.77$ and a standard deviation of, $\sigma = 0.04$. The relative agreement between each method of evaluation suggests that any of them could serve as valuable means of evaluating the coherence of topics. [35]

Given a topic $T = w_1, \dots, w_n$, the coherence of that topic is defined as the mean similarity of each possible word pairing $Sim(w_i, w_j)$ where $w_i, w_j \in T$. The similarity of two words could be one of many measures, including the cosine of the word vectors, or the jaccard coefficient. [36]

$$Coherence_{sim}(T) = \frac{\sum_{\substack{1 \leq i < j \leq n \\ i+1 \leq j \leq n}} Sim(w_i, w_j)}{\binom{n}{2}}$$

2.5.3 Comparing Models

When evaluating topic models, it is often done from the perspective selecting the best topic model. This could be done through the use of interpretability and coherence tests, or it could be done through a model's performance at some secondary task. For exploratory purposes, visual means of comparison could also be used to explore the differences between models and methods.

When regarding topics as a k – *dimensional* model, visualising the topic space can be challenging. One approach to comparing topic models is to look at the similarities between each topic in a pair of models. In the form of a bipartite graph, it is clear which topics in two models are conceptually similar. The bipartite graph in figure 2 shows that even though two methods do have different conceptual topics, there are some topics in both that are somewhat similar. The similarities and differences between the two models are seen further if relationships between documents are plotted. The noticeable difference between the two topic modelling methods made visible in the visualisation is that LDA appears to form relations between the document clusters that LSA kept distinct. [37]

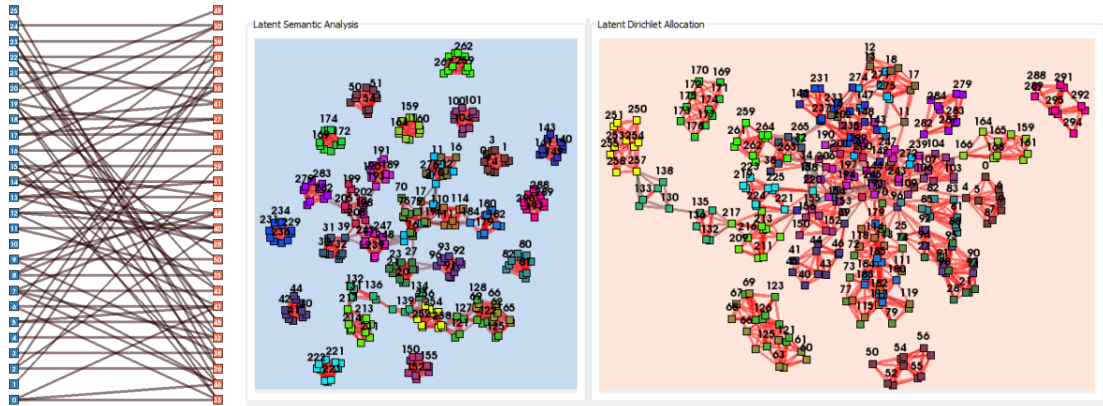


Figure 2: (Left) Bipartite topic conceptual similarity graph. (Middle, Right) Document relationships formed by LSA and LDA topic modelling methods

These visualisations might be effective at explaining why a particular model performs in a certain way at some secondary task, or it might be useful for identifying traits of models that are valuable for particular tasks. However, the visualisations don't provide any quantitative insight into how one model may perform when compared to another.

2.6 Entity Modeling

Entity Modelling exists as a small field of study aside from topic modelling and could be regarded as an extension of the NER domain. It relates to the extraction of entity descriptions, as well as the linking of distinct entities. Many of the challenges in this field stem from the high levels of ambiguity present in natural language. Attempts have been made to utilise Topic Models in the task of entity resolution as well as forming links and relationships between entities within a document. Whilst semantic resources have been investigated as tools for the automatic tagging of extracted entities.

2.6.1 Entity Linking

Entity linking is the task identifying how entities in a corpus relate to each other. It involves the resolution of entity mentions that refer to the same entity in order to perform optimally.

Topic models have been utilised in the task of entity linking through the application and development of specialised Entity-Topic Models. One such investigation involved finding mentions of concrete entities, such as '*iPod*', '*SteveJobs*', and *iPhone* and forming a unifying '*Apple Inc*' topic. [38] The approach added an additional step to the generative LDA algorithm which sampled topic, and entity assignments for each entity. The focus of the paper utilising the topical context within which each entity was mentioned, to try and resolve entities with multiple mentions through relationships that may form between entities in the model.

The Entity-Topic Models suggested are statistically sound, and perform well in predictive tasks oriented around wikipedia information. The key drawback for this study is the use of short documents to train, test, and evaluate derived models. Models were trained on a dataset of news articles from the TAC 2009 dataset [39], which being short in length often benefit the assumption made in the study regarding how entities relate to the topics within the document. For longer texts and particularly narratives, this assumption is not necessarily valid; wide ranges of topics may be present, with entities only being expressed in the context of some of these topics. The limitation of this entity-topic modelling method could come from the base LDA method.

As a bag-of-words approach,

2.6.2 Entity Topic Models

Chapter 3

New Ideas

3.1 Introduction

In chapter 2 many different topic modelling and wider NLP techniques were considered, with only those specific to entity-topic modelling being close to meeting the aims that this project has. The LDA Topic modelling method provides a base statistical method from which this investigation can be launched.

The main drawback with specific entity-topic modelling methods relate to the assumptions made regarding the relationships between entities and topics in a document. Not necessarily applicable to larger documents and written narratives it is clear that some adaptation is required for general LDA based entity-topic modelling methods. Additionally, there is clearly no publically available tools that allow non-technical individuals to explore and investigate entity-topic models, resulting in the field remaining relevant only in literature.

An obvious adaptation to entity-topic modelling methods is to incorporate some weighting metric to associated terms and topics, there is potential for this step to be done as a pre-processing method. Further improvements relate to the clustering of similar entity topic models as latent entities. This provides some estimation of different types of entities that may be present within a corpus

3.2 Gaussian Entity-Term Matrix

The Gaussian Entity-term matrix (GETM) is a means of segmenting a given document into entities and associated terms. The matrix for a given document consists of strengths of as-

sociations between a entity-word pairing. It functions under an assumption that the more separation there is between two words in a document, the weaker the association is between those two words. By modelling documents not as a bag of words, but as sequences it is possible to capture some of the relationships between different words prior to applying additional statistical methods.

Given a set of entities $\underline{E} = \{e_0, \dots, e_i\}$ where $e_i = \{e_{i,0}, \dots, e_{i,m}\}$, and a set of words $\underline{W} = \{w_0, \dots, w_j\}$, then a matrix D when indexed by i, j can be defined as:

$$D_{i,j} = \sum_{m=0}^{|\underline{e}_i|} \sum_{n=0}^{|\underline{w}_j|} \sqrt{\frac{a}{\pi}} \cdot e^{-a \cdot (e_{i,m} - w_{j,n})}$$

$$D = \begin{matrix} & e_0 & e_1 & \dots & e_i \\ \begin{matrix} w_0 \\ w_1 \\ \vdots \\ w_j \end{matrix} & \begin{pmatrix} D_{0,0} & D_{0,1} & \dots & D_{0,i} \\ D_{1,0} & D_{1,1} & \dots & D_{1,i} \\ \vdots & \vdots & \ddots & \vdots \\ D_{j,0} & D_{j,1} & \dots & D_{j,i} \end{pmatrix} \end{matrix}$$

3.3 Latent Entities

A Latent Entity is an abstract statistical model that approximately describes some entities within a corpus. Identifying Latent Entities is possible through clustering entity topic models extracted from a corpus. In the case of k-means clustering, each entity in a GETM matrix is expressed as distributions of entity topics and turned into an Entity-Topic Model. Each Entity-Topic Model can then be clustered in n-dimensional space on topic distributions. The center of each cluster extracted through k-means can then be labelled as a latent entity which abstractly represents entities within that cluster.

As each latent entity remains a probability distribution of entity topics, it is possible to state how similar a given entity is to a latent entity. Each similarity metric of entity to latent entity can then be used to express the entity in terms of similarities to latent entities.

Harry		Hermione		Ron	
Term	Strength	Term	Strength	Term	Strength
Scar	0.1	Magic	0.11	Drooble	0.07
Magic	0.09	Wand	0.09	Magic	0.06
Muggles	0.09	Superiority	0.04	Snigger	0.05
Expelliarmus	0.05	Unkempt	0.04	Sherbert	0.03

Figure 3: Example entity-term associations between three entities

Topics	Harry	Hermione	Ron	Latent Entity
Topic 1	0.6	0.5	0.4	0.5
Topic 2	0.1	0.2	0.1	0.13
Topic 3	0.1	0.2	0.2	0.17
Topic 4	0.2	0.1	0.3	0.2

Figure 4: Example entity-Topic distributions for entities in Figure 3 and the associated Latent Entity.

Given a set of entities (like Harry, Ron, and Hermione), a Latent Entity for this set could be derived by computing the average entity topic model. Conceptually this latent entity could describe entities that can be classified as 'wizards' or 'students'. The exact interpretation of what the latent entity represents is up for debate, in a similar way to how the topic expressed in a topic model is open to interpretation.

3.4 Entity Topic Model

By using a GETM prior to deriving topic models words are already associated with entities, this makes extending LDA topic models a more simple task. As LDA is a bag-of-words method, each entity and its associated terms within the GETM could be regarded as a distinct document. Using GETMs add an additional set of pre-processing steps to the standard LDA algorithm.

3.4.1 Evaluating ETMs

The evaluation of topic models is a challenging task, many techniques commonly used are inherently flawed or difficult to carry out. Many of the flaws present in topic modelling evaluation methods may be made even more visible if used on Entity Topic Models. An adaptation of the intruding word test applicable to ETMs proposed here is the intruding Entity Test.

The intruding entity test involves taking an individual topic entities that are most strongly associated with that topic, and then introducing an entity not all associated with that topic. Should the topics

3.5 Analysis Pipeline

The Analysis Pipeline is a series of steps necessary for the task of extracting Latent Entities from a Corpus. Each step can be viewed as a transformation which can be applied to the results of the previous step.

The pipeline is essential to the implementation of a system which can be used to extract Latent Entity information. Without consideration of what the steps would be, the value of any developed system could be severely diminished.

Starting with raw documents, the start of the process is the tokenization of each document in the corpus. The extraction of POS tag information and Entities within the corpus can be carried out in parallel, it is necessary that both tasks are performed prior to the calculation of the Gaussian Entity-Term Matrix. From the GETM, Entity-Topics and subsequently Latent Entities can be extracted.

Tokenization - NER & POS - GETM - Entity Topics - Latent Entity Clustering

Figure 5: Analysis Pipeline for the Extraction of Latent Entity Information.

3.6 Topic Model Explorer

As machine representations of words most probably related topically, topic models are not always easy to interpret. Presenting only a set of terms within a topic might not provide the

best insight into what that topic represents in a corpus. The challenges of interpretation are further compounded as layers of abstractions are added to a topic model. Entity Topic Models and subsequently Latent Entities would be easier to interpret when presented alongside other derived models. As such, a topic model explorer is proposed as a tool allowing the comparison and exploration of topic models, entity topic models, and latent entities.

Humans often try to understand a topic by applying a single word or small phrase as a title of any given topic. Latent Entities are abstractions of multiple Entity Topic Models, with each entity topic model being varying distributions of topics and words. In order to accommodate for the human need to label models, the Topic Model Explorer can accommodate for this by demonstrating how different clusters of entity topic models are similar, extracting words most significant across members of the cluster.

As well as allowing comparisons of entity topic models, by allowing the explorer to accommodate for the comparison of other levels of topic model extractions, it can allow individuals to evaluate the quality of each model used to describe a corpus.

Chapter 4

Implementation

4.1 Introduction

Implementation for this project involves the development of 3 key elements:

- Analysis Toolkit
- Web Server Architecture
- Latent Entity Explorer

The Analysis Toolkit involves the derivation and implementation of methods necessary to process and analyse corpora. Web Server Architecture involves the design and implementation of a database, Server, and API capable of allowing access to processed data and results from analysis. The Latent Entity Explorer entails the creation of a set of visualisations and tools allowing for the exploration of Entity-Topic models and Latent Entities

4.2 Architecture Overview

The system principally follows a standard 3 layer approach to web-server architecture, requiring a database (DB), web server and API, and a client application. This implementation uses PostgreSQL as its DB of choice, as well as following RESTful practices in the API design. The web server is implemented in python, with a DB repository layer accompanying the server and API implementation. The client application (Latent Entity Explorer) is implemented using TypeScript and the Angular 5 framework.

Alongside the server is a set of Analysis Tools that were developed as distinct methods with utility in their own right.

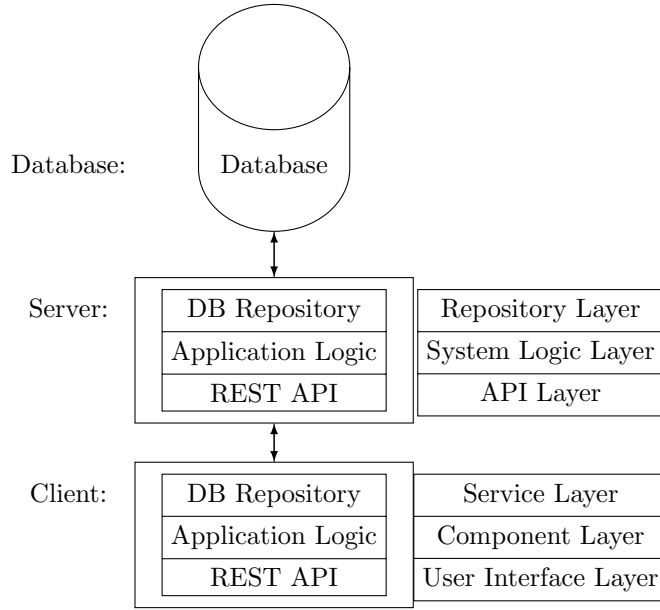


Figure 6: Architecture overview of the Latent Entity Model Analysis and Exploration System.

4.3 NLP Toolkit

The NLP Toolkit consists of a series of methods for processing a corpus of text. A Functional approach was used, with each method applying some transformation to a corpus of documents. The Analysis Pipeline defined in chapter 3 involves tasks starting with the Tokenization of plain text, and ending with the use of clustering methods to extract Latent Entities.

Development followed the steps defined in the Analysis Pipeline. Starting with the early stages, before progressing onto the more complex stages.

4.4 Corpus Pre-Processing

Corpus Pre-Processing involves applying a set of transformations necessary to begin the analysis and creation of descriptive models. It is assumed that the corpus provided to the system will contain only the text that is to be processed. In many documents there is additional preamble, publishing information, and other elements that can be considered 'noise'. This system will not concern itself with the removal of noise within the corpus.

Tokenization is the first transformation to be applied to the corpus. NLTK provides a suite

of tools capable of kickstarting any NLP project, including methods for tokenizing documents. With the option to tokenize a document into sentences, and then words, there are different approaches that can be taken to tokenization. This project uses the default tokenization methods which use XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX. It is important that the documents are tokenized into sentences, and then each sentence is tokenized into words, as this structure allows for more effective POS tagging.

Once each document in the corpus has sentences and words tokenized, the next stage in the processes is the application of POS tagging methods to determine the gramatical structure of each sentence. A Viterbi Decoded Hidden Markov Model (HMM) POS tagger was implemented in python, utilising the NLP toolkit developed for this project. POS Tagging this document means that it is possible to apply filters to the corpus before deriving topic models from the GETM. By removing certain words tagged as certain POS (like determiners, or prepositions) it is possible to derive topics using only descriptive words, which may provide more cohesive and interpretable models.

4.4.1 HMM POS Tagger

A POS tagger was implemented for this project to serve as a benchmark from which standard tools could be comapred against. Hidden Markove Models is an effective baysian method of predicting a most likely given sequence of states given a sequence of inputs.

4.5 GETM Derivation

4.6 Entity-Topic Derivation

4.7 Latent Entity Clustering

4.8 Latent Entity Explorer

Chapter 5

Results and Evaluation

5.1 Introduction

5.2 Gaussian Entity-Term Matrix

Whilst the Gaussian Entity-Term Matrix (GETM) was defined only as a means it is necessary to evaluate the quality of associations formed within the matrix. Should the quality of underlying associations be poor, then the Entity-Topic models and Latent Entity Models will never be able to accurately capture entity structures within narratives.

5.2.1 Experiment Setup

5.2.2 Results

5.3 Entity Topic-Models

The use of a GETM was proposed as a means removing an assumption made in literature regarding the relationships which entities have with the topics within a document. It is therefore necessary to evaluate how the adapted approach performs.

5.3.1 Experiment Setup

5.3.2 Results

5.4 Latent Entities

Chapter 6

Conclusion

6.1 Introduction

Bibliography

- [1] Web server survey — netcraft. <https://news.netcraft.com/archives/category/web-server-survey/>. Accessed: 2017-10-16.
- [2] Eric Brill. Transformation-based error-driven learning and natural language processing: A case study in part-of-speech tagging. *Comput. Linguist.*, 21(4):543–565, December 1995.
- [3] Eric Brill. A simple rule-based part of speech tagger. In *Proceedings of the Workshop on Speech and Natural Language*, HLT '91, pages 112–116, Stroudsburg, PA, USA, 1992. Association for Computational Linguistics.
- [4] H Huang and X Zhang. Part-of-speech tagger based on maximum entropy model. In *2009 2nd IEEE International Conference on Computer Science and Information Technology*, pages 26–29, August 2009.
- [5] Doug Cutting, Julian Kupiec, Jan Pedersen, and Penelope Sibun. A practical part-of-speech tagger. In *Proceedings of the Third Conference on Applied Natural Language Processing*, ANLC '92, pages 133–140, Stroudsburg, PA, USA, 1992. Association for Computational Linguistics.
- [6] Geoffrey Leech, Roger Garside, and Michael Bryant. CLAWS4: The tagging of the british national corpus. In *Proceedings of the 15th Conference on Computational Linguistics - Volume 1*, COLING '94, pages 622–628, Stroudsburg, PA, USA, 1994. Association for Computational Linguistics.
- [7] A Ratnaparkhi. A maximum entropy model for part-of-speech tagging. *of the conference on empirical methods in . . .*, 1996.
- [8] G Malecha and I Smith. Maximum entropy part-of-speech tagging in NLTK. *unpublished course-related report*, 2010.
- [9] David Nadeau and Satoshi Sekine. A survey of named entity recognition and classification. *Linguisticæ Investigationes*, 30(1):3–26, January 2007.
- [10] Lev Ratinov and Dan Roth. Design challenges and misconceptions in named entity recognition. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning*, CoNLL '09, pages 147–155, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics.
- [11] Guodong Zhou and Jian Su. Named entity recognition using an HMM-based chunk tagger. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, pages 473–480, Stroudsburg, PA, USA, 2002. Association for Computational Linguistics.
- [12] A Borthwick and R Grishman. A maximum entropy approach to named entity recognition. 1999.
- [13] Oliver Bender, Franz Josef Och, and Hermann Ney. Maximum entropy models for named entity recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003 - Volume 4*, CONLL '03, pages 148–151, Stroudsburg, PA, USA, 2003. Association for Computational Linguistics.
- [14] Erik Tjong Kim and Fien De Meulder. Introduction to the CoNLL-2003 shared task: Language-Independent named entity recognition. June 2003.
- [15] John D Lafferty, Andrew McCallum, and Fernando C N Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning*, pages 282–289. Morgan Kaufmann Publishers Inc., June 2001.

- [16] Andrew McCallum and Wei Li. Early results for named entity recognition with conditional random fields, feature induction and web-enhanced lexicons. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003 - Volume 4*, CONLL '03, pages 188–191, Stroudsburg, PA, USA, 2003. Association for Computational Linguistics.
- [17] Samet Atdağ and Vincent Labetut. A comparison of named entity recognition tools applied to biographical texts. August 2013.
- [18] J R Finkel, T Grenager, and C Manning. Incorporating non-local information into information extraction systems by gibbs sampling. *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL 2005)*, pages 363–370, 2005.
- [19] The stanford natural language processing group. <https://nlp.stanford.edu/projects/project-ner.shtml>. Accessed: 2018-1-10.
- [20] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3(Jan):993–1022, 2003.
- [21] Scott C Deerwester, Susan T Dumais, George W Furnas, Richard A Harshman, Thomas K Landauer, Karen E Lochbaum, and Lynn A Streeter. Computer information retrieval using latent semantic structure, June 1989.
- [22] Thomas K Landauer, Peter W Foltz, and Darrell Laham. An introduction to latent semantic analysis. *Discourse Process.*, 25(2-3):259–284, January 1998.
- [23] Cherukuri Aswani Kumar and Suripeddi Srinivas. On the performance of latent semantic indexing-based information retrieval. *CIT. Journal of Computing and Information Technology*, 17(3):259–264, October 2004.
- [24] Thomas Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22Nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '99, pages 50–57, New York, NY, USA, 1999. ACM.
- [25] István Bíró, Jácint Szabó, and András A Benczúr. Latent dirichlet allocation in web spam filtering. In *Proceedings of the 4th International Workshop on Adversarial Information Retrieval on the Web*, AIRWeb '08, pages 29–32, New York, NY, USA, 2008. ACM.
- [26] Yansong Feng and Mirella Lapata. Topic models for image annotation and text illustration. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, HLT '10, pages 831–839, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.
- [27] Tong Zhang, Zhe-Ming Lu, Kap Luk Chan, and Zhen Li. Automatic image annotation and retrieval using the latent dirichlet allocation model. *IJCSSES International Journal of Computer Sciences and Engineering Systems*, 5(1), 2011.
- [28] Michael A Newton and Adrian E Raftery. Approximate bayesian inference with the weighted likelihood bootstrap. *J. R. Stat. Soc. Series B Stat. Methodol.*, 56(1):3–48, 1994.
- [29] Hanna Megan Wallach. *Structured topic models for language*. PhD thesis, University of Cambridge, 2008.
- [30] Siddhartha Chib. Marginal likelihood from the gibbs output. *J. Am. Stat. Assoc.*, 90(432):1313–1321, 1995.
- [31] Hanna M Wallach, Iain Murray, Ruslan Salakhutdinov, and David Mimno. Evaluation methods for topic models. In *Proceedings of the 26th Annual International Conference on Machine Learning*, ICML '09, pages 1105–1112, New York, NY, USA, 2009. ACM.
- [32] Pengtao Xie, Diyi Yang, and Eric Xing. Incorporating word correlation knowledge into topic modeling. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 725–734, 2015.

- [33] Jonathan Chang, Sean Gerrish, Chong Wang, Jordan L Boyd-graber, and David M Blei. Reading tea leaves: How humans interpret topic models. In Y Bengio, D Schuurmans, J D Lafferty, C K I Williams, and A Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 288–296. Curran Associates, Inc., 2009.
- [34] Fred Morstatter and Huan Liu. A novel measure for coherence in statistical topic models. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 543–548, 2016.
- [35] Nikolaos Aletras and Mark Stevenson. Evaluating topic coherence using distributional semantics. In *Proceedings of the 10th International Conference on Computational Semantics (IWCS 2013)–Long Papers*, pages 13–22. aclweb.org, 2013.
- [36] David Newman, Jey Han Lau, Karl Grieser, and Timothy Baldwin. Automatic evaluation of topic coherence. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, HLT ’10, pages 100–108, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.
- [37] P J Crossno, A T Wilson, T M Shead, and D M Dunlavy. TopicView: Visually comparing topic models of text collections. In *2011 IEEE 23rd International Conference on Tools with Artificial Intelligence*, pages 936–943, November 2011.
- [38] Xianpei Han and Le Sun. An entity-topic model for entity linking. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, EMNLP-CoNLL ’12, pages 105–115, Stroudsburg, PA, USA, 2012. Association for Computational Linguistics.
- [39] Paul McNamee and Hoa T Dang. Overview of the tac 2009 knowledge base population track, 01 2009.