

UNIVERSITY OF ARIZONA

---

# Application of Logistic Regression in Automatic Reading Comprehension and Detection of Grammatical Structures Based on Difficulty

---

*Author:*

Ragheb AL-GHEZI

*Supervisor:*

Dr. Robert HENDERSON

*A report submitted in fulfillment of the requirements  
for the degree of Masters of Science*

*in the*

Human Language Technology  
Department of Linguistics

March 5, 2019



UNIVERSITY OF ARIZONA

*Abstract*Faculty Name  
Department of Linguistics

Masters of Science

**Application of Logistic Regression in Automatic Reading Comprehension and  
Detection of Grammatical Structures Based on Difficulty**

by Ragheb AL-GHEZI

In this study we examine the use of a probabilistic classification algorithm called Logistic Regression on two novel Natural Language Processing (NLP) tasks: Automatic Reading Comprehension and Detection of Grammatical Difficulty in English Sentences. Logistic Regression has been used extensively in machine learning and NLP research over the past decades, but it gains popularity as it is considered the building block of Artificial Neural Network (ANN). Both tasks studied in this report show a unique use to multinomial logistic. In automatic reading comprehension task, given a reading passage and a question, the goal is to extract an answer from the passage to the question if there exists one, or abstain from answering if the passage does not have an answer. For this purpose, a multinomial logistic regressor is used on two stages: a) find the sentence that most likely contains the answer (null-answer is treated as a sentence at this stage). b) within the best candidate sentence, another logistic regressor is used to extract the constituent that most likely represent the answer. The experiment results show that using simply designed features a logistic regressor score 0.71 F1 score and 0.40 F1 in stage 1 and 2 respectively, making it as a potential competitor to other state-of-the-art advanced neural network architectures. In the second task, detection of grammatical difficulty, a multinomial logistic regressor is applied in a semi-supervised manner on a low-resourced data to categorize the difficulty of English grammatical structures according to Common European Framework Reference (CEFR) guidelines. The method used shows a significant improvement from 0.67 F1 to 0.79 F1.



## *Acknowledgements*

The acknowledgments and the people to thank go here, don't forget to include your project advisor...



# Contents

<b>Abstract</b>	<b>iii</b>
<b>Acknowledgements</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background	1
1.2 Goals	1
1.3 Outline	1
<b>2 Theory</b>	<b>3</b>
2.1 Introduction	3
2.2 Classification as a Machine Learning Task	3
2.3 Feature Extraction	4
2.3.1 Frequency-based Methods	4
2.3.2 Dimensionality reduction	6
2.3.3 Distributional Method	8
2.4 Logistic Regression	8
2.4.1 Classification Functions	9
2.4.2 Objective Functions	10
2.4.3 Optimization	10
2.5 Evaluation Metrics	11
2.5.1 Confusion Matrix	11
2.5.2 Precision, Recall and F-Score	11
2.6 Conclusion	12
<b>3 The Use of Semi-supervised Learning in Classification of English Grammatical Structures</b>	<b>13</b>
3.1 Introduction	13
3.2 Related Works	14
3.3 Method	15
3.3.1 Dataset: English Grammar Profile	15
3.3.2 Multinomial Logistic Regression	15
3.3.3 Feature Design	16
3.4 Experimental Results	17
3.4.1 First Phase	17
3.4.2 Second Phase	18
3.5 Conclusion	19

<b>4</b>	<b>Two-stage Classification Method for Automatic Reading Comprehension</b>	<b>21</b>
4.1	Introduction . . . . .	21
4.2	Related Work . . . . .	22
4.3	Baseline Implementation . . . . .	23
4.3.1	Feature Design . . . . .	23
4.3.2	Training and Result . . . . .	25
4.4	Stage Two: Predicting the Answer Span . . . . .	26
4.4.1	Iteration 1 . . . . .	26
4.4.2	Error Analysis: Iteration 1 . . . . .	27
4.4.3	Iteration 2 . . . . .	28
4.5	Conclusion . . . . .	28
	<b>Bibliography</b>	<b>31</b>



# List of Figures

2.1	Figure Showing Workflow of Machine Learning . . . . .	4
2.2	Confusion Mtrix Example . . . . .	11
3.1	2d projection of the three classes A, B and C using K-means Clustering and PCA . . . . .	16
3.2	Confusion matrix illustrating the performance of LR-BoW-Word model with class weight to deal with data imbalance . . . . .	17
3.3	Confusion matrix illustrating the performance of LR-BoW-Word model after augmentation . . . . .	19
4.1	Flowchart illustrating the two-stage classification approach . . . . .	24
4.2	Confusion matrix shows which classes were predicted correctly . . . . .	25
4.3	Confusion Matrix illustrating the first iteration of Stage2 LR. . . . .	27
4.4	Confusion Matrix illustrating the Second iteration of Stage2 LR. . . . .	29



# List of Tables

2.1	Bag-of-words feature example . . . . .	5
3.1	An excerpt of English sentences and their corresponding CEFR difficulty levels . . . . .	14
3.2	Comparison of Precision-Recall - Phase One . . . . .	18
4.1	This table shows the results of running a multinomial regularized logistic regression. The class column represents the index of sentences within the paragraph, and -1 represents the unanswerable question case. Unpredicted classes are removed. . . . .	26



*For/Dedicated to/To my...*



## Chapter 1

# Introduction

### 1.1 Background

Lum is a company developing customized machine reading technology using machine learning and natural language processing. Their technology is served to a variety of medical, agricultural and educational domains and sectors. During the internship, the researcher was assigned two main NLP tasks that contribute to solving problems in teaching foreign languages using online distance learning. The two NLP tasks are automatic reading comprehension and detecting grammatical difficulty of English sentences according to CEFR. The goal of the first task is to design and implement a system that, given a reading passage and a set of questions, can automatically find answers to those questions from within the reading passage. The motivation for such a system is that it can help human language instructors create automatically graded assignment and quiz activities in virtual learning environments. The purpose of the second task is to create a system that can automatically categorize a given English text into three levels of grammatical difficulty: Elementary (A), Intermediate (B), and Advanced (C). In this report, the researcher investigates the possibility to use NLP and ML algorithms to solve these two problems. The report also covers the entire process: data description and analysis, literature survey, feature design, and engineering and evaluating the performance of the algorithm.

### 1.2 Goals

One goal of this internship report is to examine the suitability of using a machine learning linear classification algorithm called logistic regression on two novel educational tasks: automatic reading comprehension and detection of grammatical difficulty. Another goal is to explore and compare different methods of feature extraction and feature engineering to each of these tasks. Furthermore, examine the use of logistic regression in a semi-supervised fashion to bootstrap the scarcity of training data and improve the overall quality of generalization.

### 1.3 Outline

The first chapter gives a general background about the internship that the researcher undertakes as well as the goals intended to achieve through this internship report.

The second chapter is dedicated to reviewing the theory of machine learning in the context of natural language processing. It also defines classification as a supervised machine learning task and logistic regression as one popular linear classification algorithm. The chapter also includes the conventional methods of computational representation of human languages such as bag-of-words, term-frequency-inverse-document-frequency, and word embeddings. Finally, it concludes by showing the metrics used to evaluate classification systems.

The third chapter introduces the task of detecting grammatical errors in English sentence as a classification problem and explains the use of multinomial logistic regression in a semi-supervised fashion to solve it. The chapter also briefly surveys the most important studies in the literature that tackled similar textual classification tasks. Finally, it shows the experiments, the results of applying the algorithm of the data collected as well as the conclusions.

The fourth chapter introduces the task of automatic reading comprehension and the data set used for this purpose. It also discusses the common approaches tackling these tasks in the literature. Next, the chapter uses the different feature engineering used to create a numerical representation of text. Finally, the chapter is concluded with the result and evaluation as well as sections dedicated to error analysis and future work respectively.

The fifth chapter contains the overall conclusion of the internship report.



## Chapter 2

# Theory

### 2.1 Introduction

In this chapter, we explain the fundamentals of machine learning theory, with emphasis on supervised classification task. We start by describing the three primary stages of the machine learning process: feature extraction, classifier training, and testing and evaluation. We also discuss two types of feature engineering methods: frequency-based and distributional in addition to dimensionality reduction methods. Next, we introduce a linear classification technique called logistic regression and explain its mathematical foundations. Finally, we discuss the metrics used to evaluate the performance of supervised classification tasks.

### 2.2 Classification as a Machine Learning Task

Generally, there are two types of supervised learning tasks according to the data type of the predicted output: regression and classification. When the predicted output is a continuous value such as predicting the price of a house or the temperature, the task is called regression. A classification task, on the other hand, is when the predicted output is a discrete value such as predicting the sentiment of a document (negative or positive) or predicting whether it will rain (rainy, sunny, cloudy).

The goal of classification task is to select a correct class for a given input, or more generally a task of “assigning objects from a universe to two or more classes or categories” (Manning, Manning, and Schütze, 1999). A classifier is an algorithm that quantitatively models the relationship between a set of inputs and their associated outputs such that it learns to label new data. In the task of spam filtering, for instance, input data are texts and the output are binary labels (0 or 1) representing whether or not a document is spam.

Any classification task has two main phases: training and prediction. During training, a feature extractor is used to convert each input into a set of features. These features are designed to capture the basic information about each input. Pairs of feature sets and their corresponding labels are, then, fed into the machine learning algorithm to generate a model. In the prediction phase, the trained model is used on unseen data to predict the labels.

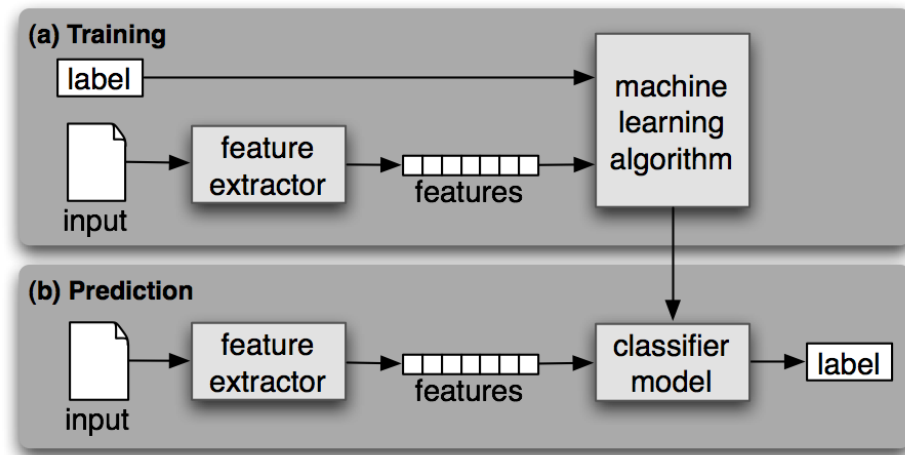


FIGURE 2.1: Figure Showing Workflow of Machine Learning

## 2.3 Feature Extraction

The first step in any machine learning task is feature extraction in which a text input is transformed into a set of attributes that characterize it. An example of a feature of a negative sentiment document may be the presence of the word *terrible*. Choosing what features to use when representing a text or a document not only has a major impact on the entire classification process, but it is also arguably the most difficult phase in any natural language processing task. The reasons go to several fundamental questions on how to represent language computationally. For example, on what level should we represent a text: *character-wise*, *word-wise* or *sentence-wise*, or even more specifically morphologically, etymologically or phonologically? Do we treat a word e.g. *try* and all its morphological variations e.g. *tries*, *trying*, *tried* similarly or differently? How to represent and model the semantics of human language computationally? Questions like these and many others are the core of natural language processing. Nevertheless, several automatic feature extraction methods have developed over the past decades that can be categorized into frequency-based or distributional, and we briefly explore some of the most common ones.

### 2.3.1 Frequency-based Methods

Frequency-based methods such as Bag-of-Words (BOW) and tfidf rely on number of occurrence of tokens in the text. They have been used extensively in a wide range of NLP applications, and each of them has advantages and disadvantages.

#### Bag-of-Words Representation

The bag-of-words (BOW) representation is the most straightforward and most intuitive representation. It represents each document by a vector whose component corresponds to the number of occurrences of a word in the document. To illustrate, the BOW vector representations of a data set of two documents (1) and (2) are shown in table 2.1:

1. The apple is better than the orange.

TABLE 2.1: Bag-of-words feature example

	the	apple	is	better	than	orange	summar	winter
Vector 1	2	1	1	1	1	1	0	0
Vector 1	0	1	1	1	0	1	1	1

2. The summar is better than winter.

Despite the simplicity of BOW method, it has two main disadvantages. One major disadvantage is that BOW method pays no respect to the order of tokens in the document, and this can give similar representations to semantically different documents. For instance, the sentences *the wolf eats the goat* and *the goat eats the wolf* have similar vector representation as they have exact same words. Another downside to BOW is that features with high frequency are more important than those of less frequency. We can see in 2.1 that the component corresponding to the first dimension (the) of vector 1 is bigger than in that of vector 2 because document 1 has two occurrence of the word *the* while the document 2 has only one occurrence.

To address the problem of order, several studies have attempted to count the occurrence of sentences instead of words (Fuhr and Buckley, 1991) (Tzeras and Hartmann, 1993). Although the sentences offer the advantage of better preserving the semantic relationship of the words, their use as features failed in practice. According to Lewis (Lewis, 1992), this representation is penalized by a large number of possible combinations which leads to low and too random frequencies. A more practical approach is representation by n-grams. It consists of breaking the text into moving sequences of  $n$  consecutive tokens. For example, the trigram segmentation of the word "language" gives the following 3-grams: lan, ang, ngu, gua, uag, age. The notion of n-grams was introduced by (Shannon, 1948); who was interested in predicting the appearance of certain characters according to the other characters. Since then, n-grams have been used in many areas such as speech recognition and information retrieval. Using n-gram on a character-level or word-level gives a contextual representation of the text. A modification to BOW to overcome the bias of high frequency tokens result in devising a method called Term Frequency Inverse Term Frequency (tfidf).

### TFIDF Representation

TFIDF has two main parts. Term Frequency is proportional to the frequency of the term in the document (local weighting). It can be used as is or in several variations (Singhal, Mitra, and Buckley, 1997) (Sable and Church, 2001).

$$tf_{ij} = f(t_i, d_j)$$

$$tf_{ij} = 1 + \log(f(t_i, d_j))$$

$$tf_{ij} = 0.5 + 0.5 \frac{f(t_i, d_j)}{\max_{t_i \in d_j} f(t_i, d_j)}$$

where  $f(t_i, d_j)$  is the term frequency of document  $j$ . Variations include logarithmic scaling of frequency counting or normalization.

Inverse Document Frequency measures the importance of a term throughout the collection (overall weighting). A term that often appears in the document base should not have the same impact as a less frequent term. Indeed, the terms that appear in the majority of documents do not have any discriminating power to distinguish documents from each other and must, therefore, have low weightings. The IDF weighting is inversely proportional to the number of documents containing the term to be weighted. Thus, the more the term appears in several documents, the less it is discriminating and is assigned a low weighting. The IDF weighting is generally expressed as follows:

$$idf(t_i) = \log \left( \frac{N}{df(t_i)} \right)$$

where  $df(t_i)$  is the term frequency of feature  $i$  and  $N$  is the number of documents in the corpus.

The TFIDF weighting combines the two weightings TF and IDF in order to provide a better approximation of the importance of a term in a document. According to this weighting, for a term to be important in a document, it must frequently appear in the document and rarely in other documents. This weighting is given by the product of the local weighting of the term in the document by its overall weighting in all the documents of the corpus.

$$tfidf(t_i, d_j) = tf_{ij} \times \log \left( \frac{N}{df(t_i)} \right)$$

Frequency-based text representation methods explored so far, despite their usefulness in a variety of tasks in natural language processing, share sparsity as one problem in common. The number of features generated using these methods can easily exceed the tens of thousands, and that does not only negatively influence the categorization process, but it is also very computationally expensive in terms of hardware resources. In addition, the higher the dimensions of the features are, the weaker the features become. This problem is known as the *curse of dimensionality* (Bellman, 2015). To remedy this issue, techniques, borrowed from the field of information theory and linear algebra, have been developed to reduce the feature space dimensionality without losing much information.

### 2.3.2 Dimensionality reduction

There are number of linguistic and non-linguistic ways to reduce the number of features in text data. Depending on the task, we can aggregate words of certain linguistic relation in concepts. For example, we replace co-hyponyms *dog* and *wolf* with their hypernym *canine*; or *vehicle* and *car* with *automobile*. Concepts, defined as units of knowledge, can be used as features to solve the ambiguity problem as well as the problem of synonymy. Indeed, each concept represents a unique meaning that can be expressed by several synonymous words. Similarly, a word with many meanings (senses) is found mapped to several concepts. Thus, a document containing the word *vehicle* may be indexed by other words such as *car* or *automobile*. The transition from a word representation to a concept representation requires the use of semantic resources external to the content of

documents such as semantic networks, thesauri, and ontologies. As a result, the performance of such a representation crucially depends on the semantic richness of the resources used in terms of the number of concepts and relationships between these concepts.

In some cases, we may need to treat morphologically variants words similarly. For example, words such as play, player, players, plays, and played will be replaced by *play*. Lemmatization and stemming are the two techniques used to find the canonical form of a word. Lemmatization uses a knowledge base containing the different inflected forms corresponding to the different possible lemmas. Thus, the inflected forms of a noun will be replaced by the singular masculine form while the infinitive form will replace the different inflected forms of a verb. Lemmatization requires the use of a dictionary of inflected forms of language as well as a grammar labeler. Stemming uses a knowledge base of syntactic and morphological rules and to transform words into their roots. One of the most well-known stemming algorithms for the English language is Porter's algorithm (Porter, 1980). Lemmatization is more complicated to implement since it depends on the grammatical labelers. In addition, it is more sensitive to misspellings than stemming.

LISTING 2.1: An example of simple Stemmer Written in Python

```
def stem(word):
    Suffix_list = ["ing", "ly", "ed", "ious", \
                  "ies", "ive", "es", "s", "ment"]
    for suffix in Suffix_list:
        if word.endswith(suffix):
            return word[: -len(suffix)]
    return word
```

Dimensionality reduction can also be done using probabilistic and linear algebraic techniques such as principal component analysis (PCA) and linear discriminant analysis (LDA) that aim to project highly dimensional data into lower dimension space. But we limit the discussion to only two methods: clustering and Latent Semantic Allocation.

Clustering (Baker and McCallum, 1998) (Sable and Church, 2001) (Slonim and Tishby, 2001) can be used to reduce dimensionality. It consists of representing the documents in a new representation space other than the original one — each dimension of the new representation space groups terms that share the same meaning. Thus, the documents will no longer be represented by terms but rather by groupings of terms representing semantic concepts. This new space of representation offers the advantage of managing the synonymy since the synonymous terms will appear in the same groupings. Likewise, the fact that a term can be included in several groupings also makes it possible to manage the polysemy. Another exciting method to reduce dimensionality, proposed by (Deerwester et al., 1990) is Latent Semantic Allocation (LSA). It uses singular value decomposition of the document  $x$  term matrix to change the representation by keeping only the  $k$  axes of strongest singular values. However, LSA is very expensive in terms of calculation time during learning as it relies on a matrix factorization method called Singular Value Decomposition (SVD) which runs computationally in cubic size. Likewise, with each new document, it is necessary to redo the whole grouping process.

### 2.3.3 Distributional Method

An alternative approach for text representation in natural language processing is the distributional similarity. The basic idea of the distribution of similarity is to represent the meaning of a word using words which co-occurs in the same context. "You shall know a word by the company it keeps" (Firth, 1957). For example, the words 'dog' and 'cat' have the same meaning in the following two sentences: I have a cat, I have a dog. While this idea might not be adequate to capture or account for the complexity or the sophistication of human language, it has succeeded tremendously in a variety of lexicalized NLP tasks. It has become the de facto representation method of text. The need for such a method does not only come from the need for more accurate numerical representation, but its density compared to other methods makes it less computationally expensive. Word2Vec (Mikolov et al., 2013) and GloVe (Pennington, Socher, and Manning, 2014) are two famous algorithms for creating word embedding.

The intuition of word2vec is to train a classifier on a binary prediction task: "Is word  $w$  likely to show up near  $X$ ?", where  $X$  is the word to which we want to find embeddings. Then, we use the weights learned by the classifier as embeddings. Skip-gram algorithm (Mikolov et al., 2013), first, initializes embeddings randomly. Then, it iteratively updates the embeddings of each word  $w$  to be equal to the words they occur within the same context. Finally, it injects  $k$  number non-neighbor words as negative examples.

## 2.4 Logistic Regression

Logistic regression is a machine learning classifier that is used in a variety of learning tasks in natural language processing and computer vision. It belongs to a group of probabilistic classifiers called discriminative models which, unlike its generative counterparts like Naive Bayes classifiers, tries to distinguish the classes rather than learn to generate them. More formally, given a document  $x$  and a class  $y$ , logistic regression computes the conditional probability of a class  $y_i$  given its input  $x_i$   $P(y|x)$ . There are two types of logistic regression classifier: binary and multinomial. The binary logistic regression is used when there are two classes of inputs to be predicted such as spam or not spam; positive or negative; malignant or benign. Multinomial or multi-class logistic regression is used when the number of predicted classes is more than two such as positive, neutral or negative.

According to (Jurafsky and Martin, 2014), classification using logistic regression, like any other supervised classification algorithm, has four main components:

- A feature extractor: to numerically represent language data (characters, words and sentences, etc) input  $X$  as feature vector  $[x_1, x_2, \dots, x_n]^T$  where  $n$  represents the number of features in the input data.
- A classification function to compute an estimation to class  $\hat{y}$  via  $p(y|x)$ . The use of classification function determines whether the classifier is binary or multinomial.
- An objective function to minimize the error predicted label  $\hat{y}$  and actual label  $y$  in training examples (cross entropy)

- An optimizer to help find the minimum of the objective function (stochastic gradient descent).

Given a dataset of  $m$  number of observations  $x$  and classes  $y$   $\mathcal{D} = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1}^m$ , the goal is to learn the set of optimal weights  $\hat{w}$  that maximizes the log probability of predicting a class  $y$  given an observation  $x$

$$\hat{w} = \underset{w}{\operatorname{argmax}} \left[ \sum_{i=1}^m \log P \left( y^{(i)} | x^{(i)} \right) \right] \quad (2.1)$$

The distribution type of  $P \left( y^{(i)} | x^{(i)} \right)$  determines whether the type of learning is binary or multiclass. The probability distribution of binary output is Bernouli, while the probability distribution of multiclass output is multinomial.

### 2.4.1 Classification Functions

What we have discussed so far gives us a method to numerically translate language from a form that is uniquely comprehensible to humans into a representation that a computer can understand, and can somehow capture the characteristics of human language. However, in order for the logistic regression to learn to classify, a classification function should be utilized. Depending on whether the classification task is binary or multinomial, there are two main functions used: Sigmoid and Softmax. Sigmoid function, used for binary classification, outputs 1 if an observation input (feature vector) is a member of a particular class e.g. (spam), and 0 otherwise. In other words, it calculates  $P(y = 1|x)$  for spam text and  $P(y = 0|x)$  for non-spam text. The mathematical form of the Sigmoid function or (logistic function, hence comes the name of the classifier) is as follow:

$$\operatorname{Sig}(x) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

where  $x$  represent the input vector,  $w$  is the learned weights, and  $b$  is the bias term or the intercept of the linear equation. One mathematical property of the Sigmoid function is that it maps any input between 0 and 1, and to make it as probability we need to make the sum of  $P(y = 1)$  and  $P(y = 0)$  equals to 1.

$$P(y = 1) = \sigma(w \cdot x + b)$$

$$P(y = 0) = 1 - \sigma(w \cdot x + b)$$

Now to make a decision, the classifier declares an observation input as *Yes* (if it is a member of class spam) if the probability is greater than a certain threshold, say 0.5, and declares an observation input as *No* otherwise.

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1|x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

In the case of multinomial classification, the logic remains the same except for the classification function where Softmax is used rather than Sigmoid. Softmax works on



normalizing the prediction of an observation by the values of all other classes in order to produce valid probability distribution.

$$p(y = i|x) = \frac{e^{w_i \cdot x + b_i}}{\sum_{j=1}^k e^{w_j \cdot x + b_j}}$$

### 2.4.2 Objective Functions

During the learning process, we need to answer the question of *how correct is the estimated output of  $\hat{y}$  of classifier from the true output  $y$ ?*. Thus our objective function is to minimize the difference between the estimated output and the true one such that the weights learned during this minimization process are, hopefully, generalizable enough to correctly label observations unseen in the training phase. So we can imagine an objective function to be a difference between the  $\hat{y}$  and  $y$  i.e.  $|\hat{y} - y|$ , but for mathematical convenience we need a non-convex objective function (or also commonly referred to as loss function) that can be easily optimized. Thus, we use maximum likelihood estimation.

Binary classification can be seen as Bernoulli distribution since the outcome is either 0 or 1. More formally, the probability of predicting class  $y$  given input  $x$  is  $p(y|x) = \hat{y}^y(1 - \hat{y})^{1-y}$ . Minimizing the probability is the same as maximizing the negative log likelihood:

$$L_{CE}(\hat{y}, y) = -\log p(y|x) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})]$$

Finally, by plugging in the value of  $\hat{y}$  we get:

$$L_{CE}(w, b) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \sigma(w \cdot x^{(i)} + b) + (1 - y^{(i)}) \log (1 - \sigma(w \cdot x^{(i)} + b))$$

, where  $m$  is the number of training examples.

On the other hand, multinomial classification uses the same process except that it uses a different classification function *softmax* and hence has different (multinomial) distribution. So the loss function for a single example  $x$  is the sum of the logs of the  $K$  output classes:

$$L_{CE}(\hat{y}, y) = - \sum_{k=1}^K 1\{y = k\} \log \frac{e^{w_k \cdot x + b_k}}{\sum_{j=1}^K e^{w_j \cdot x + b_j}}$$

### 2.4.3 Optimization

The objective functions derived in the previous section is optimized using numerical methods. Several algorithms are used for this purpose such as Stochastic Gradient Descent, that uses first derivative (gradient) information to find a minimum of a function; Newton-Raphson algorithm uses Second derivative information to find a minimum of a function. Discussing the details of how these algorithms work and the mathematical derivation of gradients is beyond the scope of this work. Interested readers are referred to (Jurafsky and Martin, 2014) for more details on the derivation of Stochastic Gradient Descent.



## 2.5 Evaluation Metrics

### 2.5.1 Confusion Matrix

A confusion matrix is a method to visualize the results of a classification algorithm. For the binary classification, the algorithm can be used to predict whether a test sample is either 0, or 1. As a way to measure how well the algorithm performs, we can count four different metrics, here 1 defined as positive and 0 defined as negative:

1. True positive (TP), the algorithm classifies 1 where the correct class is 1.
2. False positive (FP), the algorithm classifies 1 where the correct class is 0.
3. True negative (TN), the algorithm classifies 0 where the correct class is 0.
4. False negative (FN), the algorithm classifies 0 where the correct class is 1.

FIGURE 2.2: Confusion Mtrix Example

		Predicted Class	
		1	0
Actual Class	1	True Positive TP	False Positive FP
	0	False Negative FN	True Negative TN

### 2.5.2 Precision, Recall and F-Score

Precision and Recall are two popular measurement to evaluate the performance of supervised classification methods. Precision is defined as:

$$\text{Precision} = \frac{\text{TruePositive}}{\text{TruePositive} + \text{FalsePositive}}$$

, recall defined as:

$$\text{Recall} = \frac{\text{TruePositive}}{\text{TruePositive} + \text{FalseNegative}}$$

and F-score is the harmonic mean of both precision and recall:

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

A classifier with high precision means it classifies almost no inputs as positive unless they are positive. A classifier with high recall, on the other hand, would mean that it misses almost no positive values.

## 2.6 Conclusion

In this chapter, we briefly introduced some fundamental topics of Machine Learning theory required for natural language processing supervised classification tasks. We discussed two commonly used numerical methods to represent text computationally and their variants, Frequency-based methods: Bag-of word models (BoW) and Term Frequency - Inverse Document Frequency (tfidf), and Distributional similarity methods. Besides, we reviewed some techniques to reduce the dimensionality of feature space to reduce the complexity of the training and save computational resources. Finally, we covered how logistic regression as a probabilistic algorithm can be used for binary or multinomial classification along with the evaluation metrics used to gauge its performance. In the next two chapters, we show the use of multinomial logistic regression classifier on two novel tasks in natural language processing: grammatical difficulty categorization of English sentences and automatic reading comprehension. We also explore how logistic regression can be used in a semi-supervised manner to work on tasks with limited data.

## Chapter 3

# The Use of Semi-supervised Learning in Classification of English Grammatical Structures

### 3.1 Introduction

In this chapter, we examine the use of multinomial logistic regression classifier in a semi-supervised manner to automatically categorize English sentences based on difficulty to second language learners using the Common European Framework Reference. We also compare some of text representation techniques introduced in the previous chapter and determine the best one to use in this task. Finally, we evaluate the overall performance of the classifier before and the semi-supervised data augmentation.

Common European Framework Reference (CEFR) is a standard for describing language achievement on a six-point scale, from A1 for beginners, up to C2 for proficient user ability on the four primary skills reading, writing, listening, speaking, and the two secondary skills: vocabulary and grammar. The latter two skills are considered the backbone for the four primary skills in terms of difficulty. For example, a frequently uncommon lexical item like *sagacious* makes the sentence in which it appears more difficult to an English learner than a sentence with a more frequent lexical item like *wise*. Similarly, grammatical structures also play a vital role in the overall difficulty of a sentence. For instance, the word *should* conveys a different meaning in *I should study for my final exams.* than in *This is your mission should you accept it..* The latter is considered more difficult to an English learner because the use of *should* as a conditional is less frequent than its use as a modal auxiliary expressing necessity.

Wide World Web and social media are becoming indispensable resources for second language learners. However, deciding the appropriateness of materials to learners is a very labor-intensive task for English teachers. Instructors spent hours sifting through materials online to determine their difficulty and whether or not they are appropriate for their learners. To this end, we leverage the power of machine learning to build a tool that can automatically determine the difficulty of a given text according to CEFR guidelines. CEFR provides 1200 criteria ([English Profile - EGP Online](#)) for the difficulty of grammatical structures 3.1.

We treat this task as a multinomial classification problem. In our first attempt, we train a logistic regression classifier optimized by Newton-Raphson method on a dataset

TABLE 3.1: An excerpt of English sentences and their corresponding CEFR difficulty levels

Sentence	Level	
I go there every year with my friends	A1	
I think swimming is good for my body.	A2	
Tomorrow I'm expecting a delivery of our latest catalogues.	B1	
Do not hesitate to contact me should you need further information.	B2	
Living in Greece, I have had a chance to realise how much tourism can affect one's life.	C1	
There were no photographs of him in Ann's mother's albums.	C2	

of around 3000 examples provided by [englishprofile.com], described in details in section 2. We get an F1 score of 0.67. In our second attempt, we use the trained classifier from our first attempt to predict the difficulty of sentences in an unlabeled corpus to fetch more training. After that, we merge the newly fetched data with the original dataset (obtained from Englishprofile.com) and re-train another multinomial logistic regression classifier from scratch to get an F1 score of 0.79.

## 3.2 Related Works

Sentence classification is the most common task in natural language processing. A considerable number of studies conducted on lexicalized tasks like sentiment analysis, spam filtering, news categorization, etc. Rarely do we see work on classification sentences based on their syntactic structures. Thus, to the best of our knowledge, no work has ever tackled the task of classifying the grammatical difficulty of English sentences according to CEFR guidelines. Therefore, we will briefly survey techniques used to solve general sentence classification problems.

Proximity-based Algorithms such as Rocchio's algorithm (Rocchio, 1971) and K-nearest neighbor (Tam, Santoso, and Setiono, 2002) build vector for each class using a training set of a document by measuring the similarity, such as Euclidean distance or cosine similarity, between the documents. Some studies (Bang, Yang, and Yang, 2006) incorporated dictionary-based methods to construct a conceptual similarity with KNN classifier, while (Chang and Poon, 2009) combine two KNN classifier with a Naive Bayes one using TF-IDF over phrases to categorize large collections of emails. While proximity-based methods perform well on document classification, yet using them on a large training set is not feasible as computing similarities across documents is computationally and resource-wise expensive. Another disadvantage is that noise, and irrelevant data can severely degrade the performance of the classification.

Another family of classifiers that work well in text categories tasks is decision trees. Classification using this method is done through automatic creating of "if-then" rules. Their use in tasks like spam filtering is widespread even with the advent of deep neural network methods (Wu, 2009). Another common powerful classifier is Naive Bayes that based on Baye's rule. It perform surprisingly well for many real-world classification

applications under some specific conditions (“A Comparison of Event Models for Naive Bayes Text Classification”) (Rish, Hellerstein, and Thathachar, 2001). While Naive Bayes does not often outperform discriminative classifiers like Support-Vector Machine, it has the advantage of functioning well with small training data. (Kim et al., 2006) and (Isa et al., 2008) show good results by selecting Naive Bayes with SVM for text classification and clustering the documents. Using a Poisson Naive Bayes for text classification model also yield excellent results (Isa et al., 2008).

The SVM classification method and logistic regression have shown outstanding results when used in text classification tasks (Yang and Liu, 1999) (Brücher, Knolmayer, and Mittermayer, 2002) perhaps because of their capacity of handling high-dimensional sparse data well. However, they can be relatively complex in training and require high time and resource consumption.

### 3.3 Method

#### 3.3.1 Dataset: English Grammar Profile

The dataset used for this study, provided by English Grammar Profile, exhibits typical grammar profile for each level. It consists of 3615 examples and 1215 rules, divided as follows: class A has 1194 supporting examples; class B has 1775 examples, and class C has 646 supporting examples. We merged each the six levels into three supercategories in order to provide more data within each category.

- **Rule:** Can use ‘and’ to join a limited range of common adjectives.
- **Example:** The teachers are very nice and friendly.
- **Level:** A1

As we can see, the dataset provides some guidance of what characterizes each example. These rules are prepared to be used by human teachers, and they are not very computationally friendly. On the other hand, the small number of examples in the dataset is not enough to build a reliable probabilistic classifier to learn to distinguish the various properties of the classes. Therefore, in order to improve the performance, we are using a semi-supervised data augmentation technique similar to Yarowsky’s bootstrapping approach (Yarowsky, 1995).

#### 3.3.2 Multinomial Logistic Regression

Suppose that  $\mathbf{x}$  is the vectorized sentence in either of its two forms: BoW or TfIDF,  $y \in Y$  is the class label,  $\mathbf{w}_k$  and  $b_k$  are network parameters associated with class  $y$ . Then the probability of  $\mathbf{x}$  belonging to the class  $y$  can be defined by the Softmax function:

$$p(y|\mathbf{x}) = \frac{1}{z(\mathbf{x})} \exp(\mathbf{w}_y^T \mathbf{x} + b_y), \quad (3.1)$$

where  $z(\mathbf{x}) = \sum_{j=1}^K \exp(\mathbf{w}_j^T \mathbf{x} + b_j)$ .

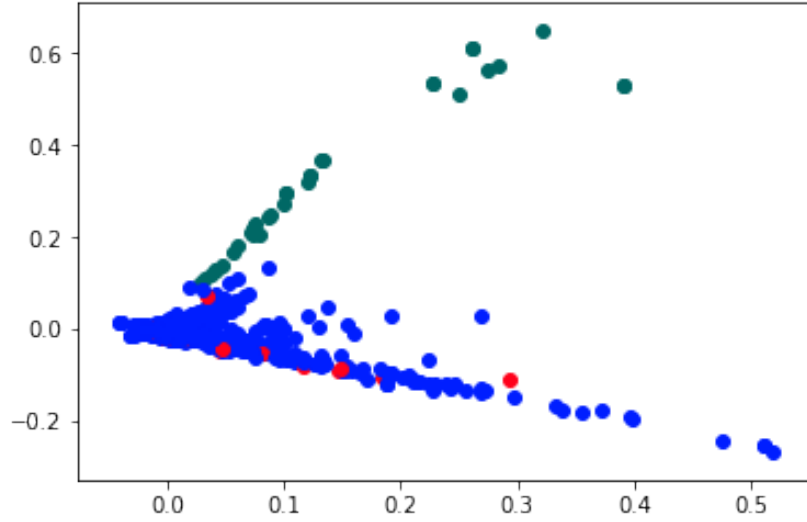


FIGURE 3.1: 2d projection of the three classes A, B and C using K-means Clustering and PCA

Let the log-likelihood function be

$$L(\beta) = \sum_{i=1}^N \log p_{g_i}(x_i; \beta)$$

$$= \sum_{i=1}^N \left[ \bar{\beta}_{g_i}^T x_i - \log \left( 1 + \sum_{l=1}^{K-1} e^{\bar{\beta}_l^T x_i} \right) \right]$$

To apply the Newton-Raphson method, we need the second derivative of the log-likelihood function

$$\frac{\partial^2 L(\beta)}{\partial \beta_{kj} \partial \beta_{mn}} = - \sum_{i=1}^N x_{ij} x_{in} p_k(x_i; \beta) [l(k = m) - p_m(x_i; \beta)]$$

The formula for updating  $\beta_{new}$  for multiclass is:

$$\beta^{new} = \beta^{old} + \left( \tilde{\mathbf{X}}^T \mathbf{W} \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T (\mathbf{y} - \mathbf{p})$$

where  $y$  is the concatenated indicator vector of dimension  $N \times (K - 1)$ ,  $p$  is the concatenated vector of fitted probabilities of dimension  $N \times (K - 1)$ ,  $\tilde{\mathbf{X}}$  is an  $N(K - 1) \times (p + 1)(K - 1)$  matrix; and Matrix  $\mathbf{W}$  is an  $N(K - 1) \times N(K - 1)$  square matrix.

### 3.3.3 Feature Design

For this task, we are comparing two standard methods of feature generation in natural language processing and information retrieval: bag-of-word model and term-frequency inverse-document frequency (TFIDF)

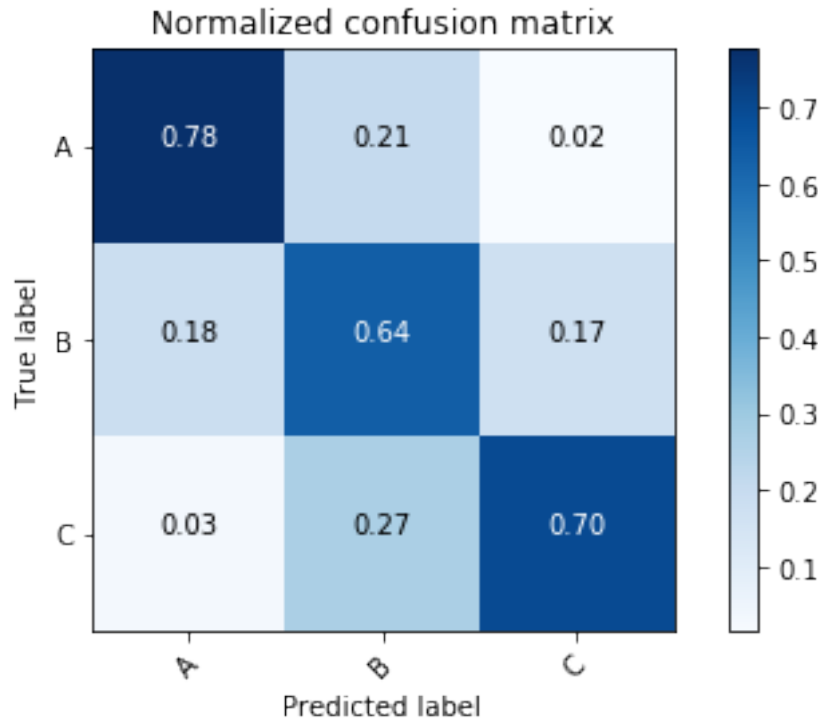


FIGURE 3.2: Confusion matrix illustrating the performance of LR-BoW-Word model with class weight to deal with data imbalance

## 3.4 Experimental Results

In this section, we evaluate several variations of our method against a random model as our baseline model. In the following experiments, all methods are done using Python 3.5 and Scikit Learn Library tested on MacBook Pro laptop with Core i5 processor and 8 GB of RAM.

### 3.4.1 First Phase

In the first phase of our experiment, we train a logistic regression classifier, optimized by the Newton-Raphson method, on English Grammar Profile dataset. We also introduce a feature design method to mask word categories with their part-of-speech tags in order to preserve grammatical information and avoid the lexical influence of the word. We use a combination of feature design methods such as:

- **BoW**: Apply unigram, bigram and trigram bag-of-words model with both word tokens, and masked tokens.
- **Tf-idf**: Apply unigram, bigram and trigram tf-idf model with both word tokens, and masked tokens.

From table 2, we compare the performance of 4 variations of our method against the baseline, which is a random model. Surprisingly, we notice that our bag-of-words

TABLE 3.2: Comparison of Precision-Recall - Phase One

Model	Precision (%)	Recall (%)	F1 (%)
Random-Baseline	39 (%)	39 (%)	39 (%)
<b>BoW-LR-Words</b>	70 (%)	<b>69 (%)</b>	<b>69 (%)</b>
BoW-LR-Masked	66 (%)	62 (%)	63 (%)
Tfidf-LR-Words	<b>75 (%)</b>	66 (%)	60 (%)
Tfidf-LR-Masked	66 (%)	64 (%)	61 (%)

model with word tokens (BoW-LR-Words) outperform the one with the masked sequence in terms of precision and recall respectively. Besides, BoW-LR-Words also outperforms (tfidf-LR-words) model in terms of recall only, while the latter has higher recall. From the confusion matrix (fig. 2), we see that (BoW-LR-Words) model predicts most of the testing data correctly even though the number of training data is not equal among the classes, especially for class C. Therefore, we needed to assign specific weights to each class during the training in order to avoid this imbalance in the data.

### 3.4.2 Second Phase

The results shown in Table 2 does indicate that the linear model has learned the associations between sentences and their corresponding grammatical class reasonably well. However, both of our feature design techniques, namely bag-of-words and tfidf have their disadvantages. The first assigns similarity based on occurrence, size and mutual words. For example, the BoW model learns that longer sentences tend to be the most difficult, while shorter ones are less difficult. Similarly, while tfidf treats the drawbacks of BoW model, it still ignores the usefulness of what-so-called *Stop Words* due to their commonality. In other words, it is very hard to trust these results with this amount of training data. Therefore, we propose a non-synthetic data augmentation method to provide more training examples, and thereby improve the overall performance of the model.

Using a text-only version of the brown corpus (**citeulike: 13797746**), we use our BoW-LR-Words model to predict the grammatical difficulty labels of its sentences. Then, we collect the sentences predicted with high certainty (above 0.9 probability) to be fed to the original dataset. It is crucial to set a higher probability threshold as the new data examples will serve as seed points for our prediction in the future, and we do not want the classifier to *drift* from the original trajectory. This technique is reminiscent of Yarowsky's Bootstrapping (Yarowsky, 1995)(semi-supervised) algorithm, but with one difference is that unlike in Yarowsky's method, we apply this step only once.

Out of 38400 sentences in the unlabeled corpus, only 1428 sentences have been detected with certainty above 90%. We also removed any sentence that is longer than 30 words in order to lessen the undesirable effect of BoW technique. We apply our best model from phase one to the augmented dataset of 5043 examples under the same training conditions to get a precision of **0.80**, recall of **0.79**, and F1 score of **0.79**.



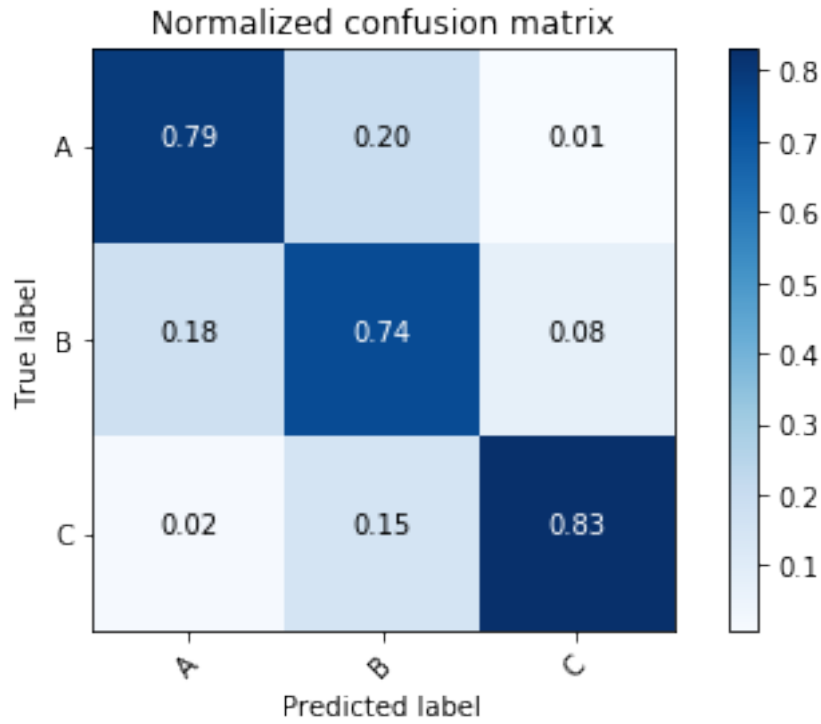


FIGURE 3.3: Confusion matrix illustrating the performance of LR-BoW-Word model after augmentation

### 3.5 Conclusion

In this chapter, we have presented a classification method based on simple multinomial logistic regression and bag-of-words model augmented with semi-supervised (bootstrapping) method to classify English sentences based on their level difficulty, A=Beginner, B=intermediate, C=Advanced according to CEFR. We have also compared several common feature design techniques to encode the sentences along with a proposed masking BoW method that replaces lexical items with their grammatical category. Finally, the model achieves an overall F1 score of 0.69, and 0.79 after augmenting it with example sentences from an unlabeled corpus.



## Chapter 4

# Two-stage Classification Method for Automatic Reading Comprehension

### 4.1 Introduction

Machine Comprehension or Automatic Reading Comprehension is the task of automatically answering comprehension questions based on information derived from short texts. While this task is relatively easy for humans to do, it is hard for the machine because it requires not only linguistic information but also cognitive, world-knowledge skills as well. Machine Comprehension has been of great interest to the NLP community because of its practical applications in education, especially in the field of automated assessment and intelligent tutoring systems. Working on similar tasks will bring us a step closer to understand how human can perform such cognitive tasks so efficiently while the machine cannot.

In this chapter, we will be using the second release of Stanford Question Answering Dataset (SQuAD) (Rajpurkar, Jia, and Liang, 2018). SQuAD 2.0 is a reading comprehension dataset, consisting of 150K questions posed by crowd workers on a set of Wikipedia articles, where the answer to every question is a segment of text, or span, from the corresponding reading passage, or the question might be unanswerable. The difference between SQuAD 2.0 and its previous release is that it has 50K unanswerable questions in addition to paragraph-supported questions. To do well on SQuAD2.0, a system must not only answer questions when possible, but also abstain from answering when no answer is supported by the paragraph.

The strategy will be used to tackle this task is mainly driven by pure linguistic intuition. Since SQUAD 2.0 task is a question answering task, the answer span lies within one of the given sentences in the paragraph, and then an answer is one of the constituents of the best candidate sentence. Thus, our goal is to, first, find the sentence containing the right answer, and we do by training a logistic regression classifier with some linguistic features. Then, we train another classifier to find the constituent best represent the right answer span within the sentence predicted in the first stage. In order to account for the unanswerability of some of the questions, we add a null-answer as a dedicated class in the first classifier along with the potential sentences. This way, if an answer is not found, the inference process stops without proceeding to the next stage, which saves time and computation.

## 4.2 Related Work

There has been a large number of studies that tackle traditional Machine Comprehension tasks, where answers are extractive and explicitly stated in the passage. The difficulty of SQUAD 2.0 task, however, lies in abstaining from answers when no answer span exists in the passage. Thus, we will focus on models that answer not only the questions but also the one that predicts the probability of questions being answered.

(Seo et al., 2016) introduced bi-directional attention flow (BiDAF) that uses a recurrent neural network to encode contextual information in both question and passage along with an attention mechanism to align parts of the question to the sentence containing the answer and vice versa. The model offers context representation at multilevel of granularity: character-level, word-level, and contextual embedding. What sets this work from others is that it does not represent the context paragraph into a fixed-length vector. Instead, it dynamically computes vectors at each time step, combines it with the one from the previous layer and allow *flow* through to the next layers. The model outputs confidence scores of start and end index of all potential answers. One problem with this model is that it is not designed to handle unanswerable questions. (Levy et al., 2017) extends the work by assigning a probability to null-answers to account for the questions whose answers do not exist in the corresponding paragraph. It achieves 59.2% EM score and 62.1% F1 score.

(Hu et al., 2018) propose a read-then-verify system that can abstain from answering when a question has no answer given the passage. They introduce two auxiliary losses to help the neural reader network focus on answer extraction and no-answer detection respectively, and then utilize an answer verifier to validate the legitimacy of the predicted answer. One essential contribution is answer-verification. The model incorporates a multi-layer transformer decoder to recognize the textual entailment that supports the answer in and passage. This model achieves an ExactMatch EM score of 71.6% and 74.23% F1 on SQuAD 2.0.

(Wang, Yan, and Wu, 2018) proposes a new hierarchical attention network that mimics the human process of answering reading comprehension tests. It gradually focuses attention on the part of the passage containing the answer to the question. The model comprises of three layers. a) **encoder layer**: builds representations to both the question and the passage using a concatenation of word-embedding representation (Pennington, Socher, and Manning, 2014) and a pre-trained neural language model (Peters et al., 2018) b) **Attention layer**: its function is to capture the relationship between the question and the passage at multiple levels using self-attention mechanisms. c) **Matching layer**: given refined representation for both question and passage, a bi-linear matching layer detects the best answer span for the question. This method achieves state-of-the-art results as of September 2018. Their single model achieves 79.2% EM and 86.6% F1 score, while their ensemble model achieves 82.4% EM and 88.6% F1 Score.

Current models that have shown significant improvement on Machine Comprehension tasks and other similar ones owe their success to a new neural architecture called transformer networks (Vaswani et al., 2017). It has become the *de facto* in recent sequential learning tasks, eschewing recurrence. This architecture creates global dependencies

between input and output using only attention mechanisms. An even more successful model is Bidirectional Encoder Representations from Transformers BERT (Devlin et al., 2018). It is a task-independent pretrained language model that uses a deep transformer network to create rich and context-specific vector representation. Using a method called Masked Language Model, the goal is to randomly mask tokens from the input and predict the vocabulary id of the masked word based only on its context. BERT has shown significant improvement on several tasks, one of which is SQUAD.

While these models achieve astonishing results, they have architecture design and implementation, and very resource-intensive. This complexity results due to lack of linguistic assumptions that can arguably constrain and direct the problem.

### 4.3 Baseline Implementation

The new version of SQuAD 2.0 task adds a new constraint to competing question-answering models. In addition to identifying the extractive answer spans, a question-answering model should abstain from answering if the passage does not contain an answer. To this end, we implement a simple multinomial logistic regression classifier to address this task. At this level, the classification task is to predict the sentence, in the paragraph, containing the right answer, or declaring that the question is unanswerable. Next, we apply a constituency parser over the sentence predicted from the first stage to get its constituents among which lies the correct answer span (see Figure 1).

#### 4.3.1 Feature Design

In order for the classifier to find the sentence containing the answer, it must determine the sentence that is most similar to the question, and by similar we mean that a good candidate sentence 1) shares more words with the question. 2) has high cosine similarity with the question. 3) shares a syntactic similarity with the question. Thus, three main features have been selected to this end:

- **Cosine Similarity:** for every sentence in the paragraph as well as the question a word vector representation is created via InferSent (Conneau et al, 2017), which is a pre-trained sentence embeddings method that provides semantic representations for English sentences. InferSent is an encoder based on a bi-directional LSTM architecture with max pooling, trained on the Stanford Natural Language Inference (SNLI) dataset. Cosine distance score is calculated for each sentence-question pair.
- **Word Overlap:** this calculates the Jaccard score between each sentence-question pair. Jaccard index is a method of computing the explicit similarity between two sets as follows:

$$J(Q, S) = \frac{|Q \cap S|}{|Q \cup S|}$$

where Q and S are sets of words in question and sentence respectively.

- **POS Overlap:** This feature computes the Jaccard score over the part-of-speech-tag representation of sentences. In other words, instead of the word tokens, it checks similarity over POS tokens. We use the default POS-tagger in the Spacy library of

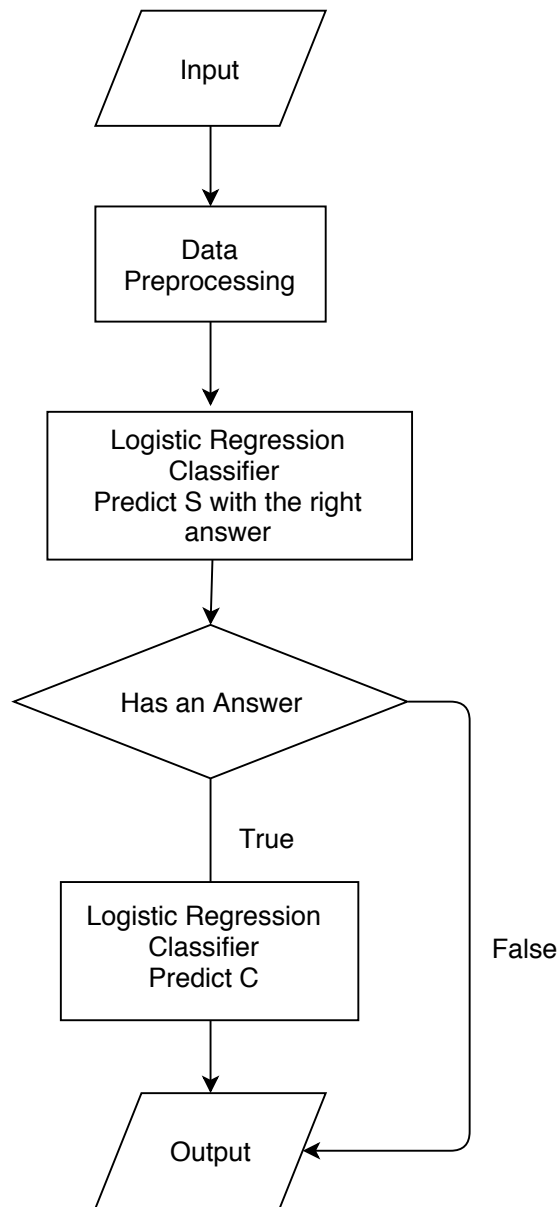


FIGURE 4.1: Flowchart illustrating the two-stage classification approach

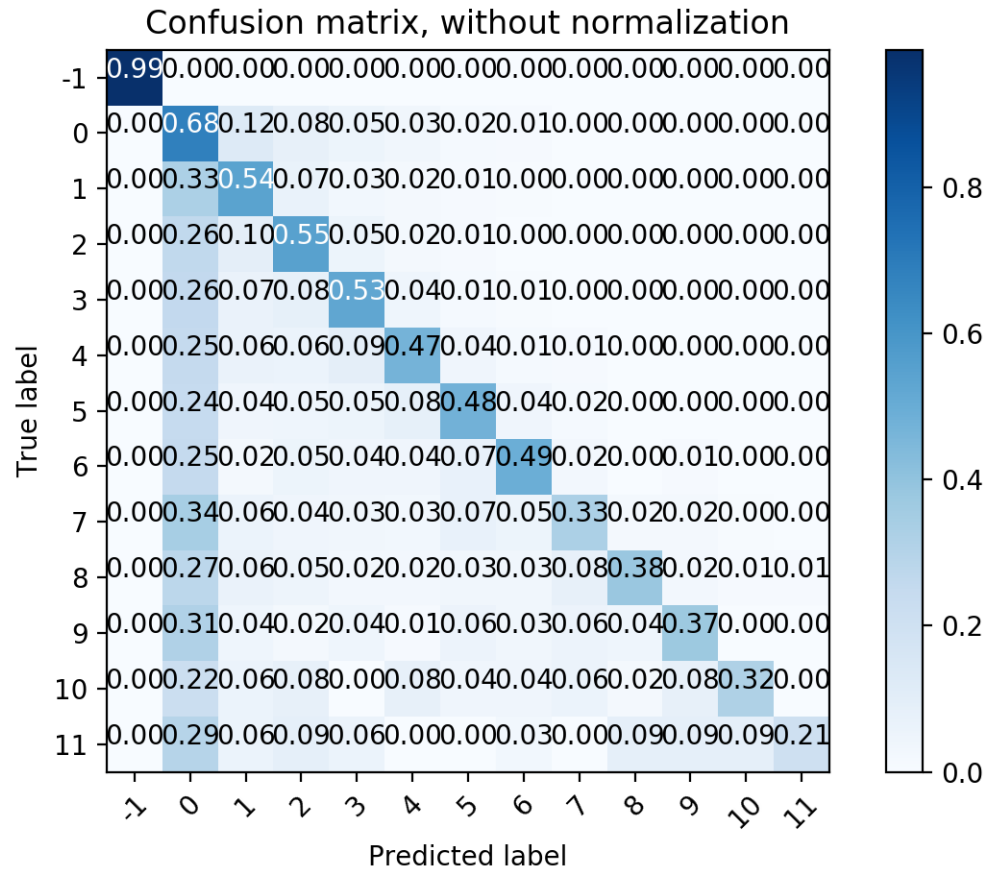


FIGURE 4.2: Confusion matrix shows which classes were predicted correctly

Python programming language to obtain the POS representation for the sentences and questions alike.

Using the three features above every question-sentence pair will have three scores and an additional binary feature indicating whether or not the question is answerable.

### 4.3.2 Training and Result

We train a logistic regression classifier with L2 regularization ('newton-cg' solver) using scikit-learn library of Python programming language, and we get the results shown in table 1. Numbers in class column represents the index of the sentence in the paragraph containing the answer, and -1 indicates that the question has no answer in the paragraph. We also limit the number of sentences to 10. The results show that with simple features we get an F1 score of 0.71.

heightclass	precision	recall	F1-score
-1	1	0.99	0.99
0	0.47	0.68	0.56
1	0.63	0.54	0.58
2	0.62	0.55	0.58
3	0.62	0.53	0.57
4	0.58	0.47	0.52
5	0.57	0.48	0.52
6	0.57	0.49	0.53
7	0.46	0.33	0.38
8	0.56	0.38	0.45
9	0.46	0.37	0.41
10	0.48	0.32	0.39
11	0.35	0.21	0.26
avg / total	0.72	0.71	0.71

TABLE 4.1: This table shows the results of running a multinomial regularized logistic regression. The class column represents the index of sentences within the paragraph, and -1 represents the unanswerable question case. Unpredicted classes are removed.

## 4.4 Stage Two: Predicting the Answer Span

### 4.4.1 Iteration 1

To select the most plausible answer span from the candidate sentence, we design a number of features:

- **Constituents:** Using a constituency parser (Kitaev and Klein, 2018), we obtain all constituents in a candidate sentence. These constituents will be the classes from which we pick the right span.
- **Contextual Overlap:** Constituents sharing context with the original question are potential candidates to be the correct answers. So we measure the cosine similarity between each constituent and the question:

$$similarity = \frac{\sum_{i=1}^n C_{|w|} Q_i}{\sqrt{\sum_{i=1}^n C_{|w|}^2} \sqrt{\sum_{i=1}^n Q_i^2}}$$

where  $w$  is the number of slide window around the candidate constituent. For our purposes features of size 2 and 3 are used.

- **Constituent Label:** Constituency parse tree label of the span combined with wh-word.

Out of 85K training example, the answer spans of nearly half of them are not within the constituents. However, answers can be part of the constituents. For example, an answer span might be *L'official* and the nearest constituent is *L'official Magazine*. So for our first attempt, we remove all data points whose answers are not explicitly found within



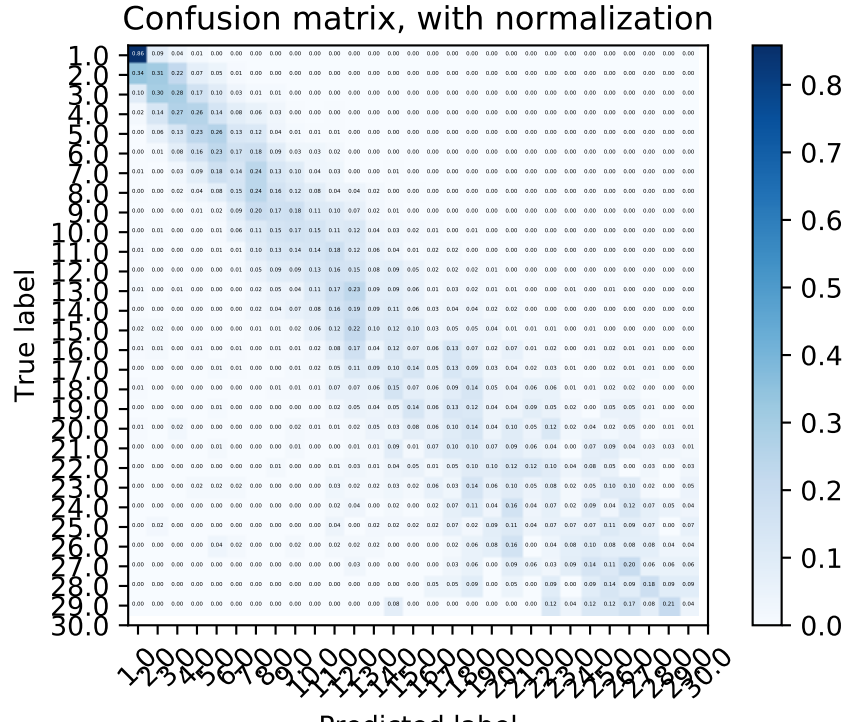


FIGURE 4.3: Confusion Matrix illustrating the first iteration of Stage2 LR.

the constituents. This results in around 45K data point. Next, we train a logistic regression classifier with L2 regularization and 100 epochs, optimized by newton method 'newton-cg' on a MacBook Pro laptop with core i5 processor and 8 GB of RAM. Python modules used are Scikit-learn and Spacy. The latter is used to drive the attentive-neural constituency parser (Kitaev, 2018). The result is 0.25 F1 score.

#### 4.4.2 Error Analysis: Iteration 1

By looking at the confusion matrix, the first thing we notice is data imbalance. Classes at the beginning have more supporting points, and therefore has less error than the others. However, class 2 has been identified as class 1 34%. For example, the answer of *At what age did Beyonce meet LaTavia Robertson?* is predicted as *At age eight* while the correct answer is *age eight*. Another example, the answer for *How much was public expenditure on the island in 2001-2002?* is predicted to be *Public expenditure* while the true answer is *£10 million*. In this case, the phrase public expenditure appears in the question, which is a stronger candidate given the features designed. Similarly, class 12 is predicted as class eleven 16%. For example, the answer to the question *Who supervised the design and implementation of the iPod user interface?* is *Steve Jobs*, but it is predicted as *Of Steve Jobs*. Another example, answer to *What roles were women recruited for in the 1950s?* is *in medicine, communication*, but it is predicted as *logistics, and administration*. Given the features, we have designed it is very difficult for the classifier to recognize this settle difference between the classes.

### 4.4.3 Iteration 2

In the first attempt, we have sacrificed valuable data because the answer span was not within the constituents of the candidate sentence. This time, we use all 85K, but we will face the same problem introduced in the previous section. The answer span can be a member of more than one constituent in the same sentence, and this will impose a very hard constraint on the performance because the inference can be close but not exact. Also because this problem requires more feature engineering efforts, we decided to leave it for future studies. Instead, we will set the target to the first constituent of which the span is a member. However, we will add three more features:

- **Distributional Distance:** We measure the distributional cosine similarity between the sum of all words in the contextual window and the question using Glove (Pennington, 2012).
- **Matching Word Frequencies:** Sum of the TF-IDF of the words that occur in both the question and the sentence containing the candidate answer.
- **Lengths:** Number of words to the left and the right of the span.

For classification, we compare the performance of a logistic regression classifier with a similar configuration as the first stage to 3-layer feed-forward neural network of 128, 64 and 32 nodes respectively. The logistic regression achieves 0.35 F1 while the neural network does 0.42 F1 using Adam optimizer and 100 epochs. Looking at the confusion matrix, it is apparent that the new features (tf-idf and distributional cosine similarity) improve the performance. However, they could not recognize the subtle difference among the phrase.

## 4.5 Conclusion

We introduced a two-step classification method for automatic reading comprehension via SQUAD 2.0 dataset. Our stage1 classifier managed to find whether or not a question is answerable within a given passage and find the sentence containing the right answer with F1 score of 0.71. Our stage2 classifier manages to detect the exact span with F1 score of 0.35 even though the predicted answer is not distant from the exact answer. In order to improve the performance of our approach, future studies should investigate the usefulness of features generated from Named Entity Recognition, Semantic Role Labeling and Dependency Parsing processes, which are expected to be potential solutions to the problems we faced in this work.

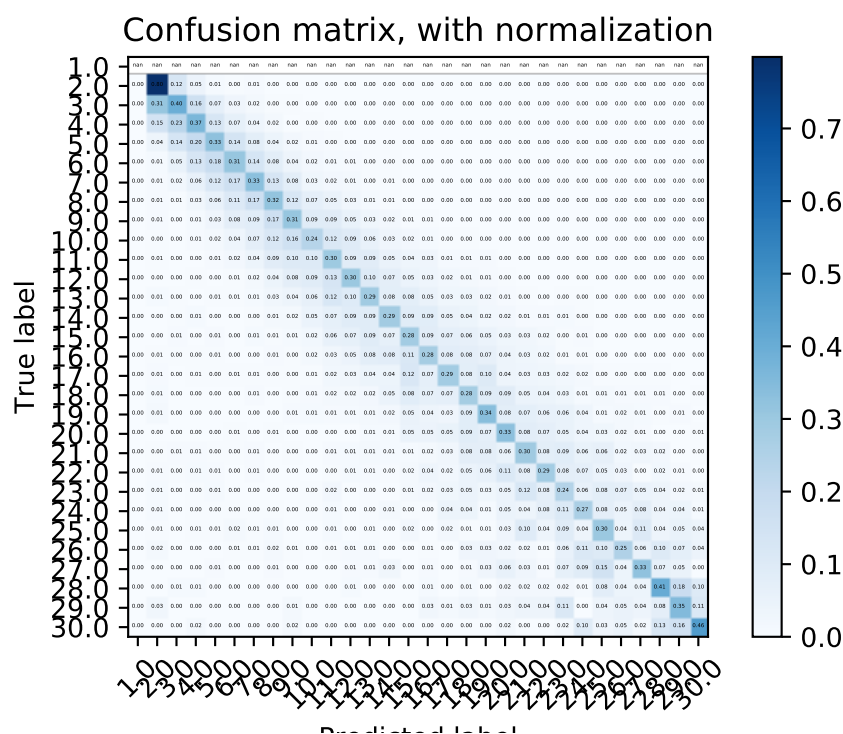


FIGURE 4.4: Confusion Matrix illustrating the Second iteration of Stage2 LR.



# Bibliography

- Baker, L Douglas and Andrew Kachites McCallum (1998). "Distributional clustering of words for text classification". In: *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, pp. 96–103.
- Bang, Sun Lee, Jae Dong Yang, and Hyung Jeong Yang (Mar. 2006). "Hierarchical document categorization with k-NN and concept-based thesauri". In: *Information Processing & Management* 42.2, pp. 387–406. ISSN: 0306-4573. DOI: 10.1016/j.ipm.2005.04.003. URL: <http://www.sciencedirect.com/science/article/pii/S0306457305000579> (visited on 12/09/2018).
- Bellman, Richard E (2015). *Adaptive control processes: a guided tour*. Vol. 2045. Princeton university press.
- Brücher, Heide, Gerhard Knolmayer, and Marc-André Mittermayer (2002). "Document classification methods for organizing explicit knowledge". In:
- Chang, Matthew and Chung Keung Poon (June 2009). "Using phrases as features in email classification". In: *Journal of Systems and Software* 82.6, pp. 1036–1045. ISSN: 0164-1212. DOI: 10.1016/j.jss.2009.01.013. URL: <http://www.sciencedirect.com/science/article/pii/S0164121209000089> (visited on 12/09/2018).
- Deerwester, Scott et al. (1990). "Indexing by latent semantic analysis". In: *Journal of the American society for information science* 41.6, pp. 391–407.
- Devlin, Jacob et al. (2018). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding". In: CoRR abs/1810.04805. arXiv: 1810.04805. URL: <http://arxiv.org/abs/1810.04805>.
- English Profile - EGP Online. URL: <http://www.englishprofile.org/english-grammar-profile/egp-online> (visited on 01/29/2019).
- Firth, John R (1957). "A synopsis of linguistic theory, 1930-1955". In: *Studies in linguistic analysis*.
- Fuhr, Norbert and Chris Buckley (1991). "A probabilistic learning approach for document indexing". In: *ACM Transactions on Information Systems (TOIS)* 9.3, pp. 223–248.
- Hu, Minghao et al. (2018). "Read + Verify: Machine Reading Comprehension with Unanswerable Questions". In: CoRR abs/1808.05759. arXiv: 1808.05759. URL: <http://arxiv.org/abs/1808.05759>.
- Isa, D. et al. (Sept. 2008). "Text Document Preprocessing with the Bayes Formula for Classification Using the Support Vector Machine". In: *IEEE Transactions on Knowledge and Data Engineering* 20.9, pp. 1264–1272. ISSN: 1041-4347. DOI: 10.1109/TKDE.2008.76.
- Jurafsky, Dan and James H Martin (2014). *Speech and language processing*. Vol. 3. Pearson London.
- Kim, Sang-Bum et al. (Nov. 2006). "Some Effective Techniques for Naive Bayes Text Classification". In: *IEEE Transactions on Knowledge and Data Engineering* 18.11, pp. 1457–1466. ISSN: 1041-4347. DOI: 10.1109/TKDE.2006.180.

- Kitaev, Nikita and Dan Klein (2018). "Constituency Parsing with a Self-Attentive Encoder". In: *arXiv preprint arXiv:1805.01052*.
- Levy, Omer et al. (2017). "Zero-Shot Relation Extraction via Reading Comprehension". In: *CoRR abs/1706.04115*. arXiv: 1706.04115. URL: <http://arxiv.org/abs/1706.04115>.
- Lewis, David Dolan (1992). "Representation and learning in information retrieval". PhD thesis. University of Massachusetts at Amherst.
- Manning, Christopher D, Christopher D Manning, and Hinrich Schütze (1999). *Foundations of statistical natural language processing*. MIT press.
- McCallum, Andrew and Kamal Nigam. "A Comparison of Event Models for Naive Bayes Text Classification". en. In: (), p. 8.
- Mikolov, Tomas et al. (2013). "Distributed representations of words and phrases and their compositionality". In: *Advances in neural information processing systems*, pp. 3111–3119.
- Pennington, Jeffrey, Richard Socher, and Christopher Manning (2014). "Glove: Global vectors for word representation". In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543.
- Peters, Matthew E. et al. (2018). "Deep contextualized word representations". In: *Proc. of NAACL*.
- Porter, Martin F (1980). "An algorithm for suffix stripping". In: *Program* 14.3, pp. 130–137.
- Rajpurkar, Pranav, Robin Jia, and Percy Liang (2018). "Know What You Don't Know: Unanswerable Questions for SQuAD". In: *arXiv preprint arXiv:1806.03822*.
- Rish, Irina, Joseph Hellerstein, and Jayram Thathachar (2001). *An analysis of data characteristics that affect naive Bayes performance*. Tech. rep.
- Rocchio, Joseph John (1971). "Relevance feedback in information retrieval". In: *The SMART retrieval system: experiments in automatic document processing*, pp. 313–323.
- Sable, Carl and Kenneth W Church (2001). "Using bins to empirically estimate term weights for text categorization". In: *Proceedings of the 2001 Conference on Empirical Methods in Natural Language Processing*.
- Seo, Min Joon et al. (2016). "Bidirectional Attention Flow for Machine Comprehension". In: *CoRR abs/1611.01603*. arXiv: 1611.01603. URL: <http://arxiv.org/abs/1611.01603>.
- Shannon, Claude Elwood (1948). "A mathematical theory of communication". In: *Bell system technical journal* 27.3, pp. 379–423.
- Singhal, Amit, Mandar Mitra, and Chris Buckley (1997). "Learning routing queries in a query zone". In: *ACM SIGIR Forum*. Vol. 31. SI. ACM, pp. 25–32.
- Slonim, Noam and Naftali Tishby (2001). "The power of word clusters for text classification". In: *23rd European Colloquium on Information Retrieval Research*. Vol. 1, p. 200.
- Tam, V., A. Santoso, and R. Setiono (Aug. 2002). "A comparative study of centroid-based, neighborhood-based and statistical approaches for effective document categorization". In: *Object recognition supported by user interaction for service robots*. Vol. 4, 235–238 vol.4. DOI: 10.1109/ICPR.2002.1047440.
- Tzeras, Kostas and Stephan Hartmann (1993). "Automatic indexing based on Bayesian inference networks". In: *Proceedings of the 16th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, pp. 22–35.

- Vaswani, Ashish et al. (2017). "Attention Is All You Need". In: *CoRR* abs/1706.03762. arXiv: 1706.03762. URL: <http://arxiv.org/abs/1706.03762>.
- Wang, Wei, Ming Yan, and Chen Wu (2018). "Multi-granularity hierarchical attention fusion networks for reading comprehension and question answering". In: *CoRR* abs/1811.11934. arXiv: 1811.11934. URL: <http://arxiv.org/abs/1811.11934>.
- Wu, Chih-Hung (Apr. 2009). "Behavior-based spam detection using a hybrid method of rule-based techniques and neural networks". In: *Expert Systems with Applications* 36.3, Part 1, pp. 4321–4330. ISSN: 0957-4174. DOI: 10.1016/j.eswa.2008.03.002. URL: <http://www.sciencedirect.com/science/article/pii/S0957417408001772> (visited on 12/09/2018).
- Yang, Yiming and Xin Liu (1999). "A Re-Examination of Text Categorization Methods". In: *SIGIR*.
- Yarowsky, David (June 1995). "Unsupervised word sense disambiguation rivaling supervised methods". In: *Proceedings of the 33rd Annual Meeting of the Association for Computational Linguistics*. Cambridge, Massachusetts, USA: Association for Computational Linguistics, pp. 189–196. DOI: 10.3115/981658.981684. URL: <http://www.aclweb.org/anthology/P95-1026> (visited on 12/09/2018).