Word Extraction

Bruno Zecchi

For this assignment we were tasked with extracting a list of CEO names, company names, and percentages from a record of news articles published each day of 2013 and 2014. We took given training data, built negative samples, and developed a logistic regression classification model that assigned a token to one of four categories: CEO name, company name, percentage, or other. Seven features were used that analyzed the characteristics of the token's characters, as well as checking whether a token was inside one of the prebuilt lists of the three entities. We checked the model on a testing dataset and determined it to be highly accurate (90% +). Lastly, we applied the model to each of the individual articles and extracted their relevant entities. The model appears to be accurate with percentages and company names, but further work is required in extracting CEO names. Further features need to be included that help the model better distinguish CEO names from regular names. A more detailed explanation of the analysis and model evaluation can be seen below.

Analysis Process

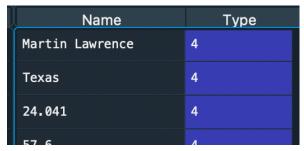
<u>Data Cleaning:</u> The given training labeled data is imported and cleaned. All duplicates and excess space characters are removed.

<u>Building Integrated Dataframe:</u> We combined the 3 different given labels into a single dataframe, and provided the label type (CEO-1, company-2, or percentage-3) into a new column called "Type".

Index	Name	Туре
1758	99 percent	3
2403	15900%	3
272	Jacobs	1
5946	Service Co Ltd	2
2191	1.75 percent	3

IEMS 308 Assignment 3 Northwestern University

<u>Form Negative Samples:</u> Data points are taken from the internet and will act as negative samples when building the classification model. Negative samples include **names of famous people, fictional company names, locations, and random English words**.

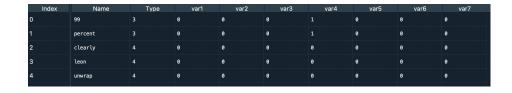


Negative Samples Dataframe

<u>Building Final Dataframe:</u> The real labels and negative samples are mixed up and combined into a final dataframe. The "Type" column is added so that testing of a model can be done later.

<u>Develop Features</u>: Each token of the final dataframe is given a 1/0 value for the seven features developed. The features are described below:

- Feature 1: Is token in prebuilt lists of CEO names?
 - o Is it a last name?
 - o Is a mix of this token and the next three tokens a full name in list?
- Feature 2: Is token in prebuilt list of percentages?
- Feature 3: Is token in prebuilt list of companies?
 - o Is a mix of this token and the next three tokens a full name in list?
- Feature 4: Is the next token "%" or "percent"?
- Feature 5: Is token's first letter capitalized?
- Feature 6: Is current token and next token capitalized?
- Feature 7: Do next three tokens contain "Co", "Co.", "Inc", etc?



<u>Build Classification Model:</u> We use a logit model that takes the seven features and the "Type" column of a training subset to build a classifier. We then test the model on a testing subset and evaluate its accuracy. The results of the evaluation are shown below:

	Percentage Predicted
Total	97%
CEO Names	97%
Company Names	90%
Other	97%
Incorrect	3%
Predictions	

It is worth noting that the performance of the model is hard to measure at this point. Testing dataset comes from the same source of the training dataset, so the first three features will be very accurate. They might not be as accurate with outside data.

Apply Model to Text Files: We build a for-loop that opens each of the articles and identifies all CEO names, company names, and percentages in it. Those values are appended to a final list for each corresponding type. Those final lists now contain those three types of labels in all articles. The three entities extracted are displayed in the 3 given txt files: finalCEO.txt, finalCompany.txt, and finalPercentage.txt.

Model Performance & Future Development

By analyzing the text files, we see that the model does a good job of extracting percentages from the articles.

10.77% 12.85% 25% 1.3% 0.7% 2.35%	Bloomberg Street Twitter "New York Times" Google HSBC	In Lu Bryant First "Federal Reserve" "Bill Gross" Johnson Gross "Janet Yellen"
Percentages Extract	Company Name Extracts	CEO Names Extracts

We also see that the model also did a good job of extracting company names from the articles. However, it seemed less accurate in extracting CEO names. Several non-CEO names appear throughout the file, as well as location and company names. This fall in performance could be linked to the fact that it is hard to distinguish a CEO name form a regular name. Many of the features will be the exact same for CEO and non-CEO name. The difference can only be seen by researching the name (something our model cannot do) or comparing a name to a prebuilt list. If a given CEO name is not in the prebuilt list or is written in a different format, the model will have a hard time understanding what it is. Further developments of the model should involve adding features that are able to distinguish CEO and non-CEO names better.