CSE255-A Assignment 2 Report

GAO, Yansong

JIANG, Tong

Zhao, Ruiwen

ABSTRACT

In this assignment, we use Amazon movie reviews data to predict the rating of a given movie. Here we mainly use two models to build the predictor – linear regression and latent factor. When choosing the features of the linear regression, we use natural language processing techniques to analyze the reviews text to get whether the reviews is positive or negative.

Keywords

data mining, natural language processing

1. EXPLANATORY ANALYSIS

In assignment 2, we use Amazon movie reviews from http://snap.stanford.edu/data/web-Movies.html.

2. PREDICTIVE TASK

Our task is to predict rating based on user, item, and review text.

3. MODEL DESCRIPTION

In this assignment, we use two different models to do the prediction: linear regression and latent factor. In linear regression, we use natural language processing techniques to analysis the sentiment of each review text.

3.1 Sentiment Analysis

The basic idea of sentiment analysis is to break the sentence into words, and see if each word is in the sentiment dictionary. Our dictionary is combined with the dictionary downloaded from https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html, and the dictionary we created from the training data.

3.1.1 Build Dictionary

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WOODSTOCK '97 El Paso, Texas USA

© 2015 ACM. ISBN 123-4567-24-567/08/06...\$15.00

DOI: 10.475/123_4

Table 1	1:	Sentiment	words	generated
---------	----	-----------	-------	-----------

refundbr	neg : pos	66.9:1.0
bolls	neg : pos	58.5:1.0
dhampir	neg : pos	52.4:1.0
crappiest	neg : pos	48.8:1.0
chucked	neg : pos	41.6:1.0
poorlywritten	neg : pos	41.6:1.0
kagan	neg : pos	40.1:1.0
whatsover	neg : pos	38.0:1.0
scientologists	neg : pos	38.0:1.0
bloodraynebr	neg : pos	38.0:1.0
stunkbr	neg : pos	38.0:1.0
craptastic	neg : pos	38.0:1.0
stinker	neg: pos	37.4:1.0

Since the dictionary downloaded from web contains only general positive and negative words, but no "movie-ish" words, we need to find out these movie-ish words on our own.

To build our own dictionary, we chose 252977 positive reviews and 46637 negative reviews (Each review with rating higher or equal to 4.0 are regarded as positive reviews, and those with rating lower or equal to 2.0 are regarded as negative reviews). We break each review into a bag of words, and label them "positive" or "negative" according to the sentiment polarity of the review they belong to. Then we use nltk.NaiveBayesClassifier to train these features.

By using Naive Bayes Classifier, we can calculate if a word is more likely to be a positive word or a negative one by calculating the following ratio:

$$\frac{P(w \text{ is positive}|w)}{P(w \text{ is negative}|w)} = \frac{P(w \in r|\text{pos review }r)P(\text{pos})}{P(w \in r|\text{neg review }r)P(\text{neg})} \quad (1)$$

And then we pick out the words with this ratio higher than 5.0 or lower than 0.2 and add them to the corresponding sentiment dictionary. Some of the sentiment words we chose through Naive Bayes Classifier are as showed in Table 1. As we can find in Table 1, many words are just movie names or characters' names. For example, Dhampir is a creature half vampire half human, and this word has high neg/pos ratio, so we assume that movies with Dhampir characters are tend to be lame movies.

3.1.2 Sentiment Analysis

After building dictionary, we can easily know whether a word is positive or negative one. We calculate the number of positive words and negative words in a review, and add

these two number as two features in our linear regression model.

RELATED LITERATURE

NLP model

We use natural language process techniques to analyze the sentiment of the review text. In order to get a more comprehensive idea of the algorithm, we mainly read the two papers below:

Mining and Summarizing Customer Reviews Opinion Observer: Analyzing and Comparing Opinions on the Web

When implementing our own algorithm to analyze the sentiment of the review text, we adopted the rough idea in the paper, and we also used the sentiment dictionary mentioned in the paper.

THE BODY OF THE PAPER

Typically, the body of a paper is organized into a hierarchical structure, with numbered or unnumbered headings for sections, subsections, sub-subsections, and even smaller sections. The command \section that precedes this paragraph is part of such a hierarchy. LATEX handles the numbering and placement of these headings for you, when you use the appropriate heading commands around the titles of the headings. If you want a sub-subsection or smaller part to be unnumbered in your output, simply append an asterisk to the command name. Examples of both numbered and unnumbered headings will appear throughout the balance of this sample document.

Because the entire article is contained in the **document** environment, you can indicate the start of a new paragraph with a blank line in your input file; that is why this sentence forms a separate paragraph.

Type Changes and *Special* Characters **5.1**

We have already seen several typeface changes in this sample. You can indicate italicized words or phrases in your text with the command \textit; emboldening with the command \textbf and typewriter-style (for instance, for computer code) with \texttt. But remember, you do not have to indicate typestyle changes when such changes are part of the structural elements of your article; for instance, the heading of this subsection will be in a sans serif² typeface, but that is handled by the document class file. Take care with the use of³ the curly braces in typeface changes; they mark the beginning and end of the text that is to be in the different typeface.

You can use whatever symbols, accented characters, or non-English characters you need anywhere in your document; you can find a complete list of what is available in the $\cancel{L}T_{FX}$ User's Guide[?].

5.2 Math Equations

You may want to display math equations in three distinct styles: inline, numbered or non-numbered display. Each of the three are discussed in the next sections.

5.2.1 Inline (In-text) Equations

A formula that appears in the running text is called an inline or in-text formula. It is produced by the math environment, which can be invoked with the usual \begin. .\end construction or with the short form \$. . .\$. You can use any of the symbols and structures, from α to ω , available in LaTeX[?]; this section will simply show a few examples of in-text equations in context. Notice how this equation: $\lim_{n\to\infty} x=0$, set here in in-line math style, looks slightly different when set in display style. (See next section).

5.2.2 Display Equations

A numbered display equation – one set off by vertical space from the text and centered horizontally – is produced by the **equation** environment. An unnumbered display equation is produced by the **displaymath** environment.

Again, in either environment, you can use any of the symbols and structures available in LATEX; this section will just give a couple of examples of display equations in context. First, consider the equation, shown as an inline equation above:

$$\lim_{n \to \infty} x = 0 \tag{2}$$

Notice how it is formatted somewhat differently in the displaymath environment. Now, we'll enter an unnumbered equation:

$$\sum_{i=0}^{\infty} x + 1$$

and follow it with another numbered equation:

$$\sum_{i=0}^{\infty} x_i = \int_0^{\pi+2} f$$
 (3)

just to demonstrate LATEX's able handling of numbering.

5.3 Citations

Citations to articles [?, ?, ?, ?], conference proceedings [?] or books [?, ?] listed in the Bibliography section of your article will occur throughout the text of your article. You should use BibTeX to automatically produce this bibliography; you simply need to insert one of several citation commands with a key of the item cited in the proper location in the .tex file [?]. The key is a short reference you invent to uniquely identify each work; in this sample document, the key is the first author's surname and a word from the title. This identifying key is included with each item in the .bib file for your article.

The details of the construction of the .bib file are beyond the scope of this sample document, but more information can be found in the Author's Guide, and exhaustive details in the \(\mathbb{H}T_{EX}\) User's Guide[?].

This article shows only the plainest form of the citation command, using \cite. This is what is stipulated in the SIGS style specifications. No other citation format is endorsed or supported.

Tables 5.4

¹This is the second footnote. It starts a series of three footnotes that add nothing informational, but just give an idea of how footnotes work and look. It is a wordy one, just so you see how a longish one plays out.

A third footnote, here. Let's make this a rather short one

to see how it looks.

³A fourth, and last, footnote.

Table 2: Frequency of Special Characters

Non-English or Math	Frequency	Comments
Ø	1 in 1,000	For Swedish names
π	1 in 5	Common in math
\$	4 in 5	Used in business
Ψ_1^2	1 in 40,000	Unexplained usage



Figure 1: A sample black and white graphic.

Because tables cannot be split across pages, the best placement for them is typically the top of the page nearest their initial cite. To ensure this proper "floating" placement of tables, use the environment **table** to enclose the table's contents and the table caption. The contents of the table itself must go in the **tabular** environment, to be aligned properly in rows and columns, with the desired horizontal and vertical rules. Again, detailed instructions on **tabular** material is found in the $\rlap/$ ETEX User's Guide.

Immediately following this sentence is the point at which Table 1 is included in the input file; compare the placement of the table here with the table in the printed dvi output of this document.

To set a wider table, which takes up the whole width of the page's live area, use the environment **table*** to enclose the table's contents and the table caption. As with a single-column table, this wide table will "float" to a location deemed more desirable. Immediately following this sentence is the point at which Table 2 is included in the input file; again, it is instructive to compare the placement of the table here with the table in the printed dvi output of this document.

5.5 Figures

Like tables, figures cannot be split across pages; the best placement for them is typically the top or the bottom of the page nearest their initial cite. To ensure this proper "floating" placement of figures, use the environment **figure** to enclose the figure and its caption.

This sample document contains examples of .eps files to be displayable with LATEX. If you work with pdfLATEX, use files in the .pdf format. Note that most modern TEX system will convert .eps to .pdf for you on the fly. More details on each of these is found in the *Author's Guide*.

As was the case with tables, you may want a figure that



Figure 2: A sample black and white graphic that has been resized with the includegraphics command.

spans two columns. To do this, and still to ensure proper "floating" placement of tables, use the environment **figure*** to enclose the figure and its caption. and don't forget to end the environment with figure*, not figure!

5.6 Theorem-like Constructs

Other common constructs that may occur in your article are the forms for logical constructs like theorems, axioms, corollaries and proofs. There are two forms, one produced by the command \newtheorem and the other by the command \newdef; perhaps the clearest and easiest way to distinguish them is to compare the two in the output of this sample document:

This uses the **theorem** environment, created by the **\newtheorem** command:

THEOREM 1. Let f be continuous on [a,b]. If G is an antiderivative for f on [a,b], then

$$\int_{a}^{b} f(t)dt = G(b) - G(a).$$

The other uses the **definition** environment, created by the **\newdef** command:

Definition 1. If z is irrational, then by e^z we mean the unique number which has logarithm z:

$$\log e^z = z$$

Two lists of constructs that use one of these forms is given in the *Author's Guidelines*.

There is one other similar construct environment, which is already set up for you; i.e. you must *not* use a **\newdef** command to create it: the **proof** environment. Here is a example of its use:

PROOF. Suppose on the contrary there exists a real number L such that

$$\lim_{x \to \infty} \frac{f(x)}{g(x)} = L.$$

Ther

$$l = \lim_{x \to c} f(x) = \lim_{x \to c} \left[gx \cdot \frac{f(x)}{g(x)} \right] = \lim_{x \to c} g(x) \cdot \lim_{x \to c} \frac{f(x)}{g(x)} = 0 \cdot L = 0,$$

which contradicts our assumption that $l \neq 0$.

Complete rules about using these environments and using the two different creation commands are in the *Author's Guide*; please consult it for more detailed instructions. If you need to use another construct, not listed therein, which you want to have the same formatting as the Theorem or the Definition[?] shown above, use the \newtheorem or the \newdef command, respectively, to create it.

A Caveat for the TEX Expert

Because you have just been given permission to use the \newdef command to create a new form, you might think you can use TEX's \def to create a new command: Please refrain from doing this! Remember that your LATEX source code is primarily intended to create camera-ready copy, but may be converted to other forms – e.g. HTML. If you inadvertently omit some or all of the \defs recompilation will be, to say the least, problematic.

Table 3: Some Typical Commands

Command	A Number	Comments				
\alignauthor	100	Author alignment				
\numberofauthors	200	Author enumeration				
\table	300	For tables				
\table*	400	For wider tables				

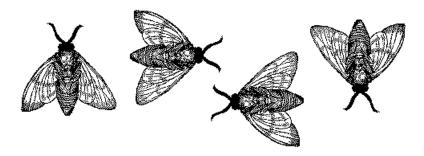


Figure 3: A sample black and white graphic that needs to span two columns of text.



Figure 4: A sample black and white graphic that has been resized with the includegraphics command.

6. CONCLUSIONS

This paragraph will end the body of this sample document. Remember that you might still have Acknowledgments or Appendices; brief samples of these follow. There is still the Bibliography to deal with; and we will make a disclaimer about that here: with the exception of the reference to the LATEX book, the citations in this paper are to articles which have nothing to do with the present subject and are used as examples only.

7. ACKNOWLEDGMENTS

This section is optional; it is a location for you to acknowledge grants, funding, editing assistance and what have you. In the present case, for example, the authors would like to thank Gerald Murray of ACM for his help in codifying this Author's Guide and the .cls and .tex files that it describes.

APPENDIX

A. HEADINGS IN APPENDICES

The rules about hierarchical headings discussed above for the body of the article are different in the appendices. In the **appendix** environment, the command **section** is used to indicate the start of each Appendix, with alphabetic order designation (i.e. the first is A, the second B, etc.) and a title (if you include one). So, if you need hierarchical structure within an Appendix, start with **subsection** as the highest level. Here is an outline of the body of this document in Appendix-appropriate form:

A.1 Introduction

A.2 The Body of the Paper

A.2.1 Type Changes and Special Characters

A.2.2 Math Equations

Inline (*In-text*) *Equations*.

Display Equations.

A.2.3 Citations

A.2.4 Tables

A.2.5 Figures

A.2.6 Theorem-like Constructs

A Caveat for the T_FX Expert

A.3 Conclusions

A.4 Acknowledgments

A.5 Additional Authors

This section is inserted by IATEX; you do not insert it. You just add the names and information in the \additionalauthors command at the start of the document.

A.6 References

Generated by bibtex from your .bib file. Run latex, then bibtex, then latex twice (to resolve references) to create the .bbl file. Insert that .bbl file into the .tex source file and comment out the command **\thebibliography**.

B. MORE HELP FOR THE HARDY

The sig-alternate.cls file itself is chock-full of succinct and helpful comments. If you consider yourself a moderately experienced to expert user of LaTeX, you may find reading it useful but please remember not to change it.