

=====

Problem Set 6

=====

Use the `business_insider_text_data.csv` and `vix_data.csv` datasets available on CCLE for this exercise. The Business Insider data is downloaded from www.kaggle.com (a really cool website for data and code, if you haven't seen it). You can and should look at the raw data using Excel before starting this exercise. The `vix_data.csv` file contains two columns where "label" indicated whether the closing value of the VIX is higher (1) or lower (0) than the closing value of the VIX for the previous trading day. The Business Insider data contains various headlines downloaded from www.businessinsider.com.

Creating a sentiment index from text data

Before you start, please download the script "PS6_data_cleaning.R" from CCLE. Run it on the data to do some pre-pre-processing and merge the two datasets.

1. Use `Corpus(VectorSource(assignment_data))` to load the corpus. `VectorSource` makes each line a document, so now each document corresponds to a different date in the dataset.
2. Pre-process the data as in the lecture notes. Feel free to use the code from Code Snippets Topic 6 on CCLE. That is, remove numbers, make all lower case, remove stopwords, stemming, etc.
3. As in the lecture note, create a `DocumentTermMatrix`, call it `dtm`. Run the line `inspect(dtm)`. Notice that the matrix is quite *sparse* (a lot of zeros).
4. As in the lecture note, create a freq matrix as the column sums of `dtm`. Show in a bar plot the frequency of words that occur more than 25 times.
5. Create a wordcloud of the 20 most frequent words. Based on this (and 4.), how would you characterize the typical headline in terms of the news subject? Are there words that, intuitively, can matter for the stock market returns that day?
6. Create the data `"y_data <- as.factor(assignment_data$Label)"` and `"x_data <- as.matrix(dtm)"`. You will try to construct an index based on the words in `dtm` that predicts the direction of stock returns.
7. Split the data into a training data-set, based on data up to and including 2016-12-31. The remaining data should be used for actual out-of-sample testing.
8. We will first let the logistic regression create the word-based index. That is, try to fit a regular logistic regression using `y_data` and `x_data` and the training dataset. Explain why this doesn't work.

9. Next, run a logistic regression with an elastic net constraint (let $\alpha = 0.5$) using cross-validation and the training dataset. Why does the regression routine work now (ie, why does it give an answer (a coefficient vector; no meltdown))? Explain.
10. Using `lambda.min`, what (if any) are the words chosen and their associated coefficients? Comment on your results.
11. Now, create instead a pre-defined sentiment word list:

```
dtm_sentiment <-  
dtm[,c("trump", "invest", "growth", "grow", "high", "strong", "lead", "good", "risk", "debt", "oil", "loss",  
      "war", "rate", "hous", "weak")]
```

Run the elastic net with `x_data_pre <- as.matrix(dtm_sentiment)` using cross-validation and the training sample. Create a bar plot with the words on the x-axis and the coefficients on the y-axis. Comment on differences and similarities to the case in 10. Again, get the coefficients using `lambda.min`.

12. Create the ROC curves for the sentiment model in 11, using `lambda.min` to get coefficient vector. Is it better than random? You likely want to use the `predict` function to get the model predictions.
13. Now, using the **test sample** and the model in 11, what is the proportion of days the model would have made the right prediction in this new sample? Is it better than random (50/50)?