## Data description

Being an avid shopper at Costco, I personally like the range of varieties it offers in all kinds of shopping. I usually suggest all my friends and family to shop at Costco. With the help of this data anlaysis, I can show them real statistics based on the reviews I searched. I conducted data analysis on costco to check how often people shop at costco and what they feel about the range of products at costco offers.

## Structured data after text processing



## Topic numbers

```
     document topic     gamma
        <chr> <int>     <dbl>
1        1        1 0.1988803
2        2        1 0.1990961
3        3        1 0.2030079
4        4        1 0.2017868
5        5        1 0.2024302
6        6        1 0.1963707
7        7        1 0.1928946
8        8        1 0.1985521
9        9        1 0.1956651
10      10        1 0.1983563
# ... with 34,030 more rows
```

## Sentiment Label

## Sentiment Score

## Hour of the day

## Test deciding on the number of topics generated





Simple approach is to analyze the metrics to find extremum.

minimization:

Arun2010

CaoJuan2009

maximization:

Deveaud2014

Griffiths2004

From this plot we made conclusion that optimal number of topics is in range 7-10.

**Top 5 words for each topic generated**

**Sentiment scores are varying over the time range**

## Distribution of topic numbers segmented by sentiment labels

## New features added with master file



## Insights on second visualization of section 4

In the first topic, we observe more positive words such as get, want, need, and take. And we just have one negative word don't. We can infer that most product are suitable for purchase through this topics.

In the second topic, we observe very positive words such as like, food,want, and date. This implies that costco has wide range of foods and people prefer buying groceries at costco.

In the third topic, we can infer that amongst the food varieties pizza is most purchased one at costco.

In the fourth topic, we can observe that cosco is a wholesale chain with chicken and pizza foods most purchased and liked.

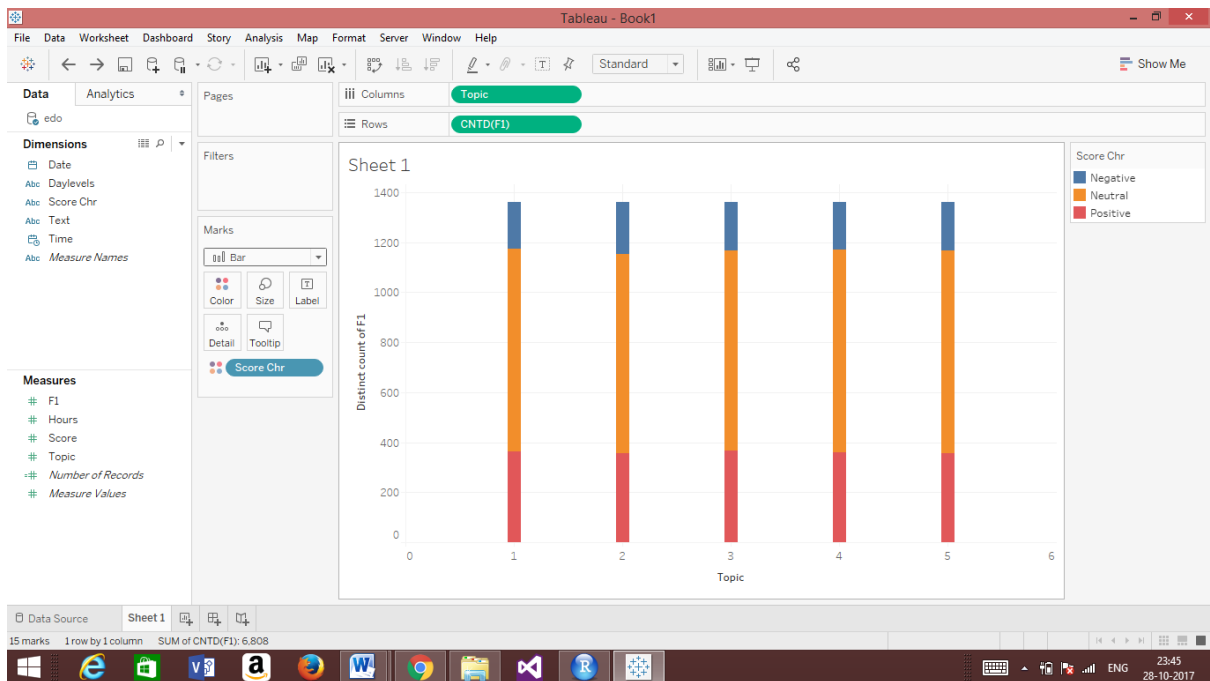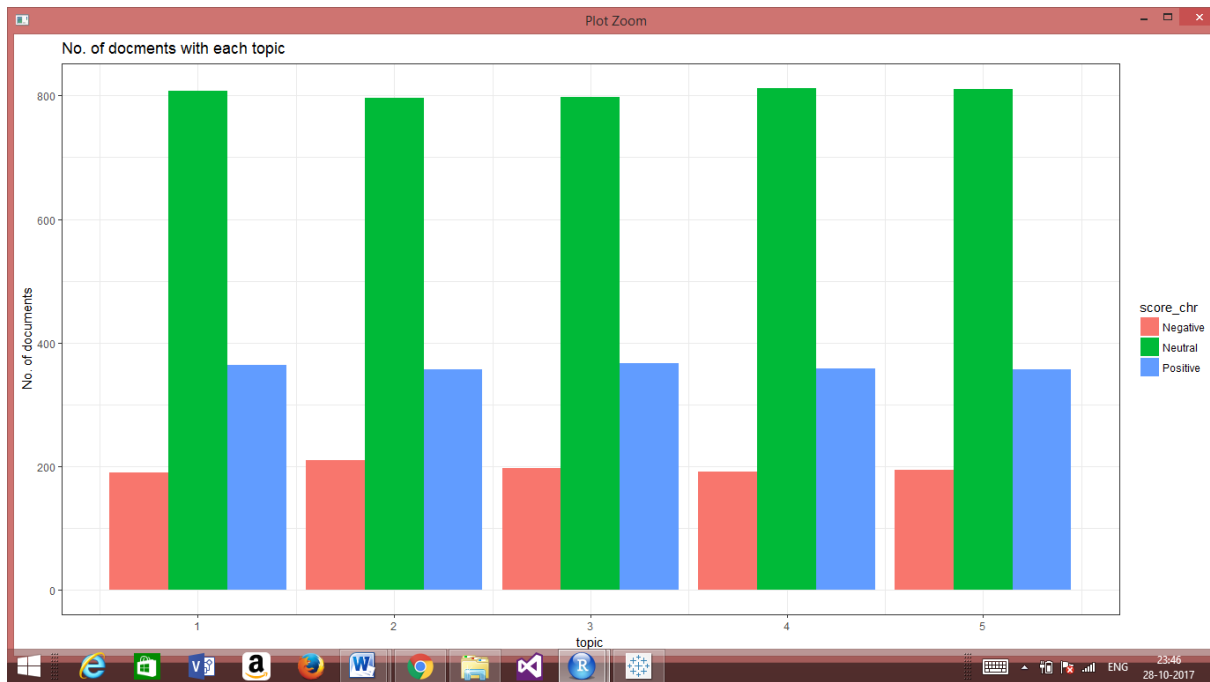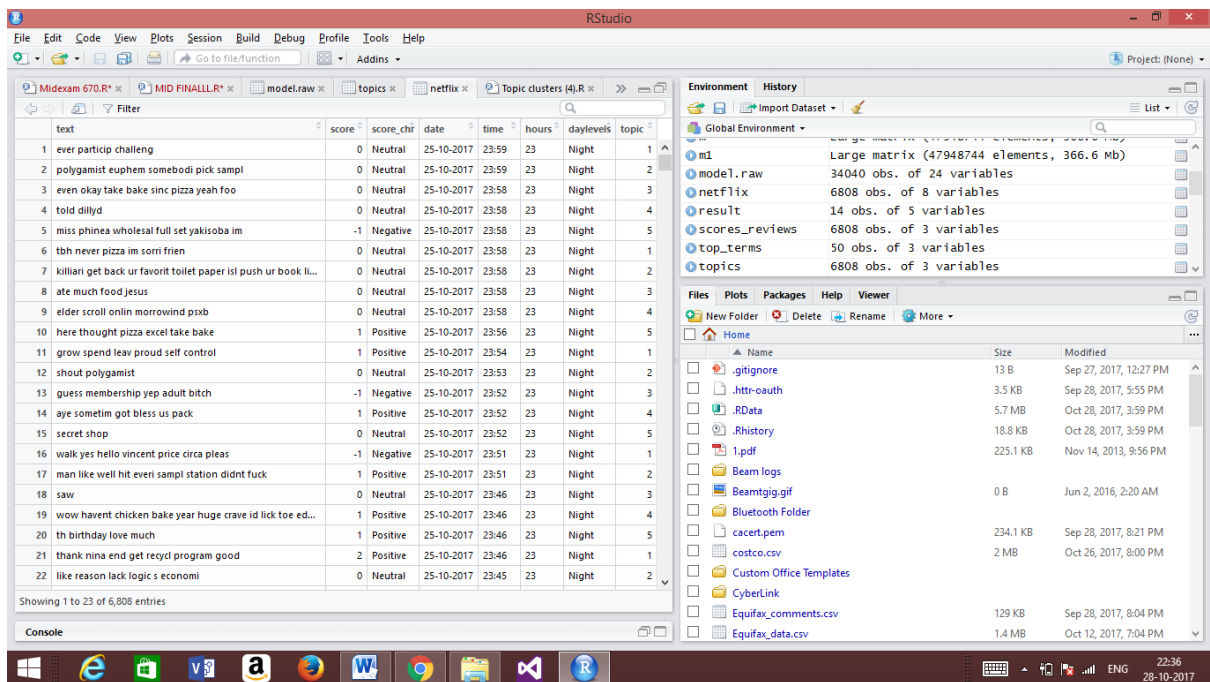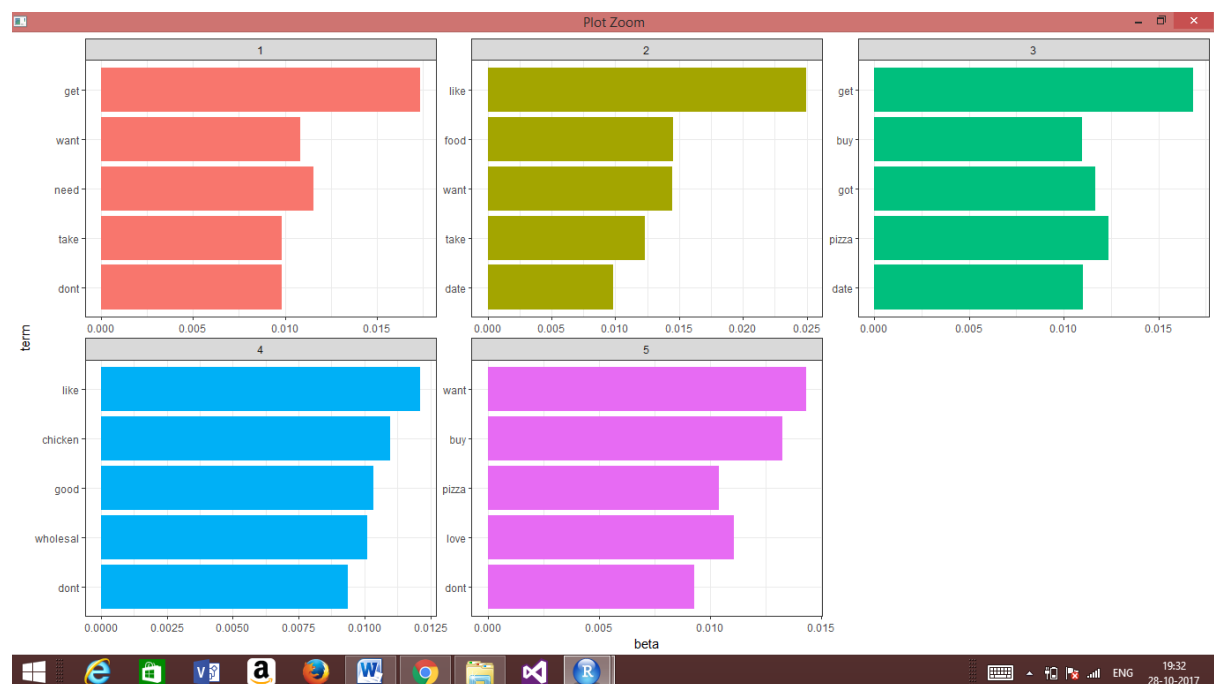In the fifth topic, we can see that pizza word is repeated in this topic. So we can infer that people prefer buying pizza at costco amongst the wide food varieties.



**Insights on third and fourth visualization of section 4**

In the third visualization, the most neutral scores were during the evenings and nights. The most positive scores were observed during evenings and nights too. The most negative scores were observed during evenings and nights. Also we can infer that most of the comments were generated during evenings and nights. However, most comments on costco are more neutral than positive in the graph.

From the fourth visualization, we calculate the number of comments per topic and also seggregate them based on sentiment label. We can observe that in topic 1, 4, and 5 the neutral comments are slightly more than the rest. In topics, 1 and 3 the positive comments are slighly higher. Whereas, the negative comments in the 2nd topic is more than the rest.

Sentiment score throughtout the day