

Louis Max Cley Slater

Sentiment Analysis of Texts on the Iraq War

Computer Science Tripos – Part II

Pembroke College

March 26, 2018

Proforma

Should be 1 page

Name:

College: **Pembroke College**

Project Title: **Sentiment Analysis of Texts on the Iraq War**

Examination: **Computer Science Tripos – Part II, July 2018**

Word Count: **Max 12,000**

Project Originator: Louis Max Cley Slater

Supervisor: Dr Tamara Polajnar

Original Aims of the Project

Due to my interest in both natural language processing and politics and a lack of previous work done in the area, I decided that I wanted to perform sentiment analysis on British political texts. After extensive research, I found a study [18] that manually assessed the biases of British newspaper articles on the Iraq war, so decided that I would use the dataset produced by the study. I aimed to develop a program to retrieve the texts of the newspaper articles specified in the study [18] and implement a classifier using this data and a bag of words model.

Work Completed

- Problem with University's licence for DowJones (and others?). Didn't cover Scraping data/API use? Be specific and add Licence agreement(s) to bibliography - Switched to Hansard and voting datasets. Made the data retrieval stage lengthier. Cite datasets. - Numbers about success of classifier - Add anything else completed?

Special Difficulties

- Mention the original dataset licence issues and change?

Declaration

I, Louis Max Cley Slater of Pembroke College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose.

Signed

Date

Contents

1	Introduction	9
1.1	Motivation	9
1.2	Corpora	10
1.2.1	Newspaper Dataset	10
1.2.2	Hansard	10
1.3	Challenges	11
1.4	Previous Work	12
2	Preparation	15
2.1	Introduction to Supervised Learning	15
2.2	Introduction to the Naïve Bayes Classifier	16
2.3	Introduction to Support Vector Machines	17
2.4	Introduction to Cross-Validation	17
2.5	Introduction to the Bag of Words Model	17
2.6	Introduction to Entailment	17
2.7	Introduction to the House of Commons	17
2.8	Spam Email Dataset	17
2.9	Requirements Analysis	17
2.10	Implementation Approach	17
2.11	Software Engineering Techniques	17
2.12	Choice of Tools	17
2.13	Starting Point	17
2.14	Summary	17
3	Implementation	19
3.1	Data Scraper	19
3.2	Database	19
3.3	Classifier	19
3.4	Entailment System	19
3.5	Summary	19
4	Evaluation	21
4.1	Unit Testing	21
4.2	Internal Evaluation	21
4.3	External Evaluation	21

4.4	Evaluation of Project Goals	21
4.5	Summary	21
5	Conclusion	23
5.1	Achievements	23
5.2	Lessons Learned	23
5.3	Future Work	23
	Bibliography	23
A	Project Proposal	27

List of Figures

Acknowledgements

Chapter 1

Introduction

This dissertation describes the development of a system that uses various natural language processing techniques to analyse texts on the Iraq war. The texts are transcripts of relevant debates in the House of Commons between 11th September 2001 and 18th March 2003. Despite underlying difficulties of the task, the system produced gives consistently good results when evaluated using a variety of different techniques. The details of the implementation and evaluation are expanded upon throughout this document.

1.1 Motivation

The initial idea for this project was due to my interest in both natural language processing and politics. While computational approaches are often applied to political texts, very few studies have ever specifically concerned a British corpus (a collection of written or spoken material stored on a computer and used to find out how language is used [4]). After reading various papers that used natural language processing techniques on political corpora, I noted the most common studies concerned sentiment analysis of short texts, such as newspaper headlines or Tweets. This motivated me to investigate the task of performing sentiment analysis on longer pieces of text, to make the project more unique still.

On a personal note, the invasion of Iraq was the first political issue I engaged in; I attended anti-war protests with my mum when I was very young and have continued to closely follow the developments since then. It was an issue for which people had very strong views that were not determined by their political leanings [15], therefore making it a particularly interesting topic for sentiment analysis. Furthermore, after extensive research I couldn't find any research that carried out computational sentiment analysis which focused on war. The project I carried out was more relevant still due to the recent publication of the Chilcot Enquiry [3] and the ongoing situation in Iraq and Syria [19].

1.2 Corpora

The most essential component of any sentiment analysis system is its corpus. For this reason, the first task I carried out was to obtain the use of a relevant corpus that I was licensed to use.

1.2.1 Newspaper Dataset

The first potential dataset I found was that produced by Robinson, P. and Goddard, P. and Brown, R. and Taylor, P.M. in which they “evaluated media performance during the 2003 Iraq War” [18]. As part of their evaluation, they manually annotated the stance of 4,893 British newspaper articles on the Iraq war. They published the resulting dataset, but it didn’t contain the body of the articles - only its headline, author, newspaper and publication date. I consequently investigated resources containing the text of the relevant articles and tried to cross-reference the data from these sources with the manually annotated stance. At the time, many newspapers published different stories online and in print, meaning that I could not rely on these. A few newspapers maintain electronic archives of their printed editions on the internet, however not enough newspapers had such archives. The final resource I looked into was Dow Jones Factiva, a “global news database” [12]. Upon inspection, this database contained the vast majority of the articles I needed and it was possible for me to cross-reference the articles in it with the labels annotated by Robinson, P. and Goddard, P. and Brown, R. and Taylor, P.M.. I initially accessed the dataset through the University of Cambridge’s subscription. I therefore (falsely) assumed that this subscription would be sufficient for use in my project, however I later discovered that an academic licence did not permit me to use the API or to carry out text-mining. I consequently contacted Dow Jones and was told that the licence I required would cost in excess of \$20,000.

1.2.2 Hansard

After exhausting all other options, I turned my attention to the House of Commons Hansard archives, which contains transcripts of debates between members of Parliament in the Commons Chamber [16]. Due to the licensing problems I encountered with Dow Jones Factiva [12], I immediately looked into the licence required to scrape data from the Hansard and found that it is covered by the Open Parliament Licence [8]. Since the Hansard archives are available under this licence, I was permitted to:

- “copy, publish, distribute and transmit the information”
- “adapt the information”
- “exploit the information commercially and non-commercially, for example, by combining it with other information, or by including it in your own product or application”

A further benefit of using the Hansard is the fact that it can be labelled using MPs’ voting records - for example, if an MP voted in favour of the invasion of Iraq, all of their

speeches on Iraq can be labelled as pro-war. This allowed me to develop the sentiment analysis system using a supervised learning model.

1.3 Challenges

There are many flaws of our ‘democracy’ in the United Kingdom. In my opinion, one major flaw is our voting system; we have a first-past-the-post system in which each member of the electorate only gets one vote. This vote goes to a candidate who is (usually [7]) a member of a political party. This party will often force an MP to vote in accordance with the party line, regardless of the MP’s own views. This subverts democracy when MPs vote with their party line, despite making contradictory election promises [5]. This issue was particularly prevalent leading up to the Iraq war, where despite protests [6] showing the public’s overwhelming opposition to the invasion of Iraq, MPs voted in favour of the invasion. Within a month of being elected as Prime Minister in 1997, Tony Blair said “Mine is the first generation able to contemplate the possibility that we may live our entire lives without going to war or sending our children to war. That is a prize beyond value.” [9], just six years before encouraging his Labour Party to vote invade Iraq. As a result, many Labour MPs who were previously opposed to the invasion voted in its favour. This hypocrisy of Members of Parliament is commonplace in British politics, meaning that as the electorate, we are frequently misled by the politicians representing us and the difficulties of party politics add to this.

At first glance, the Hansard appears to be a great resource for natural language processing, however its data is poorly presented and to use it, it’s necessary to scrape the data from the inconsistent web pages, which in itself presented a difficulty. Having done so, labelling the data is not necessarily as straightforward as I suggested in 1.2.2, as MPs do not always vote consistently with their own views.

Manually labelling speeches is cumbersome and not possible within this project (give number of relevant speeches). Manually labelling MPs’ stances is also very time consuming, as it would require a lot of research reading old newspaper articles, and even then, the stances of some of the lesser known MPs might still be ambiguous. Therefore, in this project I will have to use MPs’ voting records to label their speeches. Due to reasons I have previously mentioned, MPs will not necessarily vote consistently with their own views, meaning that there will be some degree of noise in the data. This issue is particularly great for the Iraq war, since the debate split parties and the leadership of both the Labour Party and Conservative Party were in favour of the war; before the Iraq war vote on 18th March 2003, the leaders of the two largest parties, Tony Blair (Labour Party) and Iain Duncan-Smith (Conservative Party) used party whips to persuade MPs to vote to support the invasion.

One result of this project will be to provide evidence in support of one of the following hypotheses:

- *Hypothesis 1* : MPs who feel strongly enough about a given issue to speak about it in the House of Commons will not vote on their own convictions, rather than in accordance with their Party's line.
- *Hypothesis 2* : Political Parties can influence their own MPs to vote with the party line, regardless of the views of the individual MPs.

Admittedly, viewing the issues in such a binary way is an oversimplification and how an MP's party affects their vote will differently vary on a vote-by-vote basis for different MPs. While this simplification would be naive in a Politics dissertation, overlooking nuances will allow for greater technical discussion, relevant to Computer Science dissertation. In essence, the better the performance of the classifier produced, the greater weight will be given to Hypothesis 1 over Hypothesis 2. This noise inherent in the data presents difficulties which can be mitigated through the design of the classifier.

1.4 Previous Work

While many studies into sentiment analysis have been based around political issues, since 2009 the majority of such research has concerned Tweets. The first study to use Twitter as its primary corpus was 'Twitter power: Tweets as electronic word of mouth' [11] and there have since been countless studies following suit. In 2010, Pak, Alexander and Paroubek, Patrick proposed that Twitter could be used to determine public opinion [17], which was proven true later that year when sentiment analysis of Twitter provided predictions that paralleled the results of traditional election polls for the German federal election [20].

This focus on Twitter is useful, but since most political decisions are made in Government and not on the internet, I propose that we should scrutinise the opinions of our elected politicians more than Twitter users. Part of the reason that what our MPs do in Parliament is not considered as much as it should be is due to the inaccessibility of the House of Commons; the language used by MPs when debating is unnaturally formal, making it difficult to follow and unnecessarily long-winded, whereas Tweets are inherently short and easy to interpret. A sentiment analysis system essentially summarises a text, meaning that applying sentiment analysis to Parliamentary debates would be a useful stepping stone towards summarising a debate. Unfortunately, although it makes the project more interesting, analysing longer political texts presents more challenges than analysing Tweets, in part due to the lack of guidance from similar previous work.

In the U.S., there have been a small number of papers detailing sentiment analysis on transcripts of Congress debates [1] [10]. The results of these studies indicate that determining an MP's stance on the Iraq war from their speeches in the House of Commons may be possible, however these papers use transcripts to determine a politician's political party, which is likely to be more clear-cut than their stance on a particular issue.

The lack of relevant works to this project highlights its uniqueness, which is one of the principal motivations for the project.

Chapter 2

Preparation

There were three main stages to the preparation of the project:

1. Learning about the necessary concepts and methods. This was useful as it helped me to make informed decisions about implementation decisions. This required considerable work, as most of the skills and knowledge required to undertake the project are not taught in the Cambridge BA Computer Science course and the parts that are taught are Part II courses. The time scale of the Part II project meant that I had to learn the courses ahead of the lectures. Sections 2.1 through 2.7 detail this stage of the preparation and section 2.8 is also strongly linked to this stage.
2. Defining and planning the project. A project of this scale needs clear definition of its goals and a well defined plan designed to achieve these goals. Sections 2.9 through 2.11 detail this stage of the preparation.
3. Selecting the tools to be used for implementation. Section 2.12 details this stage of the preparation.

2.1 Introduction to Supervised Learning

A supervised learning problem the task of determining the label of a given input. This is split into two phases: Learning and predicting.

In the learning phase, the system receives inputs of feature vectors and their associated labels. A feature vector of length k is usually denoted by \mathbf{x} where

$$\mathbf{x} \in \{x_1, x_2, \dots, x_k\} \quad \forall i \in \mathbb{Z}^+. \forall x_i \in \mathbb{R}. \quad (2.1)$$

A feature vector contains encodes the information necessary to predict a label. In the context of this project, there is a feature vector for each speech we consider, which contains information about the words in the speech. The label is usually denoted by y . The set of values that y can take varies depending on the context of the supervised learning problem. For example, in a regression problem, $y \in \mathbb{R}$. This project concerns binary classification, since we simplify the problem so that we consider all speeches to be

either pro-war or anti-war. Because of this, from now on, we will only consider binary classification problems, that is where $y \in \{-1, +1\}$.

The supervised learning system creates a function h that takes a feature vector as an input and outputs a label. That is

$$h(\mathbf{x}) = y. \quad (2.2)$$

This definition allows us to intuitively view each feature vector as a point in k -dimensional space. We consider each point to be either negative ($y = -1$) or positive ($y = +1$). In this analogy, h is a function that determines whether a point is negative or positive, depending on where it is in the k -dimensional space. The more points that h sees, the better its estimation of whether new unseen points are negative or positive. The prediction phase is where the system receives unlabelled feature vectors outputs estimates of the corresponding labels. The figure below shows a visualisation of our intuition of feature vectors, where $k = 2$. In this diagram, the supervised learning system learns a function to distinguish the '-' and '+' points. Given a new, previously unseen point, this function would be able to estimate whether it is a '-' or a '+'.

2.2 Introduction to the Naïve Bayes Classifier

This is one of the simplest classifiers to understand and implement. It uses the assumption that all features are independent of each other:

$$p(x_i|x_j) = p(x_i) \quad \forall i, j \in \mathbb{Z}^+. \quad (2.3)$$

We say that the classifier is 'naïve' because of this assumption. Although the assumption is very rarely true, the classifier still provides good performance [14].

In addition to this assumption, the classifier uses Bayes Theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (2.4)$$

The intuition behind the classifier is that given a set of features \mathbf{x} , we should assign it to the class that has the highest probability, given the set of features. We can compute this probability as follows:

$$p(C = y|\mathbf{x}) = \frac{p(\mathbf{x}|C = y)p(C = y)}{p(\mathbf{x})} \quad (2.5)$$

$$= \frac{p(x_1, \dots, x_k|C = y)p(C = y)}{p(\mathbf{x})} \quad (2.6)$$

$$= \frac{p(x_1, \dots, x_k, C = y)}{p(\mathbf{x})} \quad (2.7)$$

$$= \frac{p(x_1|x_2, \dots, x_k, C = y)p(x_2, \dots, x_k, C = y)}{p(\mathbf{x})} \quad (2.8)$$

2.3 Introduction to Support Vector Machines

2.4 Introduction to Cross-Validation

2.5 Introduction to the Bag of Words Model

2.6 Introduction to Entailment

2.7 Introduction to the House of Commons

2.8 Spam Email Dataset

2.9 Requirements Analysis

In my initial Project Proposal [Appendix A] I outlined the project as

2.10 Implementation Approach

2.11 Software Engineering Techniques

2.12 Choice of Tools

2.13 Starting Point

2.14 Summary

Chapter 3

Implementation

3.1 Data Scraper

3.2 Database

3.3 Classifier

3.4 Entailment System

3.5 Summary

Chapter 4

Evaluation

4.1 Unit Testing

4.2 Internal Evaluation

4.3 External Evaluation

4.4 Evaluation of Project Goals

4.5 Summary

Chapter 5

Conclusion

5.1 Achievements

5.2 Lessons Learned

5.3 Future Work

Bibliography

- [1] Maneesh Bhand, Dan Robinson, and Conal Sathi. Text classifiers for political ideologies. 2009.
- [2] C. M Bishop. *Pattern recognition and machine learning*. Information science and statistics. Springer, New York, 2006.
- [3] John Chilcot et al. The report of the iraq inquiry. *Report of a Committee of Privy Counsellors*. London, UK: House of Commons, 2016.
- [4] Cambridge Dictionary. Corpus dictionary definition.
- [5] A. Feldman, B. Harris, and J. Urban. Government tracker.
- [6] Carmen Fishwick. ‘we were ignored’: anti-war protesters remember the iraq war marches.
- [7] UK Government. Current state of the parties.
- [8] UK Government. Open parliament licence v3.0.
- [9] The Guardian. Tony blair’s key quotes.
- [10] Mohit Iyyer, Peter Enns, Jordan Boyd-Graber, and Philip Resnik. Political ideology detection using recursive neural networks. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 1113–1122, 2014.
- [11] Bernard J. Jansen, Mimi Zhang, Kate Sobel, and Abdur Chowdury. Twitter power: Tweets as electronic word of mouth. *Journal of the American Society for Information Science and Technology*, 60(11):2169–2188.
- [12] Dow Jones. Factiva.
- [13] Dan Jurafsky and James H. Martin. *Speech and language processing : an introduction to natural language processing, computational linguistics, and speech recognition / Daniel Jurafsky and James H. Martin*. Prentice Hall series in artificial intelligence. 2nd ed., international ed. edition, 2009.
- [14] Kevin P. Murphy. *Machine learning [electronic resource] : a probabilistic perspective / Kevin P. Murphy*. Adaptive computation and machine learning series. MIT Press, Cambridge, MA, 2012.

- [15] BBC News. Did your mp support the rebels?
- [16] House of Commons. Hansard archives.
- [17] Alexander Pak and Patrick Paroubek. Twitter as a corpus for sentiment analysis and opinion mining. In *LREc*, volume 10, 2010.
- [18] P. Robinson, P. Goddard, R. Brown, and P.M. Taylor. Content and framing study of united kingdom media coverage of the iraq war, 2003.
- [19] Kareem Shaheen. Us-led coalition says its strikes have killed 800 iraqi and syrian civilians.
- [20] Andranik Tumasjan, Timm Oliver Sprenger, Philipp G Sandner, and Isabell M Welp. Predicting elections with twitter: What 140 characters reveal about political sentiment. *Icwsn*, 10(1):178–185, 2010.

Appendix A

Project Proposal

Assessors like to see some sample code or example circuit diagrams, and appendices are the sensible places to include such items. Accordingly, software and hardware projects should incorporate appropriate appendices. Note that the 12,000 word limit does not include material in the appendices, but only in extremely unusual circumstances may appendices exceed 10-15 pages - if you feel that such unusual circumstances might apply to you you should ask your Director of Studies and Supervisor to apply to the Chairman of Examiners. It is quite in order to have no appendices. Appendices should appear between the bibliography and the project proposal.

Computer Science Tripos – Part II – Project Proposal

Sentiment Analysis of British Newspaper Articles on the Iraq War

Louis Slater, Pembroke College

Originator: Louis Slater

12th October 2017

Project Supervisor: Dr Tamara Polajnar

Director of Studies: Dr Anil Madhavapeddy

Project Overseers: Dr Timothy Griffin & Professor Anuj Dawar

Introduction

While there have been lots of studies involving sentiment analysis of political texts to determine their bias, none of these have uniquely involved British newspaper articles. Furthermore, after extensive research, I have not found any sentiment analysis of articles to determine their stance on a war. The purpose of this project is to develop a program that can reasonably determine the stance of any British newspaper article on the Iraq war. The core part of this project will be developing a program that achieves this using a bag-of-words method.

Starting point

In the past few years, there has been a lot of research into determining political biases of shorter segments of text, such as Tweets in the 2010 paper by Pak and Paroubek on Twitter as a Corpus for Sentiment Analysis and Opinion Mining. On the other hand, Political Ideology Detection Using Recursive Neural Networks by Iyyer, Enns, Boyd-Graber and Resnik uses a corpus containing longer texts US Congressional floor debate transcripts. Although there are clear differences between these transcripts and the newspaper articles that I plan to use (for example, the fact that the transcripts were initially spoken, whereas the articles were not), there are also many similarities (for example the length and inherently political nature of the corpora). Since this study showed that a bag-of-words method can successfully determine the bias of these transcripts with a 65% accuracy, it is justified to use a similar model to determine the bias of the newspaper articles I shall analyse.

The corpus I will use will be articles on the Iraq war from up to seven of the UKs most popular national daily newspapers and their Sunday equivalents published between 16th March 2003 and 18th April 2003. I will use a database of these articles compiled by Robinson, Goddard, Brown and Taylor in their 2003 study, Content and Framing Study of United Kingdom Media Coverage of the Iraq War, in which they manually determine the stance of 4,893 news articles from seven British daily newspapers and their Sunday equivalents (Daily Telegraph, The Times, The Guardian/The Observer, The Independent, The Daily Mail, The Mirror, The Sun/News of the World). This database does not include the articles texts, so the first part of my implementation will be to scrape this data from as many of these articles as possible. I will be able to get these texts from existing online Newspaper archives. I have already found searchable archives for The Guardian, The Observer, The Daily Telegraph, The Sunday Telegraph, The Independent, Indy on Sunday, The Times and the Sunday Times, all of which I will be able to use. Scraping textual data from the other newspapers in the database may be prove more difficult, but I will as many possibilities as I feasibly can within the scope of the project.

Resources required

In addition to the database and archives mentioned above, I will also require the use of a computer. I intend to mainly use my own computer, which has an Intel Core i7 processor and runs Windows 10. I will use the computing facilities in my college if my laptop is lost, broken or stolen. I will back up my work using both Google Drive and GitHub, which I will also use as a version control repository. I may also require the use of a server or external hard drive to store the corpus I use; however, this will be dependent on the amount of data that I scrape in the initial part of my project.

Work to be done

The project breaks down into the following sub-projects:

1. Gaining access to as many of the relevant searchable newspaper archives as possible.
2. Scraping data from as many articles as possible in the database compiled by Robinson, Goddard, Brown and Taylor.
3. Implementing a program to determine the biases of texts on the Iraq war, using the corpus I gather, along with corresponding the Reporters Tones from the database compiled by Robinson, Goddard, Brown and Taylor.
4. Running the program on the texts and comparing the results with the manually determined biases to judge the effectiveness of the program.

Success criteria

The project will be a success if I develop a program that can determine the stance of an article on the Iraq war with a greater than 50% accuracy.

Possible extensions

If I meet my success criteria early, I shall attempt one, or both, of the following extensions:

- Implementing a program that performs the same function as the initial program I develop, but using a different method, such as a recursive neural network. If I complete this extension, I will be able to compare the effectiveness of the two methods.
- Extrapolating the results using new datasets and analysing these results. Possible datasets I could use are newspaper articles from different countries, publications or times or transcripts of parliamentary debates.

Timetable

Planned starting date is the beginning of Michaelmas Week 3 (Thursday 19th October 2017).

1. **Michaelmas week 3** Gain access to as many of the relevant searchable newspaper archives as possible.
2. **Michaelmas weeks 4–5** Scrape data from as many articles as possible in the database compiled by Robinson, Goddard, Brown and Taylor, creating a database of the texts, their manually determined bias and other relevant information on them. If necessary, I will also get access to a server and store the database I compile on this server.

3. **Michaelmas weeks 6–8** Implement a program to determine the biases of texts on the Iraq war, using the corpus I gather, along with corresponding the Reporters Tones from the database compiled by Robinson, Goddard, Brown and Taylor.
4. **Michaelmas vacation** Finish the implementation, then run the program on the texts and compare the results with the manually determined biases to judge the effectiveness of the program.
5. **Lent weeks 1–2** Write the progress report and start work on possible extensions of the project.
6. **Lent weeks 3–4** Finish the extensions to the project.
7. **Lent weeks 5–6** Write the first draft of the dissertation.
8. **Lent weeks 7–8** Revise the dissertation in accordance with feedback I receive from my supervisor.
9. **Easter vacation** Finish revising the dissertation and submit the final project.