

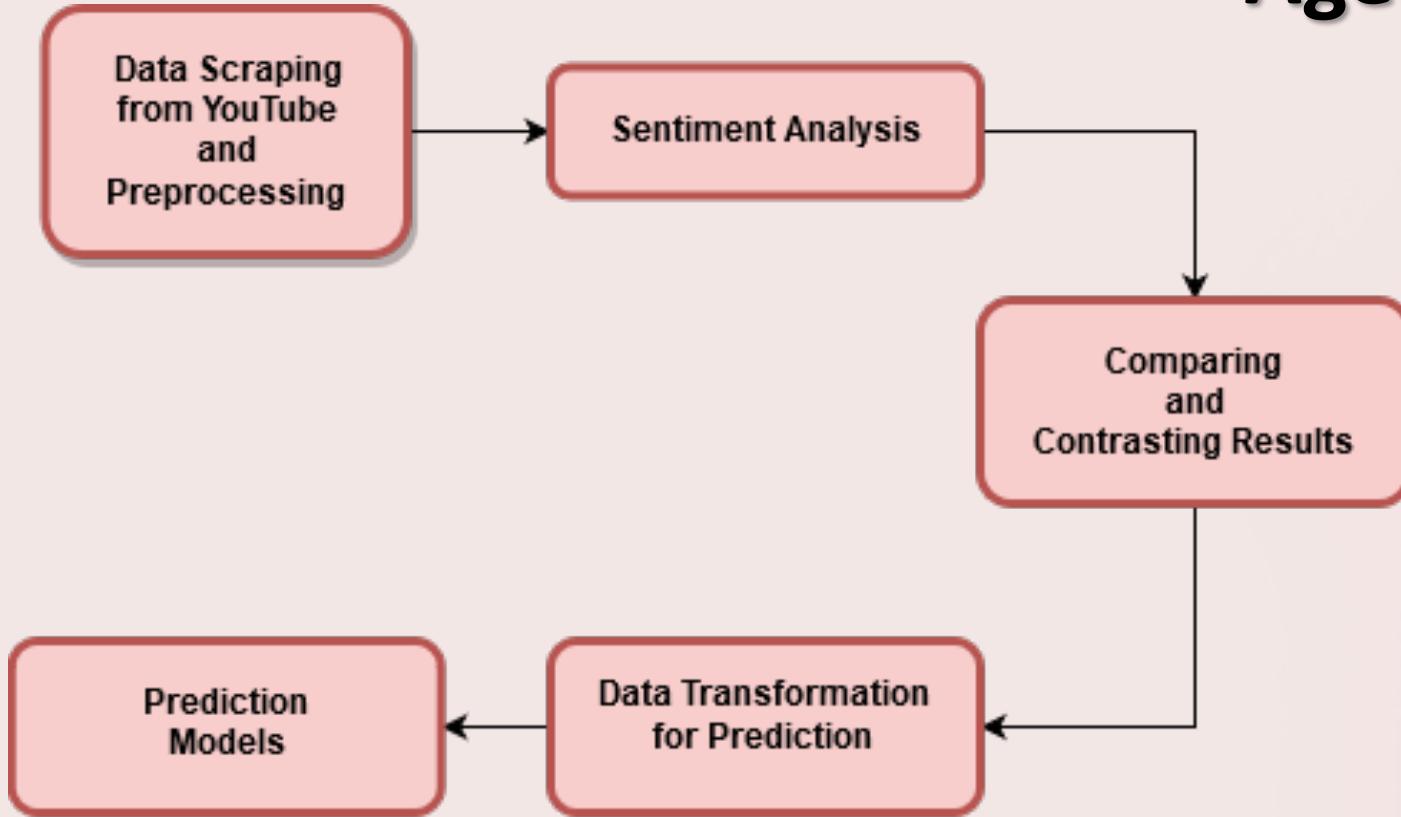


CIS-600 Social Media and Data Mining Project Presentation

YouTube Video Comment Sentiment Analysis

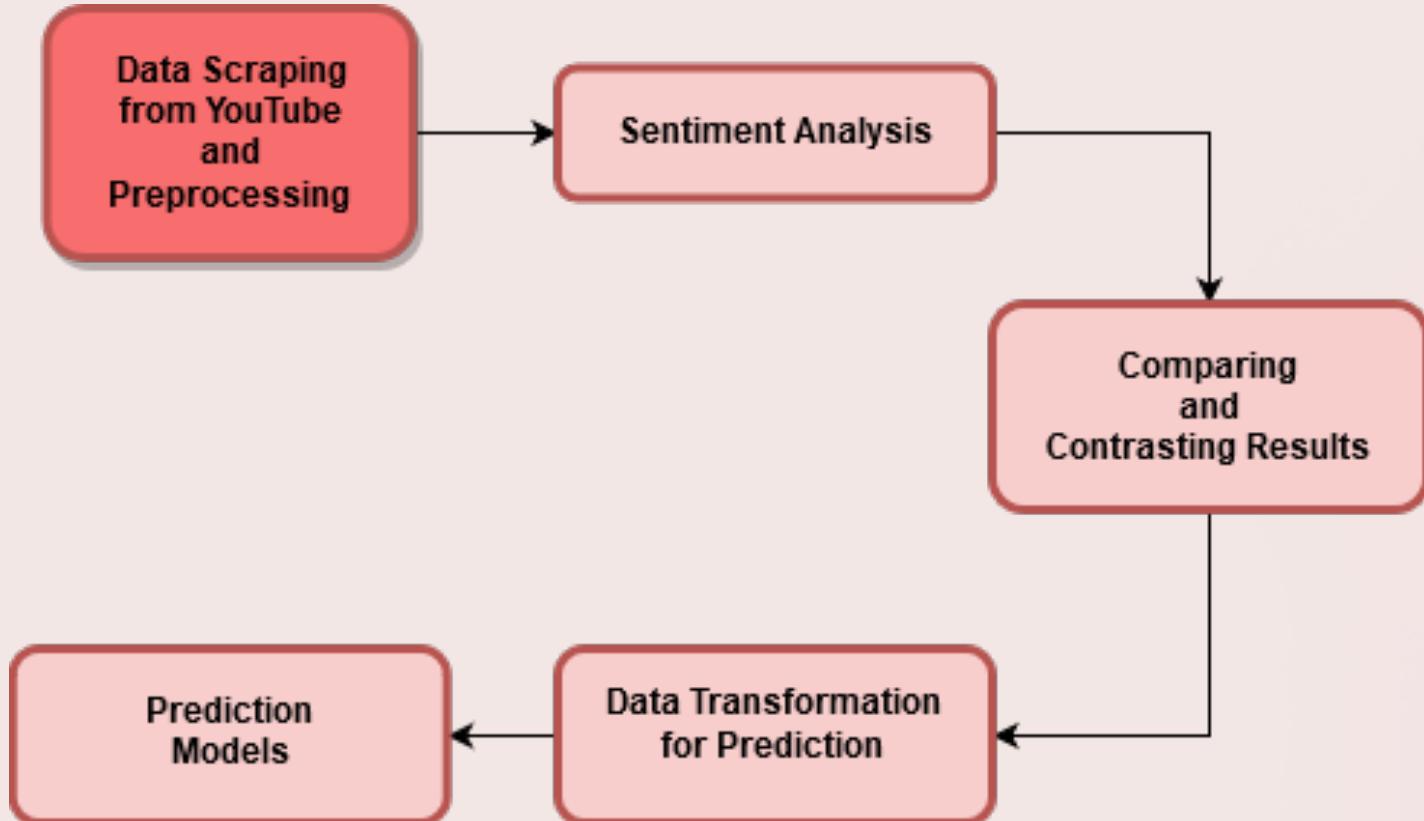
Pranay Kumar Verma
Prateek Sahu
Rohith Pattathil
Sitesh Mishra
Nitesh Nawlani

Agenda



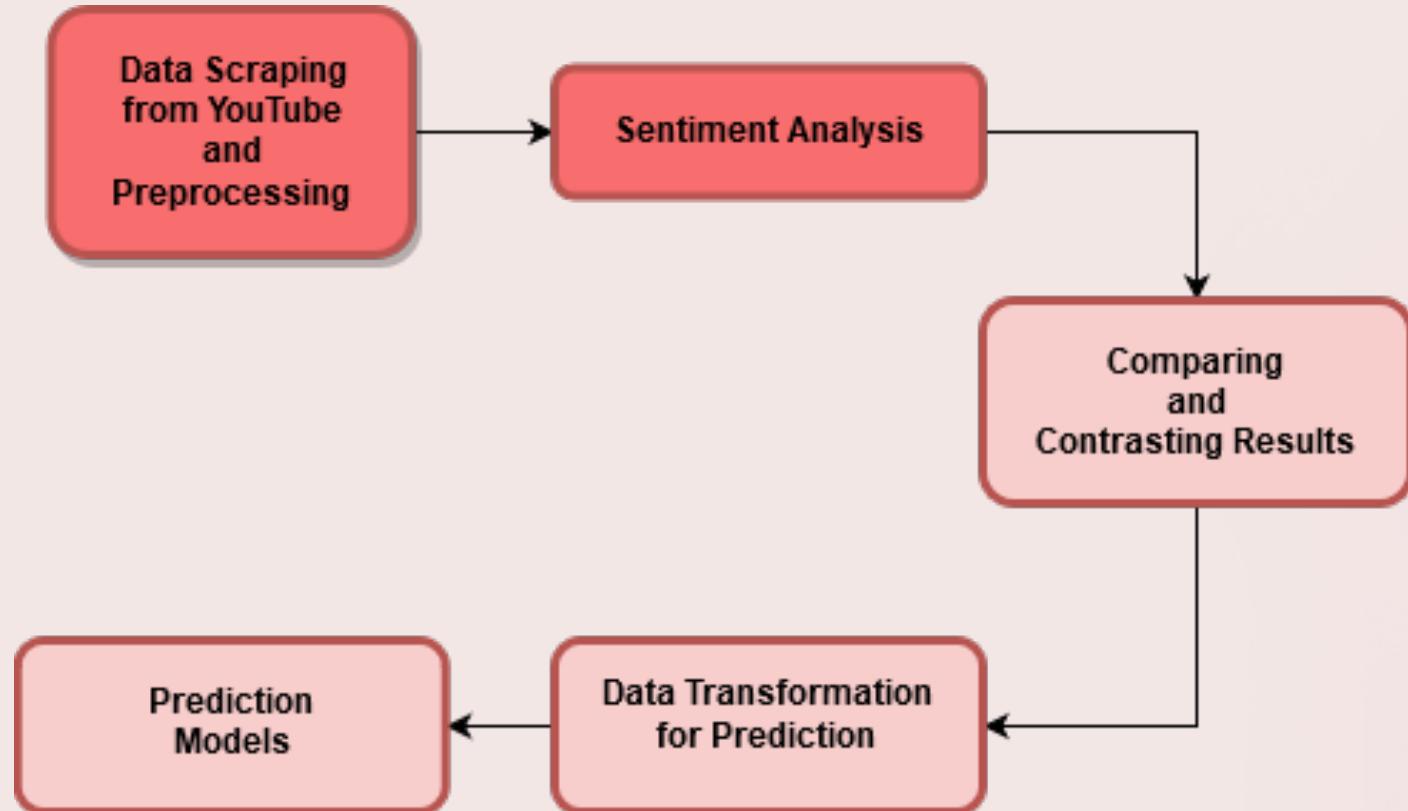
Introduction

- Founded in 2005, YouTube is the most popular platform for sharing and watching videos. The content until February 2020 amounts up to 6.5 million hours viewing (750 years).
- With the onset of the COVID-19 pandemic and restrictions in place forcing people to stay at home, YouTube witnessed a huge spike in viewing traffic, so much so that it had to reduce the default streaming quality.
- Such a huge viewing time generates exponential ad revenue. Hence, motivating many people to pursue YouTube as a career option.



All About Data

- Dataset is extracted using the YouTube Data API v3.
- The baseline data was made by extracting 1200 comments from multiple YouTube channels and manually classifying them into positive, negative and neutral comments.
- Extracted 2500 most recent comments for each of the 200 videos of a YouTube channel based on user input.
- Feature selection: Video ID, Title, Comments, their last update date, etc.
- Preprocessing the dataset includes removing punctuations, special characters, stop words, tokenizing and stemming.





Sentiment Analysis

► Models used

- **Vader** - Lexicon and rule-based sentiment analysis which is specifically attuned to sentiments expressed in social media. It works on a complete sentence analyzing text in addition to punctuations and emojis.
- **Afinn** - Builds up a score for a sentence word by word after tokenizing, removing the stop words and stemming the text.
- **NRC Lexicon** - Lexicon based approach used for its ability to grasp actual emotion rather than positive or negative.

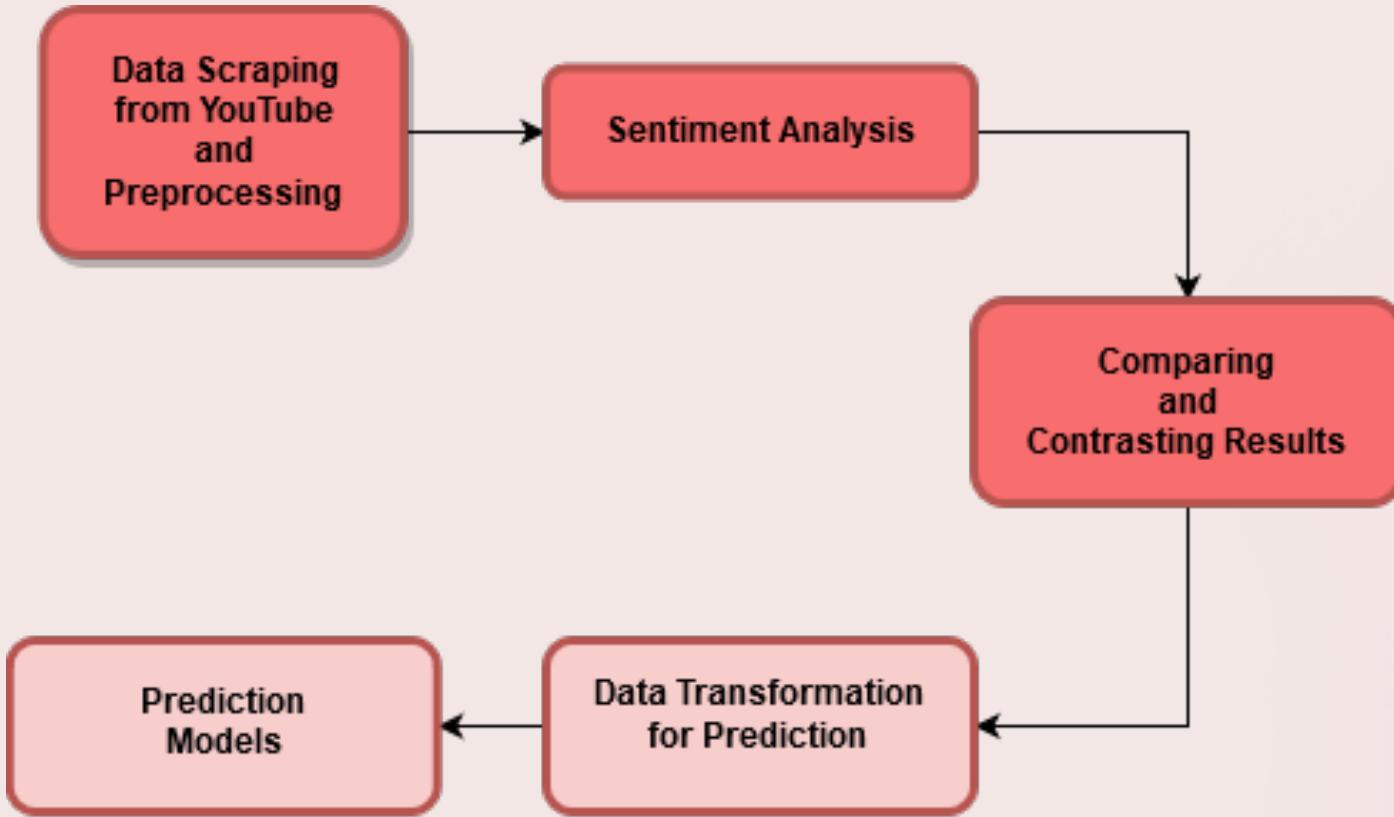


Sentiment Analysis

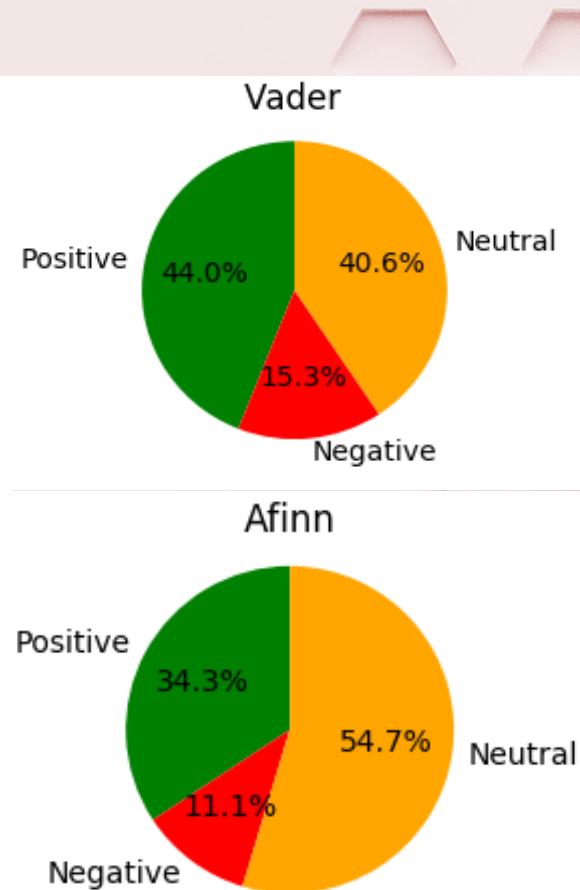
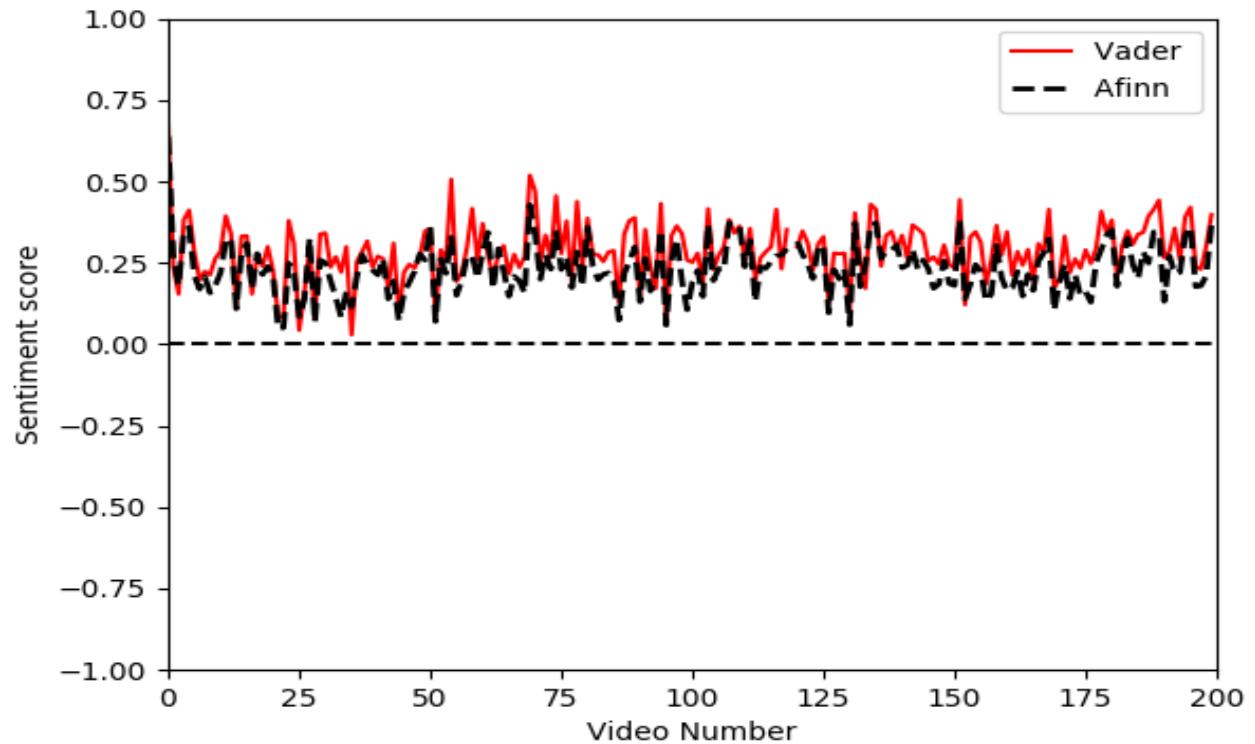
▶ Baseline accuracy:

- We selected various sentiment analysis models such as Vader and Afinn and ran them on baseline data for validation of performance.
- On analysis, Vader showed more promising results in comparison to Afinn.

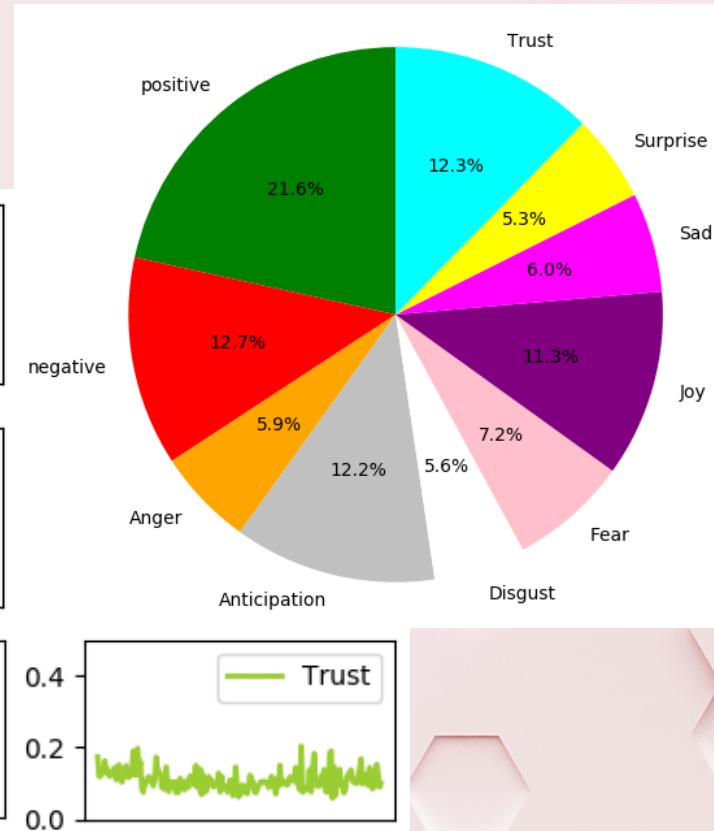
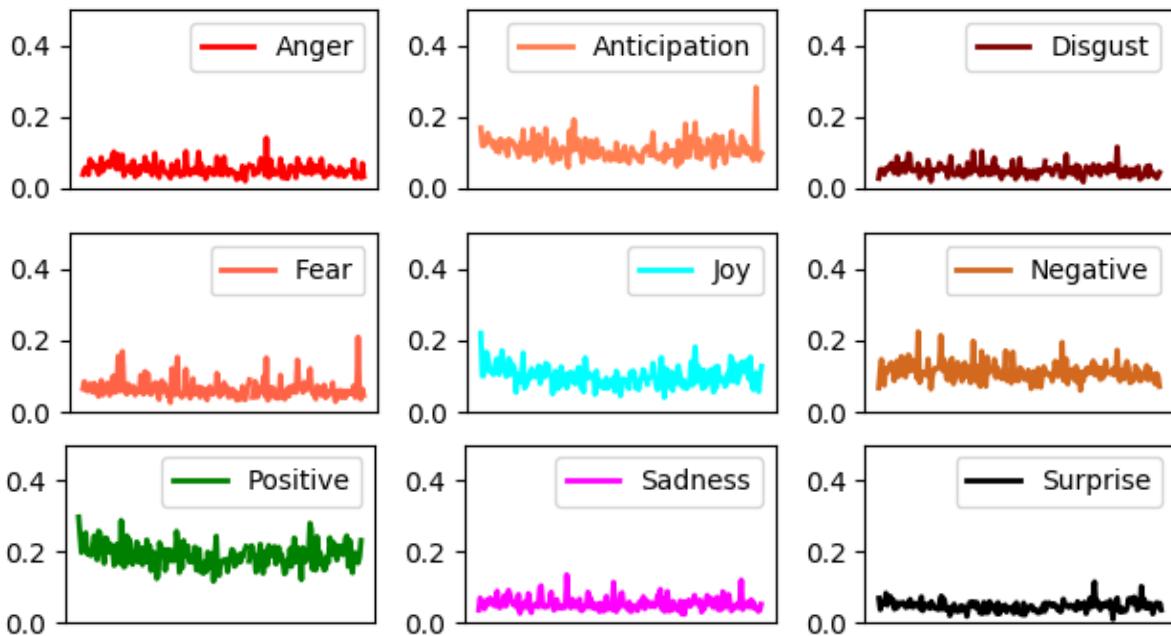
Model	Accuracy
Vader	78.79%
Afinn	70.88%



► Vader V/S Afinn



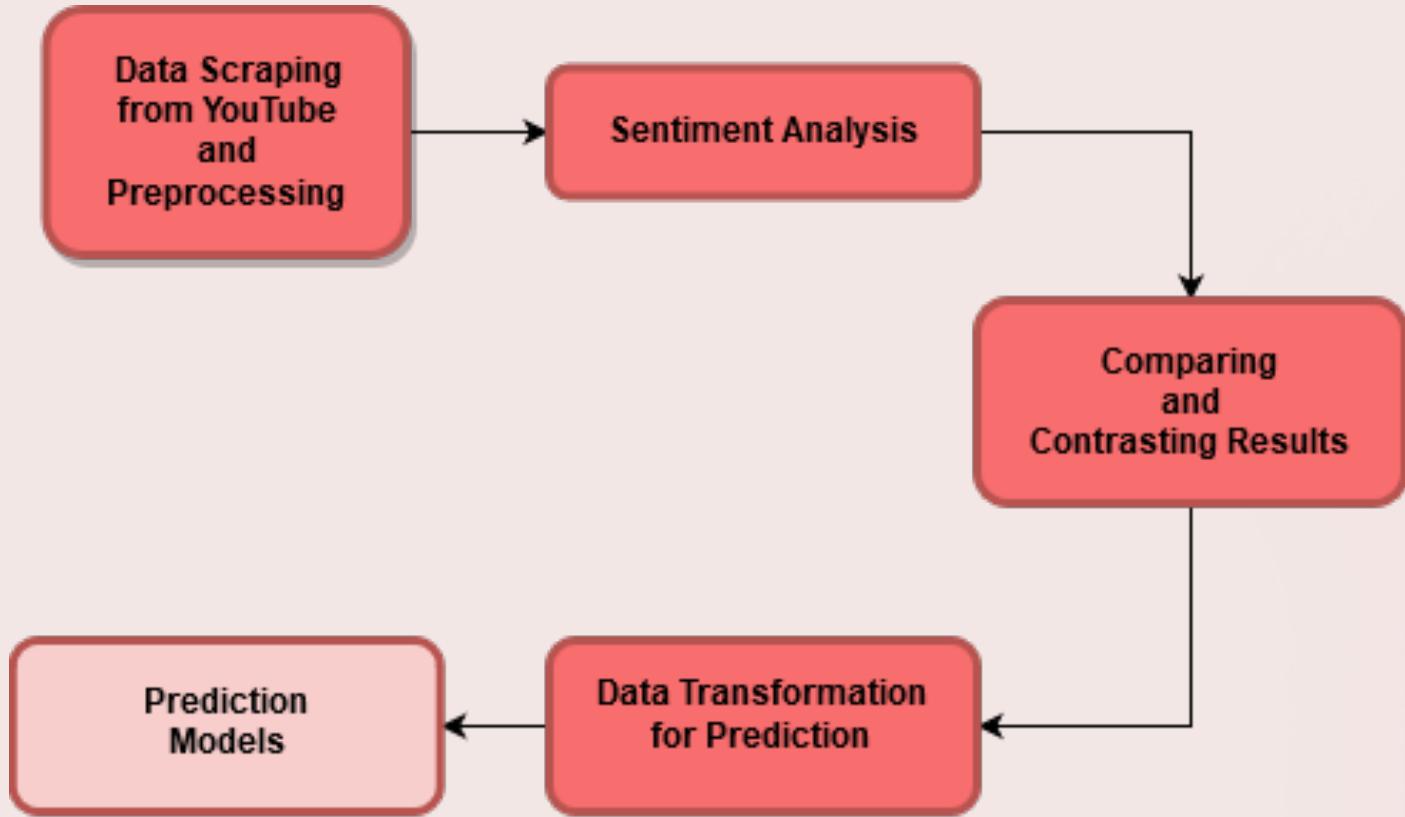
NRC Lexicon



WordCloud

The word cloud displays the following words and their approximate bounding boxes:

- look (purple)
- channel (blue)
- amazing (blue)
- awesome (blue)
- say (blue)
- Ferrari (yellow)
- time (purple)
- way (purple)
- man (green)
- something (green)
- McLaren (green)
- engine (green)
- sound (purple)
- Supra (purple)
- make (blue)
- guy (yellow)
- buy (yellow)
- house (yellow)
- seetake (blue)
- sick (green)
- watch (green)
- comment (green)
- let (green)
- please (green)
- cool (blue)
- never (green)
- come (green)
- one (purple)
- lot (green)
- said (green)
- lot (green)
- well (green)
- day (green)
- sure (green)
- drive (blue)
- best (blue)
- better (purple)
- thing (purple)
- Rolls Royce (purple)
- E63 (purple)
- BMW (purple)
- P1 (purple)
- much (blue)
- keep (green)
- video (yellow)
- know (blue)
- car (blue)
- already (green)
- people (green)
- Damn (green)
- Next (green)
- video (green)
- good (green)
- going (green)
- gon na (green)
- thought (green)
- think (green)
- got (green)
- mean (green)
- year (green)
- driving (green)
- give (green)
- go (green)
- back (green)
- still (green)
- bro (green)
- need (blue)
- go (blue)
- nice (blue)
- review (green)
- put (green)
- want (blue)
- first (blue)
- Thank (blue)
- Parker (blue)
- dude (blue)
- Lamborghini (blue)

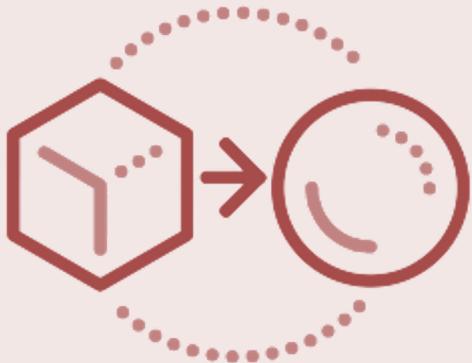


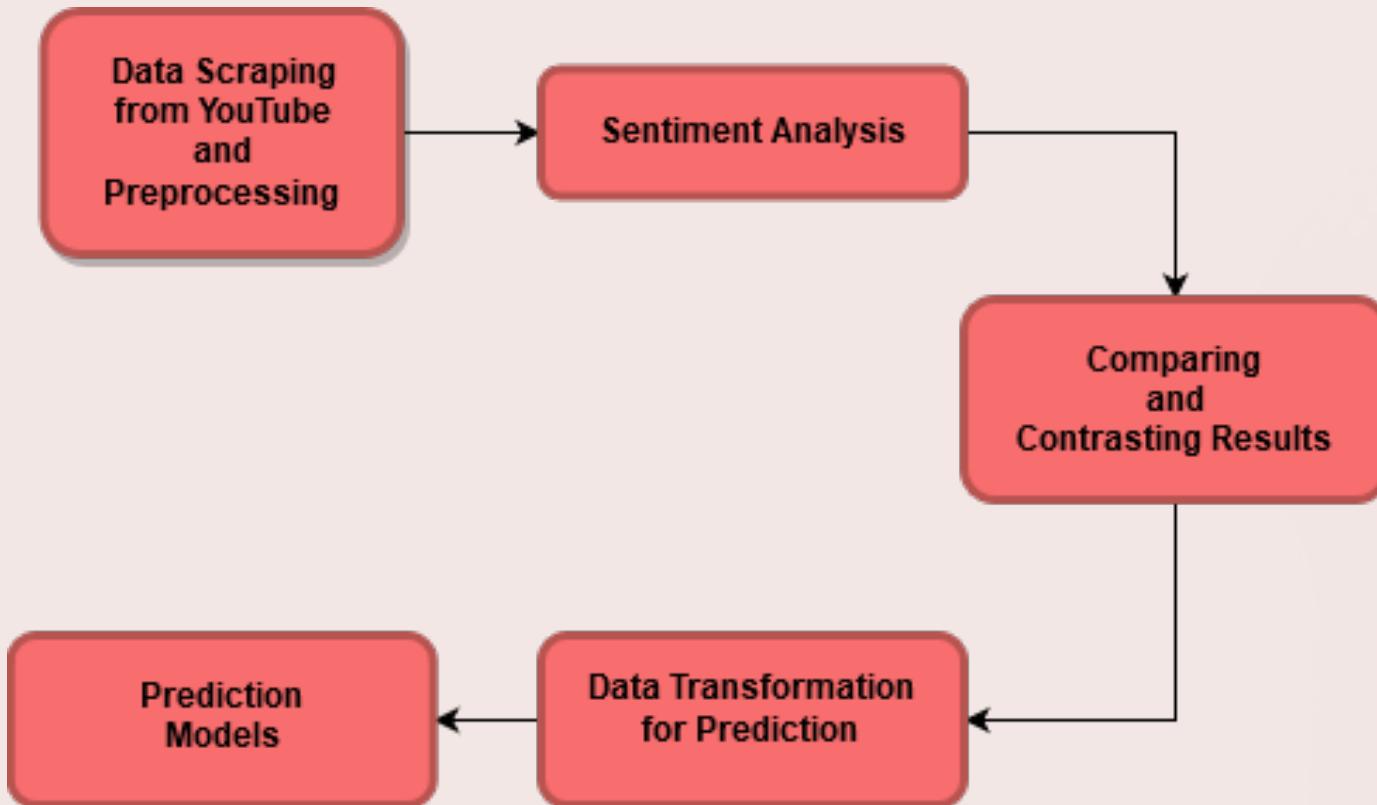


Data Transformation

Transformation to Time-Series dataset

- The comments as well as sentiment analysis scores are now transformed by grouping based on dates for creating a time-series dataset.
- This dataset is now used for performing time series prediction using simple Machine Learning algorithms as well as Neural Networks.







Prediction Models

▶ What we aim to achieve?

- The goal is to identify how the comments posted on a specific day vary in sentiment.
- Observing the past trends, we aim to make a reasonable prediction on the sentiment in the upcoming days and months.
- This in turn will help the channel owners to alter their content accordingly to keep the viewers happy.

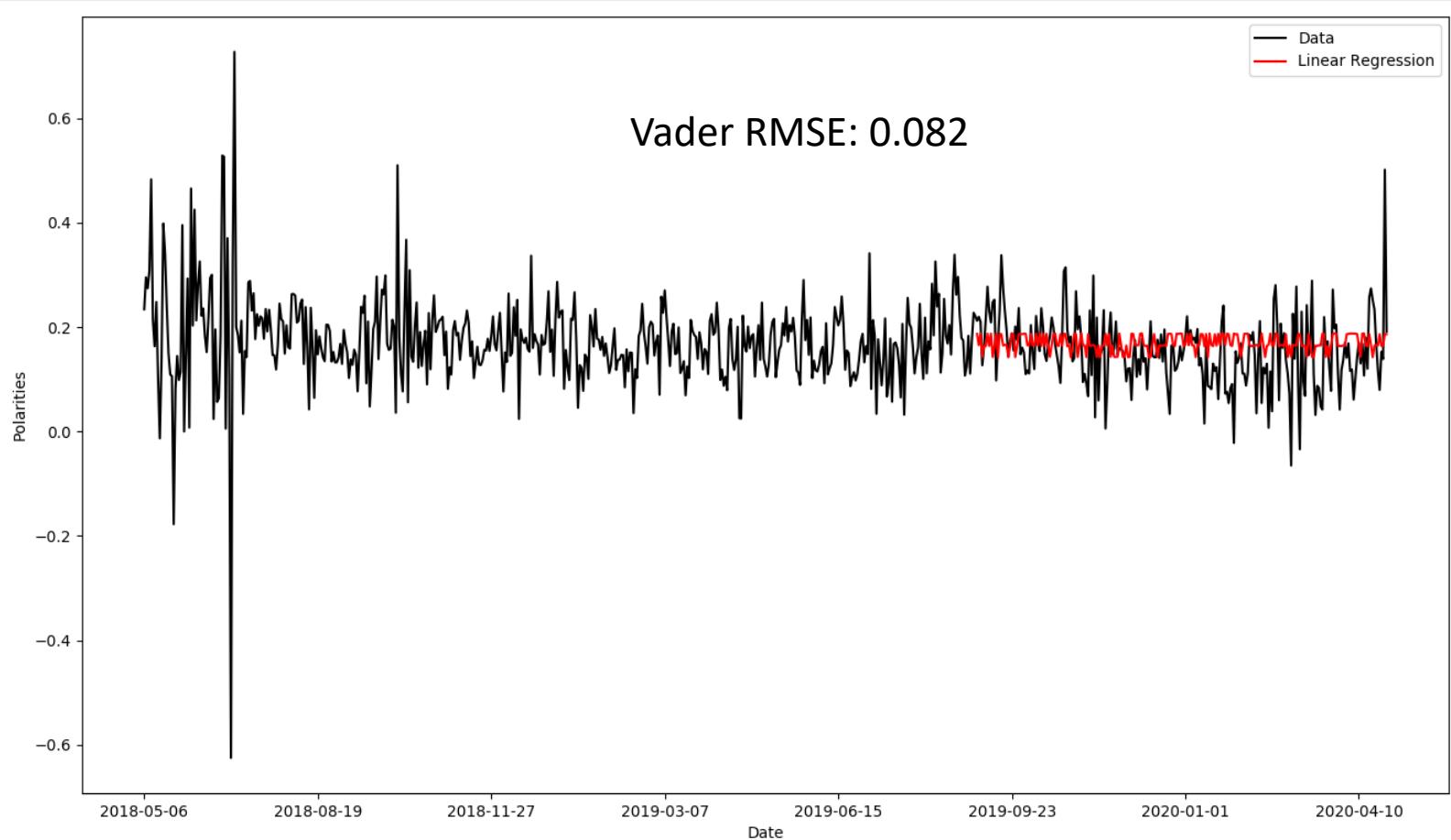


Prediction Models (Cont.)

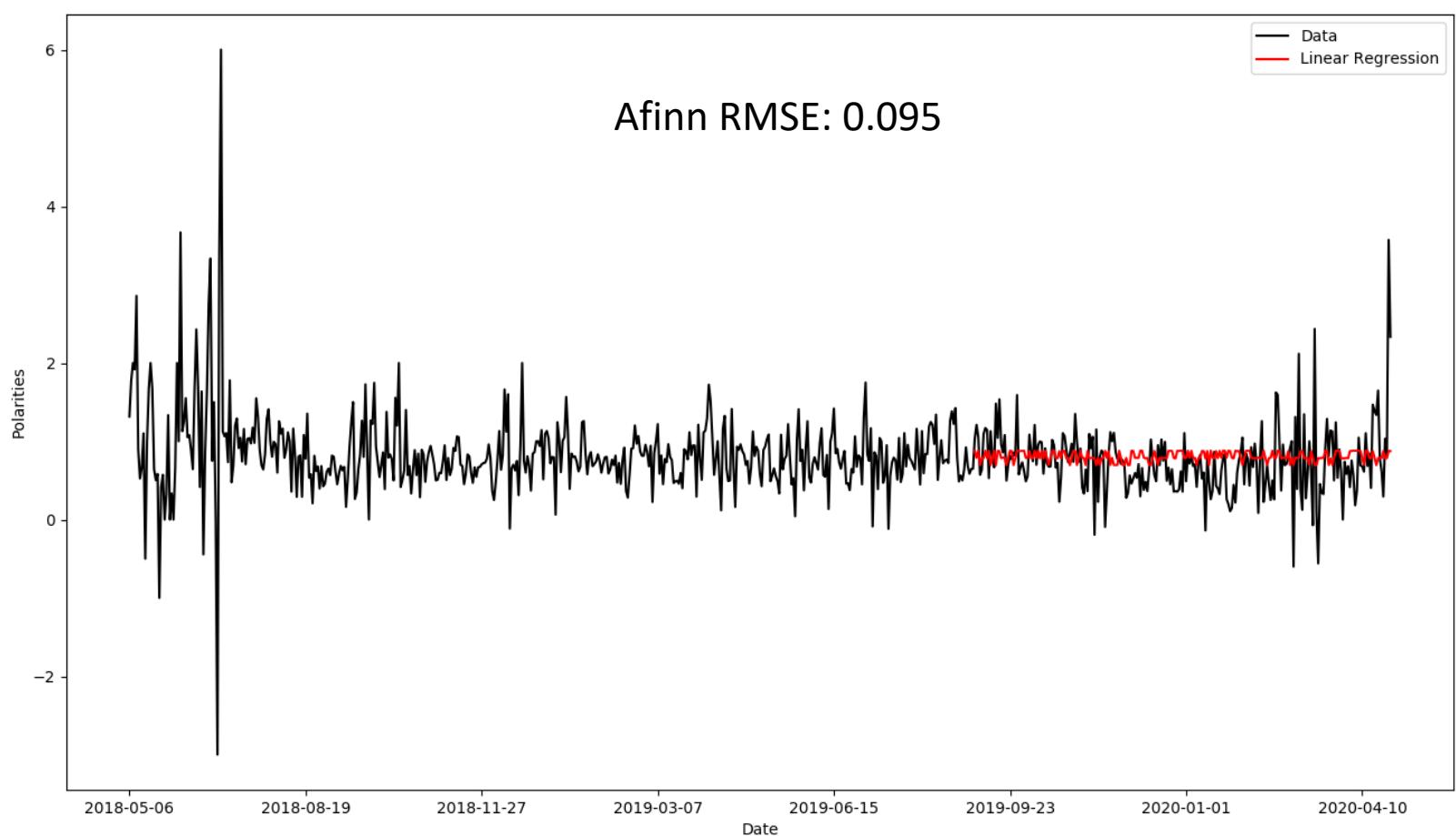
▶ Models Used

- **Linear Regression** - One of the most basic Machine Learning Algorithm, used to fit linear relations
- **Polynomial Regression** - Fits a nonlinear relationship.
- **LSTM** (Long Short Term Memory) - An advanced part of Artificial Neural Networks, which is also recurring, the previous state is preserved.

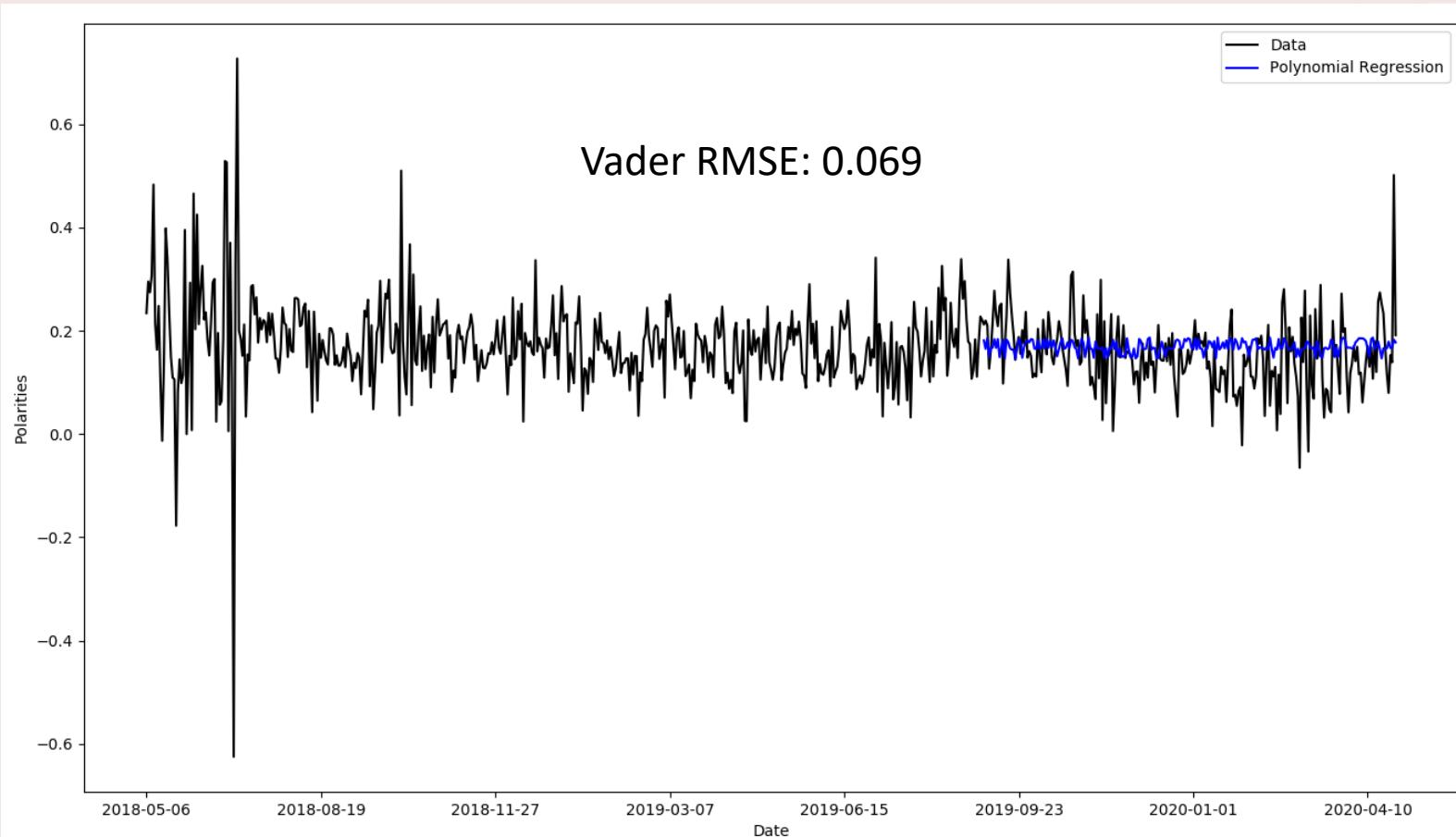
Linear Regression Vader



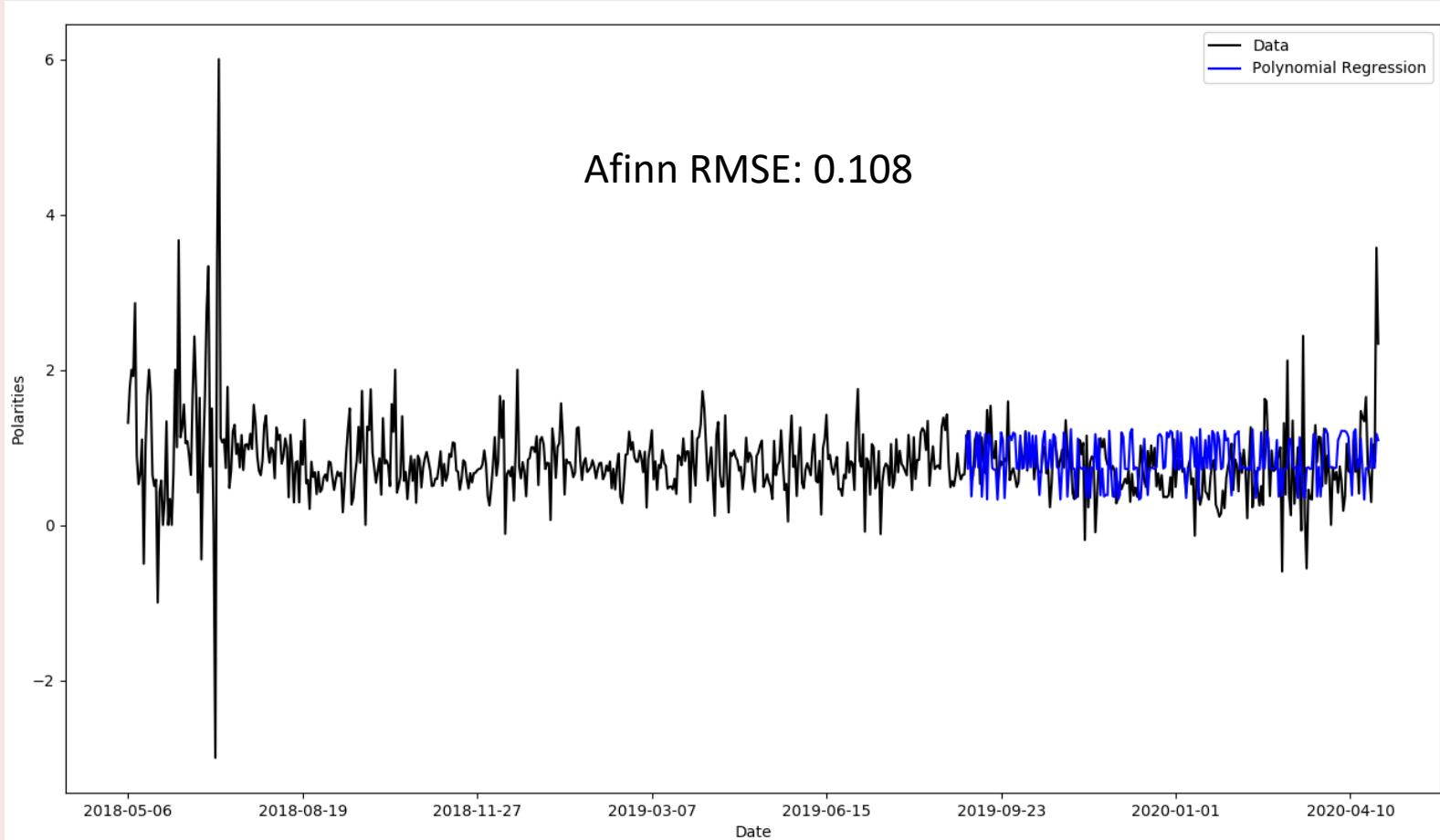
Linear Regression Afinn



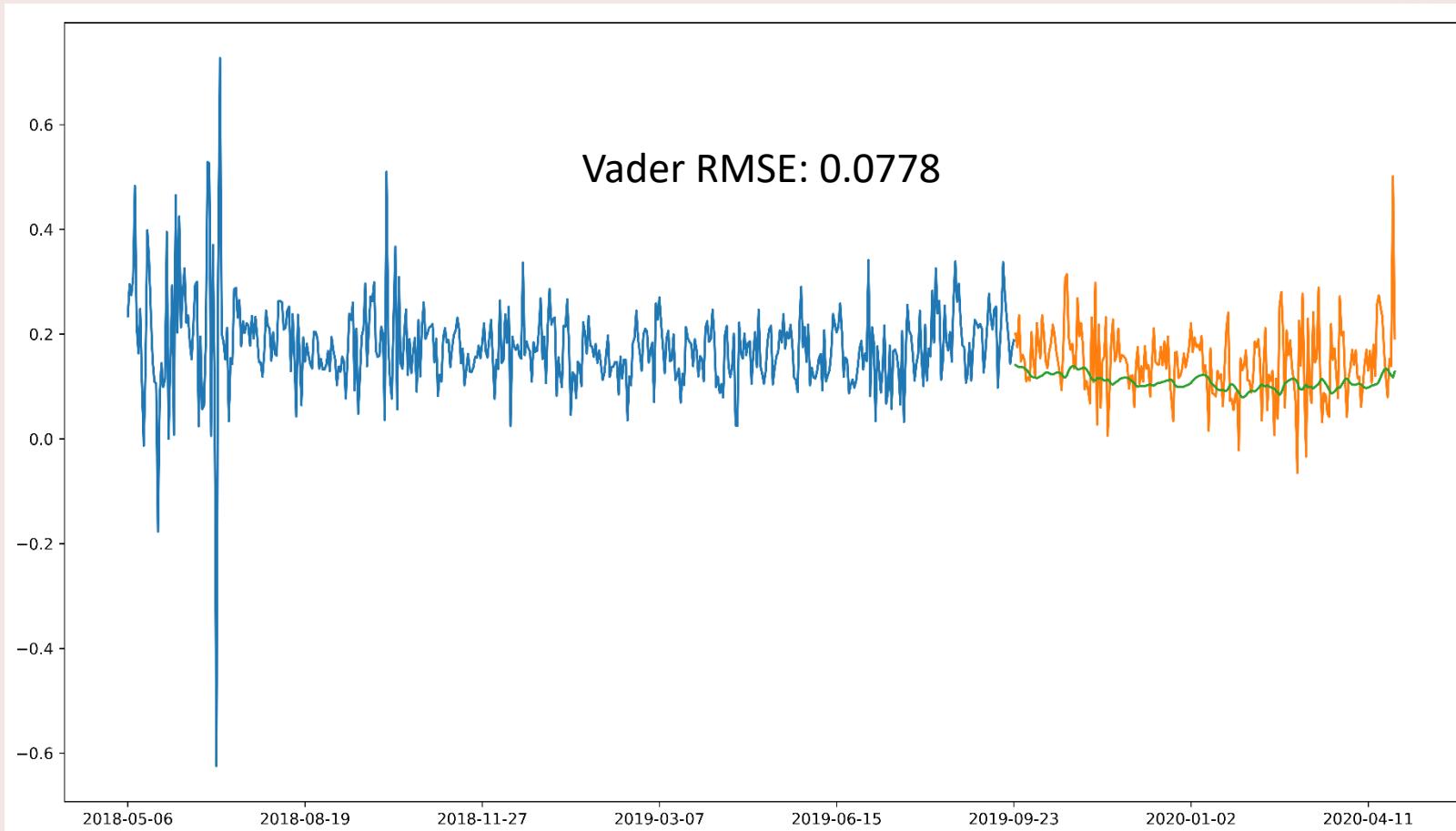
▶ Polynomial Regression Vader



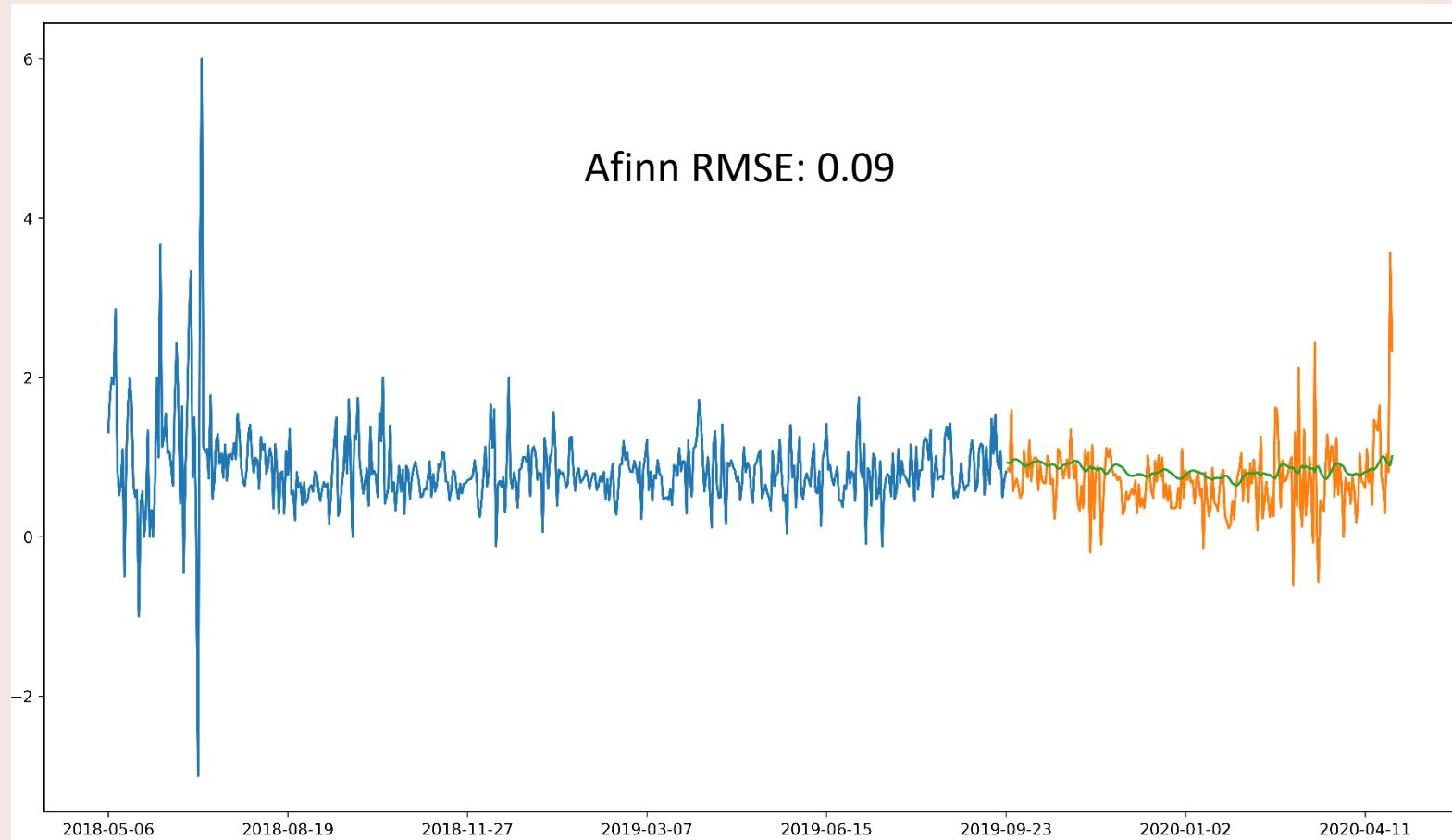
▶ Polynomial Regression Afinn



LSTM Vader



LSTM Afinn





Which Performs Best?

- As per our baseline analysis as well as the predictions performed, Vader constantly outperforms Afinn.
- The results of the model rank from best to worst are:
 1. Polynomial Regression (approx. 6.9% off the mark)
 2. LSTM (approx. 7.7% off the mark)
 3. Linear Regression (approx. 8.2% off the mark)

▶ Additional Analysis

- The common belief is that the like/dislike ratio can project how the sentiment of the comments is.
- But is this really the case?
- We used the sentiment analysis in conjunction with the view count and comment count of a video to try and predict the likes/dislikes ratio.
- What we found was interesting, in that the comments are not necessarily in line with the likes and dislikes.

▶ Future Scope

- Obtain access to YouTube Analytics API and get historical data like subscriber count, like/dislike count and average playtime which is currently not possible with data API.
- Using this, we can predict how fast the channel will grow in future.
- With average play time at our disposal we can also predict if the channel will have monetary gains.



Questions?



References

- <https://saifmohammad.com/WebPages/NRC-Emotion-Lexicon.htm>
- <https://developers.google.com/youtube/v3>
- <https://medium.com/analytics-vidhya/simplifying-social-media-sentiment-analysis-using-vader-in-python-f9e6ec6fc52f>
- <https://machinelearningmastery.com/time-series-prediction-lstm-recurrent-neural-networks-python-keras/>