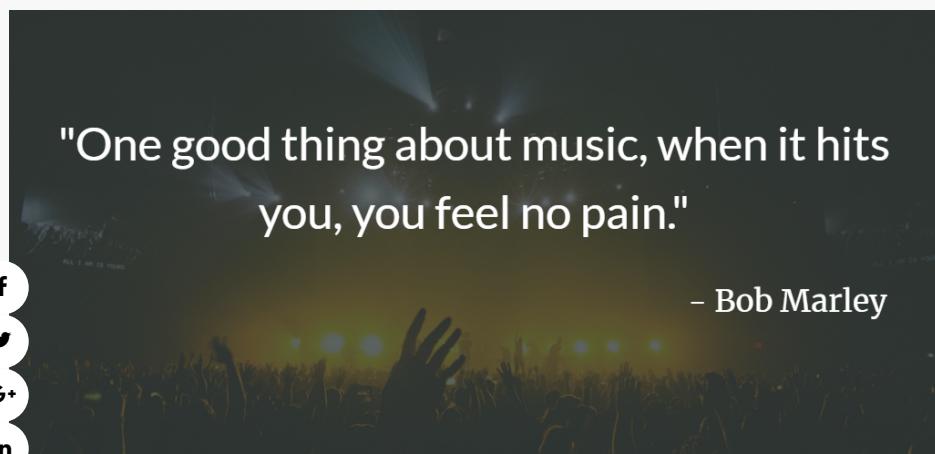


This week we're celebrating International Tourism Day -- grab 10% discount on set up fee when you crawl sites from travel domain!


[Solutions](#) [Pricing](#) [Resources](#) [Blog](#) [Company](#)
[Contact Us](#) [SUBMIT REQUIREMENT](#)

Data Visualization and Analysis of Taylor Swift's Song Lyrics

Published by Preetish on July 18, 2018



Did you know that Taylor Swift is the youngest person to single-handedly write and perform a number-one song on the Hot Country Songs chart published by Billboard magazine in the United States? Apart from that she is also the recipient of 10 Grammys, one Emmy Award, 23 Billboard Music Awards, and 12 Country Music Association Awards (enough said)! She is particularly known for infusing her personal life into her music which has received a lot of media coverage. Hence, it would be interesting to analyze her songs' content via exploratory analysis and sentiment analysis to find out various underlying themes.

Data set

Thanks to the amazing API exposed by Genius.com, we were able to extract the various data points associated with Taylor Swift's songs. We've selected only the [six albums](#) released by her and removed the duplicate tracks (acoustic, US version, pop mix, demo recording etc.). This resulted in 94 unique tracks with the following data fields:

- album name
- track title
- track number
- lyric text
- line number
- year of release of the album

Search this blog

Subscribe to our blog

[Subscribe](#)
[How to](#)
[Scrape](#)
[Data from](#)
[Twitter](#)
[Using](#)
[Python](#)


September
24, 2018

[How to](#)
[Scrape](#)
[Real](#)
[Estate](#)
[Listings](#)
[from](#)
[Trulia](#)
[using](#)
[Python](#)


September
14, 2018 **Talk to us!**

[Download the data set](#)



[Sign up for DataStock](#)

Using
Web
Scraping
to Enforce
Minimum
Advertised
Price
(MAP)



September
10, 2018

Web
scraping in
GDPR Era
– Impact
and
Opportunities



September
4, 2018

How to
Scrape
Amazon
Product
Reviews
using
Python



August
27, 2018

Goals

Our goal is to first perform exploratory analysis and then move to text mining including sentiment analysis which involves Natural Language Processing.

- Exploratory analysis

- word counts based on tracks and albums
- time series analysis of word counts
- distribution of word counts

- Text mining

- word cloud
- bigram network
- sentiment analysis (includes chord diagram)

We'll be using [R](#) and [ggplot2](#) to analyse and visualize the data. Code is also included in [this post](#), so if you download the data, you can follow along.

Exploratory analysis

in

Let's first find out the top 10 songs with the most number of words. The code snippet given below includes the packages required in this analysis and finds out the top songs in terms of length.

```

1 library(magrittr)
2 library(stringr)
3 library(dplyr)
4 library(ggplot2)
5 library(tm)
6 library(wordcloud)
7 library(syuzhet)
8 library(tidytext)
9 library(tidyr)
10 library(igraph)
11 library(ggraph)
12 library(readr)
13 library(circlize)
14 library(reshape2)
15
16 lyrics <- read.csv(file.choose())
17
18 lyrics$length <- str_count(lyrics$lyric, "\\S+")
19
20 length_df <- lyrics %>%
21   group_by(track_title) %>%
22   summarise(length = sum(length))
23
24 length_df %>%
25   arrange(desc(length)) %>%
26   slice(1:10) %>%
27   ggplot(., aes(x= reorder(track_title, -length), y=length))
28   geom_bar(stat='identity', fill="#1CCCC6") +
29   ylab("Word count") + xlab ("Track title") +

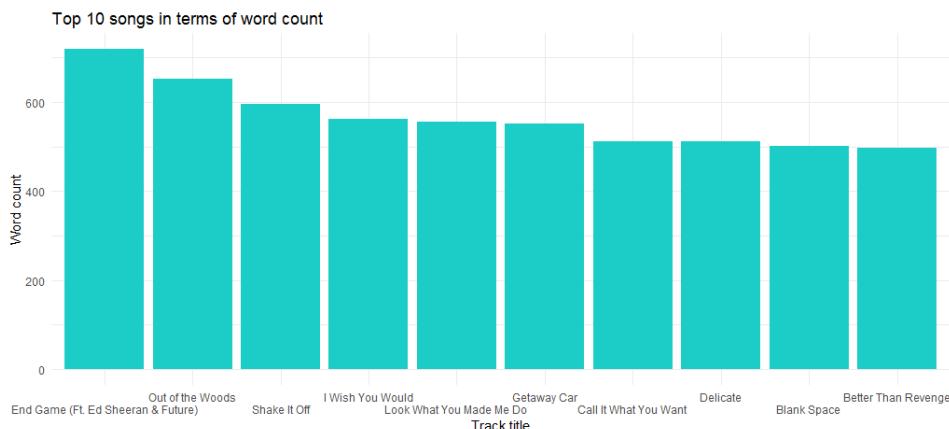
```

[Talk to us!](#)

```

30 ggtitle("Top 10 songs in terms of word count") +
31 theme_minimal() +
32 scale_x_discrete(labels = function(labels) {
33   sapply(seq_along(labels), function(i) paste0(ifelse(i %%
34 })
```

This gives us the following chart:



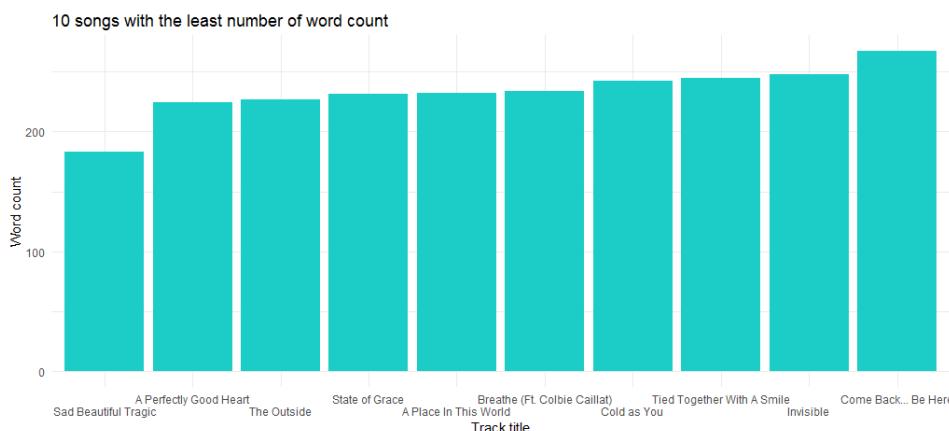
We can see that “End Game” (released in her latest album) is the song with maximum number of words and next in line is “Out of the Woods”.

Now, how about the songs with the lowest number of words? Let’s find out using the following code:

```

1 length_df %>%
2   arrange(length) %>%
3   slice(1:10) %>%
4   ggplot(., aes(x= reorder(track_title, length), y=length)) +
5     geom_bar(stat='identity', fill="#1CCCC6") +
6     ylab("Word count") + xlab ("Track title") +
7     ggtitle("10 songs with least number of word count") +
8     theme_minimal() +
9     scale_x_discrete(labels = function(labels) {
10       sapply(seq_along(labels), function(i) paste0(ifelse(i %%
```

This results in the following chart:



“Sad Beautiful Tragic” song which was released in 2012 as part of the album “Red” is the song with least number of words.

The next analysis is centered around the distribution of the number of words. Given below is the code:

```

1 ggplot(length_df, aes(x=length)) +
2   geom_histogram(bins=30, aes(fill = ..count..)) +
3   geom_vline(aes(xintercept=mean(length)),
4             color="#FFFFFF", linetype="dashed", size=1) +
```

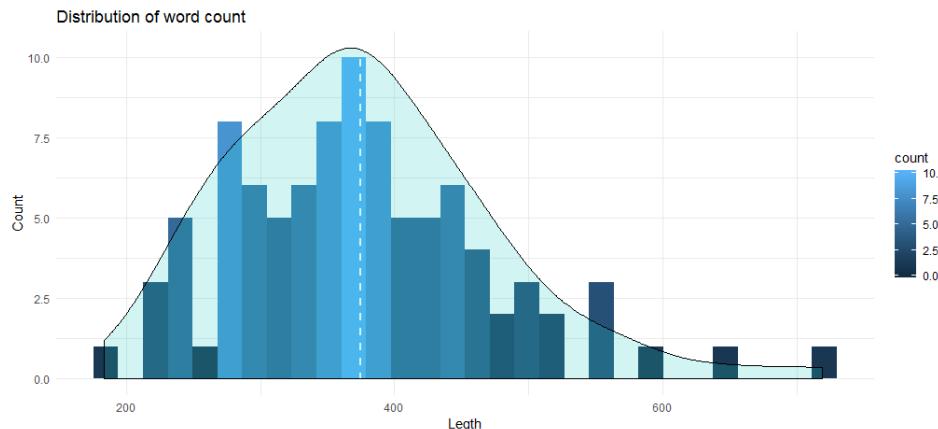
Talk to us!

```

5 |   geom_density(aes(y=25 * ..count..),alpha=.2, fill="#1CCCC6"
6 |   ylab("Count") + xlab ("Length") +
7 |   ggtitle("Distribution of word count") +
8 |   theme_minimal()

```

This code give us the following histogram along with density curve:



The average word count for the tracks stands close to 375, and chart shows that maximum number of songs fall in between 345 to 400 words.

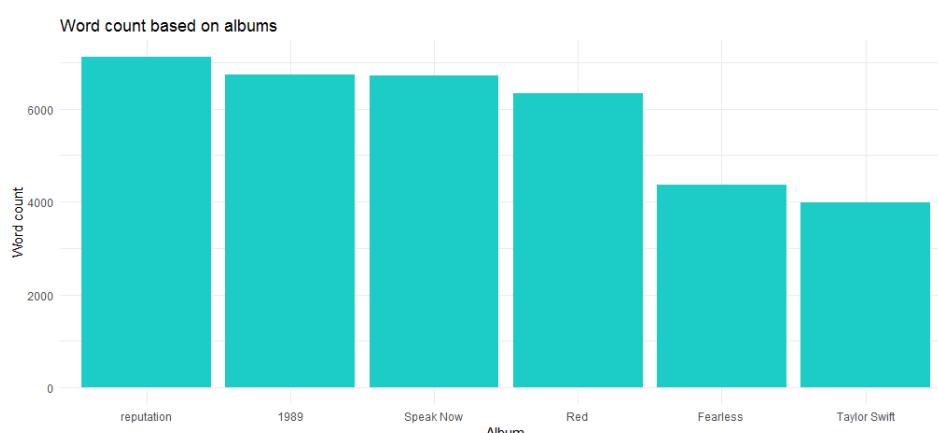
Now, we'll move to the analysis based on albums. Let's first create a data frame with word counts based on album and year of release.

```

1 | lyrics %>%
2 |   group_by(album,year) %>%
3 |   summarise(length = sum(length)) -> length_df_album
f
`ext step for us is to create a chart that will depict the length of the albums based on
` .mulative word count of the songs.
G+
1 | ggplot(length_df_album, aes(x= reorder(album, -length), y=ler
in 2 | geom_bar(stat='identity', fill="#1CCCC6") +
3 |   ylab("Word count") + xlab ("Album") +
4 |   ggtitle("Word count based on albums") +
5 |   theme_minimal()

```

The resulting chart shows that "Reputation" album which also the latest album has maximum number of words.



Now, how has the length of songs changed since the debut from 2006? The following code answers this:

```

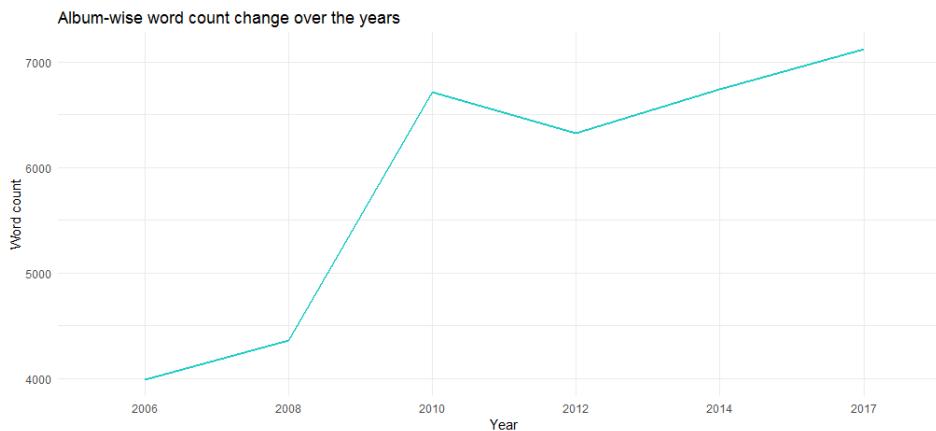
1 | length_df_album %>%
2 |   arrange(desc(year)) %>%
3 |   ggplot(., aes(x= factor(year), y=length, group = 1)) +
4 |   geom_line(colour="#1CCCC6", size=1) +
5 |   ylab("Word count") + xlab ("Year") +
6 |   ggtitle("Word count change over the years") +

```

Talk to us!

7 | theme_minimal()

The resulting chart shows that the length of the albums have increased over the years – from close to 4000 words in 2006 to more than 6700 in 2017.

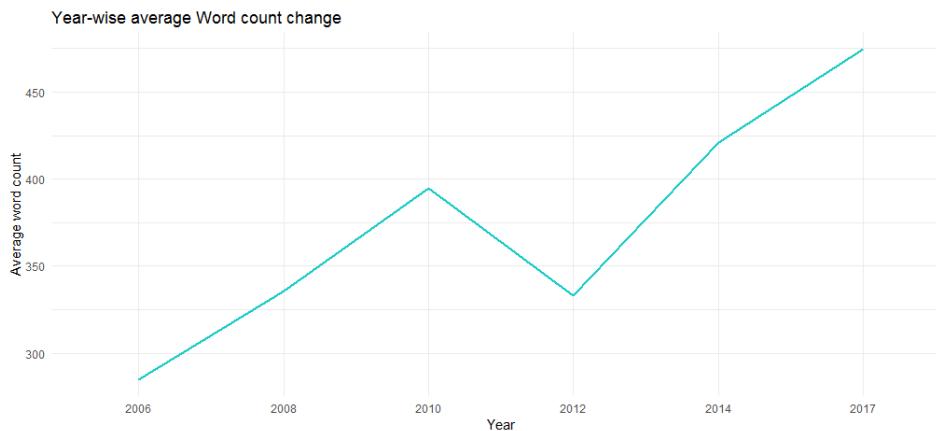


But, is that because of the number of words in individual tracks? Let's find out using the following code:

```

1 #adding year column by matching track_title
2 length_df$year <- lyrics$year[match(length_df$track_title, 1
3
4 length_df %>%
5   group_by(year) %>%
6   summarise(length = mean(length)) %>%
7     ggplot(., aes(x= factor(year), y=length, group = 1)) +
8       geom_line(colour="#1CCCC6", size=1) +
9       ylab("Average word count") + xlab ("Year") +
10      ggtile("Year-wise average Word count change") +
11      theme_minimal()
  
```

➔ The resulting chart confirms that the average word count has increased over the years (from 285 in 2006 to 475 in 2017), i.e., her songs have gradually become lengthier in terms of content.



We'll conclude exploratory analysis here and move to text mining.

Text mining

Our first activity would be to create a word cloud so that we can visualize the frequently used words in her lyrics. Execute the following code to get started:

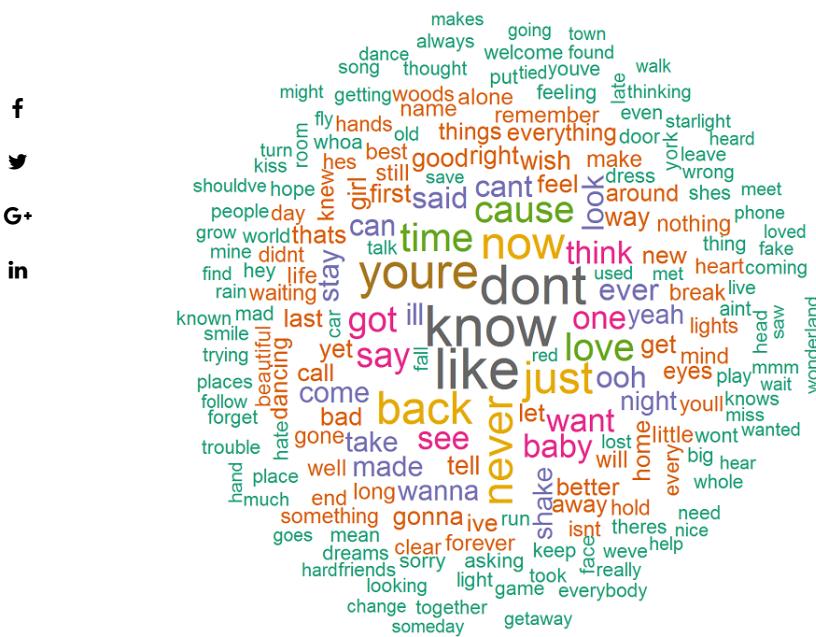
```

1 library("tm")
2 library("wordcloud")
3
4 lyrics_text <- lyrics$lyric
5 #Removing punctuations and alphanumeric content
6 lyrics_text<- gsub('[:punct:]+', '', lyrics_text)
7 lyrics_text<- gsub("([:alpha:])+", "", lyrics_text)
  
```

Talk to us!

```
8 #creating a text corpus
9 docs <- Corpus(VectorSource(lyrics_text))
10 # Converting the text to lower case
11 docs <- tm_map(docs, content_transformer(tolower))
12 # Removing english common stopwords
13 docs <- tm_map(docs, removeWords, stopwords("english"))
14 # creating term document matrix
15 tdm <- TermDocumentMatrix(docs)
16 # defining tdm as matrix
17 m <- as.matrix(tdm)
18 # getting word counts in decreasing order
19 word_freqs = sort(rowSums(m), decreasing=TRUE)
20 # creating a data frame with words and their frequencies
21 lyrics_wc_df <- data.frame(word=names(word_freqs), freq=word_freqs)
22
23 lyrics_wc_df <- lyrics_wc_df[1:300,]
24
25 # plotting wordcloud
26
27 set.seed(1234)
28 wordcloud(words = lyrics_wc_df$word, freq = lyrics_wc_df$freq,
29 min.freq = 1, scale=c(1.8,.5),
30 max.words=200, random.order=FALSE, rot.per=0.15,
31 colors=brewer.pal(8, "Dark2"))
```

The resulting word cloud shows that the most frequently used words are `know`, `like`, `don't`, `you're`, `now`, `back`. This confirms that her songs are predominantly about someone as `you're` has significant number of occurrences.



How about bigrams (word pairs that appear in conjunction)? The following code will give us the top 10 bigrams:

```
1 count_bigrams <- function(dataset) {  
2   dataset %>%  
3   unnest_tokens(bigram, lyric, token = "ngrams", n = 2) %>%  
4   separate(bigram, c("word1", "word2"), sep = " ") %>%  
5   filter(!word1 %in% stop_words$word,  
6   !word2 %in% stop_words$word) %>%  
7   count(word1, word2, sort = TRUE)  
8 }  
9  
10 lyrics_bigrams <- lyrics %>%  
11 count_bigrams()  
12  
13 head(lyrics_bigrams, 10)
```

Given below is the list of bigrams:

Talk to us!

Word 1	Word 2
ey	ey
ooh	ooh
la	la
shake	shake
stay	stay
getaway	car
ha	ha
ooh	whoa
uh	uh
ha	ah

Although we found out the word list, it doesn't divulge any insight on several relationships that exist among words. To get a visualization of the multiple relationships that can exist we will leverage network graph. Let's get started with the following:

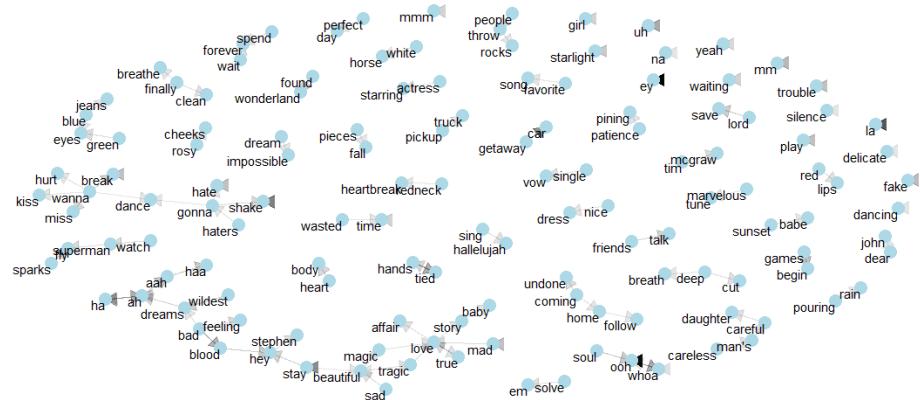
```

1 visualize_bigrams <- function(bigrams) {
2   set.seed(2016)
3   a <- grid::arrow(type = "closed", length = unit(.15, "inches"))
4
f 5 bigrams %>%
6   graph_from_data_frame() %>%
7   ggraph(layout = "fr") +
8     geom_edge_link(aes(edge_alpha = n), show.legend = FALSE, arrow = a)
9   geom_node_point(color = "lightblue", size = 5) +
10  geom_node_text(aes(label = name), vjust = 1, hjust = 1) +
11  ggtitle("Network graph of bigrams") +
12  theme_void()
13 }
14
15 lyrics_bigrams %>%
16   filter(n > 3,
17     !str_detect(word1, "\\\d"),
18     !str_detect(word2, "\\\d")) %>%
19  visualize_bigrams()

```

Check out the graph given below to see how `love` is connected with `story`, `mad`, `true`, `tragic`, `magic` and `affair`. Also, both `tragic` and `magic` are connected with `beautiful`.

Network graph of bigrams



Talk to us!

Let's now move to sentiment analysis which is a text mining technique.

Sentiment analysis

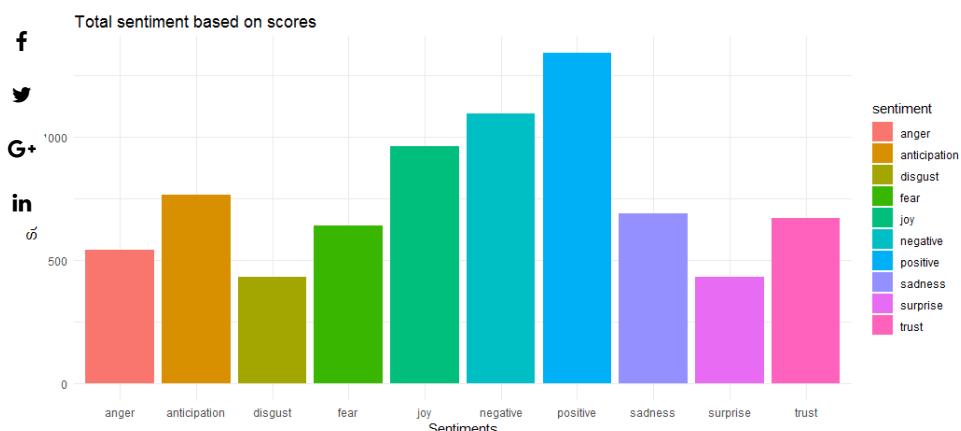
We'll first find out the overall sentiment via `nrc` method of `syuzhet` package. The following code will generate the chart of positive and negative polarity along with associated emotions.

```

1 # Getting the sentiment value for the lyrics
2 ty_sentiment <- get_nrc_sentiment(lyrics_text)
3
4 # Dataframe with cumulative value of the sentiments
5 sentimentscores<-data.frame(colSums(ty_sentiment[,]))
6
7 # Dataframe with sentiment and score as columns
8 names(sentimentscores) <- "Score"
9 sentimentscores <- cbind("sentiment"=rownames(sentimentscores),
10 rownames(sentimentscores) <- NULL
11
12 # Plot for the cumulative sentiments
13 ggplot(data=sentimentscores,aes(x=sentiment,y=Score))+ 
14 geom_bar(aes(fill=sentiment),stat = "identity")+
15 theme(legend.position="none")+
16 xlab("Sentiments") + ylab("Scores")+
17 ggtitle("Total sentiment based on scores")+
18 theme_minimal()

```

The resulting chart shows that the positive and negative sentiment scores are relatively close with 1340 and 1092 value respectively. Coming to the emotions, `joy`, `anticipation` and `trust` emerge as the top 3.



Now that we have figured out the overall sentiment scores, we should find out the top words that contribute to various emotions and positive/negative sentiment.

```

1 lyrics$lyric <- as.character(lyrics$lyric)
2
3 tidy_lyrics <- lyrics %>%
4 unnest_tokens(word,lyric)
5
6 song_wrd_count <- tidy_lyrics %>% count(track_title)
7
8 lyric_counts <- tidy_lyrics %>%
9 left_join(song_wrd_count, by = "track_title") %>%
10 rename(total_words=n)
11
12 lyric_sentiment <- tidy_lyrics %>%
13 inner_join(get_sentiments("nrc"),by="word")
14
15 lyric_sentiment %>%
16 count(word,sentiment,sort=TRUE) %>%
17 group_by(sentiment)%>%top_n(n=10) %>%
18 ungroup() %>%
19 ggplot(aes(x=reorder(word,n),y=n,fill=sentiment)) +
20 geom_col(show.legend = FALSE) +

```

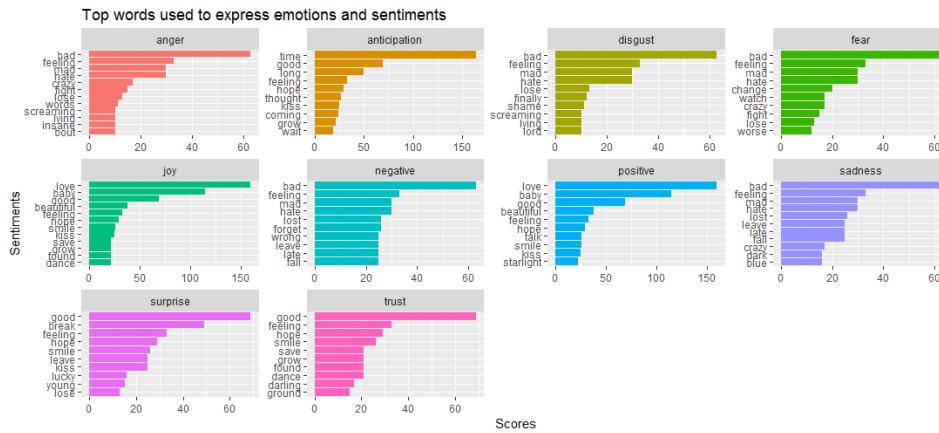
Talk to us!

```

21 | facet_wrap(~sentiment,scales="free") +
22 | xlab("Sentiments") + ylab("Scores")+
23 | ggtitle("Top words used to express emotions and sentiments",
24 | coord_flip()

```

The visualization given shows that while the word `bad` is predominant in emotions such as `anger`, `disgust`, `sadness` and `fear`, `Surprise` and `trust` are driven by the word `good`.



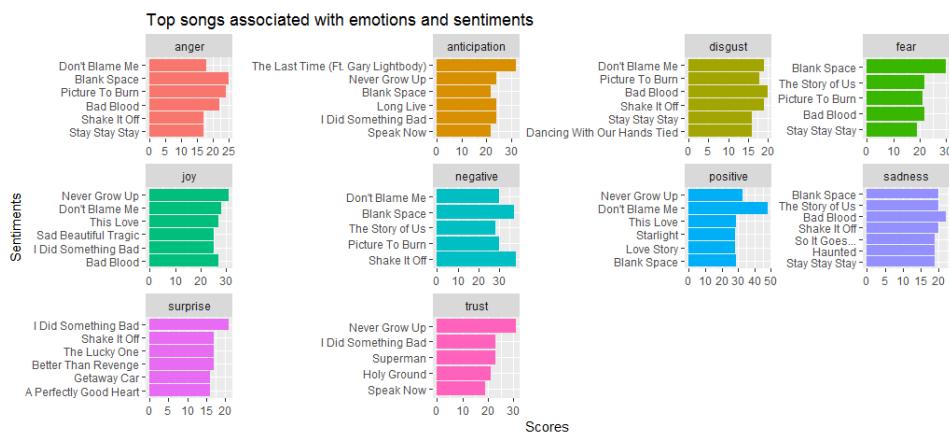
This brings to the following question – which songs are closely associated with different emotions? Let's find out via the code given below:

```

1 | lyric_sentiment %>%
2 | count(track_title,sentiment,sort=TRUE) %>%
3 | group_by(sentiment) %>%
4 | top_n(n=5) %>%
5 | ggplot(aes(x=reorder(track_title,n),y=n,fill=sentiment)) +
6 | geom_bar(stat="identity",show.legend = FALSE) +
7 | facet_wrap(~sentiment,scales="free") +
8 | xlab("Sentiments") + ylab("Scores")+
9 | ggtitle("Top songs associated with emotions and sentiments");
in 10 | coord_flip()

```

We see that the song `Black Space` has a lot of anger and fear in comparison to other songs. `Don't Blame Me` has considerable score for both `positive` and `negative` sentiment. We also see that `Shake it off` scores high for `negative` sentiment; mostly because of high frequency words such as `hate` and `fake`.



Let's now move to another sentiment analysis method, `bing` to create a comparative word cloud of positive and negative sentiment.

```

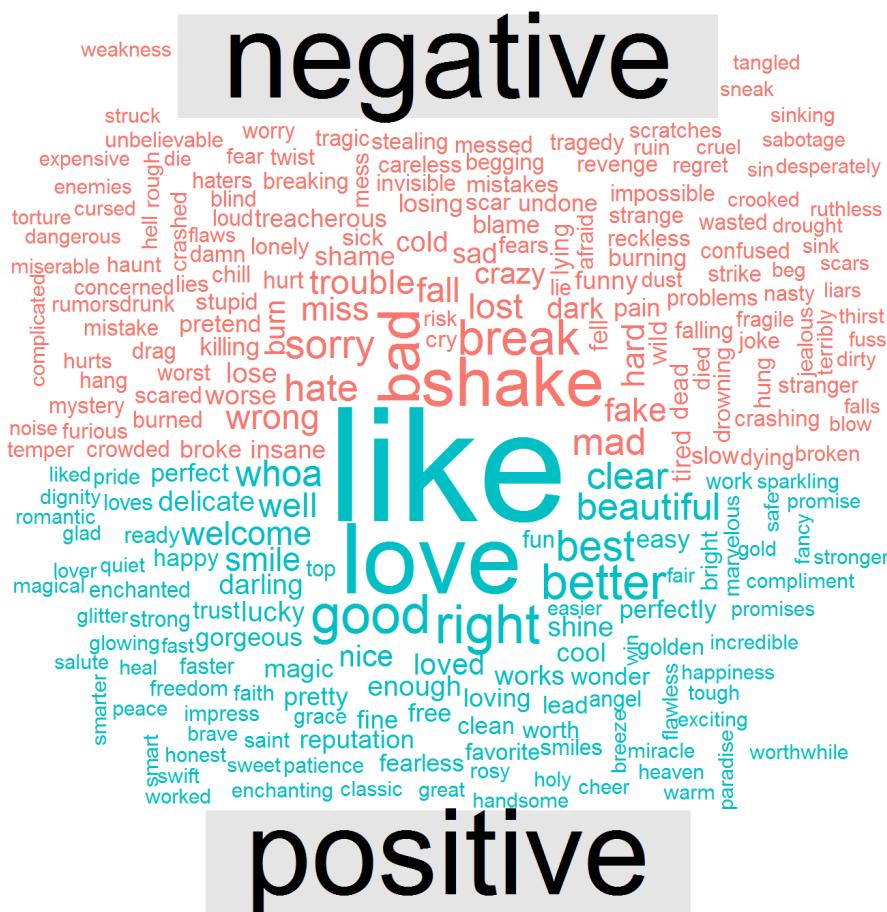
1 | bng <- get_sentiments("bing")
2 |
3 | set.seed(1234)
4 |
5 | tidy_lyrics %>%

```

Talk to us!

```
6 | inner_join(get_sentiments("bing")) %>%
7 | count(word, sentiment, sort = TRUE) %>%
8 | acast(word ~ sentiment, value.var = "n", fill = 0) %>%
9 | comparison.cloud(colors = c("#F8766D", "#00BFC4"),
10 | max.words = 250)
```

Following visualization shows that her songs have positive words such as `like`, `love`, `good`, `right` and negative words such as `bad`, `break`, `shake`, `mad`, `wrong`.



This brings to the final question – how has her sentiment and emotions changed over the years? For this particular answer we will create a visualization called [chord diagram](#).

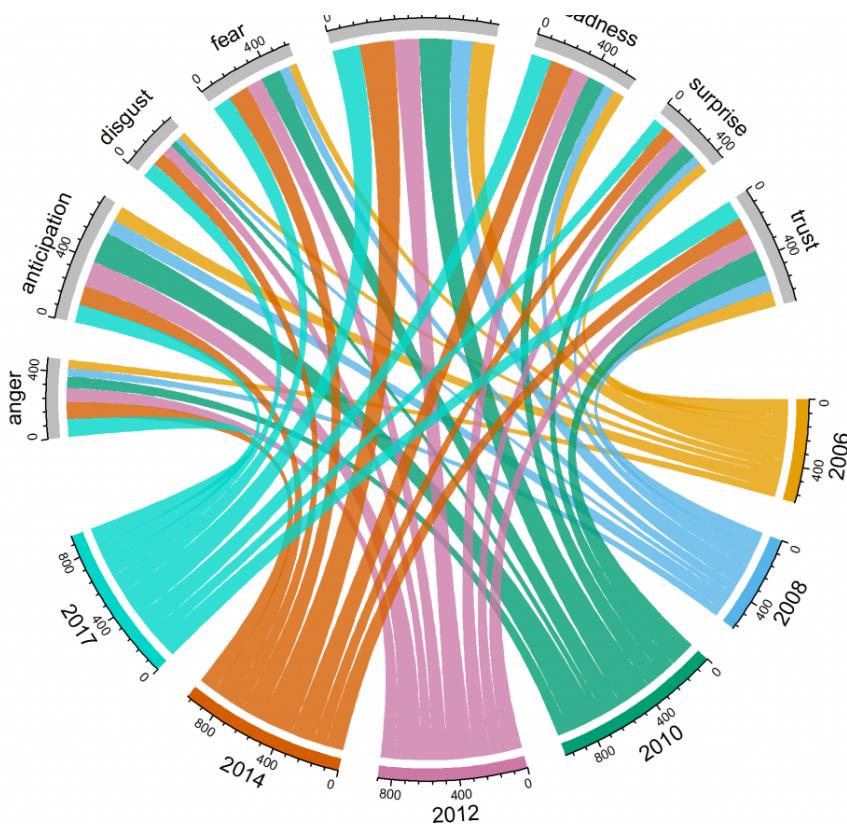
Here is the code:

```
1 grid.col = c("2006" = "#E69F00", "2008" = "#56B4E9", "2010"
2
3 year_emotion <- lyric_sentiment %>%
4 filter(!sentiment %in% c("positive", "negative")) %>%
5 count(sentiment, year) %>%
6 group_by(year, sentiment) %>%
7 summarise(sentiment_sum = sum(n)) %>%
8 ungroup()
9
10 circos.clear()
11
12 #Setting the gap size
13 circos.par(gap.after = c(rep(6, length(unique(year_emotion[1])) - 1),
14 rep(6, length(unique(year_emotion[[2]])) - 1), 15))
15
16 chordDiagram(year_emotion, grid.col = grid.col, transparency =
17 title("Relationship between emotion and song's year of release"))
```

It gives us the following visualization:

Relationship between emotion and song's year of release

Talk to us!



From the chart, we can see that **joy** has maximum share for the years 2010 and 2014. Overall, **sadness**, **surprise**, **disgust** and **anger** are the emotions with least score; however, in comparison to other years 2017 has maximum contribution for **disgust**. Coming to **anticipation**, 2010 and 2012 have higher contribution in comparison to other years.

in ver to you

In this study we performed exploratory analysis and text mining, which includes NLP for sentiment analysis. If you'd like to perform additional analyses (e.g., lexical density of lyrics and topic modeling) or simply replicate the results for learning, download the data set for free from DataStock. Simply follow the link given below and select "free" category on DataStock.



Related posts



September 24, 2018



September 14, 2018



September 10, 2018

Using Web Scraping to Enforce Minimum

Talk to us!

How to Scrape Data from Twitter Using Python

[Read more](#)

How to Scrape Real Estate Listings from Trulia using Python

[Read more](#)

Advertised Price (MAP)

[Read more](#)

Leave a Reply

Your email address will not be published. Required fields are marked *

Comment

Name *

Email *

Website

[Post Comment](#)

f



Solutions

Custom Crawling
DataStock
JobsPikr
Hosted Indexing
Live Crawls
WordPress Blog
Extraction
Price Scraping Service

Industries

Ecommerce
Travel
Jobs
Research & Analytics
Real Estate
Automobile
Finance

Resources

Blog
Case Studies
Customer Spotlight
FAQs
PromptInsights
Who we help
Web Data Acquisition
Navigator

Company

About Us
Media Center
Current Openings
Become a Partner
Wall of Fame
Contact Us



[Terms](#) [Privacy Policy](#)

© 2018 PromptCloud. All Rights Reserved.

[Talk to us!](#)