

# **Social Media Analysis of Samsung vs. Apple**

## **Final Report**

**Making Products Count:**

**Data Science for Product Managers**

**Team Member:**

Shuhan Zhang

Jingyi Ran

Ziyu Li

Xiang Li

Yiyun Hu

March 4th, 2022

# Contents

<b>Social Media Analysis of Samsung vs. Apple .....</b>	<b>1</b>
<b>1. Background .....</b>	<b>3</b>
<b>2. Data Preparation.....</b>	<b>4</b>
<b>3. Exploratory Data Analysis .....</b>	<b>5</b>
<b>4. Sentiment Analysis.....</b>	<b>7</b>
4.1 Preprocessing.....	7
4.2 Sentiment Score.....	8
4.2.1 Vader .....	8
4.2.2 TextBlob.....	9
4.3 Subjective score.....	9
4.4 Twitter, non-twitter .....	10
4.4.1 Sentiment score comparison .....	10
4.4.2 Subjective score comparison .....	10
4.4.3 Sentiment score of Samsung in each non-twitter platform .....	11
4.5 Sentiment before and after the launch.....	11
4.5.1 For Samsung.....	12
4.5.2 For Apple .....	13
<b>5. Top Attributes and Customer Response.....</b>	<b>13</b>
5.1 Top Attributes mentioned for each product.....	14
5.2 Top Attributes that are liked the most.....	14
5.3 Top Attributes that are disliked the most.....	15
<b>6. Product Recommendations .....</b>	<b>16</b>
6.1 AARRR Model .....	16
6.2 Retention.....	17
6.3 Revenue.....	17
6.4 Referral .....	17
<b>7. Future Improvement for the Project.....</b>	<b>18</b>
7.1 Models for data from non-twitter sources. ....	18
7.2 Customer segmentation. ....	18
<b>8. Contribution .....</b>	<b>18</b>
<b>9. Reference .....</b>	<b>19</b>

# 1. Background

As a product management team at Samsung, our goal is always to explore potential customer markets and promote our brand to a new level.

Here is a consumer survey which shows the relationship of Samsung and Apple that they are the two dominant companies in the worldwide electronics market. The chart below shows that Samsung has better performance and high market share in some European countries. In general, Samsung has over 30% market share in European countries like Germany, France and around 30% market share in the UK, which indicates excellent sales and strategies there. Whereas in the US market, Apple dominates Samsung as it takes up nearly 40% market share and acts as the main competitor against Samsung.

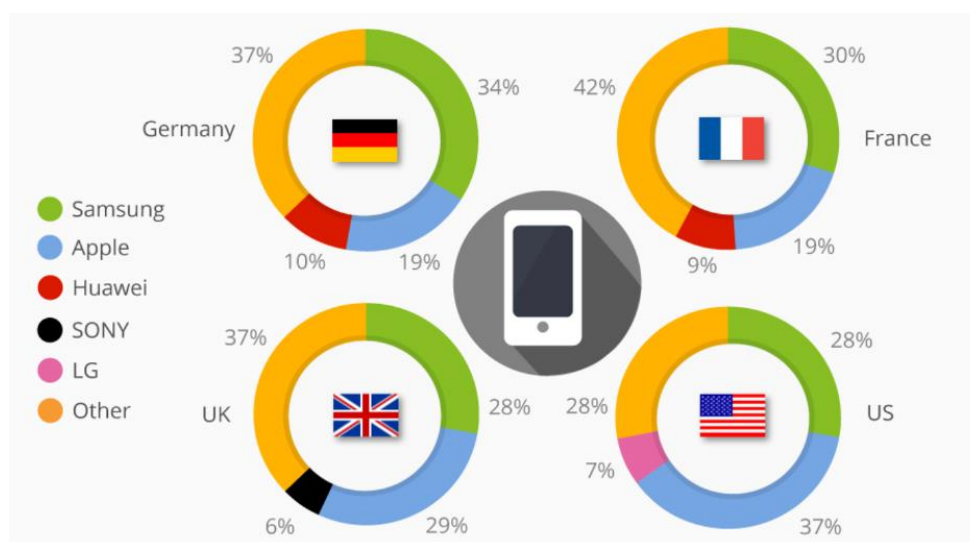


Figure 1: The market share comparison of Samsung and Apple

Source: Statista Global Consumer Survey 2017

Thus, it is obvious that there is a high profit market in the U.S. where Samsung can do better. In order to maintain the dominant position globally and take up more market share in the U.S., our team developed sentiment analysis based on social media platforms such as Twitter, Instagram, Blog, Tumblr etc. to evaluate the feedback of our newly released product Samsung Galaxy S8 and S8+, as well as comparing with iPhone8 and iPhone X which are released later in 2017 from Apple.

This final report carefully describes the end-to-end process of the social media analysis we conducted. From data preparation to product development recommendations, we mainly employed Natural Language Processing knowledge and combined technical analysis with

our product development sense. We truly believe this project can assist Samsung with refinement of its future smartphone products.

## **2. Data Preparation**

Our data engineering team provided a total of 12 datasets, among which 2 datasets are about Twitter data, while the other 10 datasets are about other non-Twitter social media data from platforms like Blog, Tumblr, Instagram etc. These data are from the beginning of March 2017 to the end of October 2017. Note that the data engineering team were not able to fetch the non-Twitter data from June 6 to September 9 and Twitter data from June 30 to September 9. To prepare our data for further analysis, the following steps are employed.

Firstly, we choose out the original post type posts, as we don't want to analyze retweets, replies or comments. We believe that original posts contain the strongest sentiment and one's original thoughts.

Secondly, we check missing values. We drop multiple columns which are full of missing values and would have no effects on EDA or sentiment analysis. In addition, we also drop rows which have missing values for the 'Post Type' column, as we want to eliminate the 10 extra sentences at the end of each excel file which describes the source of data.

Thirdly, we get rid of posts from authors with a large amount of followers. For the Twitter dataset, we believe 8000 followers will be a good threshold to use, as we don't want to delete too many valuable posts. By checking the number, we can see around 12% of the total posts are from authors with more than 8000 followers. However, for non-Twitter data, we use 1000 followers as the threshold, as the majority of non-Twitter data are from authors with no followers. In this way, we reduced the chance of analyzing posts which might be news, commercial ads, or promotions.

Fourthly, we merged these 12 datasets to 3 final datasets which are cleaned non-Twitter dataset, cleaned Twitter dataset and an overall dataset which combines Twitter and non-Twitter data together. We also dropped duplicates of posts (value of the 'Sound Bite Text' column) for these three final datasets and exported them to three new excel files for future use. There are 26334 rows for cleaned Twitter dataset, 360263 rows for cleaned non-Twitter dataset and 386597 rows for cleaned overall dataset.

### 3. Exploratory Data Analysis

The purpose of exploratory data analysis for this project is to find out potential patterns within the data, non-sense data errors, and to further prepare for the following text processing including sentiment analysis and common word extraction. The input data of EDA are the cleaned overall dataset being processed during data preparation. While we have removed the irrelevant columns with high proportion of missing values, our main concern is on the “Sound Bite Text” column. Therefore, we checked and removed the row entries with missing “Sound Bite Text”, which further reduced the number of observations by 2,118.

Although the raw data was initially splitted to twitter and non-twitter in terms of source type, we believe the posts from different social media platforms would potentially possess different characteristics regarding text length, the level of formality, and the underlying demographic groups of users. Therefore, we investigated the number of posts from each source type and the barplot is presented below for better visualization.

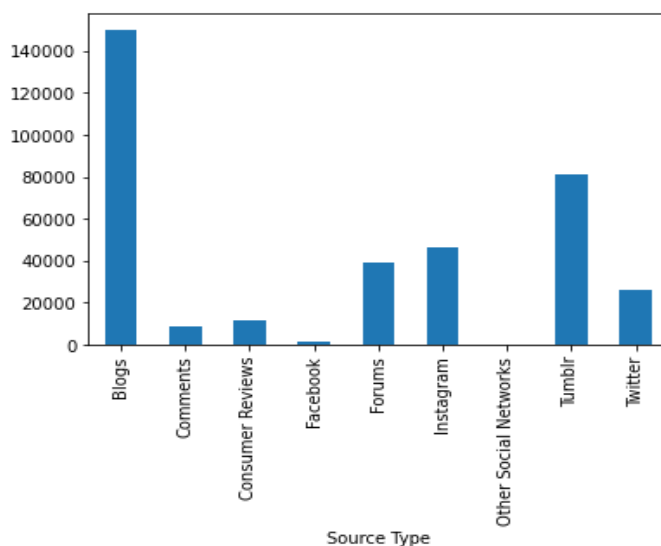


Figure 2: Number of reviews in each source

As easily observed from the bar graph, the majority of posts within our dataset originated from Blogs and Trumblr, with over 140,000 posts from Blogs and more than 80,000 posts from Trumblr. In this regard, we decided to conduct individual analysis on texts from different platforms, specifically Blogs and Trumblr, as they compose the majority of observations.

The dataset also has a numerical attribute “richness”, which, under the context of social media, measures the communication medium’s ability to reproduce the information sent

out. (Quoted from Wikipedia). Although the dataset doesn't provide detailed information about how richness is evaluated and assigned to each post, exploring the difference between richness of source type would potentially provide additional insights.

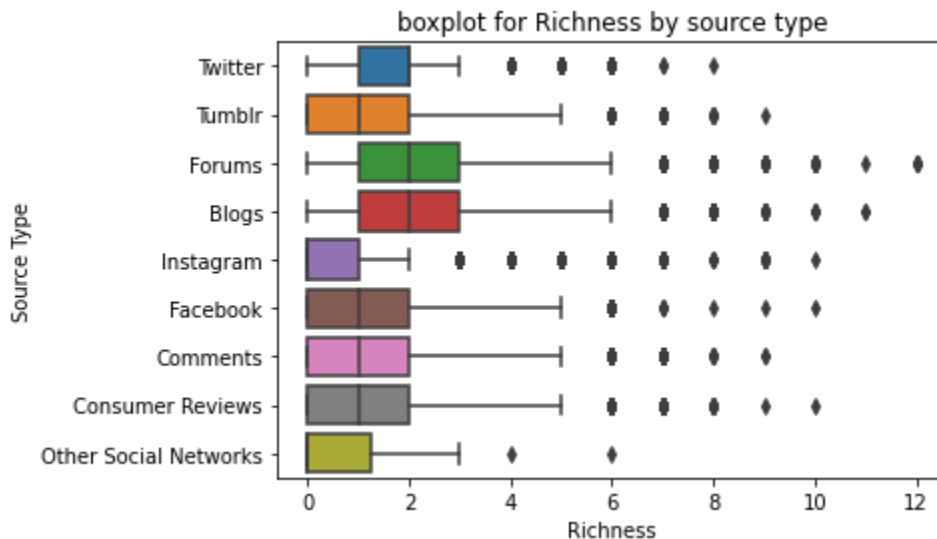


Figure 3: boxplot for Richness by source type

The boxplot shows that Forums and Blogs have higher scores on richness compared to that of the other source types, which inspired us to investigate the post length (which is the number of words per post).

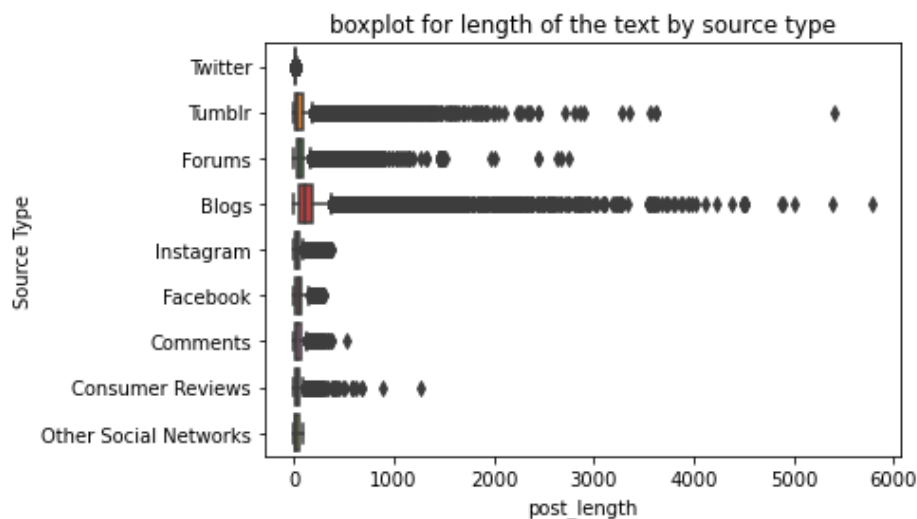


Figure 4: boxplot for length of the next by source type

It shows that Blogs, Forums, and Tumblr have plenty of extremely large posts, which were initially flagged as anomaly as we doubted whether it exceeded the maximum number of words restricted by the social media and therefore should be marked as data error. We then

found out that the character limit on Tumblr Page is 500,000, and the length limit on Blogs vary across different platforms. Therefore, we could further confirm that posts from different source types are potentially distinctive in nature, and we decide to take this factor into consideration for the following sentiment and text analysis.

## **4. Sentiment Analysis**

### **4.1 Preprocessing**

As the first step of natural language processing, we would remove the stop words that are either meaningless in terms of information gain or irrelevant to the semantics of the text and thereby improve the performance of the sentiment evaluation model. Specifically, we used both spacy and nltk Package to parse the posts into substrings of words and punctuation and remove the default stop words defined in the spacy package. Taking consideration of the distinctive nature of social media posts, we further eliminated the url link, tag sign, “.com” pattern and email within the post to remove apparent noises that would negatively affect the following sentiment analysis. To be more specific, we defined regular expressions patterns using re package and applied them to every observation.

As the previous steps had already converted each post to a list of meaningful words and punctuation, we then splitted the data into four categories based on the occurrence of specific keywords in the list of substrings per post. The four categories consist of: 1) Apple-related, 2) Samsung-related, 3) both, and 4) Neither.

The detailed decision rules are as follows:

- If the post contains either one of ['iphone', 'iphone x', 'iphone 8'], we classified it as “Apple-related”;
- If it contains either one of ['galaxy', 's8', 'samsung'], we classified it as “Samsung-related”;
- If it contains words from both list, we remove it from previous two categories which is originally assigned and reclassified it as “Both”;
- If it doesn’t contain any of the above words, we classified it as “None”.

For the following sentiment analysis, we only used the post labeled as “apple-related” and “samsung-related” and gently ignored those assigned as “both” since it’s difficult to determine the sentiment score evaluated to the post that mentioned both “Apple” and “Samsung” towards a specific one of them.

## 4.2 Sentiment Score

We conducted sentiment score calculation using both Vader and TextBlob.

### 4.2.1 Vader

Vader uses a list of lexical features which are labeled as positive or negative according to their semantic orientation to calculate the text sentiment. Vader sentiment returns the probability of a given input sentence to be positive, negative or and neutral (Analytics Vidhya, 2021).

The detailed decision rules for a review to be positive, negative or neutral are as follows:

- If the sentiment score  $> 0$ , we classified it as “positive”
- If the sentiment score  $= 0$ , we classified it as “neutral”
- If the sentiment score  $< 0$ , we classified it as “negative”

The summary statistics result is shown by the pie chart below:

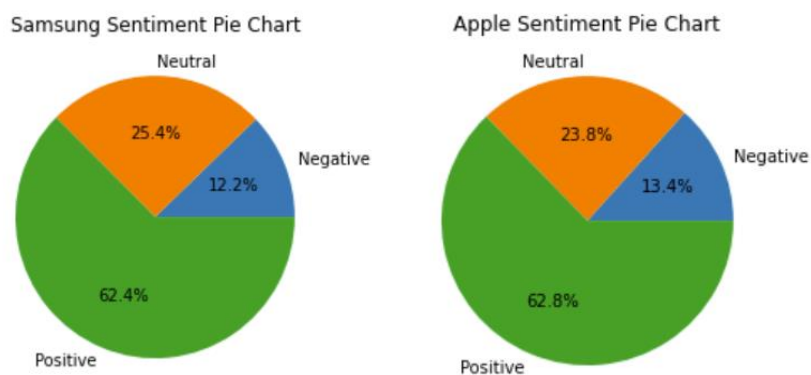


Figure 5: Vader sentiment scores comparison of Samsung and Apple

From the result we can see most reviews are positive for both Samsung and Apple, but the percentage of negative reviews of Apple is more than that of Samsung.



### 4.2.2 TextBlob

TextBlob can return the popularity of a sentence. Popularity is a float that lies in the range of  $[0,1]$ . Textblob will disregard words which they don't have any acquaintance with. In addition, TextBlob is a Lexicon-based sentiment analyzer (Analytics Vidhya, 2021).

The detailed decision rules in TextBlob for a review is the same as that in the Vader part.

The summary statistics result is shown by the pie chart below:

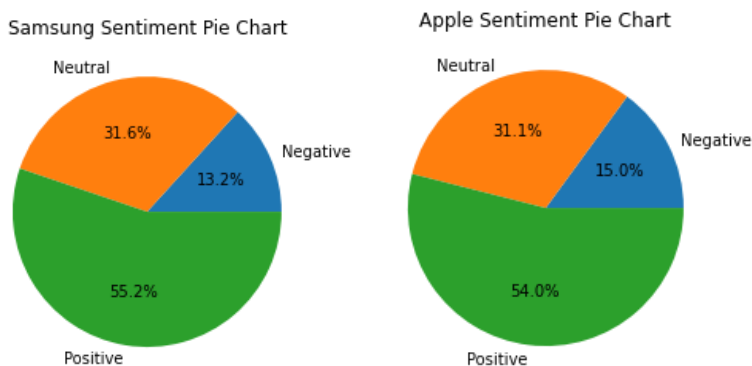


Figure 6: TextBlob sentiment scores comparison of Samsung and Apple

Same conclusion could be drawn: the percentage of negative reviews of Apple is more than that of Samsung.

## 4.3 Subjective score

The output of TextBlob is polarity and subjectivity. Subjective score shows the amount of personal opinion. If the subjective is high (close to 1), it means the text information contains more personal opinion than factual information (Analytics Vidhya, 2021).

We calculated the subjective score of each review and more information could be found in comparison of Twitter and non-Twitter part

## 4.4 Twitter, non-twitter

### 4.4.1 Sentiment score comparison

From the histograms below, we can see reviews from non-twitter platforms are more positive than those from twitter platforms. Reviews from twitter platforms tend to be neutral.

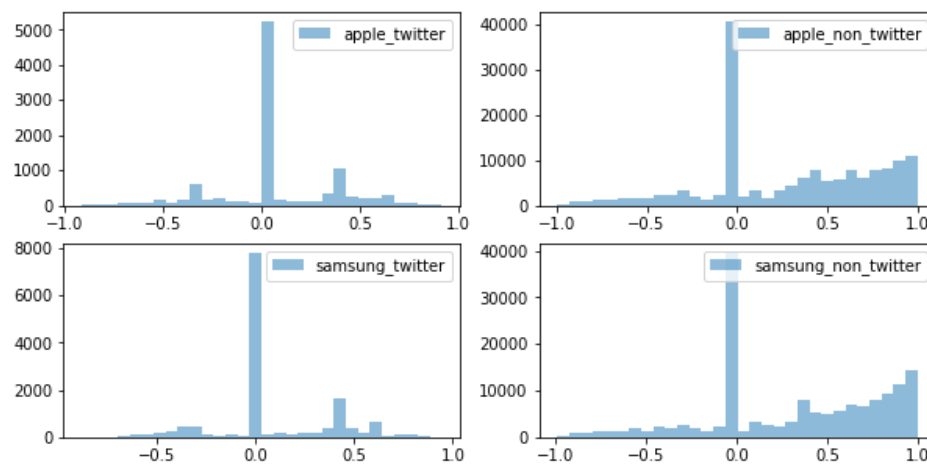


Figure 7: Sentiment comparison between twitter and non-twitter platforms

### 4.4.2 Subjective score comparison

From the histograms, it seems that reviews from non-twitter platforms are more subjective than those from twitter platforms. Reviews from twitter platforms tend to be more factual.

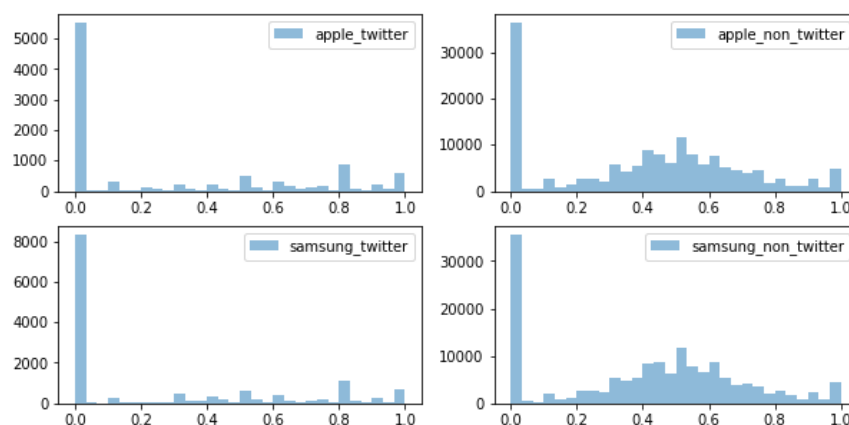


Figure 8: Subjectivity comparison between twitter and non-twitter platforms

However, the principle of subjectivity score is analyzing the amount and level of strong words like “very”, “too”, “so” and etc, which is strongly related to the post

length, so this result is also an indication of post length and richness. We don't think it's reliable enough to draw the conclusion that reviews from non-twitter platforms are more subjective.

#### 4.4.3 Sentiment score of Samsung in each non-twitter platform

Since there are different platforms in non-twitter data, we analyzed the sentiment score of reviews of Samsung in each platform to see customers' feelings about Samsung products.

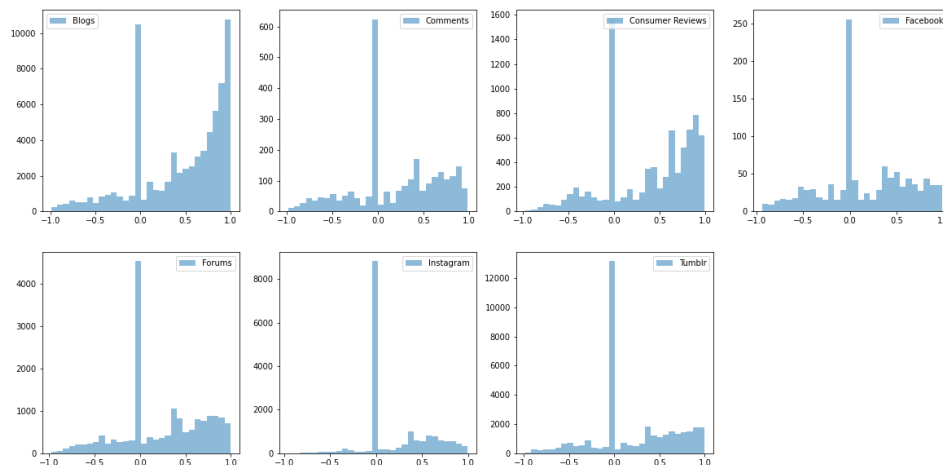


Figure 9: Sentiment score of reviews of Samsung in each non-twitter platforms

From the result, we can see reviews from blogs are relatively more positive. Combine the result with the conclusion of EDA, which shows reviews from blogs have higher richness and longer posts, we can preliminarily conclude that customers hold a positive attitude towards Samsung products.

## 4.5 Sentiment before and after the launch

The comparison between sentiment before and after the launch of the product (specifically phone in this case) could provide meaningful insights about customer's expectation before the debut and how customers react to the product after the press. Therefore, we calculated the daily average sentiment score for "Samsung" and "Apple" separately based on the post date and conducted a time series analysis. Note that the launch date for Samsung Galaxy S8 and S8+ is March 29, 2017. And the launch date for iPhone 8 and iPhone X is September 9, 2017.

### 4.5.1 For Samsung:

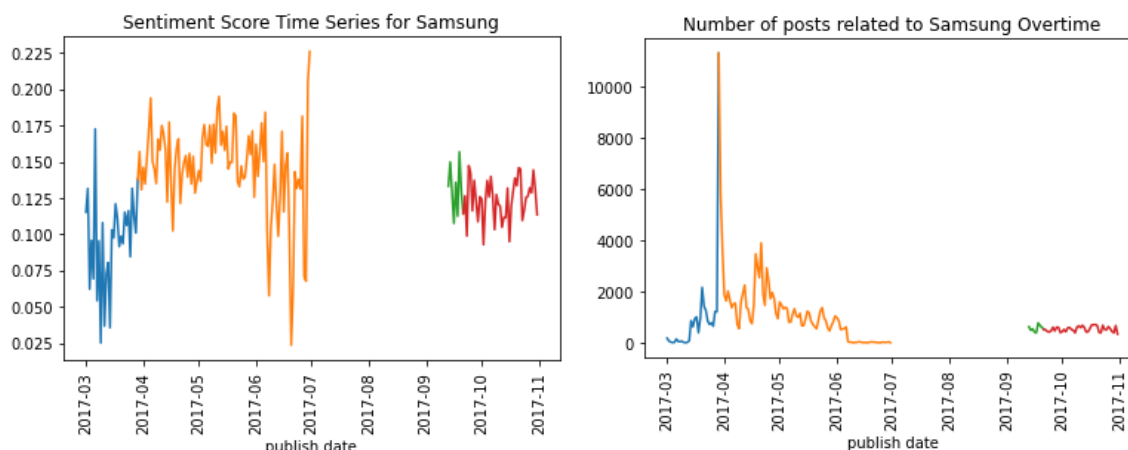


Figure 10: Time series of Samsung

The time series plot on the left above is the daily average sentiment score for Samsung over time. The blue portion represents the post before the debut of Galaxy S8. We could observe an obvious increase in the mean sentiment score, indicating that people's expectations towards Galaxy S8 are generally positive. The orange part is the post published after the launch, an increasing trend could also be found out from the graph, which demonstrates that users on social media express positive feedback on the launch of Samsung Galaxy S8. Note that the data on some dates experiences dramatic fluctuations, which is due to the fact that we don't have non-twitter data for these dates.

The plot on the left-hand side shows the corresponding number of posts related to Samsung over the same time horizon. We could see that the number of posts reached the peak and gradually decreased as time passed by. Therefore, the fluctuations on the sentiment score could be due to the imbalance data (meaning the number of posts per day fluctuate a lot across the time).

### 4.5.2 For Apple:

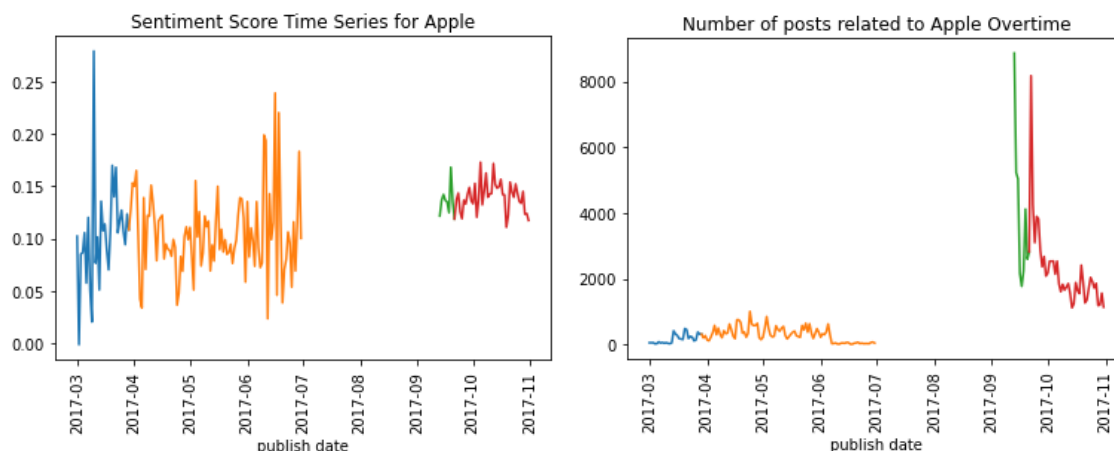


Figure 11: Time series of Apple

Similarly, the time series plot for sentiment scores of Apple-related posts are shown above. We could observe that the mean sentiment score for Apple experienced a decreasing trend right after Samsung launched Galaxy S8. For the Samsung product management team, that could be regarded as a positive sign as the public's attitudes towards the main competitor turned negative after our brand's product was introduced to the market.

The sentiment score for Apple before they launched doesn't show a clear trend, and the sentiment after the launch of iPhone X indicates a positive reaction towards iPhone X and Apple in general. However, the trend is not as obvious as that of Samsung's increase after the Galaxy S8 was introduced.

## 5. Top Attributes and Customer Response

In order to extract the most frequently discussed attributes for Samsung Galaxy S8/S8+ and iPhone 8/iPhone X, the first thing we did is to separate out the Samsung data and Apple data based on the 'Apple' and 'Samsung' column we created in sentiment analysis.

Then, we separated out negative and positive posts based on the sentiment score calculated in sentiment analysis for both Samsung and Apple data, so that we can find out the most important attributes people like or dislike the most.

After we had these four datasets which are negative Samsung posts, positive Samsung posts, negative Apple posts and positive Apple posts, we started to do text analysis. The first thing is to remove all the stop words and irrelevant information such as urls, hashtags etc, in each post. Then, we concatenated all the cleaned sound bite texts to a huge string for each dataset, in order to draw one word cloud graph per dataset to show what kind of words are mentioned the most for each product and for each sentiment polarity. To validate the word frequency, we also utilized dict subclass Counter from collections class to generate and display the most frequent adjective and noun words for each dataset. The final results are super informative.

## 5.1 Top Attributes mentioned for each product

We believe that the attributes which appear in both negative and positive sentiment are the top attributes for each product. The result shows that for both Samsung and Apple, color, screen, camera, wireless charging and facial recognition are the features mentioned the most. It is pretty reasonable, as it is the very first time that facial recognition and wireless charging were introduced to both Samsung S8 and iPhone X.

We also noticed that for Samsung, people talked about some other features like virtual assistant and storage. While for Apple, people talked more about the design and look. In order to distinguish how social media posts feel about each highly mentioned feature for each product, let's look at these features closely.

## 5.2 Top Attributes that are liked the most



Figure12: Word Cloud for Positive Samsung Posts Word Cloud for Positive Apple Posts

For positive Samsung posts, we find that people are more satisfied with the new features which were added to Galaxy S8 and S8+. For example, people talked a lot about the virtual

assistant called Bixby which is basically like Siri on the iPhone, the infinity display and the iris scanner. While for Apple, people are more satisfied with updates on hardwares such as the OLED display, portrait lighting for photo taking and A11 bionic chip.

One of the most important findings is that our customers are truly satisfied with a new unlocking feature added to Galaxy S8 and S8+ called Iris Scanning. This is a brand new technology first introduced to Samsung S8 and S8+, which basically uses the biometric information of users' pupils to identify them and enables them to unlock phones using only their eyes (Samsung Galaxy S8 | S8+, 2017). This feature is super useful especially during the current situation that we are all in the middle of the pandemic and have to wear masks all the time. Thus, our product team is determined to maintain and keep refining this feature in our future product launchment.

### 5.3 Top Attributes that are disliked the most

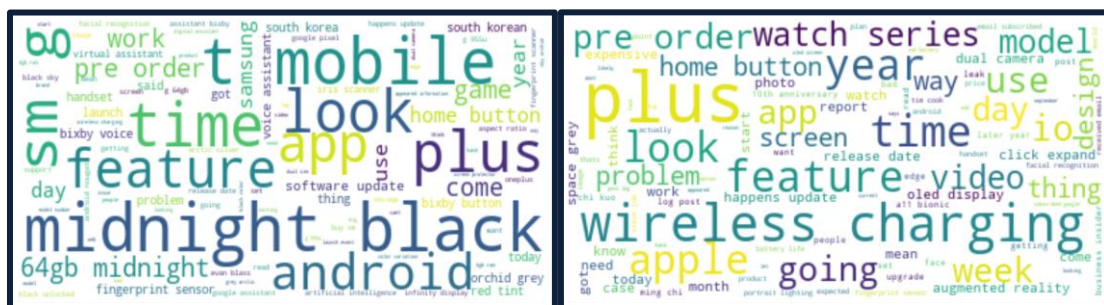


Figure 13: Word Cloud for Negative Samsung Posts      Word Cloud for Negative Apple Posts

For negative Samsung posts, we observed that people talked more about color, storage, screen issue, and fingerprint sensor. While for negative Apple posts, people talked about wireless charging, Apple Watch series and the design. Let's discuss more about the negative attributes around Samsung Galaxy S8 and S8+.

The first attribute mentioned the most is the new midnight black color for Galaxy S8. In this new launching, there are five available colors for customers to choose: midnight black, arctic silver, maple gold, coral blue, and orchid gray. Based on the word frequency result generated by Counter, we found that midnight black is not as successful as we thought it would be, compared to the new color gold offered in iPhone 8 which got a lot of positive feedback. Thus, our product team believes that the appearance of products is an important aspect to be focused on in the future product launches.

The second negative attribute is about the display. Although “Infinity Display” is one of the top attributes mentioned by people and got a lot of positive responses, some Galaxy S8 and S8+ users complained about the occurrence of the display with odd red tints, which stood out even more when being compared to the finely-tuned screen and OLED display of iPhones. This red tint issue needs to be solved as soon as possible. Along with making breakthrough innovations on the display design, our product team realized that screen quality should also be put in the same priority.

The third negative attribute is the fingerprint sensor. Samsung Galaxy S8 and S8+ actually have three unlocking ways, iris scanner, facial recognition and fingerprint. However, some users complained that the fingerprint sensor on the back of Galaxy S8 is not working properly. Our product team should have further discussion about whether to keep this feature or not in the future product design.

The fourth important attribute that got disliked a lot is the storage. For Galaxy S8 and Galaxy S8+, they are mainly provided in a storage of 64 GB (Samsung S8 | S8+ Specifications, 2017). Even if there is a 128GB edition for S8+, they are only sold in China, South Korea and India. On the other hand, the iPhone provides pretty flexible storage options including 64GB, 128GB and 256GB. Thus, in order for us to gain more customers in the U.S. market, we should think about adding more storage options to our customers, so that we can cover a wider range of customer demand.

## **6. Product Recommendations**

### **6.1 AARRR Model**

In order to make future improvements on our product and business strategy, we are going to give our suggestions based on the AARRR funnel model. AARRR model framework is an acronym for a set of five user-behavior metrics of product-led growth business should be tracking, including Acquisition, Activation, Retention, Revenue, Referral. AARRR could help with the following questions and provide a thorough perspective for our company to improve our products.

Acquisition: How are people discovering our product or company?

Activation: Are these people taking the actions we want them to take?

Retention: Are our activated users continuing to engage with our product?

Referral: Do users like the product enough to tell others about it?



Revenue: Are our customers willing to pay for our products?

As Samsung is already an international company with a large scale of customers. We decided to focus on the 3R strategy, which is Retention, Revenue and Referral.

## **6.2 Retention**

We should think about how to turn customers into repeat buyers and prevent them from switching to our competitor company. In order to do that, we would suggest the business strategy team and marketing team to make promotions on complementary products. For example, we could bundle sales of Samsung mobile phones and Samsung watch series.

The complementary products would increase the lock-in effect that customers tend to stick with our products, so that our customers would be more willing to keep a long term relationship with Samsung.

## **6.3 Revenue**

Revenue is the income generated from product selling, which is the most important topic a company will care about. In order to increase Samsung galaxy sales, we would suggest the product team do more customer research to know what product features consumers are concerned about the most. From our data, we found that our customers complained a lot about the appearance and performance of Samsung Galaxy S8, like the color and storage. The takeaway from this data analysis is that our company should hear more feedback from real customers in order to develop our new products better to satisfy our customers and increase the product sales.

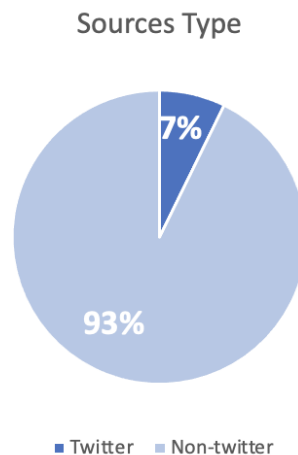
## **6.4 Referral**

As we all know, customer referrals are the most powerful selling and marketing tools. The best source of a new business is usually from a satisfied customer. To make our existing customers satisfied and be willing to recommend our products to new customers, we should make our product more attractive and unique compared to the products from our competitors. According to our analysis, we found that our customers love the flexible unlock ways in Samsung Galaxy S8. This function is our special selling point which Apple doesn't have right now. Plus, it's a very practical function especially under the pandemic situation. Therefore, we would suggest keeping the flexible unlock ways and trying to explore more special and attractive product features in the future.

## 7. Future Improvement for the Project

### 7.1 Models for data from non-twitter sources.

Based on our data, we noticed that about 93% data is from non-twitter sources. The content and form from these sources are very unlike twitter, but we haven't built up models for the data from these sources yet. We are going to develop more sentiment analysis models better fitting for different sources if we have more time to improve our project.



### 7.2 Customer segmentation.

Do analysis on different types of customers according to their gender, country, job etc. In this way, we could form a customer portrait to better understand how our customers are like and what kind of products they love the most.

## 8. Contribution

Shuhan Zhang - Data Cleaning/Preparation, Attribute Analysis

Jingyi Ran - Data Cleaning/Preparation, non-Twitter vs. Twitter Analysis

Ziyou Li - EDA, Sentiment Score Computing, Before After Launch Analysis

Xiang Li - Data Preprocessing, Attribute Analysis, Recommendations

Yiyun Hu - Data Preprocessing, Background Research, Slides Preparation

## 9. Reference

1. AARRR Pirate Metrics Framework, <https://www.productplan.com/glossary/aarr-framework/>
2. Analytics Vidhya, 2021, [analyticsvidhya.com/blog/2021/10/sentiment-analysis-with-textblob-and-vader/](https://analyticsvidhya.com/blog/2021/10/sentiment-analysis-with-textblob-and-vader/)
3. How to get customer referrals, <https://www.inc.com/guides/2010/08/how-to-get-customer-referrals.html#:~:text=Customer%20referrals%20are%20one%20of,shoppers%20read%20reviews%20before%20buying.>
4. Samsung Galaxy S8 | S8+, 2017, [https://www.samsung.com/latin\\_en/smartphones/galaxy-s8/security/](https://www.samsung.com/latin_en/smartphones/galaxy-s8/security/)
5. Social Media Richness Theory, [https://en.wikipedia.org/wiki/Media\\_richness\\_theory](https://en.wikipedia.org/wiki/Media_richness_theory)
6. Samsung S8 | S8+ Specifications, 2017, <https://www.samsung.com/global/galaxy/galaxy-s8/specs/>