

# Technical Report: Final Project DS 5110: Introduction to Data Management and Processing

Team Members: Nicholas Nehemia, Rebecca Bronfeld  
Khoury College of Computer Sciences  
Data Science Program  
`nehemia.n@northeastern.edu`, `bronfeld.r@northeastern.edu`

December 10, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Literature Review</b>	<b>2</b>
<b>3</b>	<b>Methodology</b>	<b>2</b>
3.1	Data Collection . . . . .	2
3.2	Data Preprocessing . . . . .	2
3.3	Analysis Techniques . . . . .	3
<b>4</b>	<b>Results</b>	<b>3</b>
<b>5</b>	<b>Discussion</b>	<b>4</b>
<b>6</b>	<b>Conclusion</b>	<b>4</b>
<b>7</b>	<b>References</b>	<b>4</b>
<b>A</b>	<b>Appendix A: Code</b>	<b>4</b>
<b>B</b>	<b>Appendix B: Tableau Analysis</b>	<b>7</b>

# 1 Introduction

In today's age, many individuals express the opinions of current events using various social media platforms. In our final project assignment, we will utilize social media, specifically Reddit, to gather comments from posts discussing iOS 15, 16, 17 and 18 upgrades, and conduct a sentiment analysis based on the comments and replies found.

# 2 Literature Review

# 3 Methodology

The Reddit API and PRAW Package allowed us to obtain top-level comments and their first reply, then gather that data. The transformers package allowed us to pass the data and begin labeling and classification.

Once data collection and preprocessing was complete, we had outputted two .CSV files for analysis: one related to the comment data, and the second related to the word count from the comments, grouped by iOS upgrades. The two .CSV files were then uploaded to a Tableau Desktop workbook. A horizontal bar-chart comparison, a word cloud, and a sentiment line graph was created, to better understand the types of comments that are within the negative/positive sentiments, which words were found most frequently for each iOS update, and the sentiment surrounding them.

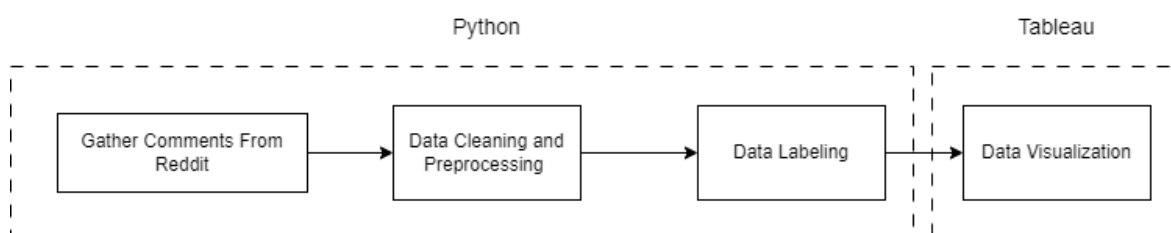


Figure 1: Process Flow of Project

## 3.1 Data Collection

Explain how the data was collected, including the sources and tools used.

In order to collect the data, we will manually find the reddit post we would like to extract comments from , and we will add to a .CSV stored in our system. From there , the application will read into the CSV and extract the data for all links in the CSV, and obtain all top level comments and the first level reply. (see App. A listing 1)

Data was collected using Reddit's API feature. We located posts relating to iOS 15-18 within Reddit, and applied their URLs to the API. Once we had the comments available, we then began processing and cleaning the comments.

## 3.2 Data Preprocessing

Our dataset contains 2 level of comments , the main comment and the 1st level replies to the 1st level comment. We only take the first reply higher level replies might go off topic.

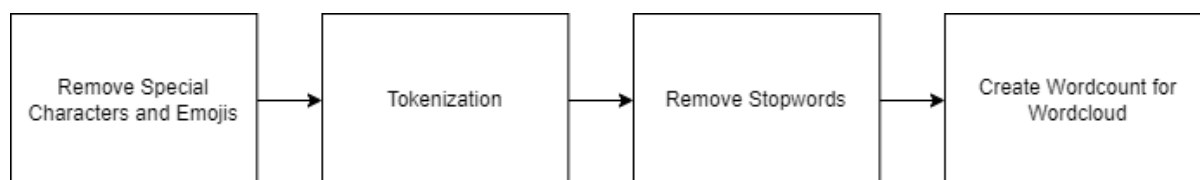


Figure 2: Data Cleaning and preprocessing for analysis

The notebook script will then find all parent comments ,and join the child comments , so we can have the parent and child comments in one row for easier identification and visualization purposes. (See App. A listing 3)

To generate the wordcloud , we also need to have a wordcount table. First , we will clean the data by removing the stopwords and emojis, and doing tokenization. We will then do a word count grouped by iOS version and classification label. (See App. A listing 4)

### 3.3 Analysis Techniques

In this project , we use the transformers package to label the sentiment of all the comments. Using the HuggingFace Transformers pipeline will simplify this process as we can just pass through the comment, and it will run through a premade pipeline and return the sentiment. (See App. A listing 3) Although we do not need specific preprocessing for the labeling of data, we still need it for the word count and analysis, as the stopwords will appear on the wordcloud and give no insight.

We will also need to perform transformations on the Data Received, as the data structure of the comments is similar to a tree, where we have to link the child comment to the parent comment in order to gain insight in the visualization more easily. To accommodate this , we have transformed the data to have both the child comment and parent comment in one row through some transformation.

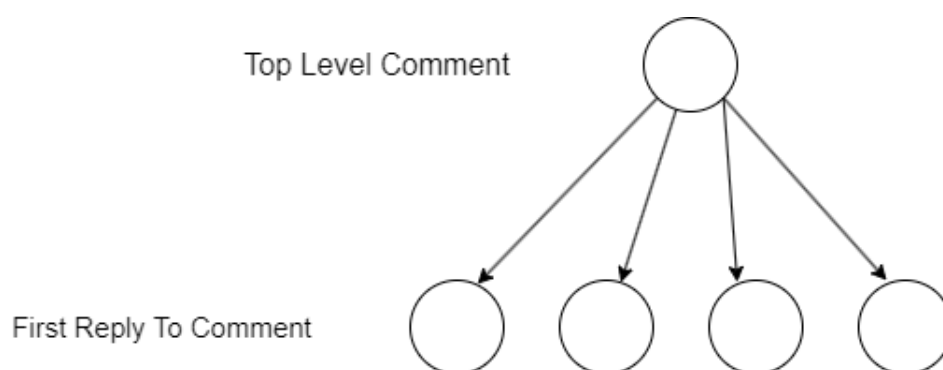


Figure 3: Reddit Comment Data Structure

## 4 Results

Overall, each iOS update has a majority of negative sentiment, across all comments. More than 60 percent of sentiment was negative per update (see Figures 4-7).

By analyzing the word cloud by each update, we could see the most frequently referenced

enhancements included notifications, the control center, and lock screen (see Figures 4-7). Once we had the two .CSV files structured for visualization, we imported the files to a Tableau Workbook and began analysis. Figures 8-10 display the backend of creating the charts used for analysis.

## 5 Discussion

Our results from the Tableau dashboard indicate a need for the business to consider fixing current enhancements as a priority, over net-new features. If a user is not satisfied with their experience, customer loyalty will diminish, and therefore there will be less users to appreciate future updates.

## 6 Conclusion

A key finding from our sentiment analysis is that users displayed negative sentiment for major iOS features, including the battery life, notification bar, and control center. Out of all the updates, the worst sentiment came from iOS 15, with only 10 percent sentiment being positive.

There were a few limitations from end-to-end process. First, we had removed emojis from final analysis, as there was a limitation in Tableau for visualizing these character types. This limited our sentiment to only be shown as text.

The model we had used also does not detect sarcasm, or have consistent accuracy with questions. This could skew the data towards positive sentiment.

We also had limited data in our analysis, as it was only a sample of reviews users had towards this updates. If we had a larger sample size, our sentiment could be viewed differently.

Finally, sampling bias could have occurred, as we were only reviewing iOS reactions from a social media site. Users tend to be vocal online when they have a topic to complain about, therefore skewing the sentiment towards negative.

## 7 References

### References

Remove symbols : [saturncloud.io](https://saturncloud.io)

Source for how to scrape using PRAW : [GeeksForGeeks](https://www.geeksforgeeks.org/pRAW/)

PRAW Documentation : [praw.readthedocs.io](https://praw.readthedocs.io)

Source for obtaining word counts : [earthdatascience.org](https://earthdatascience.org)

Text Analysis Tableau Reference : [Tableau Public](https://public.tableau.com/)

Remove Emojis : [stackoverflow](https://stackoverflow.com/questions/49232925/removing-emojis-from-a-string-in-python)

## A Appendix A: Code

```
1 def scrape_reddit(urlCSVpath):  
2
```

```

3 # Define the post URL or ID
4 urlcsv=pd.read_csv(urlCSVpath) # get URL to scrape
5 urlList=urlcsv['urlList'].to_list() #convert to list
6 all_comments_data = []
7 for val in urlList:
8     post_url = val
9     submission = reddit.submission(url=post_url)
10    # Fetch top-level comments with their scores and store in a
    list
11    submission.comments.replace_more(limit=0) # Removes "More
    Comments" placeholders
12    for top_level_comment in submission.comments:
13        all_comments_data.append({
14            "subreddit": submission.subreddit.display_name,
15            "post_title": submission.title,
16            "post_url": val,
17            "comment_id": top_level_comment.id,
18            "parent_id": submission.id, # Replies to the post
    itself
19            "comment": top_level_comment.body,
20            "comment_score": top_level_comment.score,
21            "comment_timestamp": datetime.fromtimestamp(
top_level_comment.created_utc),
22            "is_top_level": True # Flag for top-level comment
23        })
24
25    # Check if there's a first reply and add it
26    for first_reply in top_level_comment.replies:
27        all_comments_data.append({
28            "subreddit": submission.subreddit.display_name,
29            "post_title": submission.title,
30            "post_url": val,
31            "comment_id": first_reply.id,
32            "parent_id": top_level_comment.id, # Points to the
    top-level comment ID
33            "comment": first_reply.body,
34            "comment_score": first_reply.score,
35            "comment_timestamp": datetime.fromtimestamp(
first_reply.created_utc),
36            "is_top_level": False # Flag for first reply
37        })
38
39    # Convert to DataFrame
40    df_comments = pd.DataFrame(all_comments_data)
41    return(df_comments)
42 1

```

Listing 1: getting all top level comments and first level replies.

```

1 classifier=pipeline("sentiment-analysis")
2
3 def classifyComment(comment):
4     sentiment=classifier(comment)[0]
5     return sentiment['label'],sentiment['score']
6
7 df_commentsAll[['classification_result','classification_score']]=
    df_commentsAll['comment'] \
8     .apply(lambda x :pd.Series(classifyComment(x)))

```

```
9 df_commentsAll
```

Listing 2: Function to generate the sentiment using HuggingFace

```
1 classifier=pipeline("sentiment-analysis")
2 structured_data = []
3
4 # Loop through each top-level comment and its replies
5 for idx, comment in df_commentsAll.iterrows():
6     if comment["is_top_level"]: # Process top-level comments
7         structured_data.append({
8             "subreddit": comment["subreddit"],
9             "post_title": comment["post_title"],
10            "post_url": comment["post_url"],
11            "parent_comment_id": comment["comment_id"], # Top-level
comment's ID
12            "parent_comment": comment["comment"],
13            "parent_comment_score": comment["comment_score"],
14            "parent_comment_sentiment": comment["classification_result"]
15        ],
16            "parent_adjusted_classification_score" : comment["
adjusted_classification_score"],
17            "parent_comment_timestamp": comment["comment_timestamp"],
18            "reply_comment_id": None, # No reply for this row, as it's
the top-level comment
19            "reply_comment": None,
20            "reply_comment_score": None,
21            "reply_comment_sentiment": None,
22            "reply_comment_timestamp": None,
23            "reply_adjusted_classification_score" : None,
24            "IOS":comment["IOS"]
25        })
26
27 # Add each reply to the top-level comment as a new row
28 replies = df_commentsAll[(df_commentsAll["parent_id"] ==
comment["comment_id"]) & (df_commentsAll["is_top_level"] == False)]
29 for _, reply in replies.iterrows():
30     structured_data.append({
31         "subreddit": comment["subreddit"],
32         "post_title": comment["post_title"],
33         "post_url": comment["post_url"],
34         "parent_comment_id": comment["comment_id"], # ID of
the top-level comment
35         "parent_comment": comment["comment"],
36         "parent_comment_score": comment["comment_score"],
37         "parent_comment_sentiment": comment["
classification_result"],
38         "parent_adjusted_classification_score" : comment["
adjusted_classification_score"],
39         "parent_comment_timestamp": comment["comment_timestamp"]
40     ],
41         "reply_comment_id": reply["comment_id"], # ID of the
reply
42         "reply_comment": reply["comment"],
43         "reply_comment_score": reply["comment_score"],
44         "reply_comment_sentiment": reply["classification_result"]
45     ],
46         "reply_comment_timestamp": reply["comment_timestamp"],
```

```
44         "reply_adjusted_classification_score" : reply["
adjusted_classification_score"],
45         "IOS":reply['IOS']
46     })
47
48 # Create a DataFrame from structured data
49 df_structured_comments = pd.DataFrame(structured_data)
```

Listing 3: Generate data structure for visualization

```
1 def cleanCount(df,ios,sentiment):
2     stop = stopwords.words('english')
3     cleaned_stopwords=df[(df["IOS"]==ios )& (df["classification_result"
]==sentiment)][["unigrams"].apply(lambda x: [item for item in x if
item not in stop])]
4     all_words_nsw = list(itertools.chain(*cleaned_stopwords))
5     counts_nsw = collections.Counter(all_words_nsw)
6     final_counts_nsw=counts_nsw.most_common()
7     dfCounts=pd.DataFrame(final_counts_nsw,columns=["words","count"])
8     dfCounts["ios"]=ios
9     dfCounts["Sentiment"]=sentiment
10    return dfCounts
```

Listing 4: Function to Remove stop words and generate word count

## B Appendix B: Tableau Analysis

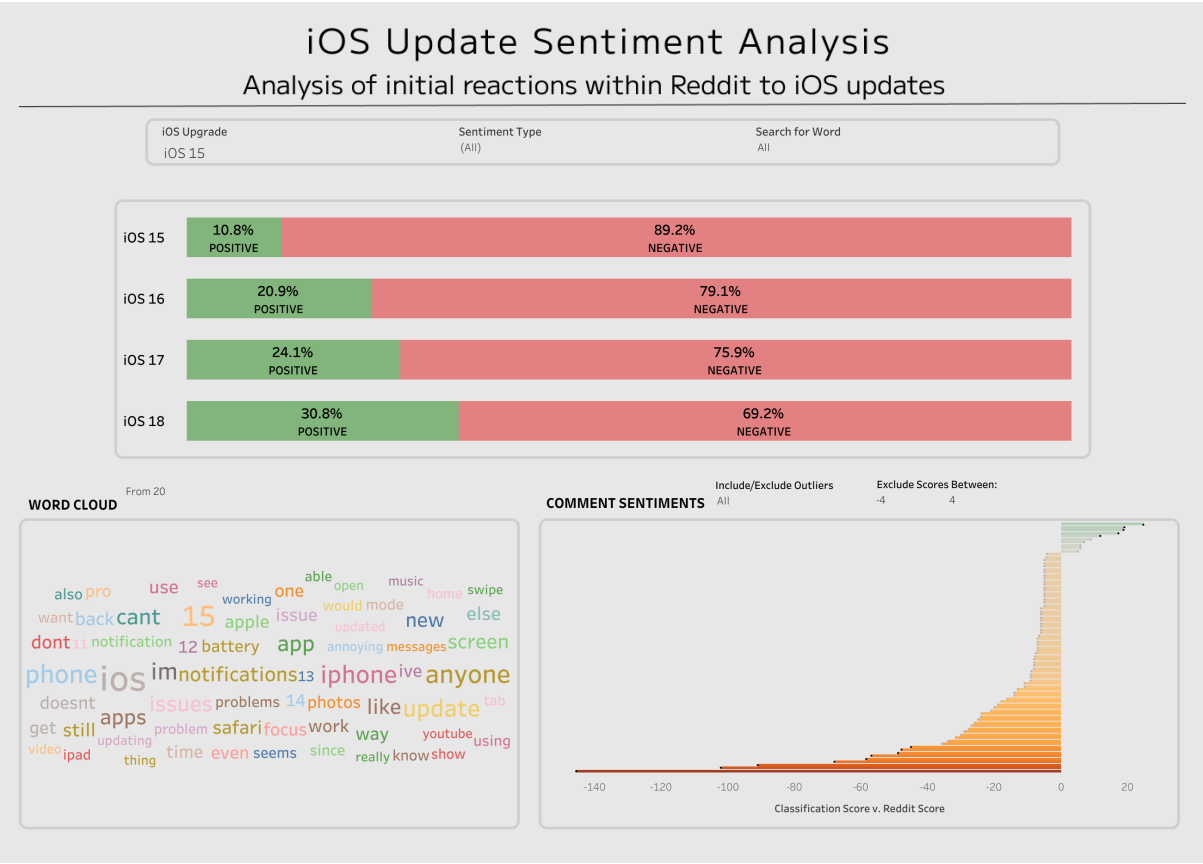


Figure 4: iOS Dashboard Filtered to iOS 15



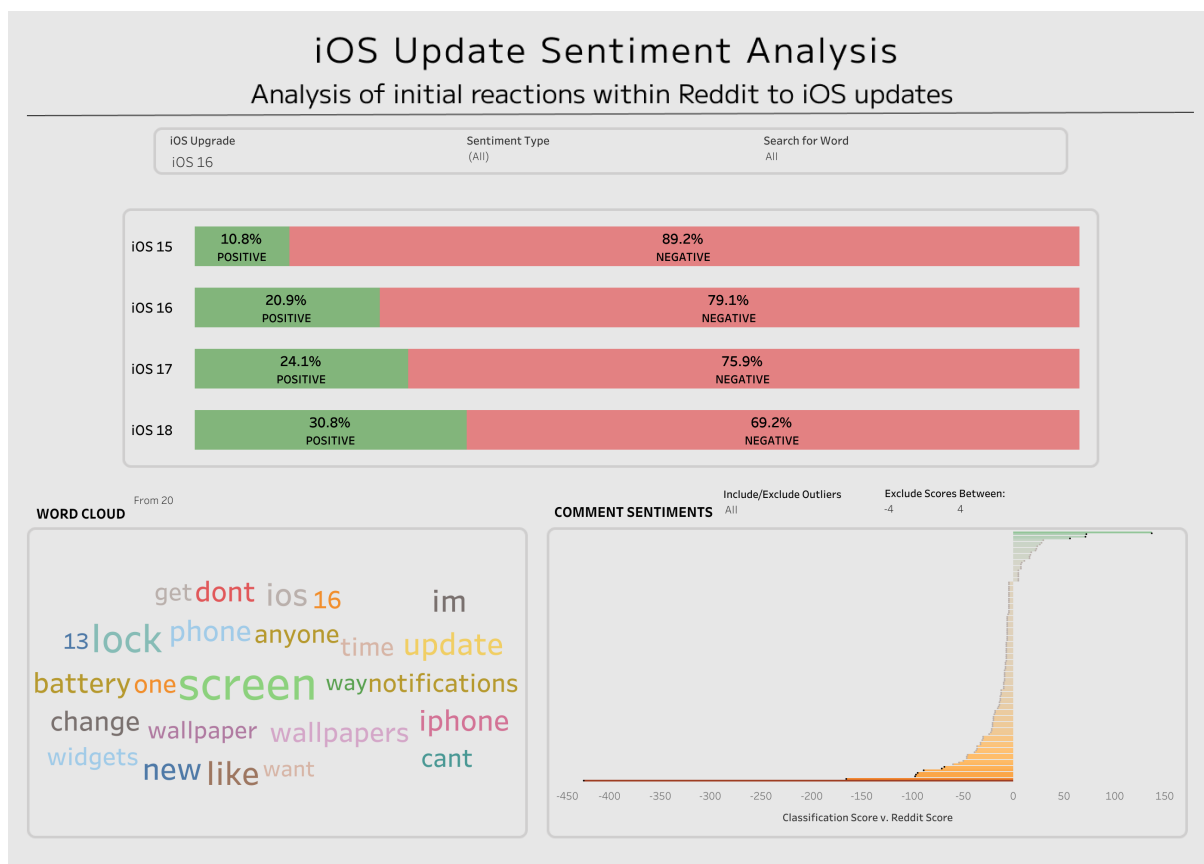


Figure 5: iOS Dashboard Filtered to iOS 16

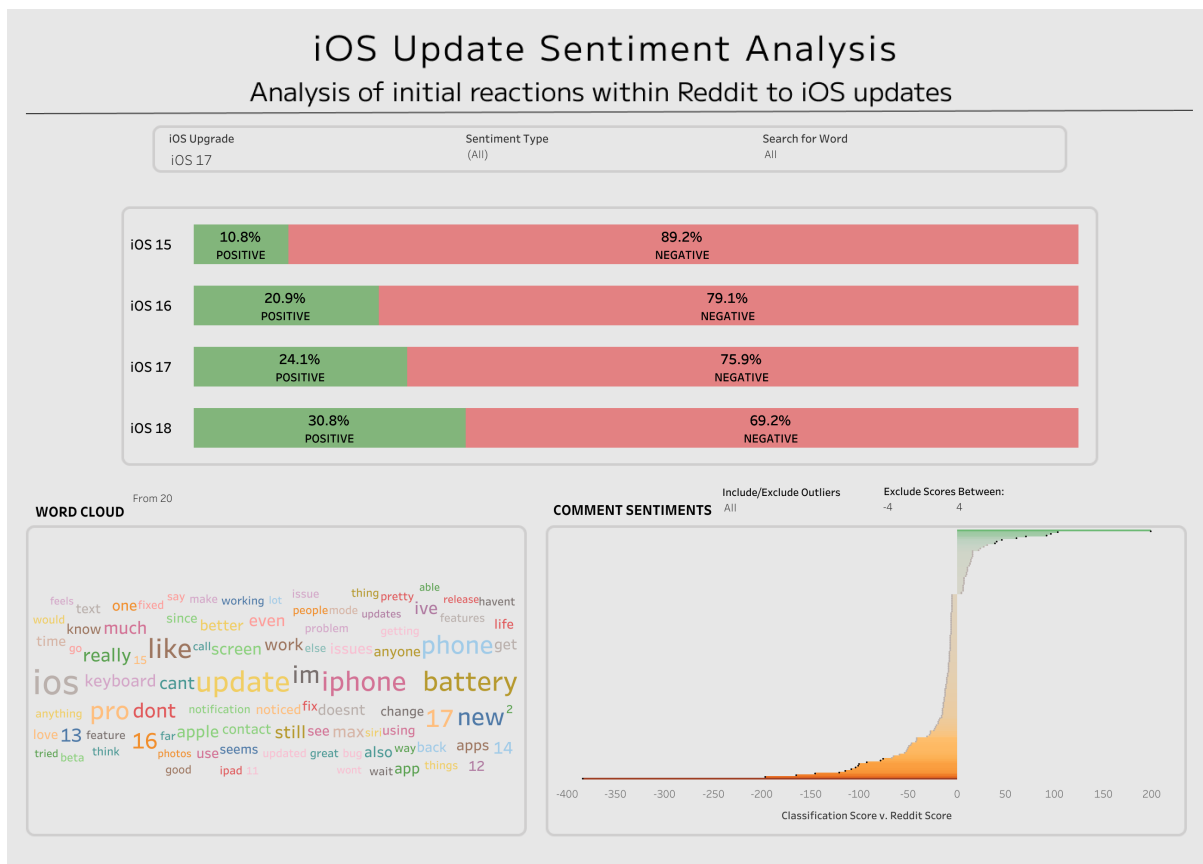


Figure 6: iOS Dashboard Filtered to iOS 17

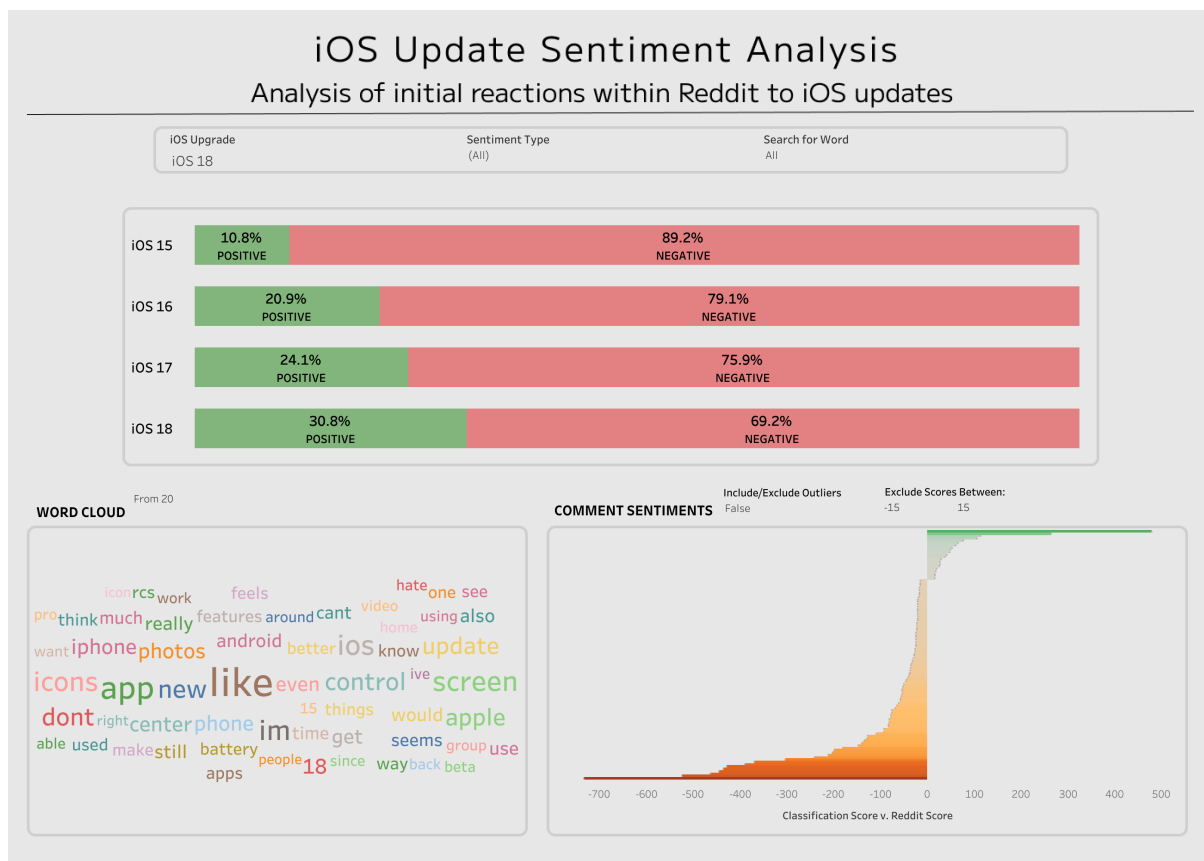


Figure 7: iOS Dashboard Filtered to iOS 18

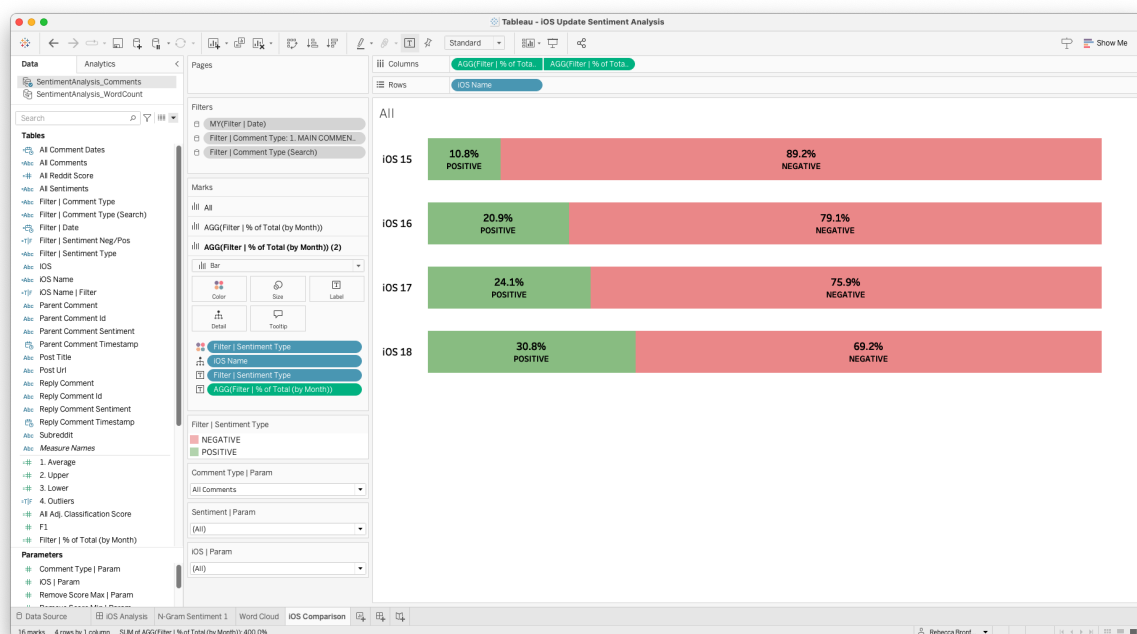


Figure 8: Tableau Backend - Creating the Comparison Bar Chart

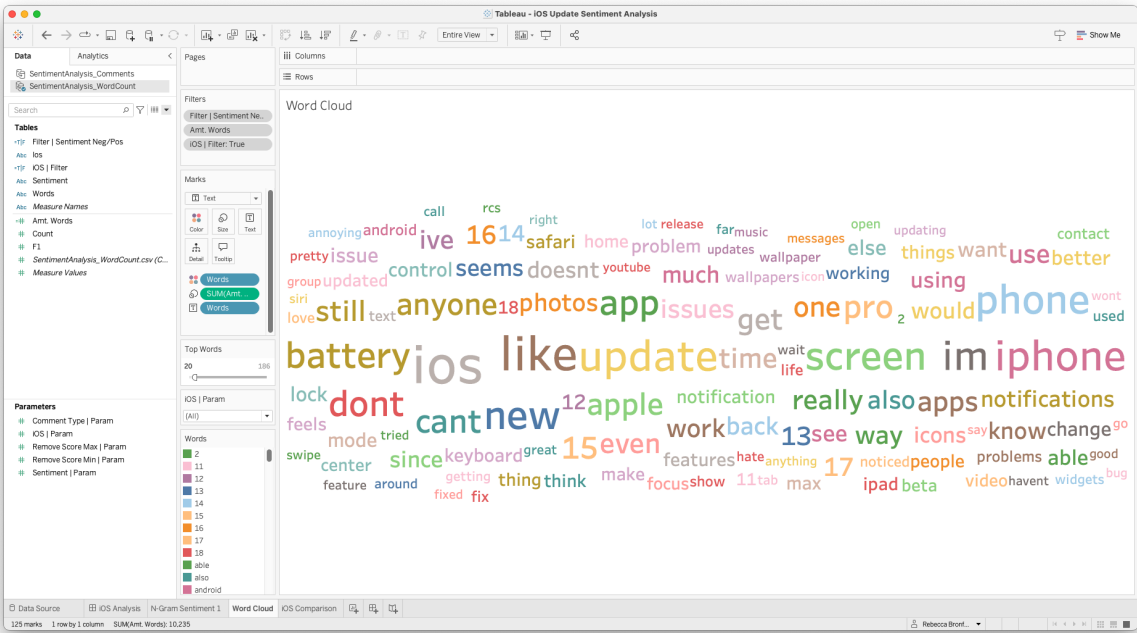


Figure 9: Tableau Backend - Creating the Word Cloud

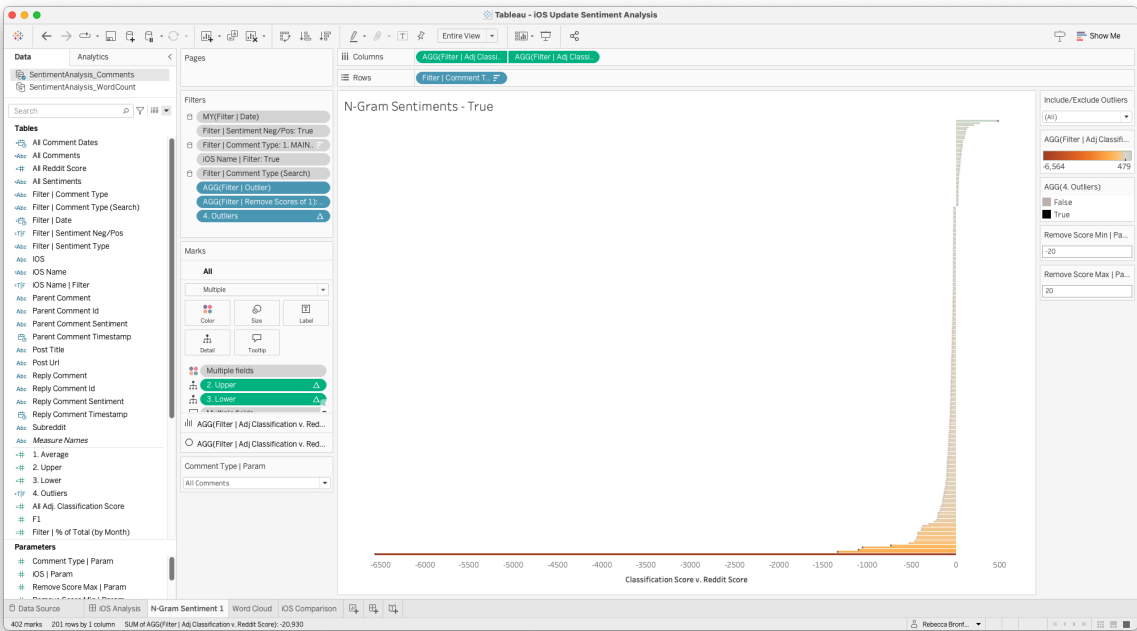


Figure 10: Tableau Backend - Creating the Comment-Level Sentiment graph