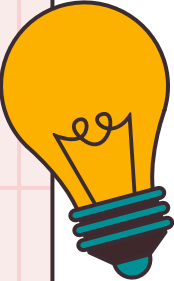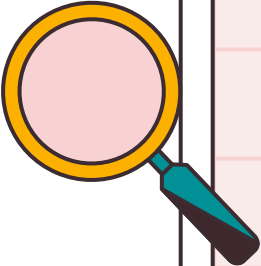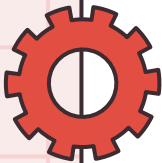Data 400

# Social Media Sentiment Analysis

By Abhik and Amanda

# Contents

1. Research Topic
2. Data set
3. Exploratory Data Analysis
4. Modelling
5. Implications for Stakeholders
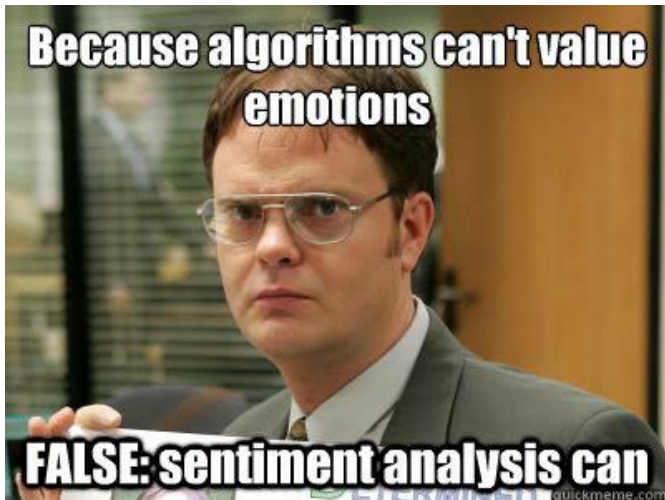6. Ethical, Legal, Societal Implications

# 1. Research Topic

Topic: Social Media Sentiment Analysis

Objective: Finding out people's feeling about a brand or product at scale

Task: Conduct sentiment analysis through models like Naive Bayes Classifier, Logistic Regression, Random Forest, and KNN

# 2. Data set

The data set is from the internet and contains 732 records and includes the following columns:

**Text**: Content of social media posts.
**Sentiment**: Sentiment classification of the text
**Timestamp**: Date and time when the post was made.
**User**: Username or identifier of the poster.
**Platform**: Social media platform used (e.g., Twitter, Instagram, Facebook).
**Hashtags**: Hashtags included in the posts.
**Retweets**: Number of retweets the post received.
**Likes**: Number of likes the post received.
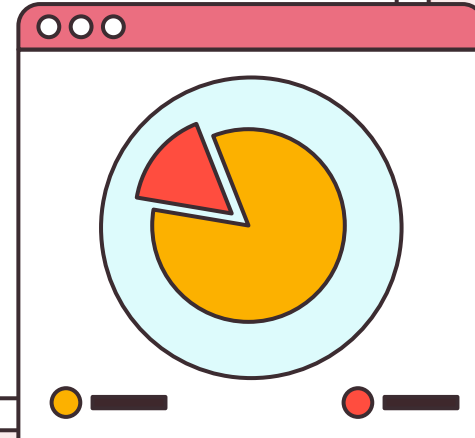**Country**: Country associated with the post.
**Year, Month, Day, Hour**: Time-related columns extracted from the timestamp for analysis.

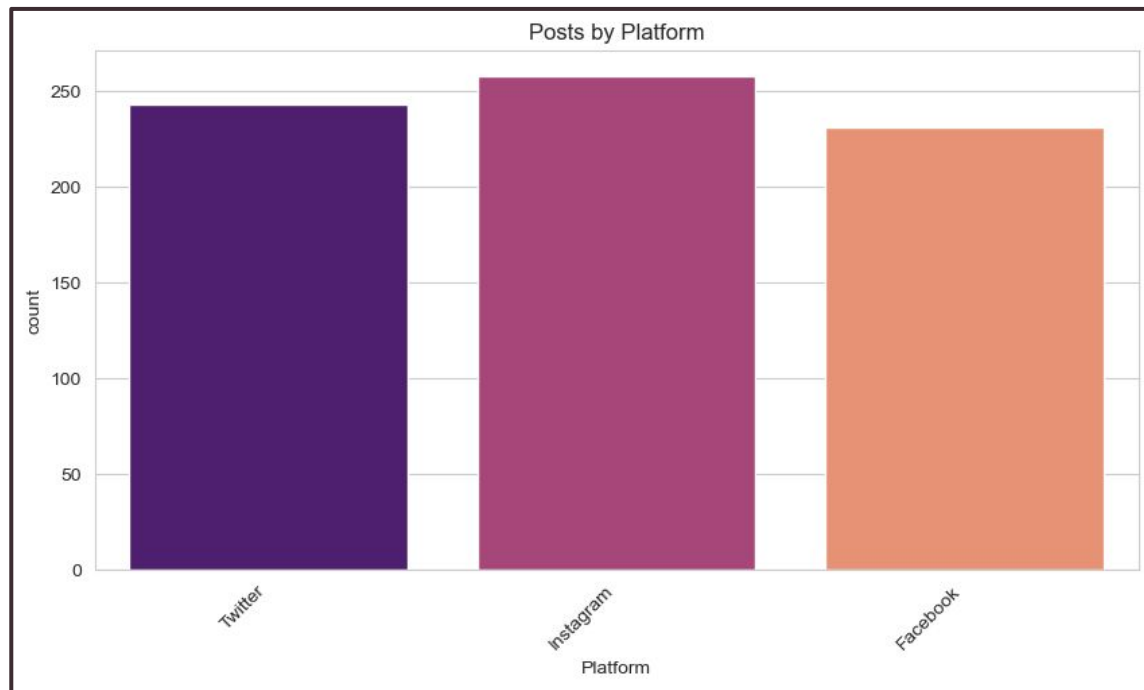| | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Unnamed: 0 | Text | Sentiment | Timestamp | User | Platform | Hashtags | Retweets | Likes | Country | Year | Month | Day | Hour |
| 2 | 0 | Enjoying a be | Positive | 1/15/23 12:30 | User123 | Twitter | #Nature #Pai | 15 | 30 | USA | 2023 | 1 | 15 | 12 |
| 3 | 1 | Traffic was te | Negative | 1/15/23 8:45 | CommuterX | Twitter | #Traffic #Moi | 5 | 10 | Canada | 2023 | 1 | 15 | 8 |
| 4 | 2 | Just finished | Positive | 1/15/23 15:45 | FitnessFan | Instagram | #Fitness #Wo | 20 | 40 | USA | 2023 | 1 | 15 | 15 |
| 5 | 3 | Excited abou | Positive | 1/15/23 18:20 | AdventureX | Facebook | #Travel #Adv | 8 | 15 | UK | 2023 | 1 | 15 | 18 |
| 6 | 4 | Trying out a r | Neutral | 1/15/23 19:55 | ChefCook | Instagram | #Cooking #F | 12 | 25 | Australia | 2023 | 1 | 15 | 19 |
| 7 | 5 | Feeling grate | Positive | 1/16/23 9:10 | GratitudeNo | Twitter | #Gratitude #I | 25 | 50 | India | 2023 | 1 | 16 | 9 |
| 8 | 6 | Rainy days ca | Positive | 1/16/23 14:45 | RainyDays | Facebook | #RainyDays | 10 | 20 | Canada | 2023 | 1 | 16 | 14 |
| 9 | 7 | The new mov | Positive | 1/16/23 19:30 | MovieBuff | Instagram | #MovieNight | 15 | 30 | USA | 2023 | 1 | 16 | 19 |
| 10 | 8 | Political disc | Negative | 1/17/23 8:00 | DebateTalk | Twitter | #Politics #De | 30 | 60 | USA | 2023 | 1 | 17 | 8 |
| 11 | 9 | Missing sum | Neutral | 1/17/23 12:20 | BeachLover | Facebook | #Summer #B | 18 | 35 | Australia | 2023 | 1 | 17 | 12 |
| 12 | 10 | Just publishe | Positive | 1/17/23 15:15 | BloggerX | Instagram | #Blogging #N | 22 | 45 | USA | 2023 | 1 | 17 | 15 |
| 13 | 11 | Feeling a bit | Negative | 1/18/23 10:30 | WellnessCh | Twitter | #SickDay #H | 7 | 15 | Canada | 2023 | 1 | 18 | 10 |
| 14 | 12 | Exploring the | Positive | 1/18/23 14:50 | UrbanExplor | Facebook | #CityExplore | 12 | 25 | UK | 2023 | 1 | 18 | 14 |

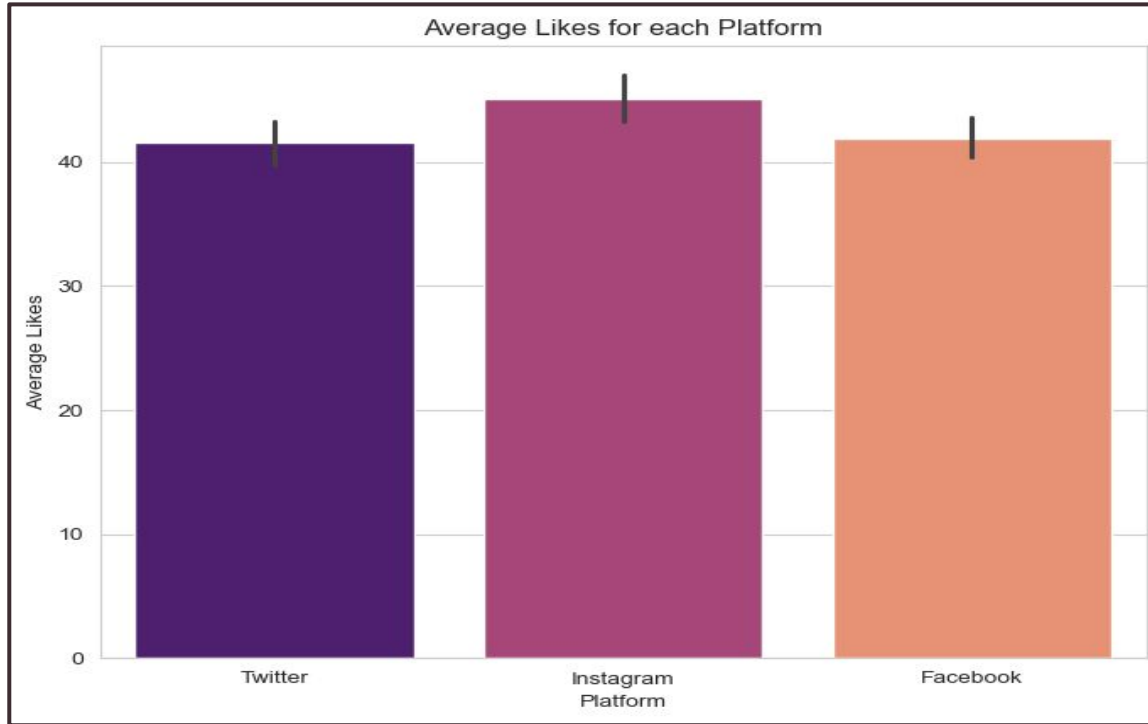# 3. Exploratory Data Analysis

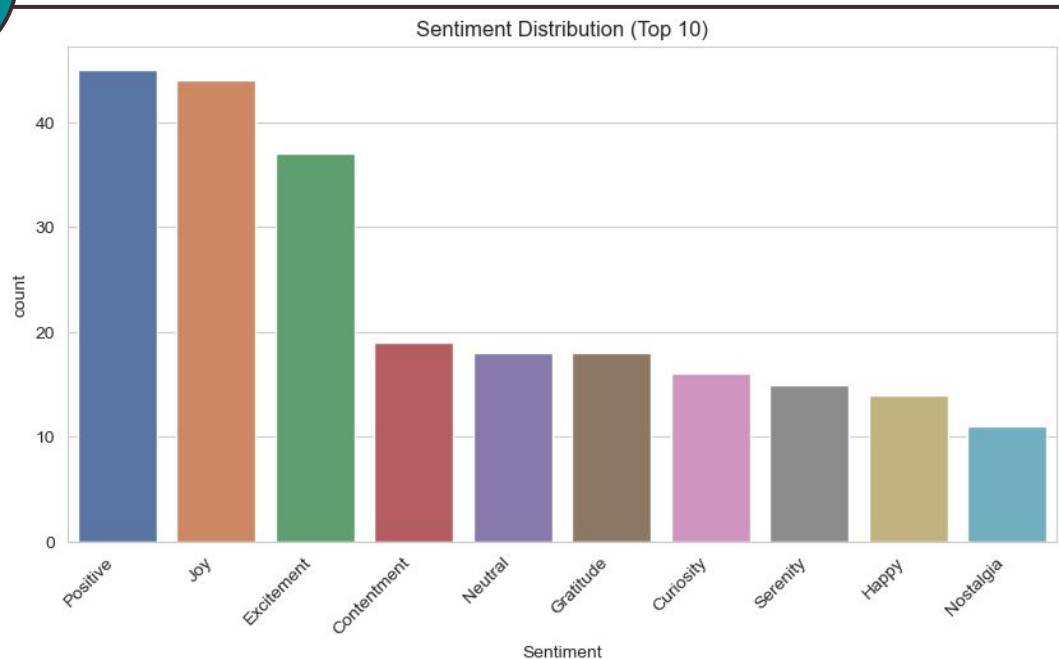yay

# Number of Posts by Platform

Posts by Platform



- Instagram has the most posts from this data set.
- Twitter has second most, then Facebook.
- All around the same amount.
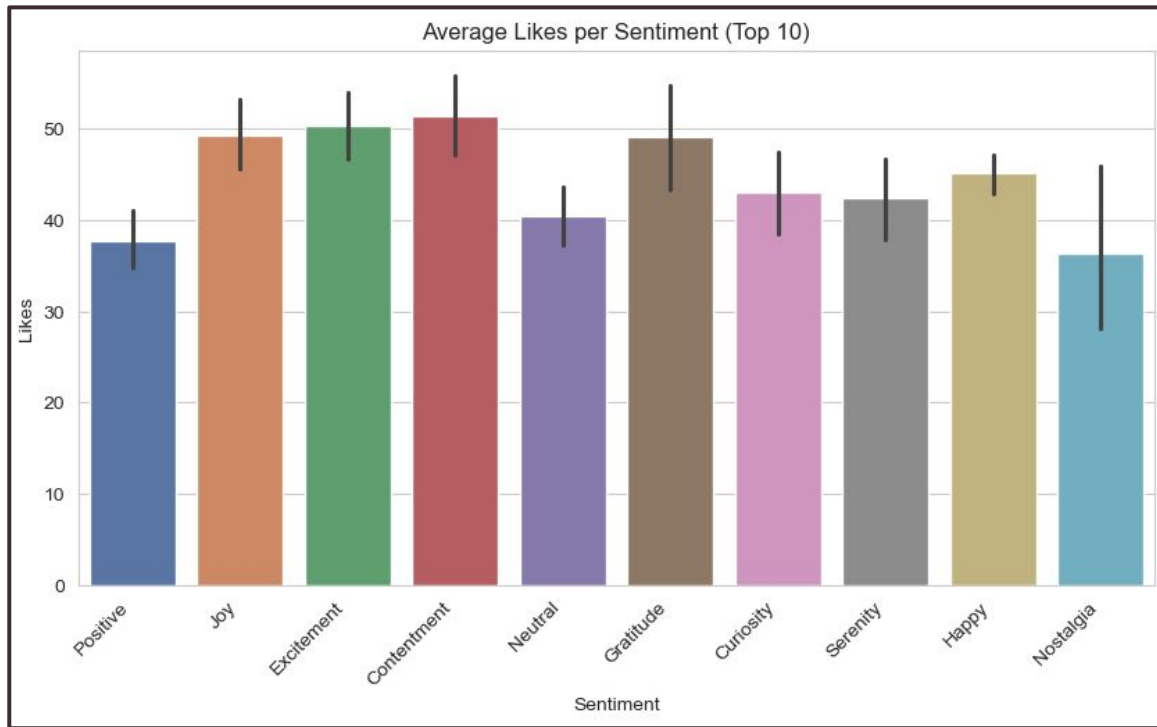
# Average Likes for each Platform



Average Likes for each Platform

- Very similar to last graph with instagram leading.
- Facebook is higher than twitter in likes.
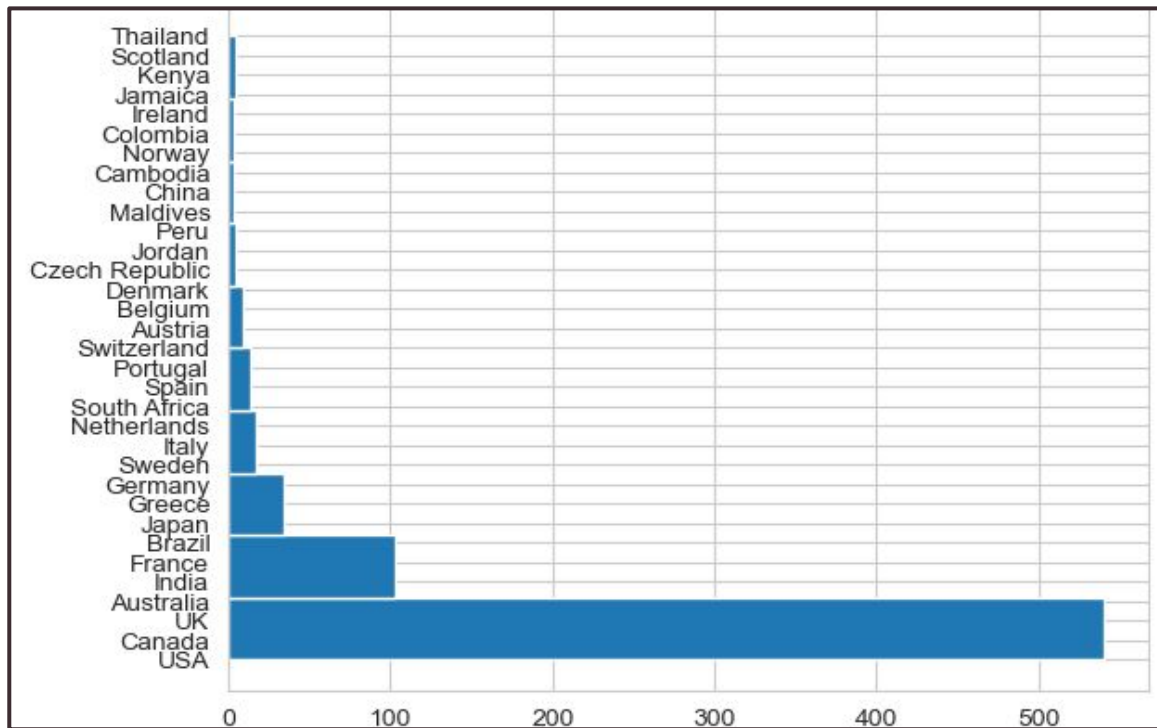
# Top 10 Sentiment Distribution



Sentiment Distribution (Top 10)

Top Ten: Positive, Joy,
Excitement,
Contentment, Neutral,
Gratitude, Curiosity,
Serenity, Happy,
Nostalgia

# Average Likes per Sentiment



Average Likes per Sentiment (Top 10)

- Positive is the top sentiment yet has the lowest likes among the top ten.
- Contentment has the highest average likes per sentiment.
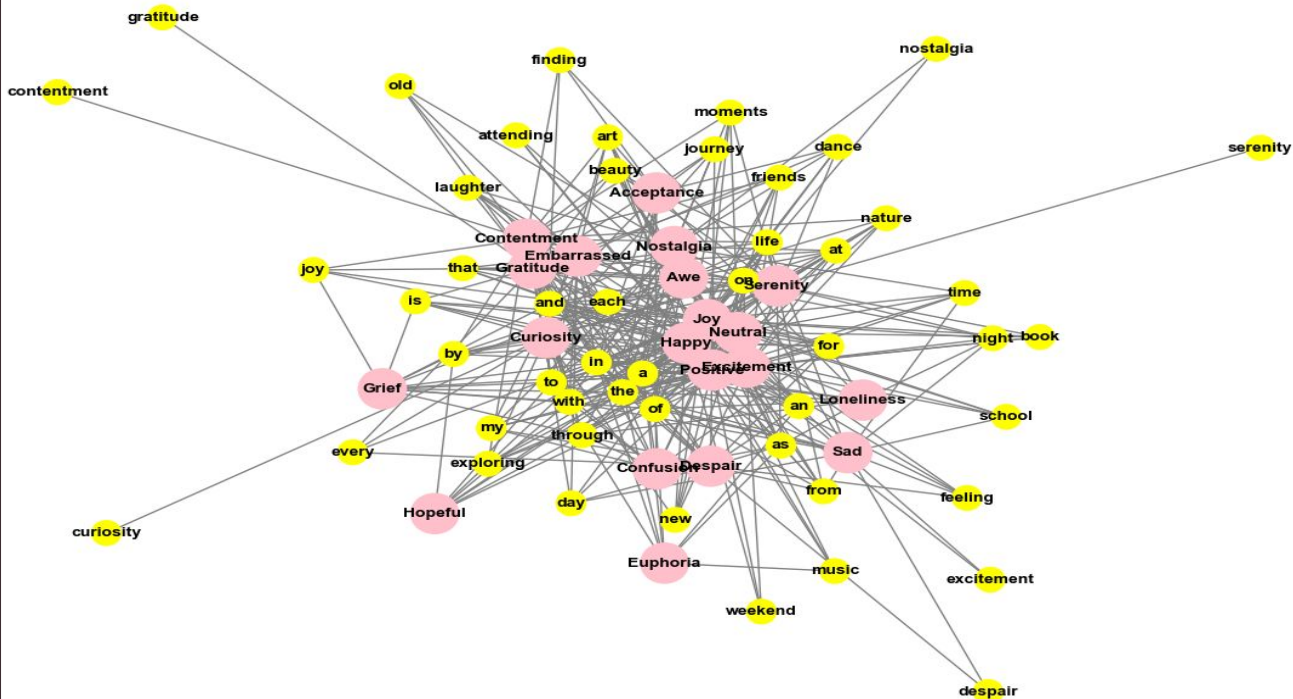
# Data Collected: Countries



- Most of the data comes from:
  - USA
  - Canada
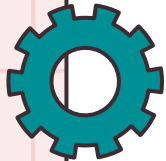  - UK
  - Australia
  - India

# Sentiment to Word Relationship
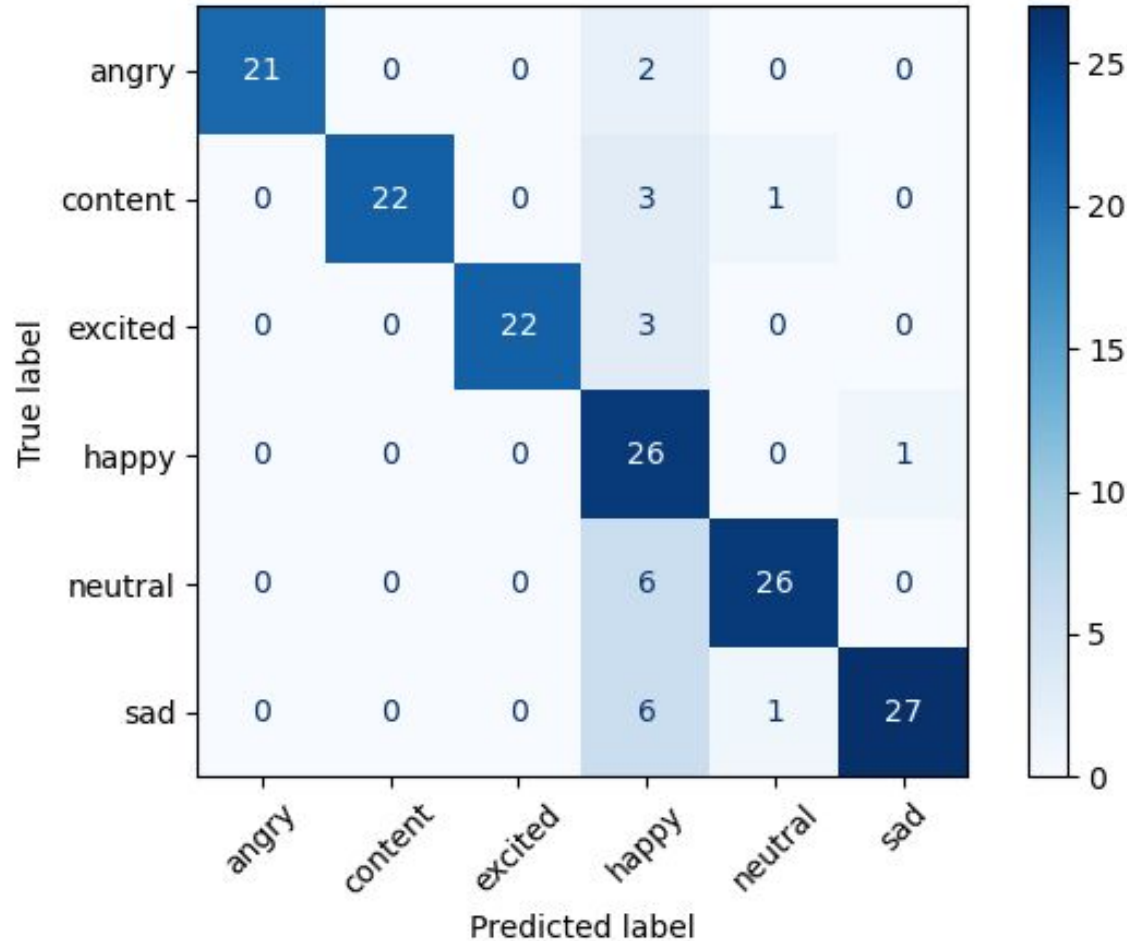


Word-Sentiment Relationship Graph

- Pink → Top 20 Sentiments
- Yellow → Top 50 frequently used words
- Some close connections: new + confusion and despair, laughter + contentment and embarrassed, school + sad and loneliness (haha)
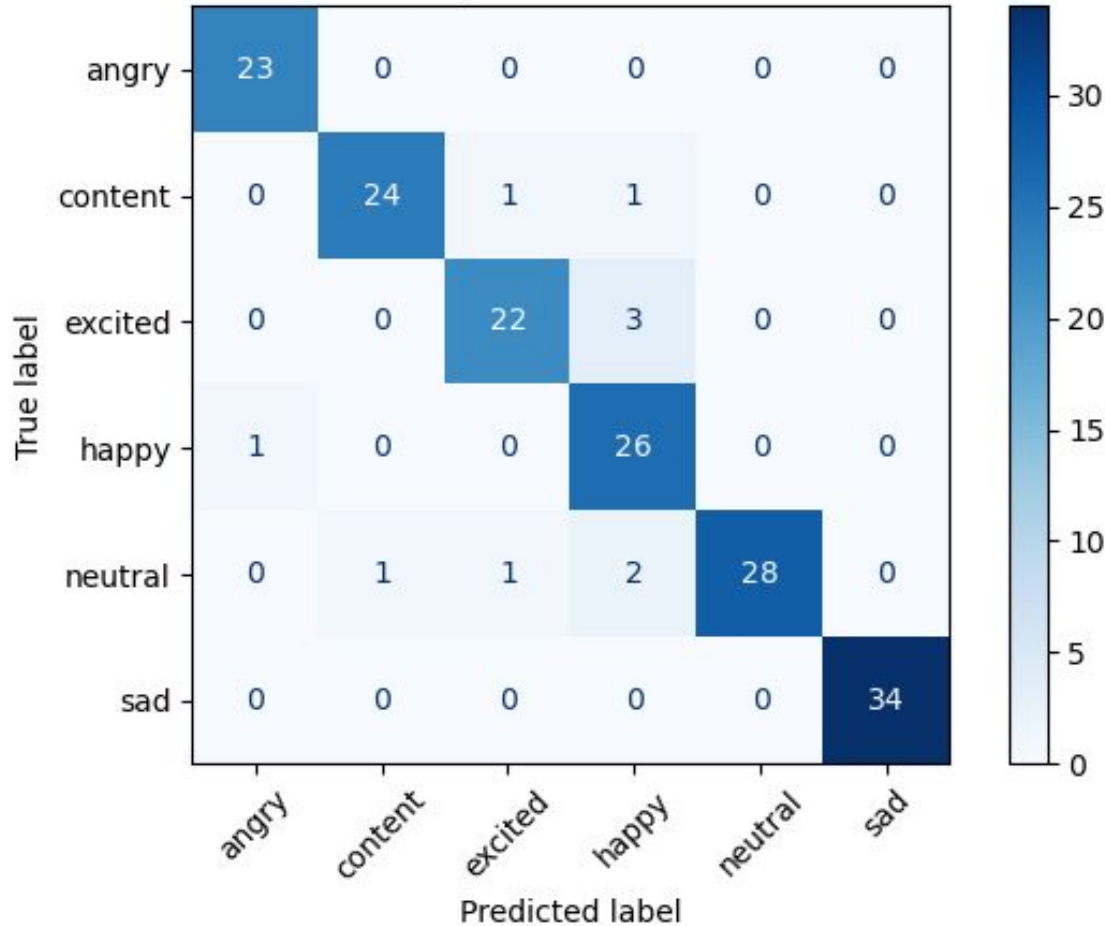
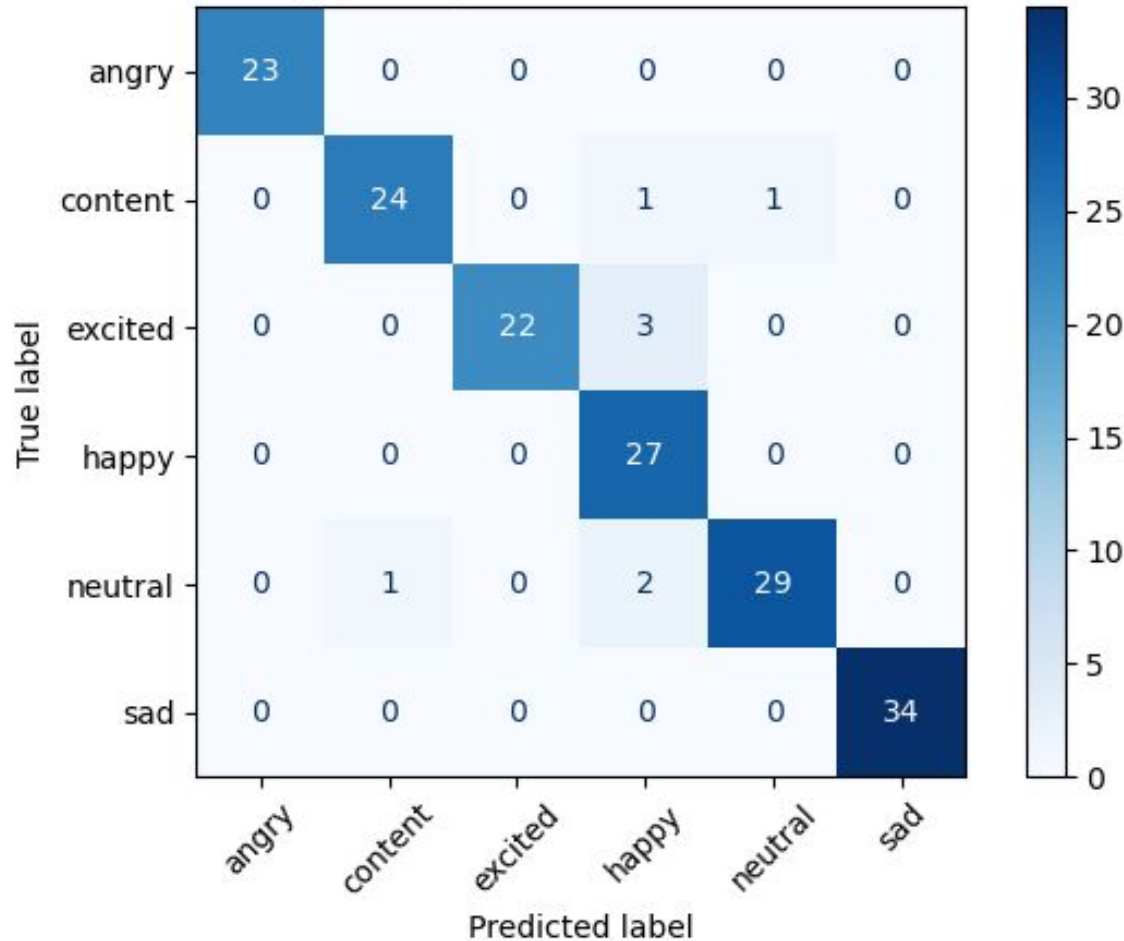# 4. Modelling

KNN Confusion Matrix

- 2 neighbors
- Most errors are between similar tones like *happy/sad/neutral*
- Some confusion between *happy* and other categories
- Because this model looks at the 'closest' past examples to make a prediction, it can struggle when sentiments are semantically similar
- Accuracy Rate: 84%
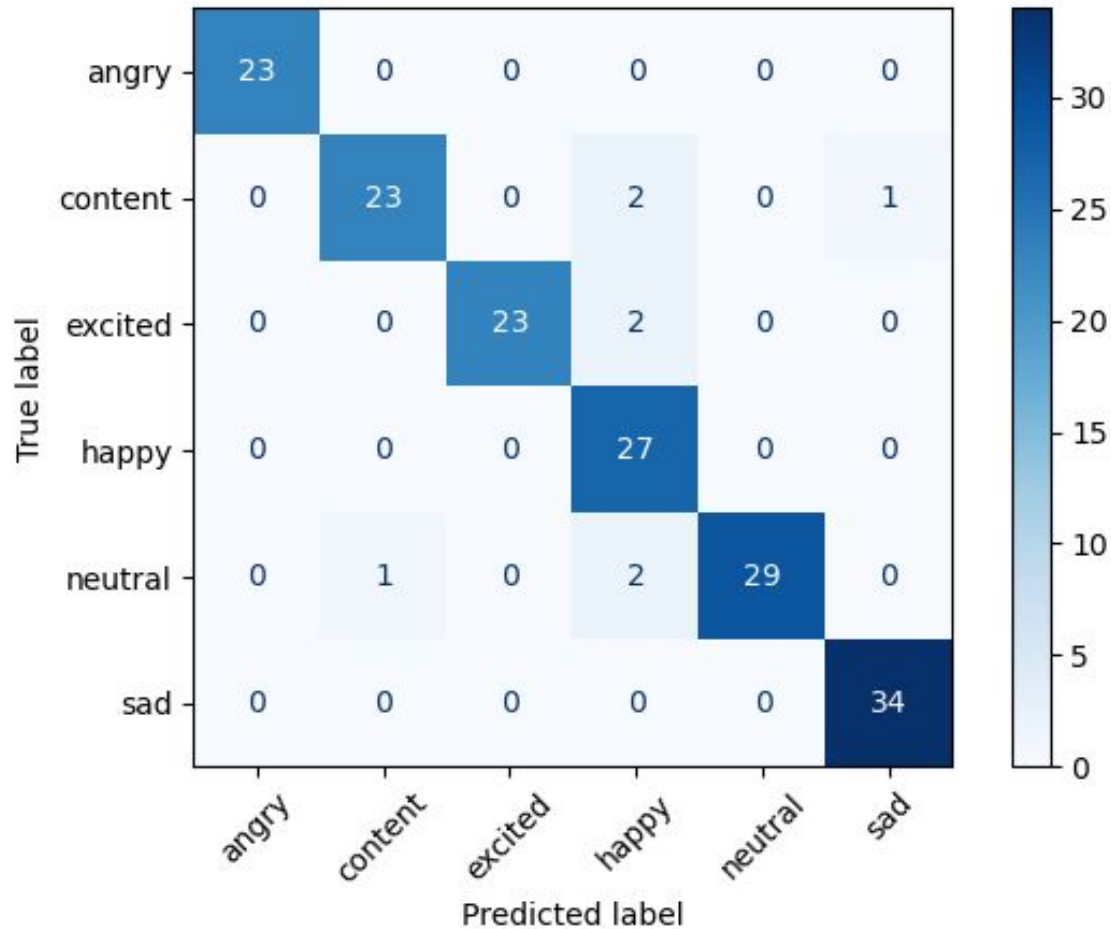
## Naive Bayes Confusion Matrix

- Naive Bayes is fast and handles text well, but it can be overly confident
- It tends to predict 'sad' more often, probably because certain negative words show up more often in training.
- Accuracy Rate: 94%
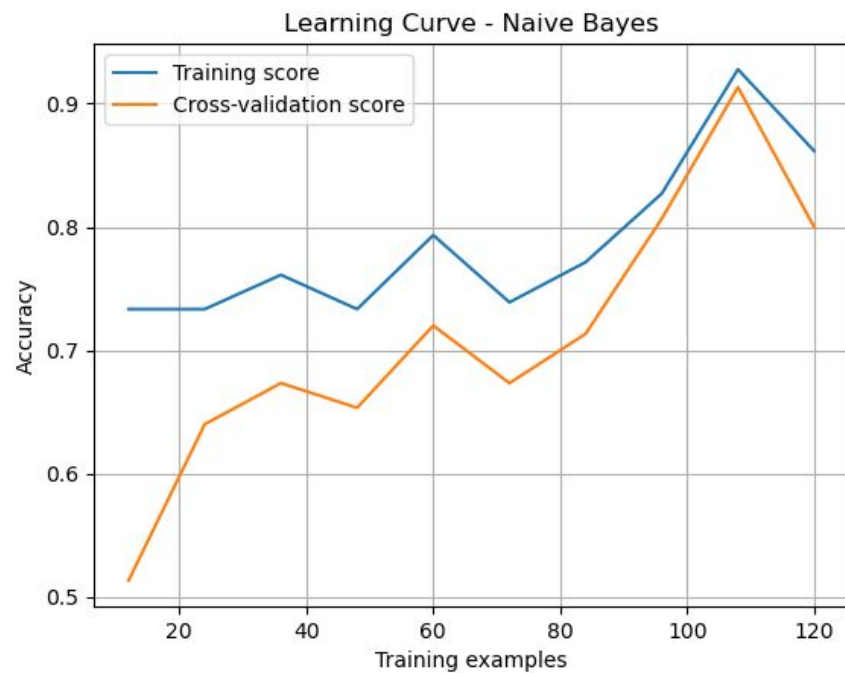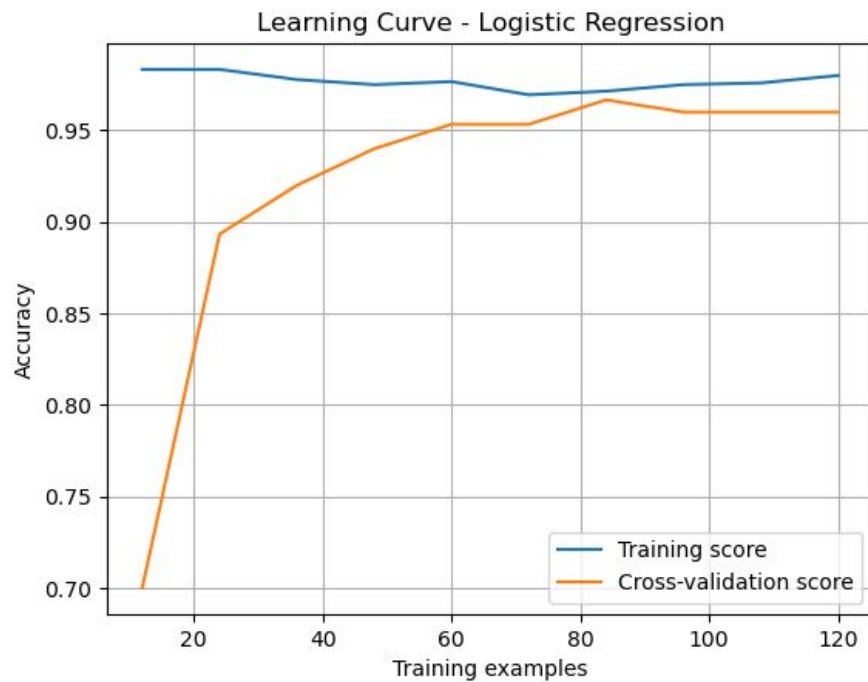
Logistic Regression Confusion Matrix

- This model performed cleanly with only a few misclassifications
- It uses probabilities to decide the most likely sentiment and showed strong precision and recall for most labels
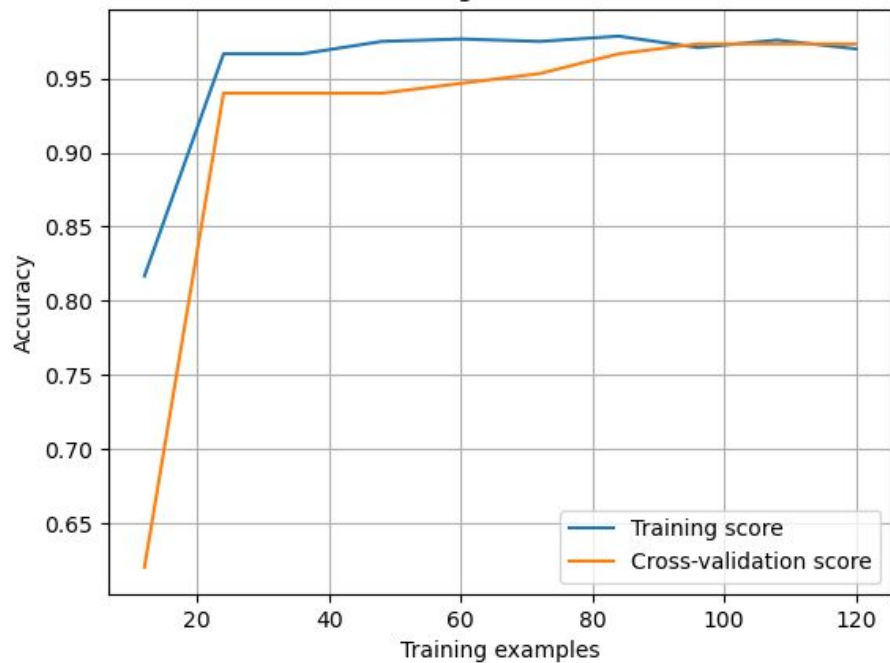- Accuracy Rate: 95%

## Random Forest Confusion Matrix
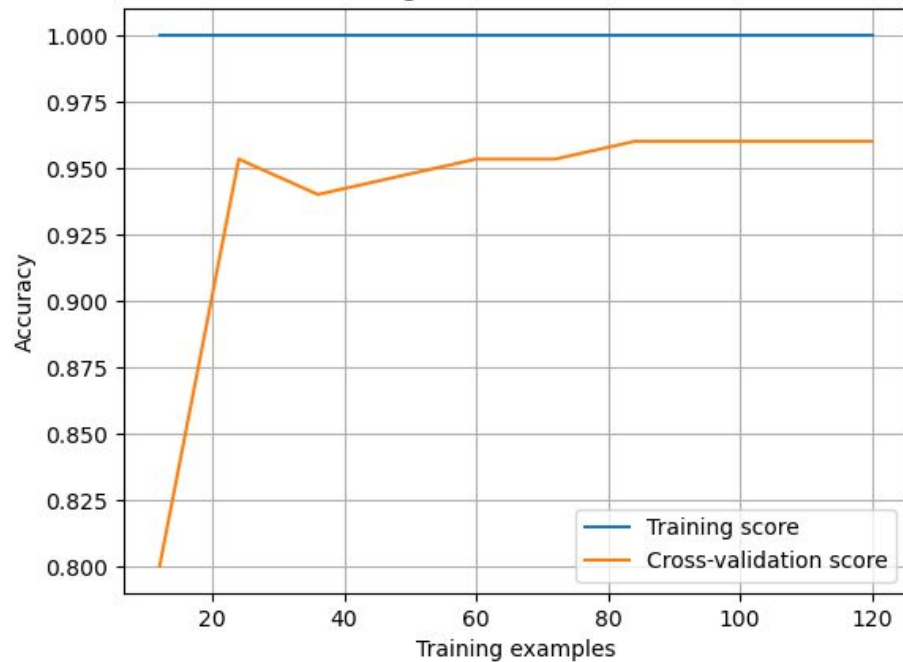
- Random Forest combines multiple decision trees and was the most accurate model
- It had the fewest mistakes and handled all sentiment classes well
- Accuracy Rate: 95%

**Learning Curve - Logistic Regression**

Accuracy vs Training examples

- Training score
- Cross-validation score

**Learning Curve - Naive Bayes**

Accuracy vs Training examples

- Training score
- Cross-validation score

Learning Curve - KNN

Learning Curve - Random Forest

# 5. Implications for Stakeholders

**Brand & Marketing Teams:**
- Real-time sentiment dashboards by country/platform
- Early warning on campaign misfires or going viral

**Customer Support / Community Managers:**
- Automatically flag negative posts for fast response
- Prioritize high-impact complaints as necessary

**Product & R&D:**
- Identify trending feature requests or areas that need to be addressed
- Fuel topic-modeling pipelines

**Policy Makers / Public Affairs:**
- Monitor public mood around social issues
- Track regional shifts in discourse for targeted outreach

# 6. Ethical, Legal, Societal Implications

**Privacy & Consent:**
- Social-media data may include private or personally identifying info
- Must respect platform TOS and, where required, anonymize or aggregate

**Bias & Fairness:**
- VADER-based labeling can misclassify sarcasm or non-English idioms
- Uneven country/platform representation may skew results
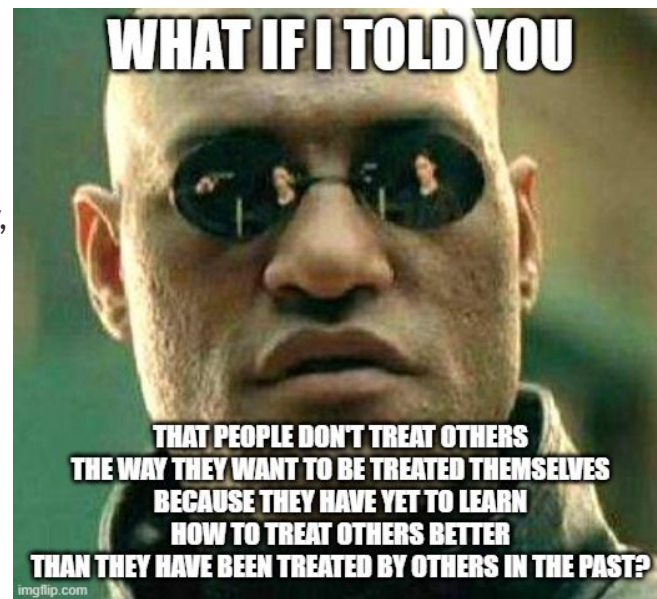
**Transparency & Accountability:**
- Stakeholders need clear documentation of model limitations
- Publish performance metrics by subgroup (country, language)

**Potential for Misuse:**
- Sentiment targeting could be weaponized for manipulation or "astroturfing"
- Require governance guardrails on automated outreach

**Legal Considerations:**
- Data-protection regulations around cross-border data flows
- Intellectual-property constraints on scraping vs. API usage

# THANK YOU!