**Assignment 2**
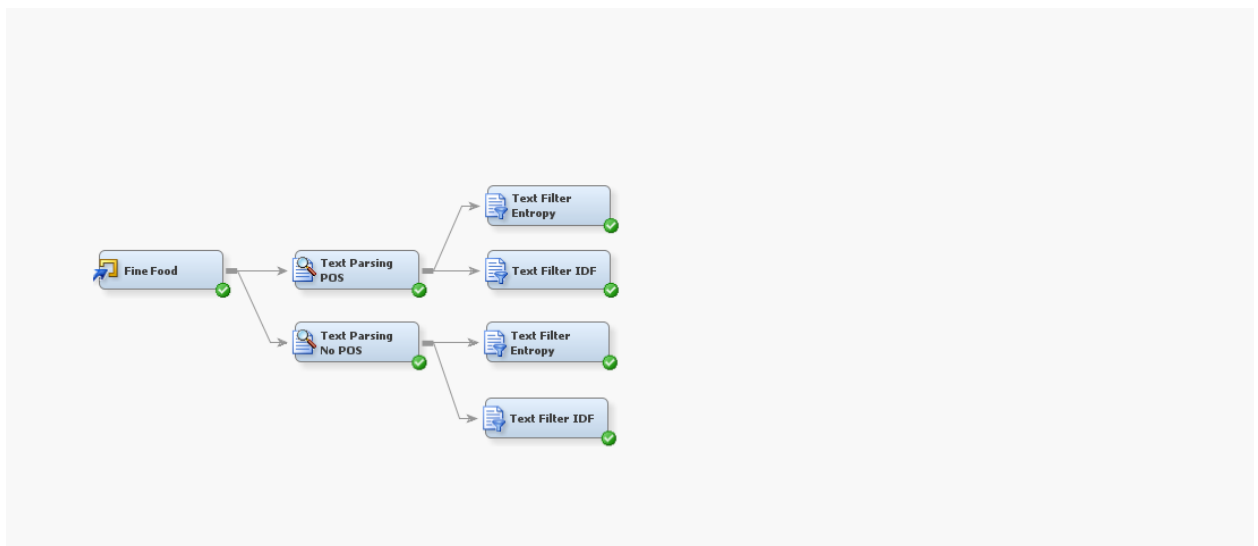
The goal of this assignment is to apply text parsing and filtering in real context. To do so do the following:

1- Go to the dataset: https://www.kaggle.com/snap/amazon-fine-food-reviews/downloads/amazon-fine-food-reviews.zip/2

(You may need to register to Kaggle.com to download the dataset)

2- Import the data file to SAS enterprise miner

3- Apply text parsing and filtering on the dataset.

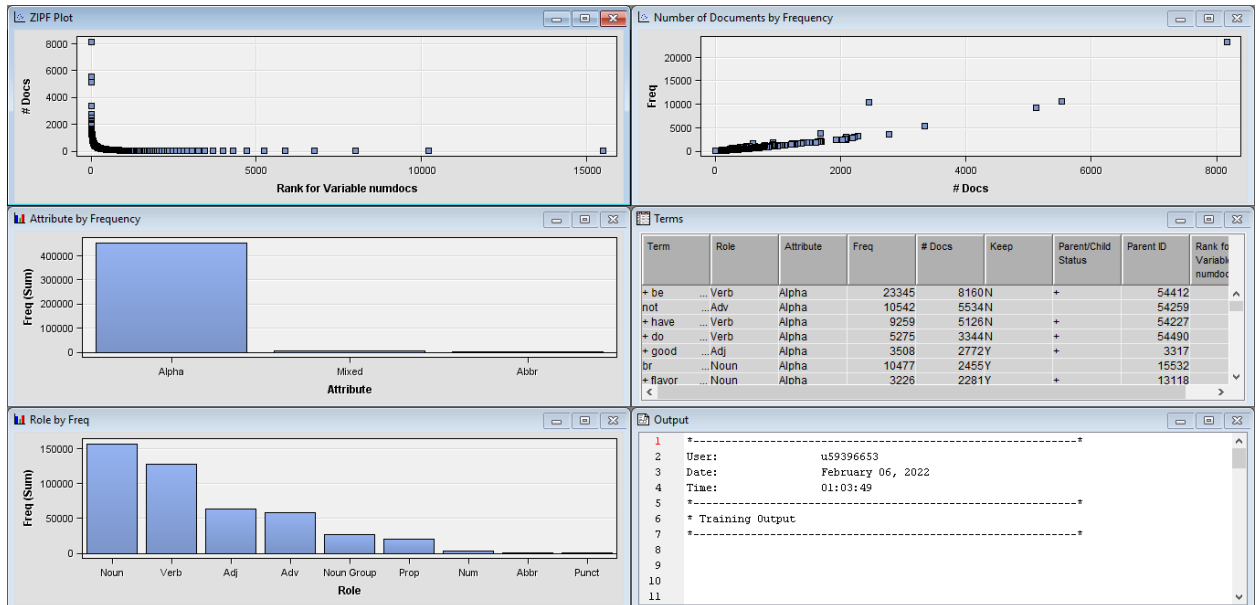4- Take screen captures of the outcome of your graphs and findings
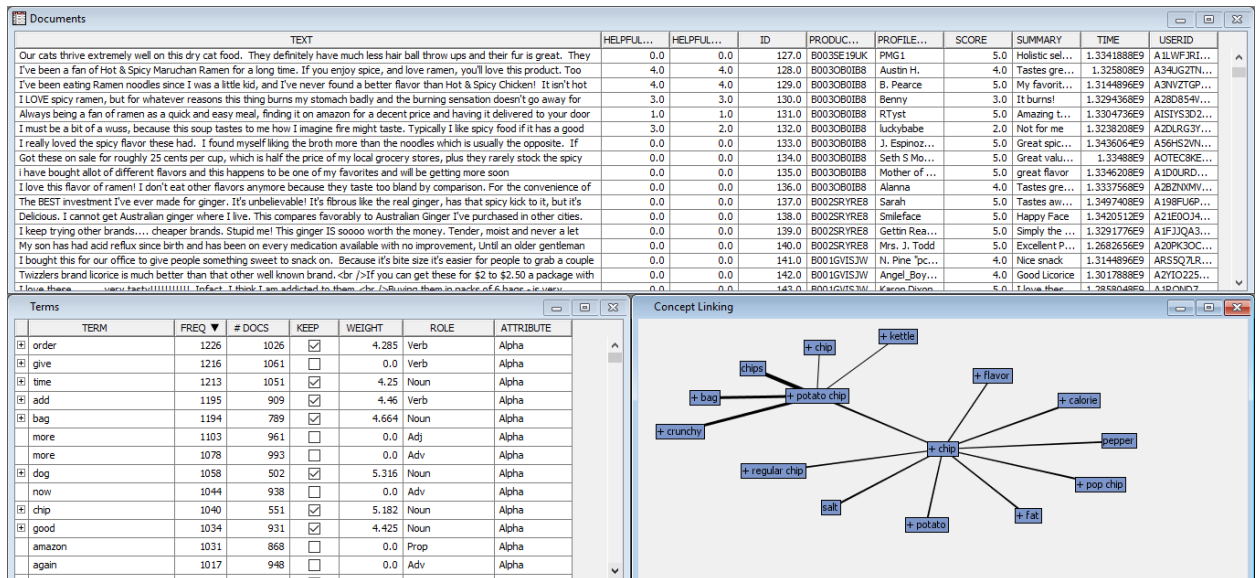


*Figure 1 Text Parsing POS Results*



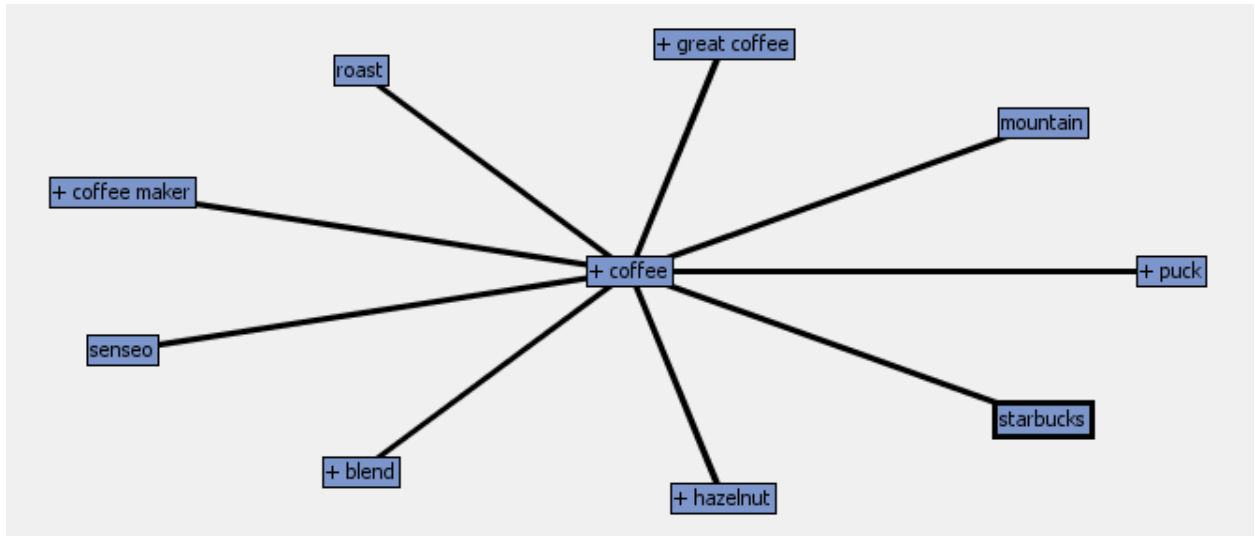*Figure 2 Text Filter Interactive View*

*Figure 3 Identifying the word Coffee and words associated with it*

5- Identify synonyms in your data and combine them together using interactive text explorer

6- What words/word-combinations have you combined?

- br occurred 10477 and it is code for line break, Amazon (Prop & noun) Occurred 497 times and it adds no meaning therefore they were dropped.
- Purchase & buy gives the same meaning in the data set.
- Great & Excellent has the same level in sentiment and depending on the analysis other words can be added as synonyms like good.
- Cat & Dog can be added if we are referring to pets.
- Coffee, Tea, Cocoa (Cocoa & Hot Cocoa are synonyms) can be combined as hot beverages.

7- Give some examples of the entities identified by the text parser

Starbuck, coffee, tea, bag, sugar, chocolate, juice, soda, cookie, potato chip

8- Repeat the previous steps while disabling POS. how are your results affected?

While using parts of speech words like great and excellent had to be matched with the right POS before treating them as synonyms, if POS was not used then it's not a problem.
By comparing the two methods some words changed in position in terms of frequency, but I don't think it make much of a difference while looking at the top 100 words for example.

| Term | Role | Attribute | Freq ▼ |
|------|------|-----------|--------|
| + be | ... | Alpha | 23526 |
| not | ... | Alpha | 10554 |
| br | ... | Alpha | 10477 |
| + have | ... | Alpha | 9261 |
| + good | ... | Alpha | 6184 |
| + do | ... | Alpha | 5430 |
| + taste | ... | Alpha | 4565 |
| + flavor | ... | Alpha | 4047 |
| + coffee | ... | Alpha | 3947 |
| + make | ... | Alpha | 3104 |
| + try | ... | Alpha | 3037 |
| + product | ... | Alpha | 3010 |
| + love | ... | Alpha | 2903 |
| + great | ... | Alpha | 2831 |
| very | ... | Alpha | 2822 |
| just | ... | Alpha | 2813 |
| + get | ... | Alpha | 2741 |
| + buy | ... | Alpha | 2673 |
| + use | ... | Alpha | 2626 |
| + find | ... | Alpha | 2228 |
| more | ... | Alpha | 2183 |
| + other | ... | Alpha | 2138 |
| + food | ... | Alpha | 2074 |
| + like | ... | Alpha | 2058 |
| + little | ... | Alpha | 2005 |
| + no | ... | Alpha | 1818 |
| so | ... | Alpha | 1799 |
| + one | ... | Alpha | 1792 |
| + tea | ... | Alpha | 1790 |
| really | ... | Alpha | 1777 |
| + eat | ... | Alpha | 1708 |
| + drink | ... | Alpha | 1708 |
| + cup | ... | Alpha | 1610 |
| too | ... | Alpha | 1598 |
| + order | ... | Alpha | 1597 |

*Figure 4 No POS*

| Term | Role | Attribute | Freq ▼ |
|------|------|-----------|--------|
| + be | ... Verb | Alpha | 23345 |
| not | ... Adv | Alpha | 10542 |
| br | ... Noun | Alpha | 10477 |
| + have | ... Verb | Alpha | 9259 |
| + do | ... Verb | Alpha | 5275 |
| + coffee | ... Noun | Alpha | 3755 |
| + good | ... Adj | Alpha | 3508 |
| + flavor | ... Noun | Alpha | 3226 |
| + make | ... Verb | Alpha | 3015 |
| + product | ... Noun | Alpha | 2999 |
| very | ... Adv | Alpha | 2812 |
| + get | ... Verb | Alpha | 2710 |
| + buy | ... Verb | Alpha | 2513 |
| + love | ... Verb | Alpha | 2484 |
| + taste | ... Noun | Alpha | 2449 |
| + try | ... Verb | Alpha | 2446 |
| + use | ... Verb | Alpha | 2283 |
| + taste | ... Verb | Alpha | 2097 |
| + find | ... Verb | Alpha | 2039 |
| + like | ... Verb | Alpha | 1970 |
| just | ... Adv | Alpha | 1903 |
| so | ... Adv | Alpha | 1789 |
| really | ... Adv | Alpha | 1775 |
| no | ... Adv | Alpha | 1767 |
| + food | ... Noun | Alpha | 1763 |
| + eat | ... Verb | Alpha | 1693 |
| + tea | ... Noun | Alpha | 1686 |
| too | ... Adv | Alpha | 1596 |
| + great | ... Adj | Alpha | 1561 |
| + little | ... Adj | Alpha | 1542 |
| other | ... Adj | Alpha | 1480 |
| one | ... Num | Alpha | 1475 |
| also | ... Adv | Alpha | 1441 |
| + go | ... Verb | Alpha | 1440 |
| + cup | ... Noun | Alpha | 1440 |

*Figure 5 POS*