

*Alex Merryman*

First, the text was imported as Unicode to remove non-Unicode characters (likely vestiges from CSS formatting). Then spaCy was used to tokenize the corpus, identifying separate words and characters. Since the data is already plain-text and now in Unicode, there is minimal preprocessing to be done.

Since all percent-related numbers to be extracted are succeeded by “%,” “percent,” “percentage,” “percentile,” or some variation thereof, Regex can be used. Three classes of varying complexity of the desired percentage numbers were determined:

- Additionally, each of these classes can be used to specify a range of percentages, for example 20-35% / 20-35 percent / twenty to thirty-five percent.

The second class required a bit more complex Regex pattern:

Finally, the third class required an algorithm:

- This captures all desired percentages.

To identify CEOs, the algorithm again iterates across all tokens. Whenever the token lemma is also one found in the CEO keywords list, or two neighboring tokens are proper nouns, the algorithm searches for and identifies proper nouns (based on POS tagging) within 5 tokens of the CEO keyword.

### Company Identification

Similar to CEO identification, the algorithm searches for and identifies proper nouns within 7 tokens on either side of a token that is also found within the company keywords list.

My model identified:

- 77523 percentages
- 467255 CEOs
- 471223 companies

And it processed all 730 text files in less than 30 minutes.