

GloDETC : Global Document Embeddings for Multi-label Text Classification

*A Thesis Submitted
in Partial Fulfilment of the Requirements
for the Degree of*

B.Tech. - M.Tech. (Dual Degree)

by

Nitish Gupta

Roll No. : 10327461

under the guidance of
Prof. Harish Karnick
Prof. Rajesh M. Hegde



Department of Electrical Engineering
Indian Institute of Technology Kanpur
April, 2015

CERTIFICATE

It is certified that the work contained in this thesis entitled "*Merging Word Senses*", by *Sumit Bhagwani* (Roll No. Y8127515), has been carried out under my supervision and that this work has not been submitted elsewhere for a degree.



(Prof. Harish Karnick)
Department of Computer Science and Engineering,
Indian Institute of Technology Kanpur
Kanpur-208016

June, 2013



Abstract

Abstract goes here

*Dedicated to
My Family*

Acknowledgement

I wouldn't like to express my sincere gratitude towards my thesis supervisor, Prof. Harish Karnick, for his constant support and encouragement. I am grateful for his patient guidance and advice in giving a proper direction to my efforts. I am also grateful to Prof. Rajesh M. Hegde for giving me the freedom to work on a topic of my liking.

Last, but not the least, I would like to thank my parents and brother for their love and encouragement. Without their support and patience this work would not have been possible.

Nitish Gupta

Contents

Abstract	i
List of Tables	vii
List of Figures	ix
List of Algorithms	xi
1 Related Work	1
1.1 Text Representation	1
1.1.1 Bag of Words	2
1.1.2 Dimensionality Reduction / Feature Selection	3
1.2 Learning Algorithms	5
1.2.1 With Multiple Binary Classifiers	6
1.2.2 With Single Joint Classifier	9
2 Distributed Document Embeddings	13
2.1 Motivation	13
2.2 Background on Word Embeddings	14
2.2.1 Neural Probabilistic Language Model (NPLM)	15
2.2.2 Log-Linear Models : word2vec	17
2.2.2.1 Continuous Bag-of-Words (CBOW)	17
2.2.2.2 Continuous Skip-gram	18
2.2.2.3 Dependency-based Word Embeddings	19
2.3 Document Embeddings	20

2.3.1	Problem Setup	21
2.3.2	Our Model	22
2.3.2.1	Projection Layer (Context Representation)	23
2.3.2.2	Estimating Prediction Probability	24
2.3.2.3	Training Objective	24
2.3.2.4	Noise Contrastive Estimation (NCE)	25
2.3.2.5	New Training Objective	27
2.3.2.6	Parameter Estimation	27
2.3.2.7	Hyper-parameters	30
3	Multi-Label Text Categorization	33
3.1	Logisitic Regression (LR) for Multi-label Document Categorization .	33
3.1.1	Training Data	34
3.1.2	Logistic Regression Model	35
3.1.2.1	Training Objective	36
3.1.2.2	Parameter Estimation	36
	Bibliography	39

List of Tables

List of Figures

2.1	Bengio’s Neural Network Architechture for Neural Probabilistic Lan- guage Model	16
2.2	Continuous Bag-of-Words Model (CBOW) TODO: Add ref?	18
2.3	Continuous Skip-gram Model TODO: Add ref?	19
2.4	Dependency-based context extraction example TODO: Add ref? . .	20
2.5	GloDETC : Neural Network Archtitecture TODO: Change figure . .	23

List of Algorithms

Chapter 1

Related Work

The task of text classification, i.e. classification of documents into a fixed number of predefined categories has been long studied in-depth for many years now. This multi-class classification problem has further evolved into a multi-label text classification task where each document can belong to multiple, exactly one or no category at all.

Supervised machine learning techniques that learn classifiers to perform this category assignment task can be broken down into two main components, namely, text representation and learning algorithm. Text representation involves converting the documents, that are usually strings of characters, into numerical vectors that are suitable inputs to the learning algorithm while the learning algorithm uses pairs of labeled input text representations and the categories it belongs in, to learn a model so as to classify new documents into categories.

1.1 Text Representation

Any text-based classification system requires the documents to be represented in an appropriate manner dictated by the task being performed [Lewis, 1992a]. Moreover, [Quinlan, 1983] showed that the accuracy of the classification task depends as much on the document representation as on the learning algorithm being employed. Different from the data mining task, which deals with structured documents, text classification deals with unstructured documents that need to be appropriately trans-

formed into numerical vectors, i.e. the need for text representation. In this section we introduce the most effective and widely-used techniques to represent documents for text classification.

1.1.1 Bag of Words

It is found in information retrieval research that word stems work well as representations units for documents and that their ordering in a document is of minor importance for many tasks. This is attributed by the fact that the most widely-used model to represent documents for the classification task is the *Vector Space Model (VSM)* [Salton and Yang, 1973].

In the Vector Space Model, a document d is represented as a vector in the term/word space, $d = (w_1, w_2, \dots, w_{|V|})$ where $|V|$ is the size of the vocabulary. Each of the $w_i \in [0, 1]$, represents the weightage of the term i in the document d . This is called the *bag-of-words* model as it ignores word ordering and each document is reduced to a bag of words that it contains or not.

An important requirement of such a representation is that, the terms that help in defining the semantic content of the document and play an important role in classification be given higher weightage than the others. Over the years, there has been much research in the information retrieval field on term weighting schemes. The most important term-weighting techniques are described below :

1. **One Hot Representation** : This is the most trivial representation, where each document is represented by a vector that is size of the vocabulary. Each element in the vector is either a 0 or a 1 to denote the absence or presence of a specific term in the document.
2. **Term Frequency (tf)** : The term frequency representation weighs the terms present in the document relative to their occurrence frequency in the document. Hence a document d is represented as, $d = (w_1, w_2, \dots, w_{|V|})$, where, w_k is the number of times the term k appears in the document d .

3. **Inverse Document Frequency (idf)** : Though using tf as a term weighting scheme is a good starting point, it faces a challenge when high frequency terms are not concentrated in a few particular documents but are prevalent in the whole collection. Those terms then stop being characteristic of the semantic content of a few documents and need not be given high weightage. To overcome this problem, Salton and Buckley [1988] suggested a new term weighting called the inverse document frequency (idf). The idf weight of a term varies inversely with the number of documents n it belongs to in a collection of total N documents. A typical idf vector can be computed as

$$w_k = \log \frac{N}{n} \quad (1.1)$$

4. **Term Frequency Inverse Document Frequency (tf-idf)** : Given the above two term weighing schemes, it is clear that an important term in a document should have high tf but a low overall collection frequency (idf). This suggests that a reasonable measure for term importance may be then obtained by the tf and the idf ($tf \times idf$). As we will see in the results section, the $tf-idf$ weighed bag-of-words document representation gives one of the best accuracies in the multi-label text classification task.

1.1.2 Dimensionality Reduction / Feature Selection

The bag-of-words representation scheme has several drawbacks but the most important drawback it suffers from is that document vectors are very sparse and high dimensional. Typical vocabulary sizes of a moderate-sized document collection ranges from tens to hundreds of thousands of terms which is prohibitively high for many learning algorithms. To overcome this issue of high-dimensional bag-of-words document representations, automatic feature selection is performed that removes uninformative terms according to corpus statistics and constructs new orthogonal features by combining several lower level features (terms/words). Several

techniques used in practice are discussed below,

1. **Information Gain** : Information Gain is widely used as a term-goodness criterion in the field of machine learning, mainly in decision trees [Quinlan, 1986] and also in text classification [Lewis and Ringuette, 1994], [Moulinier et al., 1996]. It is a feature space pruning technique that measures the number of bits of information obtained (entropy) for category prediction by knowing the presence or absence of a term in a document. For terms where the information gain was below some predefined threshold are not considered in the document vector representation. The information gain of a term t is defined as

$$G(t) = - \sum_{i=1}^{|C|} P(c_i) \log P(c_i) + P(t) \sum_{i=1}^{|C|} P(c_i|t) \log P(c_i|t) + P(\neg t) \sum_{i=1}^{|C|} P(c_i|\neg t) \log P(c_i|\neg t) \quad (1.2)$$

2. **Mutual Information** : Similar to the Information Gain scheme, Mutual Information estimates the information shared between a term and a category and prunes terms that are below a specific threshold. The mutual information between a term t and a category c is estimated in the following fashion,

$$I(t, c) = \log \frac{P(t \wedge c)}{P(t) \times P(c)} \quad (1.3)$$

To measure the goodness of a term in global feature selection, the category specific scores of a term are combined using,

$$I_{avg}(t) = \sum_{i=1}^{|C|} P(c_i) I(t, c_i) \quad (1.4)$$

3. **χ^2 Statistic** : The χ^2 statistic measures the lack of independence a term t and a category c and can be compared to the χ^2 distribution with one degree of freedom. The term-goodness factor is calculated for each term-category pair and is averaged as above. The major difference between Mutual Information and χ^2 statistic is that the later is a normalized value and the goodness factors

across terms are comparable for the same category.

4. **Latent Semantic Indexing (LSI)** : LSI first introduced by Deerwester et al. [1990], is a popular linear algebraic dimensionality reduction technique that uses the term co-occurrence statistics to capture the latent semantic structure of the documents and represent them using low-dimensional vectors. It is an efficient technique to deal with synonymy and polysemy. LSI aims to find the best subspace approximation to the original document bag-of-word vector space using Singular Value Decomposition. Given a term-document matrix $X = [x_1, x_2, \dots, x_{|D|}] \in \mathbb{R}^{|V|}$, its k -rank approximation as found using SVD, can be expressed as,

$$X = TSD^T \tag{1.5}$$

where, $T \in \mathbb{R}^{|V| \times k}$ and $D \in \mathbb{R}^{|D| \times k}$ are orthonormal matrices called the left and right singular vectors respectively. The matrix $S \in \mathbb{R}^{k \times k}$ is a diagonal matrix of singular values arranged in descending order. The k -dimensional rows of the matrix D contain the dimensionality reduced representations of the $|D|$ documents in the collection. The representations obtained using LSI alleviate the issue of data sparsity and high-dimensionality in bag-of-words representations and also helps unfold the latent semantic structure of the documents.

1.2 Learning Algorithms

Multi-label text classification has seen growing number of statistical learning methods being applied to it. Over the years, various learning algorithms like, Regression models ([Cooper et al., 1994], [Fuhr et al., 1991]), Conditional Random Field ([Ghamrawi and McCallum, 2005]), Nearest Neighbour techniques ([Yang, 1994], [Zhang and Zhou, 2005], [Zhang and Zhou, 2007]), Bayesian classifier and topic modelling ([Lewis and Ringuette, 1994], [McCallum, 1999], [Nigam et al., 2000], [Rubin et al., 2012], [Nigam et al., 1999], [Ueda and Saito, 2002]), SVM ([Joachims, 1998], [Elisseeff and Weston, 2001]), Neural Networks ([Wiener et al., 1995], [Ng

et al., 1997]), Decision Trees ([Tong and Appelbaum, 1994]), Online learning algorithms ([Lewis et al., 1996], [Crammer and Singer, 2002]), Non-negative Matrix Factorization ([Liu et al., 2006]) etc. have been used or developed for Multi-label document categorization.

Earlier learning algorithms reduced the problem of multi-label classification into multiple binary classification problems and independently learned binary classifiers for each category. While these algorithms performed well, their drawback of considering correlation among categories led to the development of algorithms that learn a single classifier and jointly classify each document.

Multi-label classification problems can be also be classified into classification-based and ranking-based approaches, where the former assigns each test instance a $|L|$ -sized label vector of ones and zeros indicating the presence and absence of labels. In the case of a ranking-based approach, the ranking system outputs the list of labels arranged in the increasing order of a ranking score which is then thresholded at an optimum and the top labels are considered appropriate label assignments for test instances.

Below we describe some of the famous learning algorithms for multi-label text classification,

1.2.1 With Multiple Binary Classifiers

The most common approach of multi-label text classification treats each label independently and learns multiple binary classifiers, one for each category and then assigns to a test document all the categories for which the corresponding classifier says ‘yes’. Below we describe some of the algorithms, in the context of multi-label text classification, that learn multiple independent binary classifiers.

1. **Logistic Regression (LR)** : Introduced by [Hosmer and Lemeshow, 1989], LR is a probabilistic binary classification regression model, that, for binary text classification learns a category weight vector and estimates the probability of a document belonging to the category using dot-product and the logistic

link function. LR can be extended for multi-label document classification by learning multiple category vectors, specifically, one for each category. At test time, one would need to query all category vectors for each document to make the category assignments. In our work, we use logistic regression for multi-label text classification, the details for which are given in Sec 3.1.

2. **Support Vector Machines (SVM)** : Support Vector Machines ([Cortes and Vapnik, 1995], [Vapnik, 2000]) based on the *Structural Risk Minimization* principle, are universal learners. In their basic form, SVMs learn linear threshold functions to find linear hyperplanes in the input data space to separate data of the two different classes. In the case, where data is not linearly separable, SVMs can be plugged-in with appropriate kernel functions to learn polynomial classifiers, radial basic functions etc. For multi-label text classification, training data is treated separately for each category and maximum margin separating hyperplanes are found for each category independently [Joachims, 1998].

Elisseeff and Weston [2001] study a ranking based variant of SVM, where the positive/negative distance from the separating hyperplane of a specific category is the score assigned to the particular instance for that category. Their formulation then aims to maximize the margin between the score of a category that belongs to the document and a category that does not belong to the document. This is also called the Rank-SVM.

3. **Neural Networks (NNet)** : Classification-based, Neural Network approaches to multi-label text classification were mainly studied by Wiener et al. [1995], developed at Xerox PARC and called NNet.PARC and Ng et al. [1997], called CLASSI. Both neural networks are examples of multiple-classifier based approaches where a separate neural network was trained for each category to make binary classifications. While CLASSI used a linear perceptron approach to classify text into categories, NNet.PARC built a three-layered nonlinear neural network that extends logistic regression by modelling higher order term

interactions and hence finding non-linear decision boundaries.

4. **Naive Bayes (NB)** : Naive-bayes as studied in Lewis [1992b] and Lewis and Ringuette [1994], is one of the most effective and simple statistical model for text classification. For multi-label classification, classifiers are learnt so as to estimate $P(C_j = 1|D)$, i.e., the probability that the document, D belongs to the category C_j , for each category. This probability is estimated by estimating the probability $P(W_i = 1|C_j = 1)$, i.e. probability that a particular word appears in the document when it belongs to a particular category. Though this approach makes the assumption of word independence, experiments show that this fast-learning algorithm can yield excellent results.

Although, approaches to multi-label classification discussed above give competitive accuracies in the task, they suffer from inefficiencies due to the following reasons,

- make assumptions of category independence and learn 1-vs-All binary classifiers. It is realized that such assumption would not hold true in most real-life situations. Fine-grained categorization of texts usually involve strongly correlated category classes and information about the presence of one gives information about the presence/absence of many others. For eg. in the sentence,

*Chicago Board of trade grain traders and analysts voiced a lot of
interest in how farmers planned to handle their upcoming spring
plantings prompting sales of new crop months of corn and oats and
purchases in new crop soybeans in the futures markets*

information from words about the presence of categories like *oats*, *corn* etc. can also aid the prediction of the *agriculture* category which can be boosted using joint classification.

Apart from inefficiencies induced by ignoring category correlations, learning independent classifiers poses other drawbacks, such as, in case of millions of labels, learning millions of high-dimensional classifiers is a computationally expensive. Sec-

only, the cost of prediction for each test instance would be high as all the classifiers need to be evaluated to make a single prediction.

1.2.2 With Single Joint Classifier

To overcome the difficulties and drawback of learning multiple binary classifiers, researchers have since developed learning algorithms that jointly classify each document into categories it belongs to. Outputs of such algorithms are $|L|$ -dimensional label vectors $\mathbf{y} \in \{0, 1\}^L$, with $\mathbf{y}_l = 1$ if label l is relevant for the particular document. Below we describe algorithms for multi-label text classification that learn a single classifier for assigning all relevant labels to a document jointly.

1. **k-Nearest Neighbor (kNN)** : k-nearest neighbor classification is one of the most effective lazy learning approaches to classification. Given an arbitrary text document input, the algorithm first ranks the nearest neighbors among the training documents using some similarity measure. It then uses the category information of the top-k ranked nearest neighbors to predict the categories of the input test document. One simple approach is to take a weighted average of the label vector of the k-nearest neighbors, weights being the similarity score while estimating document distances. This yields a category ranking for the test input which can be thresholded to yield binary classifications.

Other approach as devised by Zhang and Zhou [2007] is based on the k-NN and the maximum a posteriori(MAP) principle. Their approach is, given a test instance, to first identify its k-nearest neighbors and then based on the statistical information gained from the label sets of the neighboring instances, use the MAP principle to determine the label set of the given input. The prior probability of label occurrences and the posterior probability, $P(C_l = n | l = 1)$ i.e. given a document belongs to label l , exactly n of its k neighbors also belong to the label l is determined from the training instances to utilize the MAP principle.

2. **Linear Least Squares Fit (LLSF)** : LLSF[Yang and Chute, 1992] learns a multivariate regression model automatically from a training set of documents and their categories. Documents are input as vectors in the desired representation and the corresponding output is a $|L|$ -dimensional binary label vector. By solving a linear least squares fit on the training pairs of vectors a matrix of word-category regression coefficients is learnt, which defines the mapping from an arbitrary document to a weighted category label vector. This weighted vector can be sorted to yield a ranked list of categories for the input document.

3. **Probabilistic Models** : Generative probabilistic models described in McCallum [1999], Nigam et al. [1999], Ueda and Saito [2002] etc. argue that the words in a document belonging to a multi-category class can be regarded as a mixture of characteristic words related to each of the categories. Therefore, they represent the multi-label nature of the document by specifying each document with a set of mixture weights, one for each class and also indicate that each document is generated by a mixture of word distributions, one distribution for each label. Once the word distributions are learnt using the training data, classification is performed using the Bayes Rule which selects the labels that are most likely to generate the given test document. Hence, along with giving the information on the labels responsible for generating the document, such models also fill the missing information of which labels were responsible for generating each word.

McCallum [1999] and Ueda and Saito [2002] define a multinomial distribution $\theta_l = \{\theta_{l1}, \theta_{l2}, \dots, \theta_{l|V|}\}$ over the vocabulary for each label, and the word distribution for a document for a given label vector \mathbf{y} , is computed by taking a weighted average of the word distributions of the labels that are present in the document. Therefore, if $\phi(\mathbf{y}) = \{\phi_1(\mathbf{y}), \phi_2(\mathbf{y}), \dots, \phi_2(\mathbf{y})\}$ is the required

word distribution, it can be represented by,

$$\boldsymbol{\phi}(\mathbf{y}) = \sum_{l=1}^{|L|} h_l(\mathbf{y}) \boldsymbol{\theta}_l \quad (1.6)$$

where $h_l(\mathbf{y})$'s are the mixing proportion that add upto 1. The word distributions for each label are found by maximizing the posterior in [Ueda and Saito, 2002] and by employing the Expectation-Maximization algorithm in [McCallum, 1999].

Chapter 2

Distributed Document Embeddings

In this chapter we describe the concept of distributed word and document embeddings and why distributed representations of words and documents are better than one-hot or bag-of-words representations as described in 1.1. We then give a background on different models that learn distributed representations for words in a fully unsupervised manner and finally describe in detail our proposed model for learning distributed embeddings for documents that can be used for multi-label text classification.

2.1 Motivation

TODO: Get in tune to document representations. Say words and documents suffer in the same manner with one-hot or bow representations. Express problems in docs with changing words. Give example of sentence

TODO: Can be tackled with distributed repr. Similarity measures as simple as cos-distance can be introduced in documents. Lets model joint distributions of words with continuous distributions. Words have distributed representations but not docs.

Words are regarded as atomic symbols in most rule-based and statistical natu-

ral language processing(NLP) tasks and hence need the appropriate representation to solve the NLP tasks with greater ease and accuracy. Words are traditionally expressed as one-hot vectors, i.e. as vectors of the size of the vocabulary where exactly one element is 1 and the rest all are zero. Though these representations have been widely used, one-hot representations have a plethora of drawbacks that pose problems and limit the ability of systems to perform better.

1. **Curse of Dimensionality** : One-hot representations lead word vectors to be the size of the vocabulary which often consists of tens to hundreds of thousands of words. Due to this curse of dimensionality, language modelling becomes almost impossible where the number of parameters would grow exponentially with the size of the vocabulary if the words are represented as one-hot vectors.
2. **No Word Similarity** : As words are represented by sparse orthogonal vectors, there is no notion of word similarity that can be introduced. In one-hot representation, the word “symphony” is equally close to the words “bark” and “guitar”. We would want word representations such that they capture the semantic or topical similarity between words.

Due to the problems dicussed above there is a need for more robust, low-dimensional, non-sparse vector representations for words that capture the semantic similarity between them, can be used to model language with continuous distributions and can be used as inputs for various other NLP tasks.

2.2 Background on Word Embeddings

Distributed word representations are dense fixed-sized feature vectors learnt for words in an unsupervised manner from large text corpus that capture the semantic similarity between words. Each word w_i in the corpus is represented by a vector, $v_{w_i} \in \mathbb{R}^m$, where m usually ranges from 50 – 300. These dense representations help deal with sparsity and high-dimensionality issues in ont-hot representations and also

provide provision for estimating similarities between words; which is as simple as taking the dot-product or calculating the cosine-distance between the vectors.

All of the word vector learning models make use of neural networks ([Bengio et al., 2003a], [Mnih and Kavukcuoglu, 2013], [Mikolov et al., 2013b], [Collobert et al., 2011], [Bottou, 2014], [Turian et al., 2010], [Levy and Goldberg, 2014]) but differ in their training objectives.

Below we describe in detail two models to show how models with very different learning objectives and architecture can lead to learning high-quality word vectors.

2.2.1 Neural Probabilistic Language Model (NPLM)

Introduced by Bengio et al. [2003a], their model aims to learn distributed word vectors and a probability function that uses these vectors to learn a statistical model of language. In their model, the probability of a word sequence is expressed as the product of conditional probabilities of the next word given the previous ones.

$$P(w_1^T) = \prod_{t=1}^T P(w_t | w_1^{t-1}) \quad (2.1)$$

And making the n-gram assumption,

$$P(w_t | w_1^{t-1}) \approx P(w_t | w_{t-n+1}^{t-1}) \quad (2.2)$$

i.e. the probability of the next word in the sequence is mostly affected by the local context, in this the previous n -words and not the whole past sequence.

Their model maps each word to a m -dimensional vector in a matrix $C \in \mathbb{R}^{|V| \times m}$ and estimates the probability $P(w_t = i | w_{t-n+1}^{t-1})$ i.e. the probability that the t^{th} word in the sequence is w_i . The neural network that is used to estimate this probability using the word vectors is shown in Figure 2.1 For each input sequence, the neural network outputs a vector $y \in \mathbb{R}^{|V|}$, where y_i is the unnormalized log-probability that

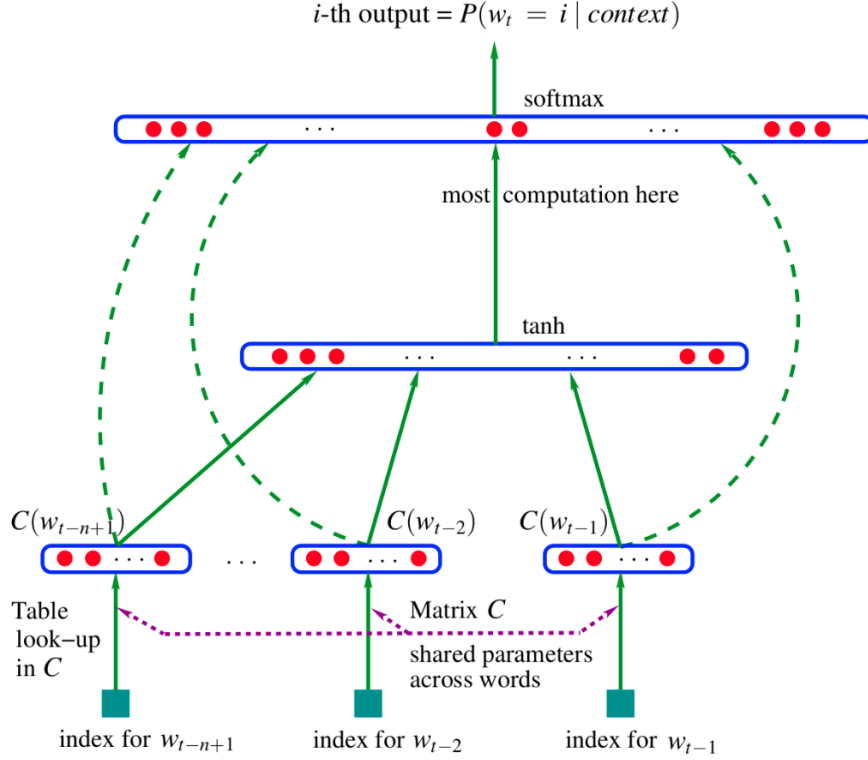


Figure 2.1: Bengio's Neural Network Architecture for Neural Probabilistic Language Model

the t^{th} word in the sequence is w_i .

$$y = b + Wx + U \tanh(d + Hx) \quad (2.3)$$

where \tanh is the hyperbolic tangent applied to introduce non-linearity and x is the word feature layer activation vector constructed by the concatenation of the context word vectors,

$$x = (C(w_{t-1}), C(w_{t-2}), \dots, C(w_{t-n+1})) \quad (2.4)$$

The unnormalized log probabilities in y are converted to positive probabilities summing to 1 by using a *softmax* output layer that computes,

$$P(w_t = i | w_{t-1}, \dots, w_{t-n+1}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}} \quad (2.5)$$

The parameters of the model (b, d, W, U, H) and the word vectors C are estimated by maximizing the log-likelihood of the training corpus.

2.2.2 Log-Linear Models : word2vec

Simple log-linear models are proposed in Mikolov et al. [2013a] as opposed to the non-linear NPLM model to bring down the training time complexity without sacrificing with the quality of the word vectors. **TODO: Also these models are based on the Distributional Hypothesis** The models bring down the complexity of learning vectors by not having a non-linear layer and matrix weighting of the input vectors that are the costliest operations in NPLM. The two models proposed in Mikolov et al. [2013a] are Continuous Bag-of-Words and Continuous Skip-Gram model, described below.

2.2.2.1 Continuous Bag-of-Words (CBOW)

This model is different from the NPLM in that the projection layer is shared for all words; i.e. all words get projected into the same hidden layer vector (their vectors are averaged). This architecture hence neglects the ordering of the words as opposed to NNLM that uses the concatenation of input vectors for the projection layer. The training criteria in this model is to to classify the current (middle) word given its context. It also uses word sequence from the future to aid this task with the relaxation that the aim is not to learn a language model. The model architecture is given in Figure 2.2. The model first computes the hidden layer vector h ,

$$h(w_{t-k}, \dots, w_{t+k}) = \frac{w_{t-k} + \dots + w_{t-1} + w_{t+1} + \dots + w_{t+k}}{2k} \quad (2.6)$$

where, w_{t-i} is the i -th previous word in the context of the middle word w_t and k is the window length. The neural network then computes a unnormalized log-probability vector y similar to Sec.2.2.1, and uses the *softmax*-classifier to estimate $P(w_t|w_{t-k}, \dots, w_{t+k})$,

$$y = b + Uh(w_{t-k}, \dots, w_{t+k}) \quad (2.7)$$

$$P(w_t|w_{t-k}, \dots, w_{t+k}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}} \quad (2.8)$$

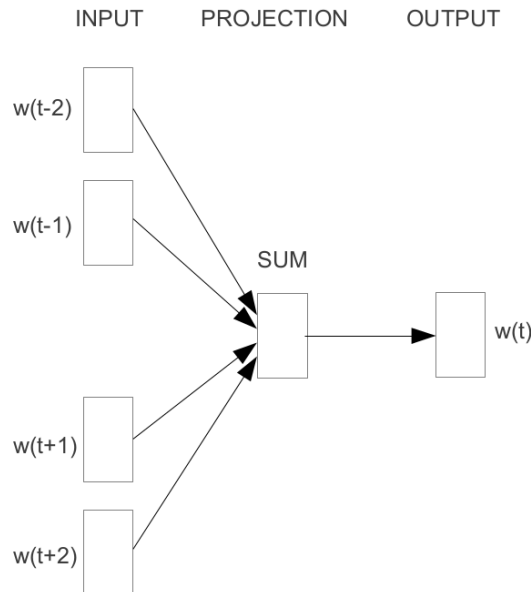


Figure 2.2: Continuous Bag-of-Words Model (CBOW) **TODO: Add ref?**

The parameters of the CBOW model, (b, U) and the word vectors (w_i) are learnt by maximizing the average log probability (Eq. 2.8) of the training corpus.

2.2.2.2 Continuous Skip-gram

This model is similar to the CBOW model, but instead of predicting the middle word based on the context, it tries to maximize the classification of a word based on another word in the context. More precisely, given each word, the skip-gram model tries to predict words within a certain range before and after the current word. The model architecture is given in Figure 2.3 Formally, given a sequence of words in a context w_{t-k}, \dots, w_{t+k} , the skip-gram model defines $P(w_{t+j}|w_t)$ using the *softmax*-classifier in the following manner,

$$P(w_{t+j}|w_t) = \frac{e^{(v_{w_t} \cdot v_{w_{t+j}})}}{\sum_i e^{(v_{w_t} \cdot v_{w_i})}} \quad (2.9)$$

The only parameters of the Skip-gram model are the word vectors (v_{w_i}) that are learnt by maximizing the average log probability (Eq. 2.9) of predicting all the context words for all the words in the training corpus.

The CBOW and the Skip-gram models use the *hierarchical softmax* [Morin and

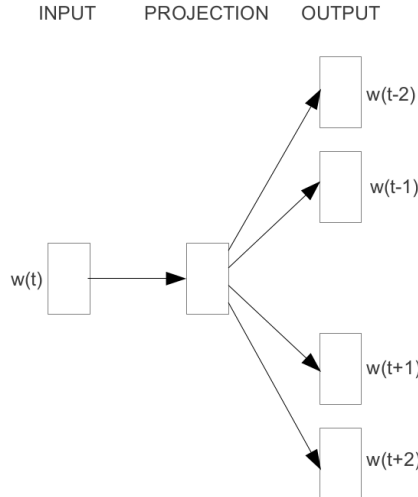


Figure 2.3: Continuous Skip-gram Model **TODO: Add ref?**

Bengio, 2005] instead of the full softmax to speed-up the learning process.

The quality of the word vectors is tested using the *Semantic-Syntactic Word Relationship test* that evaluates the model performance on retrieving semantically and syntactically similar words to the given test words. The word vectors learnt using the skip-gram model are also shown to encode many linguistic regularities and pattern [Mikolov et al., 2013c] and show additive compositionality using simple vector arithmetics. For example, the result of the vector calculation $vec(Madrid) - vec(Spain) + vec(France)$ is closest to $vec(Paris)$ than any other word vectors.

2.2.2.3 Dependency-based Word Embeddings

Instead of using bag-of-words based context as used in *NPLM* and *word2vec*, Levy and Goldberg [2014] use arbitrary contexts to investigate its effects on the word vectors and the properties they encode. The most important of their techniques is to derive the contexts based on the syntactic relations that the word participates in. For each word w and its modifiers m_1, \dots, m_k found using the parse tree of the sentence, contexts $(m_1, lbl_1, \dots, m_k, lbl_k)$ are extracted, where lbl is the type of the dependency relation between word and the modifier and lbl^{-1} is used to mark the inverse-relation. An example of the contexts extracted for a sentence is given in Figure 2.4. After extracting the contexts, their model uses the neural network architecture and the

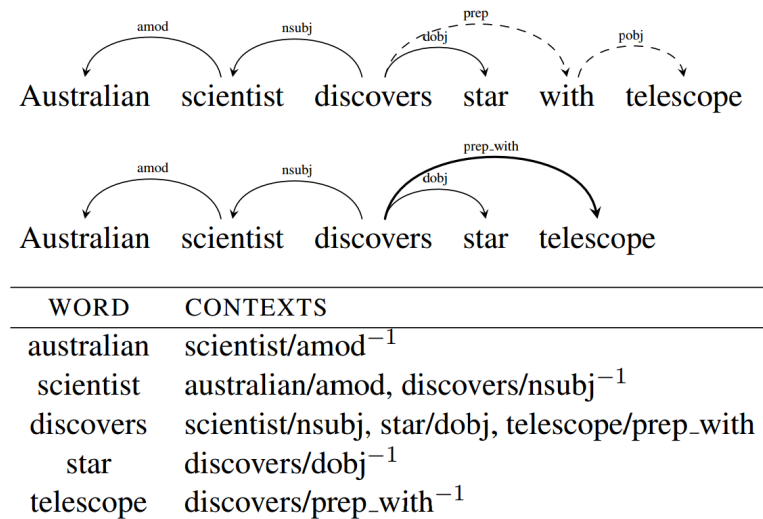


Figure 2.4: Dependency-based context extraction example **TODO: Add ref?**

training objective of the skip-gram model to learn word vectors. On comparison to the vectors learnt from the skip-gram model on the tasks of *topical similarity* and *functional similarity* estimation, it is found that the vectors learnt from this model perform better on the *functional similarity* task that expects word vectors to encode syntactic relationships better. In the task of *topical similarity* estimation, the vectors from the skip-gram model performed better as they encode semantic similarity between words because of the bag-of-words context used during training.

2.3 Document Embeddings

In the previous section we saw how distributed word embeddings that encode semantic similarity can be learnt from text. Though these semantic word spaces are very useful for a lot of tasks, their ability to capture the complexity and compositionality of human language is limited. Word embeddings cannot be directly used to represent longer phrases, sentences and documents to express their meaning. Tasks such as word sense disambiguation, sentiment analysis, text categorization etc. all require the text representation to capture the semantic content of the text for better inputs to learning algorithms as compared to a simple bag-of-words model.

Progress towards learning distributed representations for longer pieces of text,

such as phrase-level or sentence-level representations [Mitchell and Lapata [2010], Zanzotto et al. [2010], Yessenalina and Cardie [2011], Grefenstette et al. [2013], Mikolov et al. [2013b]] that capture semantic compositionality has been promising, but most models do not go beyond simple weighted average of word vectors to represent longer texts. Socher et al. [2013] proposes a more sophisticated approach using recursive tensor neural network where the dependency parse-tree of the sentence is used to compose word vectors in a bottom-up approach to represent sentences for sentiment classification of phrases and sentences. Both the techniques have weaknesses for learning document representations. The first approach is analogous to a bag-of-words approach and neglects word order while representing documents whereas the second approach considers syntactic dependencies but cannot go beyond sentences as it relies on parsing.

Below we present our model on learning universal distributed vector representations for documents and words in the corpus such that,

1. The learned vectors encode semantic and topical content of the documents and words.
2. Semantically similar documents/words have similar vector representations.

To learn vectors that satisfy 1. and 2. above, we hypothesize that document representations should be learnt such that they can aid in the prediction of words in a given word sequence from the document. In the sections below we formally introduce the problem and present our model to learn document and word vector representations.

2.3.1 Problem Setup

Given a set of documents, $\mathbf{D} = \{d_1, \dots, d_{|\mathbf{D}|}\}$ and a vocabulary of words, \mathbf{V} constructed using the set of documents, we wish to embed each document $d_i \in \mathbf{D}$ and each word in the vocabulary onto the same k -dimensional space such that the learnt vectors encode semantic content of the entities.

For every sequence of words w_{t-c}, \dots, w_{t+c} in, say document d_i , we wish to estimate the probability $p(w_t|d_i, w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c})$ of predicting the middle word in the sequence using the information about the document and the words in the context. We will estimate this probability using the vector representations for documents and words and learn vectors such that the probability of predicting the middle word in the context correctly is maximized.

2.3.2 Our Model

A document $d_i \in \mathbf{D}$, indexed by ‘ i ’, in our model is represented by a vector $\mathbf{v}_i^D \in \mathbb{R}^k$, which is also the i -th column of the matrix $\mathbf{D} = [\mathbf{v}_1^D, \dots, \mathbf{v}_{|\mathbf{D}|}^D] \in \mathbb{R}^{k \times |\mathbf{D}|}$. Similarly, a word indexed by ‘ i ’ in the vocabulary \mathbf{V} is represented by vector $\mathbf{v}_i^W \in \mathbb{R}^k$, which is also the i -th column of the matrix $\mathbf{W} = [\mathbf{v}_1^W, \dots, \mathbf{v}_{|\mathbf{V}|}^W] \in \mathbb{R}^{k \times |\mathbf{V}|}$.

Given a sequence $(w_{t-c}, \dots, w_{t+c})$ of $2c + 1$ words and the document it occurs in, our training objective is to maximize the probability of correctly predicting the middle word w_t using the surrounding context words $(w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c})$, which we now denote as *context* as a shorthand, and the information about the document in terms of their distributed vector representations. Therefore, our training objective is to maximize the probability $p(w_t|d_i, \text{context})$ of correctly predicting the middle word using the information about the surrounding words and the document the sequence occurs in.

To learn distributed word and document representations, we present a neural network model using which we,

1. Represent each word and document in the corpus by a k -dimensional distributed representation stored as vectors in the matrices \mathbf{W} and \mathbf{D} , respectively.
2. Estimate the probability of predicting the middle word in a sequence, given the document it occurs in, using the vector representation of the document and the words in the context.

3. Learn the word and document vectors simultaneously with the parameters of the function to estimate the probability.

The architecture for the proposed neural network is given in Fig. 2.5. Also note that the word vector representations learned and stored in the matrix W are universal representations and shared across all documents and contexts.

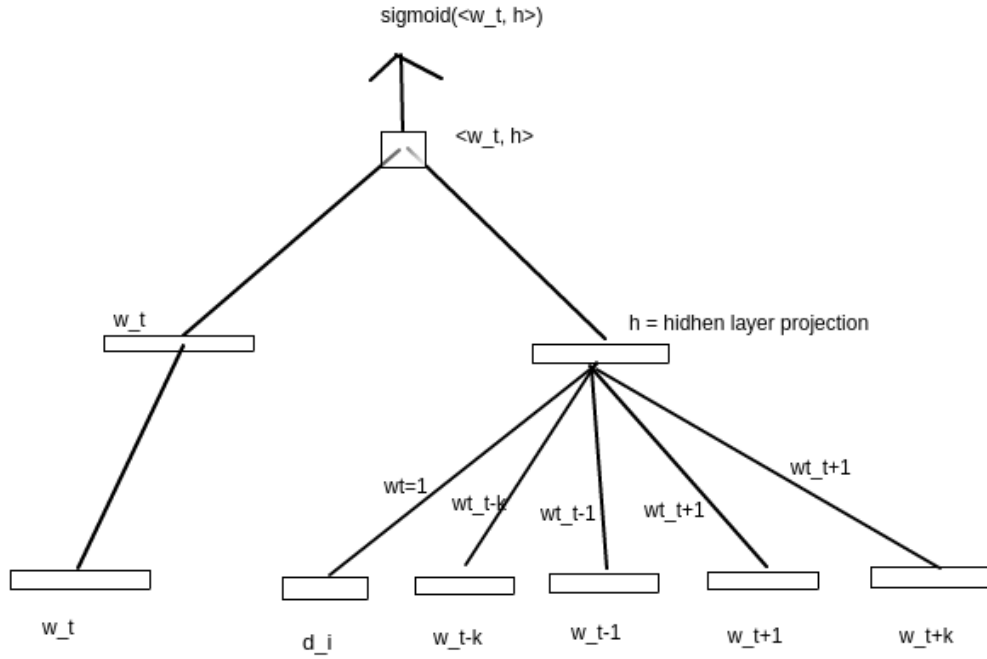


Figure 2.5: GloDETC : Neural Network Archtichture **TODO: Change figure**

2.3.2.1 Projection Layer (Context Representation)

We represent the context(words surrounding the middle word to be predicted) and the document together in the same projection layer, denoted by $h_c \in \mathbb{R}^k$, by taking a weighted sum of the corresponding vector representations. The weights for the context words $\Lambda = \{\lambda_i | i = \{t - c, \dots, t - 1, t + 1, \dots, t + c\}\}$ are kept universal for different sequences across the corpus as we expect the weights to learn some kind of syntactic quality of the language to better represent the context. Also the weight corresponding to the document vector is kept constant at 1 as we expected the document to have equal contribution to all sequences. This also gave the best

results. We also (unsuccessfully) experimented by taking matrix weights instead of scalar weights(λ_i) to learn better syntactic qualities of the language.

$$h_c = v_i^D + \lambda_{t-c} v_{t-c}^W + \dots + \lambda_{t-1} v_{t-1}^W + \lambda_{t+1} v_{t+1}^W + \lambda_{t+c} v_{t+c}^W \quad (2.10)$$

2.3.2.2 Estimating Prediction Probability

We expect in absence of any non-linearity that the projection layer vector should be aligned to the correct middle word of the sequence. Hence we estimate the probability of predicting the word w_t as the middle word in the following manner.

1. An output score $s_{w_i} \in \mathbb{R}$ for every w_i in the vocabulary is estimated by,

$$s_{w_i} = \sigma(v_{w_i}^W \cdot h_c) \quad (2.11)$$

where $\sigma(z) = \frac{1}{1+e^{-z}}$ is the standard sigmoid function.

2. After calculating the score for each of the word in the vocabulary, we use the *softmax* classifier to estimate the probability of predicting the actual correct word w_t as the middle word in the sequence,

$$p(w_t | d_i, w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c}) = \frac{e^{s_{w_t}}}{\sum_{i \in \mathbf{V}} e^{s_i}} \quad (2.12)$$

2.3.2.3 Training Objective

The training data \mathcal{T} , is composed of M training sequences each of which consists a $2c+1$ length sequence of words and the document index it belongs to. For example, $t = \{d_i^{(m)}, w_{t-c}^{(m)}, \dots, w_{t+c}^{(m)}\}$ represents the m^{th} training sequence in \mathcal{T} .

Given the training data \mathcal{T} , our objective is to learn an optimum parameter set $\Theta = (D, W, \Lambda)$ consisting of the document and word vector matrices and the projection layer weights for the context words, by maximizing the average log probability of estimating the middle word correctly in all the training word sequences where the

probability of estimating the middle word as w_i is given by Eq. 2.12. Therefore,

$$\hat{\Theta} = \arg \max_{\Theta} l(\mathcal{T}, \Theta) \quad (2.13)$$

$$l(\mathcal{T}, \Theta) = \frac{1}{M} \sum_{m=1}^M \log \left[p(w_t^{(m)} | d_i^{(m)}, w_{t-c}^{(m)}, \dots, w_{t-1}^{(m)}, w_{t+1}^{(m)}, \dots, w_{t+c}^{(m)}) \right] \quad (2.14)$$

To learn the optimize the above training objective, we can use the Stochastic Gradient Descent(SGD) algorithm to find gradient of the objective function (Eq. 2.14) w.r.t. to the individual parameters θ_i and apply the update rule as follows,

$$\theta_i^{(x)} = \theta_i^{(x-1)} + \gamma \frac{\partial l(\mathcal{T}, \Theta)}{\partial \theta_i} \quad (2.15)$$

where x is the current iteration number and γ is the learning rate. Also note that we add the gradient to $\theta_i^{(x)}$ because we wish to maximize the training objective. Updating the parameters for sufficient number of iterations should yield the optimum document and word vectors along with the weights for the neural network.

2.3.2.4 Noise Contrastive Estimation (NCE)

As we see in Eq 2.12, estimating the probability for each training word sequence requires a sweep through the whole vocabulary of size $|\mathbf{V}|$ which can be a very expensive computation given that typical vocabulary sizes range from a few tens to a few hundreds of thousand words for large datasets. Approaches to reduce this training time, such as, use of heirarchical soft-max [Morin and Bengio, 2005] and use of importance sampling to approximate the likelihood gradient [Bengio et al., 2003b], [Bengio and Senecal, 2008] have been proposed. Using heirarchical softmax reduces the training time from linear to logarithmic in vocabulary size but is considerable more involved and finding well-performing trees is not trivial. Also, though importance sampling provides substantial speedups, it suffers from stability problems.

Noise Contrastive Estimation (NCE) [Gutmann and Hyvärinen, 2012] is

method for fitting unnormalized probabilities by reducing the problem of *density estimation* to *probabilistic binary classification*. It has also been adapted to NPLM (Sec. 2.2.1) [Mnih and Teh, 2012] and learning word embeddings Mnih and Kavukcuoglu [2013] and shows significant improvements in training time with no considerable degradation in the quality of word vectors learned.

The basic idea of NCE is to train a logistic classifier to distinguish between the correct middle word in the given word sequence and corrupt samples from some “noise” distribution. Therefore, given a training sequence of the form $t = \{d_i^{(m)}, w_{t-c}^{(m)}, \dots, w_{t+c}^{(m)}\}$, our training objective now is to train a classifier such that it can distinguish between positive training sample $w_t^{(m)}$ as positive example and negative training samples w_x drawn from a noise distribution $P_n(w)$ as negative examples for the middle word given the surrounding words (context) and the document the sequence belongs to.

Our training data \mathcal{T} is now converted to a set of labeled sequences of the form $\{d_i^{(m)}, w_{t-c}^{(m)}, \dots, w_{t+c}^{(m)}, Y^{(m)} = 1\}_{m=1}^{m=M}$ where $Y = 1$ denotes that the sequence is a positive sample occurring in the corpus. For every positive training sequence t we also have n corrupt training sequences where in each of them only the middle word w_t has been replaced by a corrupt word sampled from the noise distribution $P_n(w)$ and the value of the label $Y = 0$. Therefore, for every positive training example there exists n negative training examples and the total number of training samples now in \mathcal{T} is $M + nM$. We now need to train a binary classifier such that, given a sequence of words and the document it belongs to, it can predict correctly whether the sequence is legitimate (correct value of the label indicator Y).

Given a training sequence we estimate the probability that the given sequence is positive using,

$$P(Y = 1 | d_i, w_{t-c}, \dots, w_{t+c}, \Theta) = \sigma(\mathbf{v}_{w_t}^W \cdot h_c) \quad (2.16)$$

where h_c is the projection layer vector calculated using Eq. 2.10. Similarly, the

probability of estimating that a given sequence is corrupt is given by,

$$P(Y = 0|d_i, w_{t-c}, \dots, w_{t+c}, \Theta) = 1 - \sigma(v_{w_t}^W \cdot h_c) \quad (2.17)$$

From Eq. 2.16 and Eq. 2.17 we get,

$$P(Y|d_i, w_{t-c}, \dots, w_{t+c}, \Theta) = [\sigma(v_{w_t}^W \cdot h_c)]^Y [1 - \sigma(v_{w_t}^W \cdot h_c)]^{1-Y} \quad (2.18)$$

As a shorthand notation, we would express the probability estimation in Eq. 2.18 as $P_\Theta(Y)$.

2.3.2.5 New Training Objective

Our new training objective involves maximizing the log-likelihood of observing the modified training data \mathcal{T} that includes the negative examples sampled from the noise distribution $P_n(w)$ along with the original positive training sequences.

$$\hat{\Theta} = \arg \max_{\Theta} l(\mathcal{T}, \Theta) \quad (2.19)$$

$$l(\mathcal{T}, \Theta) = \sum_{m=1}^{M+nM} \log P_\Theta(Y_m = Y^{(m)}) \quad (2.20)$$

where Y_m is the predicted label for the m -th training sequence and,

$$\log P_\Theta(Y_m = Y^{(m)}) = Y^{(m)} \log \sigma(v_{w_t^{(m)}}^W \cdot h_c^{(m)}) + (1 - Y^{(m)}) \log(1 - \sigma(v_{w_t^{(m)}}^W \cdot h_c^{(m)})) \quad (2.21)$$

2.3.2.6 Parameter Estimation

To learn the optimum parameters $\Theta = (D, W, \Lambda)$ we would maximize the log-likelihood of observing the training data given in Eq 2.19 using the Stochastic Gradient Descent(SGD) algorithm described below.

Firstly, for the SGD algorithm we would need to calculate the gradient of the log probability estimate (Eq 2.21) with respect to individual parameters $\theta \in \Theta$. The

derivative of the log probability estimate w.r.t. to a paramter $\theta \in \Theta$ is given by,

$$\frac{\partial \log P_{\Theta}(Y_m = Y^{(m)})}{\partial \theta} = \left[Y^{(m)} \frac{1}{\sigma(d^{(m)})} - (1 - Y^{(m)}) \frac{1}{(1 - \sigma(d^{(m)}))} \right] \frac{\partial \sigma(d^{(m)})}{\partial \theta} \quad (2.22)$$

$$= \left[Y^{(m)} \frac{1}{\sigma(d^{(m)})} - (1 - Y^{(m)}) \frac{1}{(1 - \sigma(d^{(m)}))} \right] [\sigma(d^{(m)})(1 - \sigma(d^{(m)}))] \frac{\partial d^{(m)}}{\partial \theta} \quad (2.23)$$

$$\frac{\partial \log P_{\Theta}(Y_m = Y^{(m)})}{\partial \theta} = [Y^{(m)} - \sigma(d^{(m)})] \frac{\partial d^{(m)}}{\partial \theta} \quad (2.24)$$

where $d^{(m)} = (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})$ is the dot-product of the projection layer vector with the word vector for the middle word. Therefore,

$$\frac{\partial \log P_{\Theta}(Y_m = Y^{(m)})}{\partial \theta} = [Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})] \frac{\partial (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})}{\partial \theta} \quad (2.25)$$

For any training sequence m , there are four types of parameters θ that need to be updated. Firstly, the document vector for $d_i^{(m)}$, the middle word vector for word $w_t^{(m)}$, word vectors for context words $w_{t+j}^{(m)}$ and the neural network weights λ_i . The derivate $\frac{\partial (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})}{\partial \theta}$ w.r.t. each of them is given by,

$$\frac{\partial (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})}{\partial \mathbf{v}_{d_i^{(m)}}^D} = \mathbf{v}_{w_t^{(m)}}^W \quad (2.26)$$

$$\frac{\partial (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})}{\partial \mathbf{v}_{w_t^{(m)}}^W} = \mathbf{h}_c^{(m)} \quad (2.27)$$

$$\frac{\partial (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})}{\partial \mathbf{v}_{w_{t+j}^{(m)}}^W} = \lambda_{t+j} * \mathbf{v}_{w_t^{(m)}}^W \quad (2.28)$$

$$\frac{\partial (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})}{\partial \lambda_{t+j}} = \mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{v}_{w_{t+j}^{(m)}}^W \quad (2.29)$$

Therefore the derivative of the log-probability estimate w.r.t. the

1. **Document Vector :**

$$\frac{\partial \log P_{\Theta}(Y_m = Y^{(m)})}{\partial \mathbf{v}_{d_i^{(m)}}^D} = \left[Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot h_c^{(m)}) \right] \mathbf{v}_{w_t^{(m)}}^W \quad (2.30)$$

2. **Middle Word :**

$$\frac{\partial \log P_{\Theta}(Y_m = Y^{(m)})}{\partial \mathbf{v}_{w_t^{(m)}}^W} = \left[Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot h_c^{(m)}) \right] h_c^{(m)} \quad (2.31)$$

3. **Context Word :**

$$\frac{\partial \log P_{\Theta}(Y_m = Y^{(m)})}{\partial \mathbf{v}_{w_{t+j}^{(m)}}^W} = \left[Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot h_c^{(m)}) \right] \lambda_{t+j} \mathbf{v}_{w_t^{(m)}}^W \quad (2.32)$$

4. **Neural Network Weight :**

$$\frac{\partial \log P_{\Theta}(Y_m = Y^{(m)})}{\partial \lambda_{t+j}} = \left[Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot h_c^{(m)}) \right] (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{v}_{w_{t+j}^{(m)}}^W) \quad (2.33)$$

According to the SGD algorithm, the update to be made to a parameter $\theta \in \Theta$ on observing a training sequence m is therefore given in Eq. 2.34. We also include L_2 regularization for the parameters as it helps in avoiding overfitting and restricts the parameters to blow up in value.

$$\theta^{(i+1)} \leftarrow \theta^{(i)} + \gamma \left[\frac{\partial \log P_{\Theta}(Y_m = Y^{(m)})}{\partial \theta^{(i)}} - \beta \theta^{(i)} \right] \quad (2.34)$$

here $\theta^{(i)}$ denotes the value of the parameter in the i -th iteration, θ^{i+1} is the value after the update, γ is the learning rate and β is the regularization constant. The update rules for the document vector, word vectors and the neural network weights are hence given by,

1. **Document Vector :**

$$(\mathbf{v}_{d_i^{(m)}}^D)^{(i+1)} = (\mathbf{v}_{d_i^{(m)}}^D)^{(i)} + \gamma \left[(Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot h_c^{(m)})) \mathbf{v}_{w_t^{(m)}}^W - \beta \mathbf{v}_{d_i^{(m)}}^D \right] \quad (2.35)$$

2. Middle Word Vector :

$$(\mathbf{v}_{w_t^{(m)}}^W)^{(i+1)} = (\mathbf{v}_{w_t^{(m)}}^W)^{(i)} + \gamma \left[(Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})) \mathbf{h}_c^{(m)} - \beta \mathbf{v}_{w_t^{(m)}}^W \right] \quad (2.36)$$

3. Context Word Vectors :

$$(\mathbf{v}_{w_{t+j}^{(m)}}^W)^{(i+1)} = (\mathbf{v}_{w_{t+j}^{(m)}}^W)^{(i)} + \gamma \left[(Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})) \lambda_{t+j} \mathbf{v}_{w_t^{(m)}}^W - \beta \mathbf{v}_{w_{t+j}^{(m)}}^W \right] \quad (2.37)$$

4. Neural Network Weights :

$$\lambda_{t+j}^{(i+1)} = \lambda_{t+j}^{(i)} + \gamma * \left[(Y^{(m)} - \sigma(\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{h}_c^{(m)})) (\mathbf{v}_{w_t^{(m)}}^W \cdot \mathbf{v}_{w_{t+j}^{(m)}}^W) - \beta \lambda_{t+j} \right] \quad (2.38)$$

To learn the vectors \mathbf{D} , \mathbf{W} and weights Λ we initialize them to small random vectors and scalars respectively and using Eqs. 2.35, 2.36, 2.37 and 2.38 we iterate through the training data, making the appropriate updates for a fixed number of epochs that is learnt using the development data. For each training sequence we make one update to the document vector, $2c + 1$ updates for the word vectors in the sequence and $2c$ updates for the neural network weights, where c is the window length we consider while training.

2.3.2.7 Hyper-parameters

Our model, like any other, has hyper-parameters that need to be tuned for optimum model performance and learning high quality document and word vectors and the model parameters for achieving the the best accuracies in the task. Below we describe the hyper-parameters in our model and the effect they have on learning document representations.

1. **Embedding Dimensionality (k)** : The most important hyper-parameter in our model is the size of the document and word embdding vectors k . The embedding dimensionality needs to be big enough such that the document

vectors can encode the different semantic topics across the corpus but shouldn't be very large so that it introduces noise in the vectors.

2. **Window Size (c)** : The length of the sequence or the window size c , that we consider as context surrounding a word plays an important role in the vectors that are learnt. While a smaller window could result in the negligence of important/similar words that surround the middle word, a large window could introduce noise in the context that can deteriorate the performance of the model.
3. **Number of Negative Samples (n)** : In NCE, the number of negative samples introduced in the training data per positive example plays an important role in deciding the trade-off between learning better word density distribution with larger n while smaller n leads to lesser training times. Hence, the number of negative samples introduced needs to be tuned using the development data.
4. **Number of Epochs** : The number of times we need to loop through the training data to learn the vectors needs to be optimized to prevent overfitting with large epochs while at the same time learn high quality representations.

Chapter 3

Multi-Label Text Categorization

TODO: Introducing multi-label text classification. Multiple algos as in related work.

TODO: Section on our model. why logisitic regression. We go ahead with logisitic regression. Similar to Matrix factorization. abides by the idea of embddings. helps learn correlations among cats using low-dimensional embeddings. **TODO:** sub1 : Dataset details. sub2 : model details. training objective updates. **TODO:** sub3 : Similarity to matrix factorization and relational learning.

In this chapter we will give an overview of the training data required for the document categorization task and present the multinomial logisitic regression algorithm in context of the multi-label text categorization, discuss its advantages and similarity to matrix factorization and relational learning.

3.1 Logisitic Regression (LR) for Multi-label Document Categorization

Introduced by [Hosmer and Lemeshow, 1989], Logisitic Regression (LR) is a probabilistic binary classification regression model that, given labeled binary data, performs regression over the data and learns weight vectors to predict whether a given data point belongs to the positive or the negative class. The probability of the data

point to belong in a class is estimated using the *logistic (sigmoid) function*, hence the name logistic regression.

Logistic Regression, though is a technique to discriminate between two categories can be easily extended to classification between multiple categories which is then referred to as Multinomial Logistic Regression. Though we use multinomial logistic regression for our task of multi-label text classification, for the sake of brevity we would refer to our algorithm as logistic regression.

In the sections below we describe the training data required for the task of multi-label document classification, the logistic regression model as modified for the task and also its similarity to relational learning.

3.1.1 Training Data

The training data \mathcal{T} is composed of a set of documents \mathbf{D} , set of categories \mathbf{C} and data about in what categories do each of the documents belong to.

Document-Category Data : Each document d_i in \mathbf{D} belongs to atleast one category from \mathbf{C} . To store this relational data between the documents and the categories, we create a database \mathcal{D} in which for the m -th training instance we store tuple of the form $\{d_i^{(m)}, c_j^{(m)}, y^{(m)}\}_{m=1}^{m=T}$ where $y^{(m)} \in \{0, 1\}$ denotes whether the document $d_i^{(m)}$ belongs the category $c_j^{(m)}$ or not.

Mostly the data about document categories is given such that it is known what categories do the documents belong to, without conclusive information about whether a document necessarily does not belong to a particular category. In such cases, if we assume the given data to be complete, then along with positive data examples of the form, $\{d_i, c_j, 1\}$, we introduce negative samples, $\{d_i, c_k, 0\}$ for every category c_k each document d_i does not belong to. If the document-category data is viewed as a matrix with documents as rows and categories as columns, then in such case, we would only observe positive examples (1) in matrix but at sparse locations. To make the training data complete in such cases, we would fill the matrix with negative examples (0) at every empty location.

Document Rerresentations : Along with the document-category data, the training data also composes of the document representations in the form of either bag-of-words representations or distributed document embeddings as learnt in Sec. 2.3. Therefore for every document $d_i \in \mathbf{D}$ indexed by i , we have a vector representation $\mathbf{v}_i^D \in \mathbb{R}^k$ of the document.

3.1.2 Logistic Regression Model

For multi-label document categorization we extend the standard logisitic regression (which is a binary classification algorithm) to,

1. Train from multi-labeled document-category data data
2. Given a document-category pair $\{d_i, c_j\}$, estimate the probability of the document d_i belonging to the category c_j .

As we explained in Sec. 2.3, we learn low-rank distributed vector representation for every document in the corpus. Similarly, for multi-label logisitic regression, we represent each category $c_i \in \mathbf{C}$ using a low-rank embedding $\mathbf{v}_{c_i}^C \in \mathbb{R}^k$ of the same dimensionality as the document embeddings, k . Similar to \mathbf{D} , we stack these category embeddings as columns in the matrix $\mathbf{C} \in \mathbb{R}^{k \times |\mathbf{C}|}$.

Given a document-category tuple of the form $\{d_i, c_j\}$, we estimate the probability of the document belonging to the category ($y = 1$) using the logisitic function as,

$$P(y = 1 | d_i, c_j, \mathbf{D}, \mathbf{C}) = \sigma(\mathbf{v}_{d_i}^D \cdot \mathbf{v}_{c_j}^C) \quad (3.1)$$

This model is similar to the standard logisitic regression (LR) as in standard LR for binary classification, we learn a universal weight vector \mathbf{w} that is used to estimate the probability in Eq. 3.1 instead of $\mathbf{v}_{c_j}^C$. Here, because we have multiple categories, we learn multiple weight vectors (category embeddings) for each category separately and hence perform multiple binary classifications.

3.1.2.1 Training Objective

As explained in Sec. 3.1.1, the training data \mathcal{T} , is composed of T tuples of the form $\{d_i^{(m)}, c_j^{(m)}, y^{(m)}\}$. Our training objective involves learning optimum category embeddings such that for any unobserved document $d_x \notin \mathbf{D}$, we should be able to predict the categories it belongs to.

For the m -th training instance $\{d_i^{(m)}, c_j^{(m)}, y^{(m)}\}$, we denote the prediction that whether the document $d_i^{(m)}$ belongs the category $c_j^{(m)}$ by y_m . Therefore, if we estimate $d_i^{(m)}$ belongs $c_j^{(m)}$, $y_m = 1$ otherwise $y_m = 0$. Using Eq. 3.1,

$$P(y_m = 1 | d_i^{(m)}, c_j^{(m)}, \mathbf{D}, \mathbf{C}) = \sigma(\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C) \quad (3.2)$$

We denote the above probability estimate as $P_{\mathbf{D}, \mathbf{C}}(y_m = 1)$ for brevity. Therefore,

$$P_{\mathbf{D}, \mathbf{C}}(y_m = 0) = 1 - \sigma(\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C) \quad (3.3)$$

$$P_{\mathbf{D}, \mathbf{C}}(y_m) = \sigma(\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C)^{y_m} (1 - \sigma(\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C))^{1-y_m} \quad (3.4)$$

To learn the optimum parameter set $\Theta = (\mathbf{C})$ consisting of the set of category embeddings, we would maximize the log-likelihood of observing the training data,

$$\hat{\Theta} = \arg \max_{\Theta} l(\mathcal{T}, \Theta) \quad (3.5)$$

$$l(\mathcal{T}, \Theta) = \sum_{m=1}^T \log P_{\mathbf{D}, \mathbf{C}}(y_m = y^{(m)}) \quad (3.6)$$

where,

$$\log P_{\mathbf{D}, \mathbf{C}}(y_m = y^{(m)}) = y^{(m)} \log \sigma(\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C) + (1 - y^{(m)}) \log (1 - \sigma(\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C)) \quad (3.7)$$

3.1.2.2 Parameter Estimation

To learn the optimum parameters $\Theta = (\mathbf{C})$ we would maximize the log-likelihood of observing the training data given in Eq 3.6 using the Stochastic Gradient De-

scent(SGD) algorithm as described earlier in Sec 2.3.2.6. We first need to calculate the gradient of the log probability estimate in Eq. ?? with respect to the category embeddings which is given by,

$$\frac{\partial \log P_{D,C}(y_m = y^{(m)})}{\partial \mathbf{v}_{c_j^{(m)}}^C} = \left[y^{(m)} \frac{1}{\sigma(s^{(m)})} - (1 - y^{(m)}) \frac{1}{(1 - \sigma(s^{(m)}))} \right] \frac{\partial \sigma(s^{(m)})}{\partial \mathbf{v}_{c_j^{(m)}}^C} \quad (3.8)$$

$$= \left[y^{(m)} \frac{1}{\sigma(s^{(m)})} - (1 - y^{(m)}) \frac{1}{(1 - \sigma(s^{(m)}))} \right] [\sigma(s^{(m)})(1 - \sigma(s^{(m)}))] \frac{\partial s^{(m)}}{\partial \mathbf{v}_{c_j^{(m)}}^C} \quad (3.9)$$

$$\frac{\partial \log P_{\Theta}(y_m = y^{(m)})}{\partial \mathbf{v}_{c_j^{(m)}}^C} = [y^{(m)} - \sigma(s^{(m)})] \frac{\partial s^{(m)}}{\partial \mathbf{v}_{c_j^{(m)}}^C} \quad (3.10)$$

where, $s^{(m)} = (\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C)$ is the *pre-sigmoid activation*. Therefore,

$$\frac{\partial \log P_{\Theta}(y_m = y^{(m)})}{\partial \mathbf{v}_{c_j^{(m)}}^C} = \left[y^{(m)} - \sigma(\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C) \right] \frac{\partial (\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C)}{\partial \mathbf{v}_{c_j^{(m)}}^C} \quad (3.11)$$

According to the SGD algorithm and Eq. 2.34, the update to be made to the category embedding on observing the m -th training instance is given by Eq. 3.12. We also include L2 regularization for the category embeddings as it helps in avoiding overfitting and restricts the embeddings to blow up in value.

$$(\mathbf{v}_{c_j^{(m)}}^C)^{(i+1)} = (\mathbf{v}_{c_j^{(m)}}^C)^{(i)} + \gamma * \left[\left[y^{(m)} - \sigma(\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C) \right] \frac{\partial (\mathbf{v}_{d_i^{(m)}}^D \cdot \mathbf{v}_{c_j^{(m)}}^C)}{\partial \mathbf{v}_{c_j^{(m)}}^C} - \beta (\mathbf{v}_{c_j^{(m)}}^C)^{(i)} \right] \quad (3.12)$$

Bibliography

- Yoshua Bengio and J-S Senecal. Adaptive importance sampling to accelerate training of a neural probabilistic language model. *Neural Networks, IEEE Transactions on*, 19(4):713–722, 2008.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *The Journal of Machine Learning Research*, 3:1137–1155, 2003a.
- Yoshua Bengio, Jean-Sébastien Senécal, et al. Quick training of probabilistic neural nets by importance sampling. In *AISTATS Conference*, 2003b.
- Léon Bottou. From machine learning to machine reasoning. *Machine learning*, 94(2):133–149, 2014.
- Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, 12:2493–2537, 2011.
- W Cooper, Aitao Chen, and F Gey. Full text retrieval based on probabilistic equations with coefficients fitted by logistic regression. *NIST SPECIAL PUBLICATION SP*, pages 57–57, 1994.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- Koby Crammer and Yoram Singer. A new family of online algorithms for category ranking. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 151–158. ACM, 2002.
- Scott C. Deerwester, Susan T Dumais, Thomas K. Landauer, George W. Furnas, and Richard A. Harshman. Indexing by latent semantic analysis. *JAsIs*, 41(6):391–407, 1990.
- André Elisseeff and Jason Weston. A kernel method for multi-labelled classification. In *Advances in neural information processing systems*, pages 681–687, 2001.
- Norbert Fuhr, Stephan Hartmann, Gerhard Lustig, Michael Schwantner, Kostas Tzeras, and Gerhard Knorz. *AIR, X: a rule based multistage indexing system for large subject fields*. Citeseer, 1991.
- Nadia Ghamrawi and Andrew McCallum. Collective multi-label classification. In *Proceedings of the 14th ACM international conference on Information and knowledge management*, pages 195–200. ACM, 2005.

- Edward Grefenstette, Georgiana Dinu, Yao-Zhong Zhang, Mehrnoosh Sadrzadeh, and Marco Baroni. Multi-step regression learning for compositional distributional semantics. *arXiv preprint arXiv:1301.6939*, 2013.
- Michael U Gutmann and Aapo Hyvärinen. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *The Journal of Machine Learning Research*, 13(1):307–361, 2012.
- David W Hosmer and Stanley Lemeshow. Applied logistic regression. 1989. *New York: Johns Wiley & Sons*, 1989.
- Thorsten Joachims. *Text categorization with support vector machines: Learning with many relevant features*. Springer, 1998.
- Omer Levy and Yoav Goldberg. Dependencybased word embeddings. *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, 2: 302–308, 2014.
- David D Lewis. Text representation for intelligent text retrieval: a classification-oriented view. *Text-based intelligent systems: current research and practice in information extraction and retrieval*, pages 179–197, 1992a.
- David D Lewis and Marc Ringuette. A comparison of two learning algorithms for text categorization. In *Third annual symposium on document analysis and information retrieval*, volume 33, pages 81–93, 1994.
- David D Lewis, Robert E Schapire, James P Callan, and Ron Papka. Training algorithms for linear text classifiers. In *Proceedings of the 19th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 298–306. ACM, 1996.
- David Dolan Lewis. *Representation and learning in information retrieval*. PhD thesis, University of Massachusetts, 1992b.
- Yi Liu, Rong Jin, and Liu Yang. Semi-supervised multi-label learning by constrained non-negative matrix factorization. In *Proceedings of the national conference on artificial intelligence*, volume 21, page 421. Menlo Park, CA; Cambridge, MA; London; AAI Press; MIT Press; 1999, 2006.
- Andrew McCallum. Multi-label text classification with a mixture model trained by em. 1999.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013a.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119, 2013b.
- Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word representations. In *HLT-NAACL*, pages 746–751, 2013c.

- Jeff Mitchell and Mirella Lapata. Composition in distributional models of semantics. *Cognitive science*, 34(8):1388–1429, 2010.
- Andriy Mnih and Koray Kavukcuoglu. Learning word embeddings efficiently with noise-contrastive estimation. In *Advances in Neural Information Processing Systems*, pages 2265–2273, 2013.
- Andriy Mnih and Yee Whye Teh. A fast and simple algorithm for training neural probabilistic language models. *arXiv preprint arXiv:1206.6426*, 2012.
- Frederic Morin and Yoshua Bengio. Hierarchical probabilistic neural network language model. In *Proceedings of the international workshop on artificial intelligence and statistics*, pages 246–252. Citeseer, 2005.
- Isabelle Moulinier, Gailius Raskinis, and J Ganascia. Text categorization: a symbolic approach. In *proceedings of the fifth annual symposium on document analysis and information retrieval*, pages 87–99, 1996.
- Hwee Tou Ng, Wei Boon Goh, and Kok Leong Low. Feature selection, perceptron learning, and a usability case study for text categorization. In *ACM SIGIR Forum*, volume 31, pages 67–73. ACM, 1997.
- Kamal Nigam, John Lafferty, and Andrew McCallum. Using maximum entropy for text classification. In *IJCAI-99 workshop on machine learning for information filtering*, volume 1, pages 61–67, 1999.
- Kamal Nigam, Andrew Kachites McCallum, Sebastian Thrun, and Tom Mitchell. Text classification from labeled and unlabeled documents using em. *Machine learning*, 39(2-3):103–134, 2000.
- J Ross Quinlan. Learning efficient classification procedures and their application to chess end games. In *Machine learning*, pages 463–482. Springer, 1983.
- J. Ross Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.
- Timothy N Rubin, America Chambers, Padhraic Smyth, and Mark Steyvers. Statistical topic models for multi-label document classification. *Machine Learning*, 88(1-2):157–208, 2012.
- Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval. *Information processing & management*, 24(5):513–523, 1988.
- Gerard Salton and Chung-Shu Yang. On the specification of term values in automatic indexing. *Journal of documentation*, 29(4):351–372, 1973.
- Richard Socher, Alex Perelygin, Jean Y Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*, volume 1631, page 1642. Citeseer, 2013.
- Richard M Tong and Lee A Appelbaum. Machine learning for knowledge-based document routing (a report on the trec-2 experiment). *NIST SPECIAL PUBLICATION SP*, pages 253–253, 1994.

- Joseph Turian, Lev Ratinov, and Yoshua Bengio. Word representations: a simple and general method for semi-supervised learning. In *Proceedings of the 48th annual meeting of the association for computational linguistics*, pages 384–394. Association for Computational Linguistics, 2010.
- Naonori Ueda and Kazumi Saito. Parametric mixture models for multi-labeled text. In *Advances in neural information processing systems*, pages 721–728, 2002.
- Vladimir Vapnik. *The nature of statistical learning theory*. Springer Science & Business Media, 2000.
- Erik Wiener, Jan O Pedersen, Andreas S Weigend, et al. A neural network approach to topic spotting. In *Proceedings of SDAIR-95, 4th annual symposium on document analysis and information retrieval*, pages 317–332. Citeseer, 1995.
- Yiming Yang. Expert network: Effective and efficient learning from human decisions in text categorization and retrieval. In *Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 13–22. Springer-Verlag New York, Inc., 1994.
- Yiming Yang and Christopher G Chute. A linear least squares fit mapping method for information retrieval from natural language texts. In *Proceedings of the 14th conference on Computational linguistics-Volume 2*, pages 447–453. Association for Computational Linguistics, 1992.
- Ainur Yessenalina and Claire Cardie. Compositional matrix-space models for sentiment analysis. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 172–182. Association for Computational Linguistics, 2011.
- Fabio Massimo Zanzotto, Ioannis Korkontzelos, Francesca Fallucchi, and Suresh Manandhar. Estimating linear models for compositional distributional semantics. In *Proceedings of the 23rd International Conference on Computational Linguistics*, pages 1263–1271. Association for Computational Linguistics, 2010.
- Min-Ling Zhang and Zhi-Hua Zhou. A k-nearest neighbor based algorithm for multi-label classification. In *Granular Computing, 2005 IEEE International Conference on*, volume 2, pages 718–721. IEEE, 2005.
- Min-Ling Zhang and Zhi-Hua Zhou. Ml-knn: A lazy learning approach to multi-label learning. *Pattern recognition*, 40(7):2038–2048, 2007.