

# Capstone Project 2: Analysis of Berkshire Hathaway Shareholder Letters Using Natural Language Processing (NLP) Techniques



Tom Halloin

SPRINGBOARD DATA SCIENCE CAREER TRACK

“Read 500 pages every day. That's how knowledge works. It builds up, like compound interest. All of you can do it, but I guarantee not many of you will do it.”

- Warren Buffett, when asked how to prepare for an investing career.

## The Problem

Warren Buffett became one of the most successful investors of all time because of his ability to comprehend thousands of pages of data. Buffett is the exception to the rule. Most people don't have the patience to read 500 pages a day, let alone comprehend it (and I am not the exception). Instead, I am going to use machine learning to summarize the equivalent of reading 500 pages, Warren Buffett's annual letters to shareholders since 1977. This capstone will be an exploration of Berkshire Hathaway's annual shareholder letters. Using Natural Language Processing, I will try to summarize the letters and create a topic model to discover hidden themes that are present across the letters. This model will try to extract meaning from these letters by identifying recurring themes or topics and their change from letter to letter.

## The Dataset

The data comes from Berkshire Hathaway's shareholder letters available on their website. The letters come in both HTML and PDF format. Using Python, I managed to scrape the text from annual letters dating back from 1977. See the code in the notebook linked here for more information. Note that Berkshire eventually caught on that I was a bot, and therefore denied me access to their reports. Fortunately, this was after I gathered all of my data, so while the code used to work, there is no way for it to work anymore without access. The letters are also saved as part of this Milestone Report.

## Data Wrangling

Before doing any analysis on the letters, I needed to do some data wrangling to clean the output from the pages and PDFs the letters came from into useable text. Before any data wrangling, the files looked like the following sample.

```
<!-- Global site tag (gtag.js) - Google Analytics -->
<script async src="https://www.googletagmanager.com/gtag/js?id=UA-136883390-1"></script>
<script>
  window.dataLayer = window.dataLayer || [];
  function gtag(){dataLayer.push(arguments);}
  gtag('js', new Date());

  gtag('config', 'UA-136883390-1');
</script>
<HTML>
<HEAD>
  <TITLE>Chairman's Letter - 1977</TITLE>
</HEAD>
<BODY>
<P ALIGN=CENTER>
<B>BERKSHIRE HATHAWAY INC.</B>
</P>
<PRE>

<I>To the Stockholders of Berkshire Hathaway Inc.:</I>

    Operating earnings in 1977 of $21,904,000, or $22.54 per
share, were moderately better than anticipated a year ago. Of
these earnings, $1.43 per share resulted from substantial
realized capital gains by Blue Chip Stamps which, to the extent
of our proportional interest in that company, are included in our
operating earnings figure. Capital gains or losses realized
directly by Berkshire Hathaway Inc. or its insurance subsidiaries
are not included in our calculation of operating earnings. While
too much attention should not be paid to the figure for any
single year, over the longer term the record regarding aggregate
capital gains or losses obviously is of significance.
```

Most of the cleaning for these files included removing HTML, anything before and after the actual contents of the letter, and non-ascii characters. A series of regular expressions removed unnecessary characters in tables and newline characters to make the text more readable. I skipped removing stop words and numbers because I want the ending summaries to read as close to actual sentences as possible.

The end result was a series of paragraphs in documents looking like this:

## Summarizing the Data

[illegible]

3

really good for doing extractive summarization, but abstractive summarization is much harder because it involves creating new text to explain the document so that humans can understand it.

There have been attempts with deep learning to create abstractive summaries, but that is beyond the scope of this project.

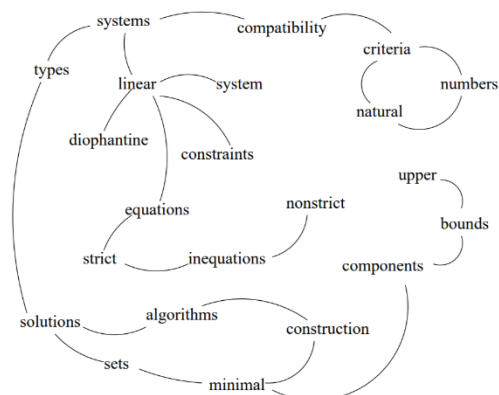
I will try three different extractive summarization methods: Textrank, Lexrank, and LSA. [Textrank](#) is a graph-based ranking model similar to Google's Pagerank algorithm. Each word links to other words based on the number of times the words appear near each other. The algorithm then assigns hidden weights to each of these links given how likely each of the linked words appears with each other.

Suppose in the example to the right I start with

“natural”. Natural is linked with “criteria” and “numbers”. Based on Textrank's weighting, however, it will determine that “natural numbers” is a better keyword than “natural criteria”. The weighting itself can be a black box – see *Textrank: Bringing Order into Texts* by Rada Mihalcea and Paul Tarau for more details on the mathematical implementation.

In the example shown on the right, connected words are good examples of keywords that can summarize documents. Examples of summarization of the text on the right include “linear Diophantine equations”, “upper bounds”, and “natural numbers”. Without any comprehension of the actual text, the graph shows that the text is at least about these three keywords. The algorithm will then look for sentences with these keywords to summarize the document.

Compatibility of systems of linear constraints over the set of natural numbers. Criteria of compatibility of a system of linear Diophantine equations, strict inequations, and nonstrict inequations are considered. Upper bounds for components of a minimal set of solutions and algorithms of construction of minimal generating sets of solutions for all types of systems are given. These criteria and the corresponding algorithms for constructing a minimal supporting set of solutions can be used in solving all the considered types systems and systems of mixed types.



**Keywords assigned by TextRank:**

linear constraints; linear diophantine equations; natural numbers; nonstrict inequations; strict inequations; upper bounds

**Keywords assigned by human annotators:**

linear constraints; linear diophantine equations; minimal generating sets; nonstrict inequations; set of natural numbers; strict inequations; upper bounds

See *Textrank: Bringing Order into Texts* by Rada Mihalcea and Paul Tarau, p. 4

[Lexrank](#) is very similar to Textrank. Both use the same graph-based weighting model, but differ in the method of generating weights. Lexrank uses a method of cosine similarity weighted by term frequency – inverse document frequency, otherwise known as TF-IDF, to weight terms that have more meaning instead of just looking at document frequencies. The end result of applying the Lexrank algorithm to a series of documents is a graph similar to the picture on the right here. Each circle is a document and the edges represent the cosine similarities between each of the documents. In this example, d1s1 (document 1, sentence 1) and d2s1 (document 2, sentence 1) have a high weighted cosine similarity and would be good candidates for summarizing the corpus of documents.

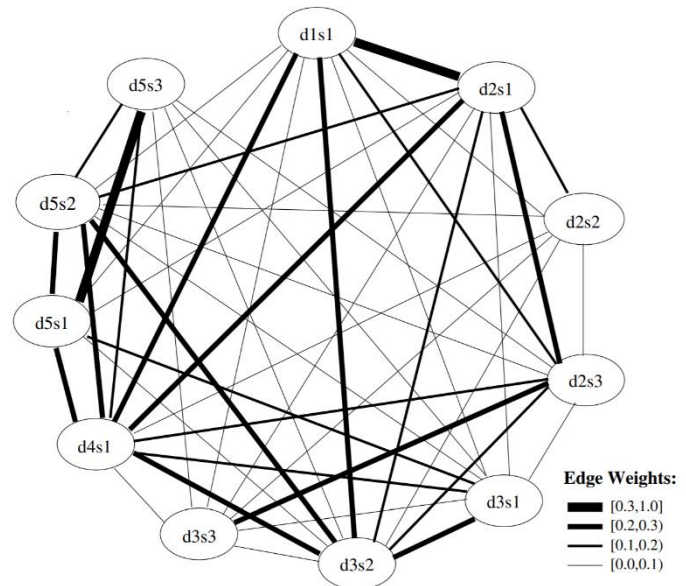


Figure 2: Weighted cosine similarity graph for the cluster in Figure 1.  
See *LexRank: Graph-based Lexical Centrality as Saliency in Text Summarization* by Gunes Erkan and Dragomir R. Radev for graph.

[LSA, or Latent Semantic Analysis](#), analyzes similarities between the meanings of a series of documents and produces a series of concepts related to the documents. Like Lexrank and Textrank, LSA creates a term-document matrix consisting of word frequencies for each term in each document. Each word is a row while each document is a column. After creating this matrix, a linear algebra technique called singular value decomposition reduces the number of words, while keeping the same meaning in each of the documents. By reducing the number of words in each document, it is easier to find similarities between different documents. Sentences that keep the most words from the reduction are considered the “most important”. These reduced words and phrases also generate topic models that will be shown later.

For each annual report, I used each of these three techniques to summarize each year in the corpus by extracting the five most meaningful sentences according to its method. Then, I used the Spacy library to create a summary score to see how similar the summaries are to each other. Spacy has a large vocabulary of words that can capture similar phrases even if the words do not match, so the techniques should produce similar summaries according to Spacy. The below graph plots these similarity scores between the three different methods. Each line is a comparison between two of the methods for each year in the annual report. The divergences, as shown in the year 1997, provide a good contrast to see the different sentences each technique deems as important.



Here are the summaries for each method in 1997:







## Milestone Report 2: Topic Modeling

### Topic Modeling

- Rationale for topic modeling, how it differs from summarizing
  - Differs in output, output in TM is clusters and latent embedding.
  - TM works on entire corpus.
  - Summarization is for each individual document.
- Different types of topic modeling, tools used, expectation for results
- Applying topic modeling to corpus
- Measuring coherence of topic models

### Recommendations