# 2022 CLIMATE CHANGE BELIEF ANALYSIS

Data Science Lifecycle

# Natural Language Processing Lifecycle

**Data Exploration (EDA)**

Data Visualization
Feature Analysis

**Data Preprocessing**

Data Cleaning

**Feature Text Engineering**

Normalizing
Vectorizing
Scaling
Imbalanced Data

**Hyperparameter Tuning**

GridSearchCV
Model Evaluation

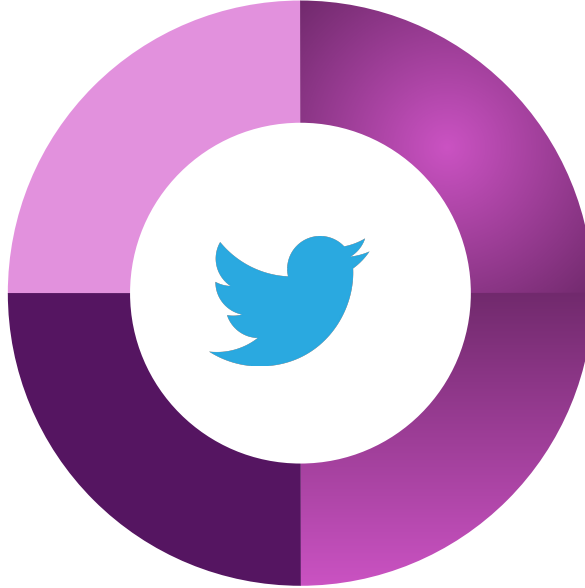**Build Model**

Visualizing
Model Performance

# Understanding Problem Statement

**Translate Language**
Machine Learning Techniques to translate language

**Climate Change**
All humans have the fundamental rights to live in a sustainable environment.

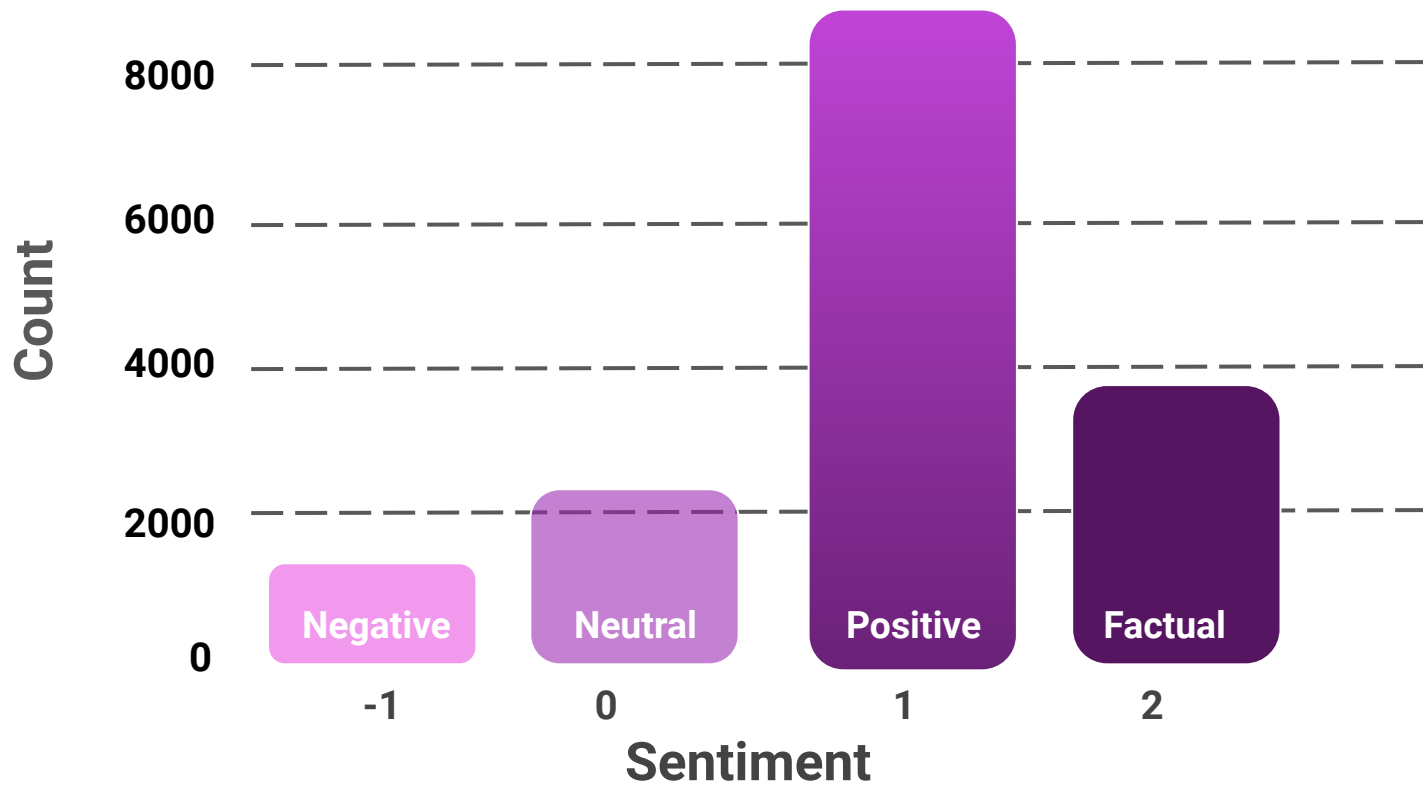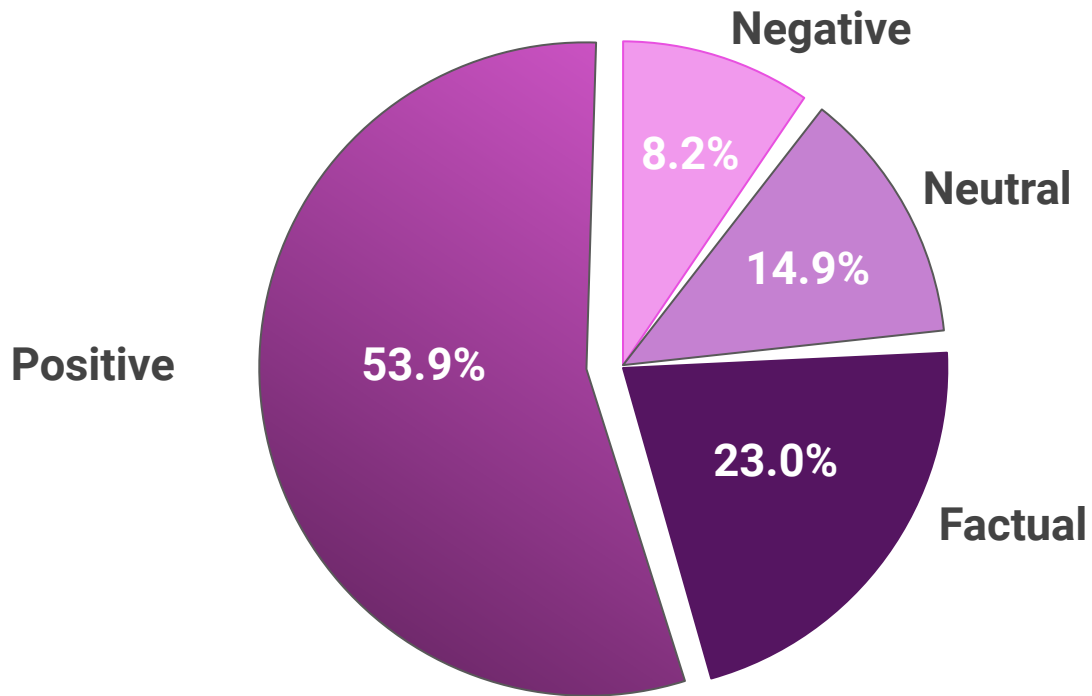**Sentimental Insights**
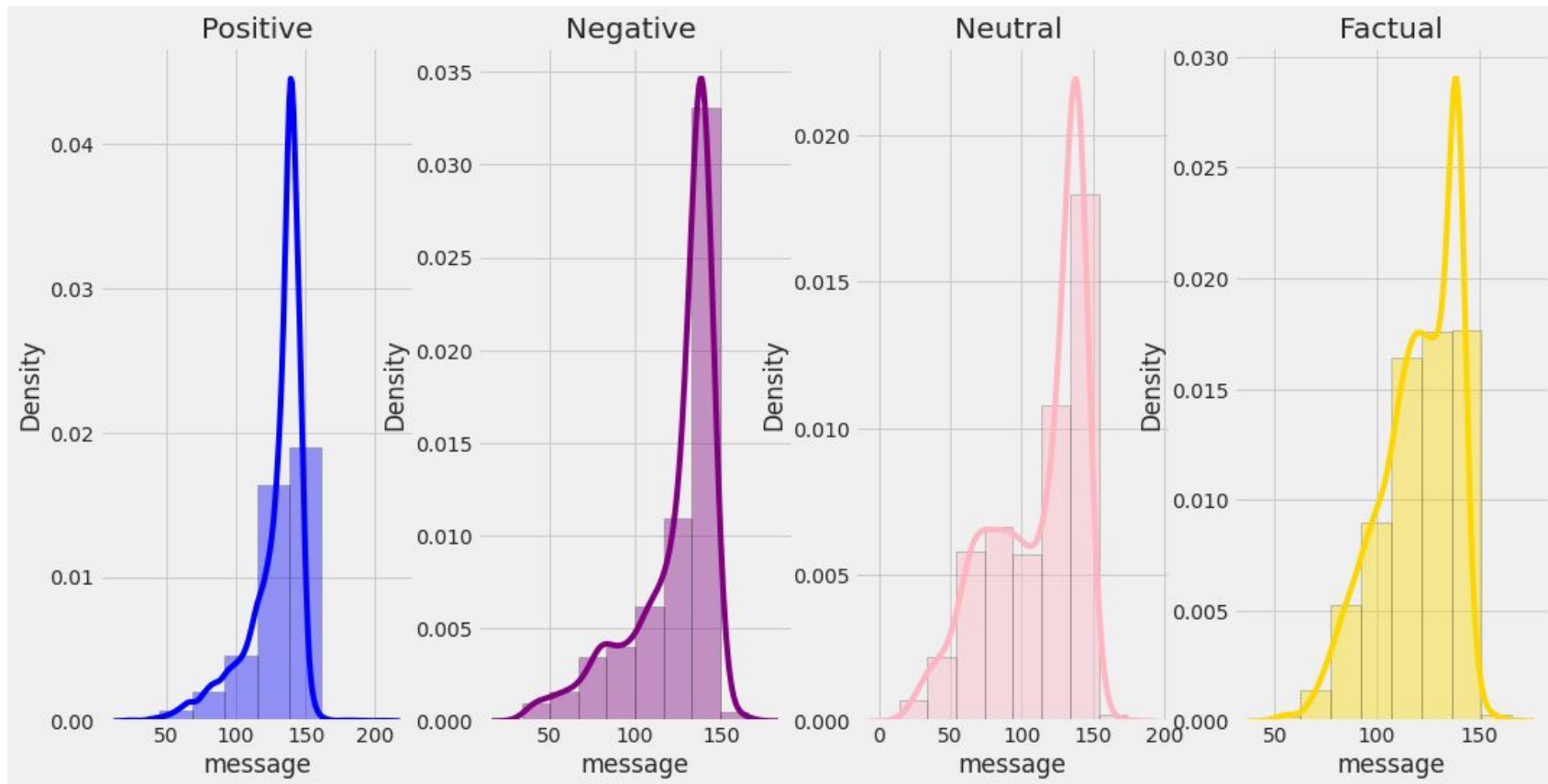Drawing insights from what people have twitted about climate change

Number of Messages Per Sentiment

# Percentage of Messages per Sentiment
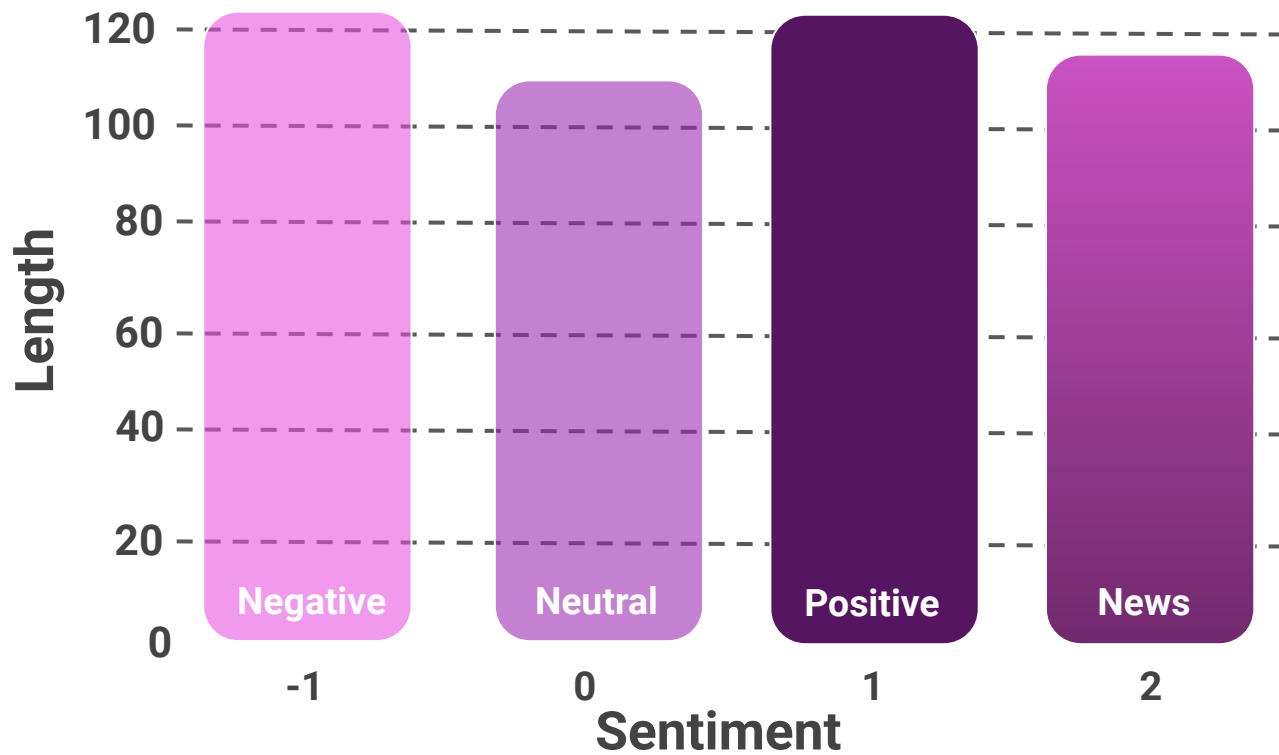
# Distribution of Length per Label

# Average Length of Messages by Sentiment



**Length**

120

100

80

60

40

20

0

Negative    Neutral    Positive    News

-1    0    1    2

**Sentiment**

# Data Preprocessing

**Data Cleaning**

**Tokenize**

**Stopwords**

**Stemming**

**Lemmitizing**

Remove Noise

Tokens of strings

English Stopwords

Roor meaning of a word

Returning the dictionary meaning of a word

# Data Preprocessing

**Characters**

Hashtags
URLs
Retweets
Mentions
Emojis

**Punctuations**

Remove [!*&%_-?\/.,<,>]

**Data Cleaning**

**Numbers**

Remove Numbers

**Lower char**

Transform to
lowercase characters
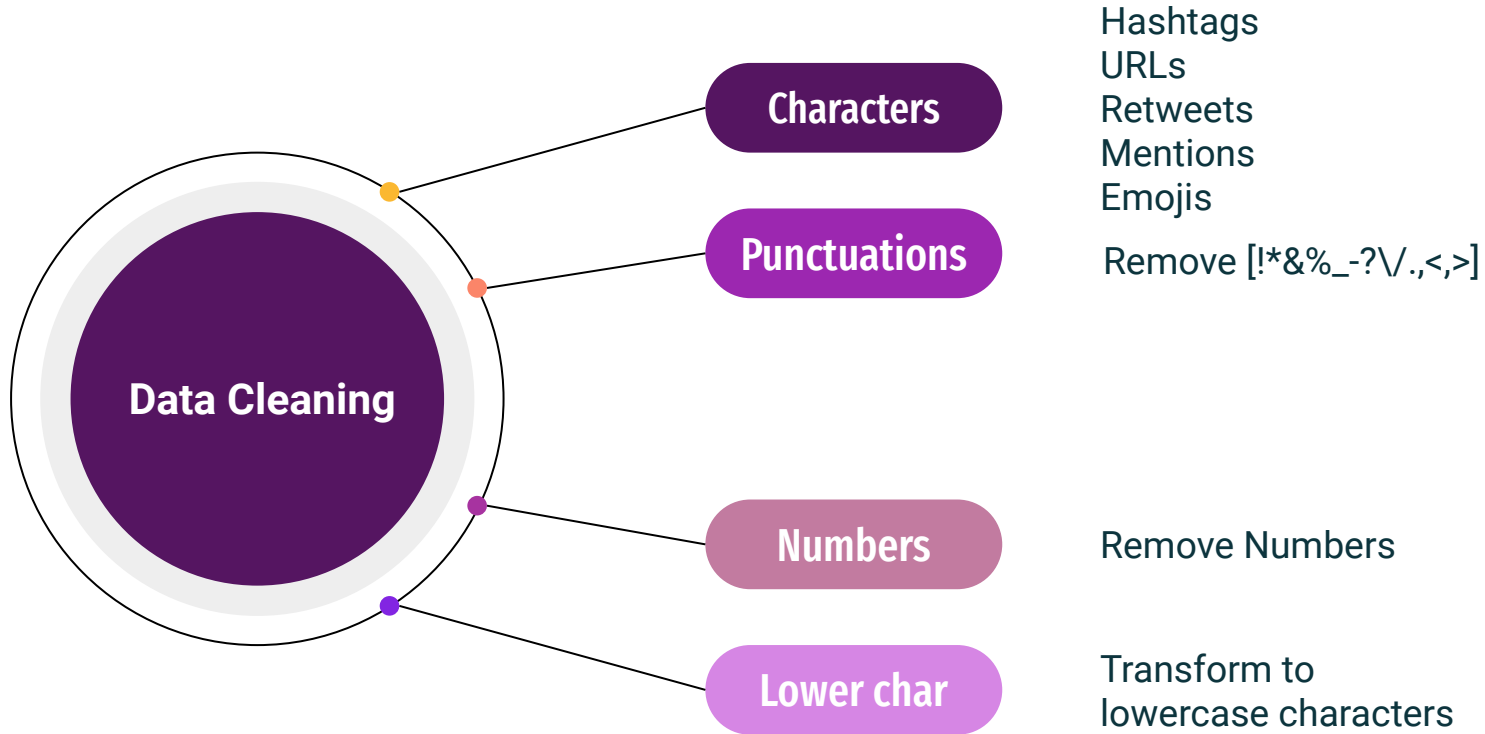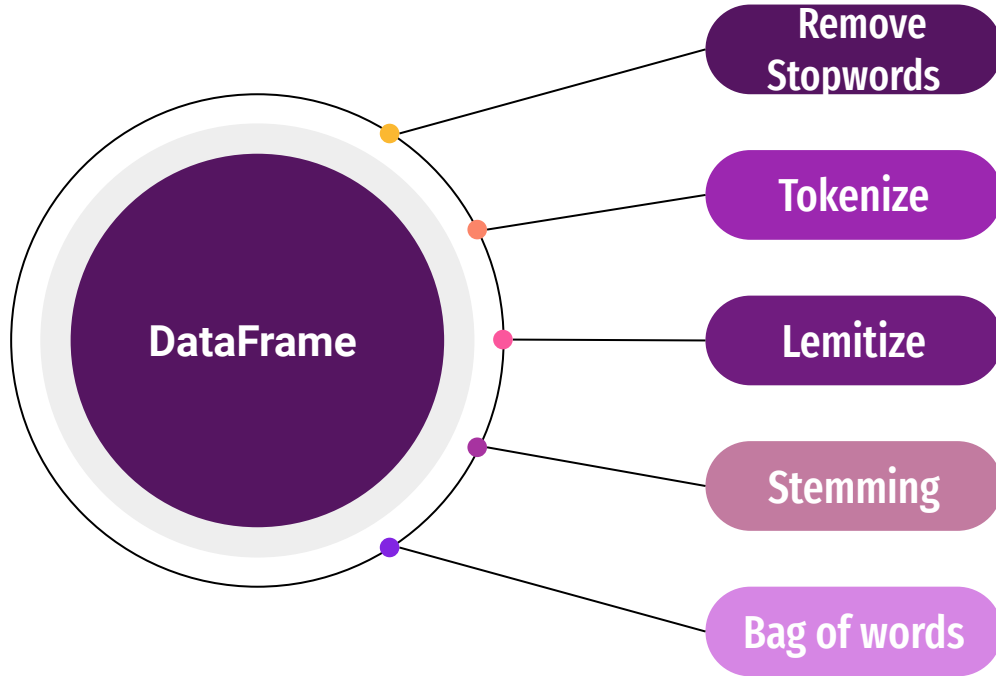
# Data Preprocessing Example

'PolySciMajor EPA chief doesn't think carbon dioxide is main cause of global warming and.. wait, what!? https://t.co/yeLvcEFXkC via @mashable'



**Data preprocessor**
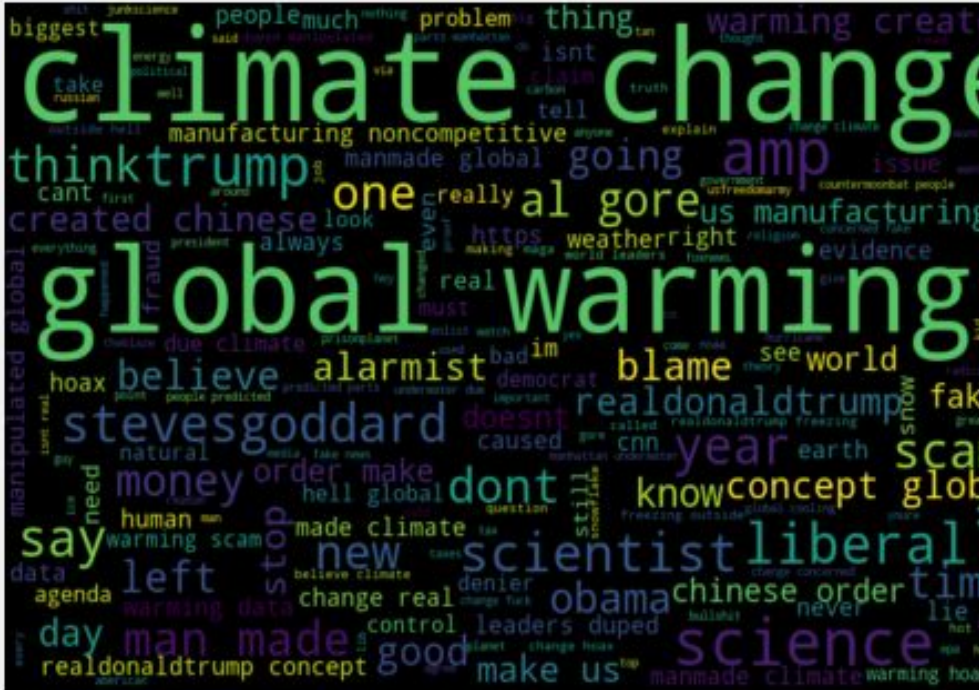
polyscimajor epa chief doesnt think carbon dioxide is main cause of global warming and wait what urlweb via mashable

Data Normalization

DataFrame

- Remove Stopwords
- Tokenize
- Lemitize
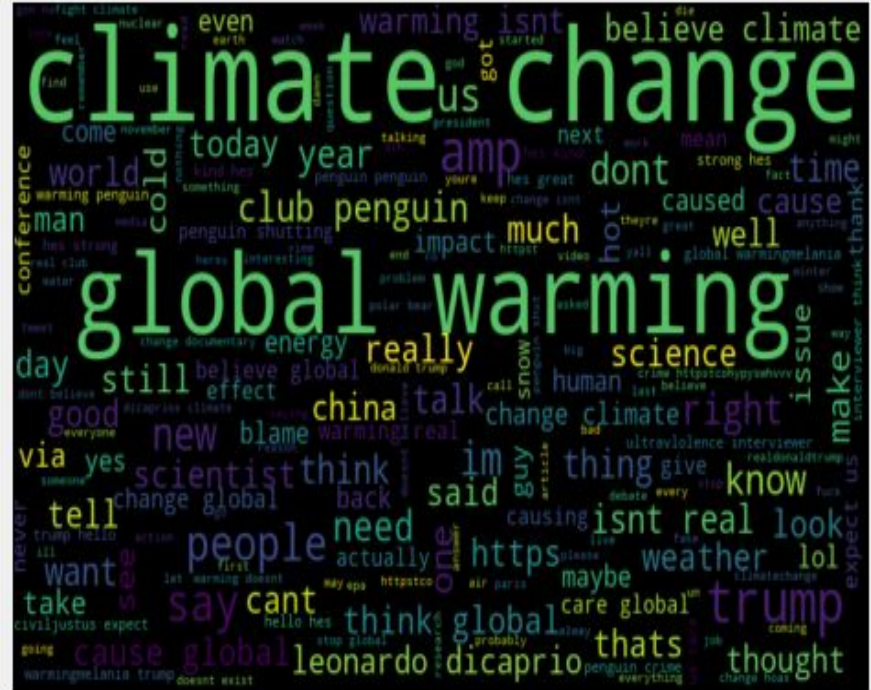- Stemming
- Bag of words

# Further Text Exploration using wordcloud

**Class -1**

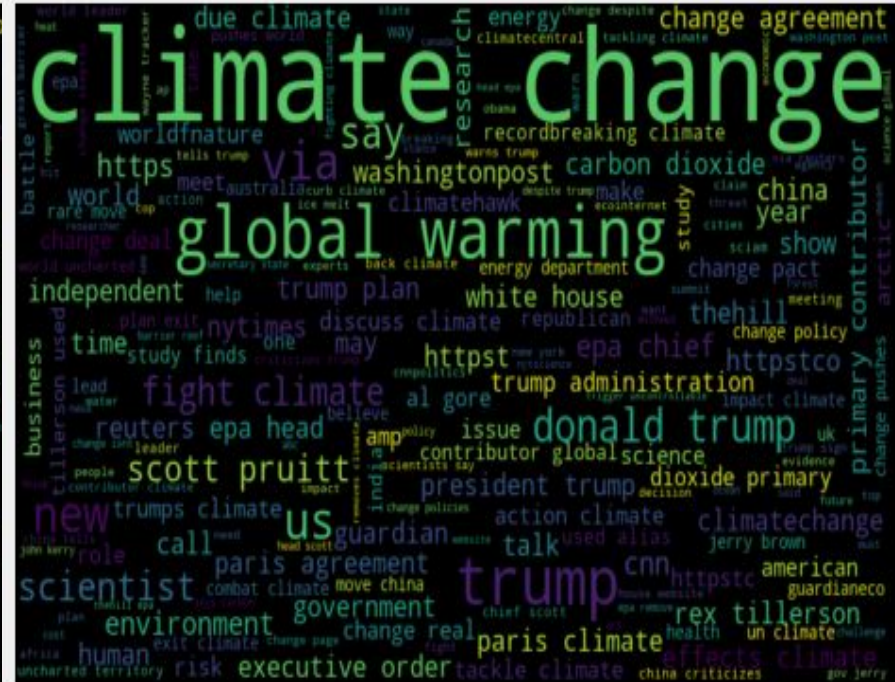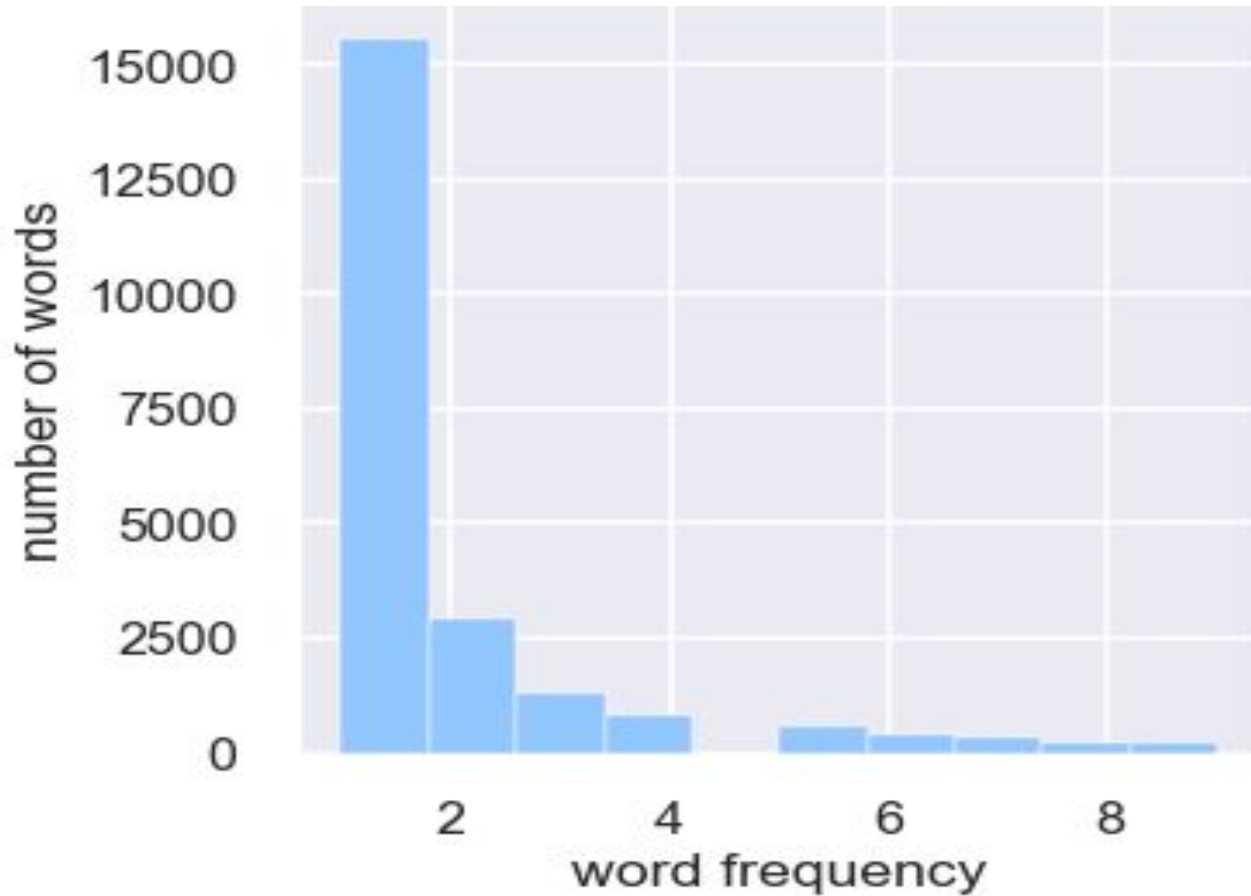**Class 0**

# Further Text Exploration using wordcloud
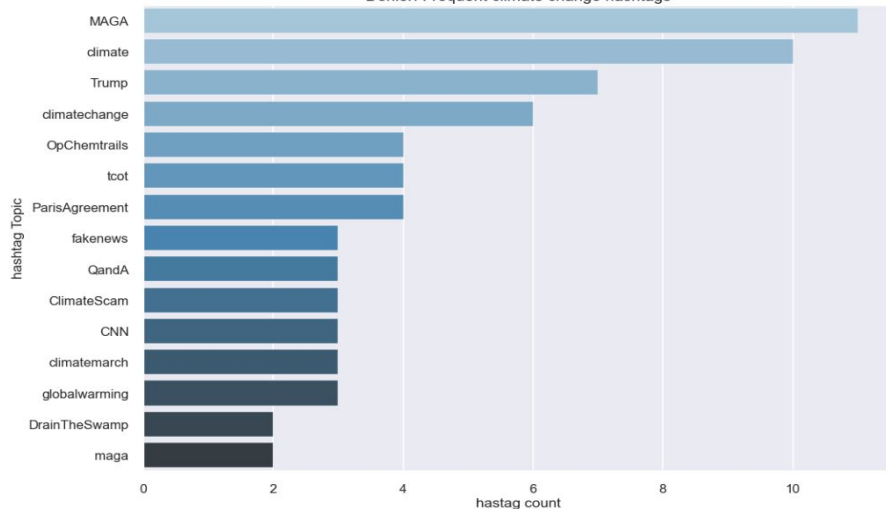
**Class 1**

**Class 2**

# Distribution of words Appearing < 10

# Frequent HashTags for all the classes

Denier: Frequent climate change hashtags

| hashtag Topic | hastag count |
|---|---|
| MAGA | 11 |
| climate | 10 |
| Trump | 7 |
| climatechange | 6 |
| OpChemtrails | 4 |
| tcot | 4 |
| ParisAgreement | 4 |
| fakenews | 3 |
| QandA | 3 |
| ClimateScam | 3 |
| CNN | 3 |
| climatemarch | 3 |
| globalwarming | 3 |
| DrainTheSwamp | 2 |
| maga | 2 |

Neutral: Frequent climate change hashtags

| hashtag Topic | hastag count |
|---|---|
| climate | 16 |
| climatechange | 11 |
| Trump | 11 |
| ClimateChange | 4 |
| amreading | 4 |
| BeforeTheFlood | 4 |
| QandA | 3 |
| ParisAccord | 3 |
| COP22 | 3 |
| Tillerson | 2 |
| ICYMI | 2 |
| China | 2 |
| firstdayofspring | 2 |
| GlobalWarming | 2 |
| AprilFoolsDay | 2 |

Believer: Frequent climate change hashtags

| hashtag Topic | hastag count |
|---|---|
| climate | 187 |
| BeforeTheFlood | 128 |
| climatechange | 94 |
| ImVotingBecause | 65 |
| COP22 | 62 |
| ParisAgreement | 50 |
| ActOnClimate | 42 |
| Ã | 35 |
| Trump | 33 |
| IVotedBecause | 32 |
| globalwarming | 23 |
| ClimateChange | 22 |
| environment | 22 |
| auspol | 22 |
| BeforetheFlood | 16 |

News: Frequent climate change hashtags

| hashtag Topic | hastag count |
|---|---|
| climate | 130 |
| environment | 44 |
| climatechange | 42 |
| Trump | 27 |
| news | 20 |
| p2 | 14 |
| COP22 | 13 |
| science | 13 |
| ClimateChange | 12 |
| GreatBarrierReef | 11 |
| News | 10 |
| ParisAgreement | 9 |
| climatemarch | 7 |
| CLIMATEchange | 6 |
| China | 6 |

# Data Preprocessing

192
Rear words

192
Short Words

# Experiments

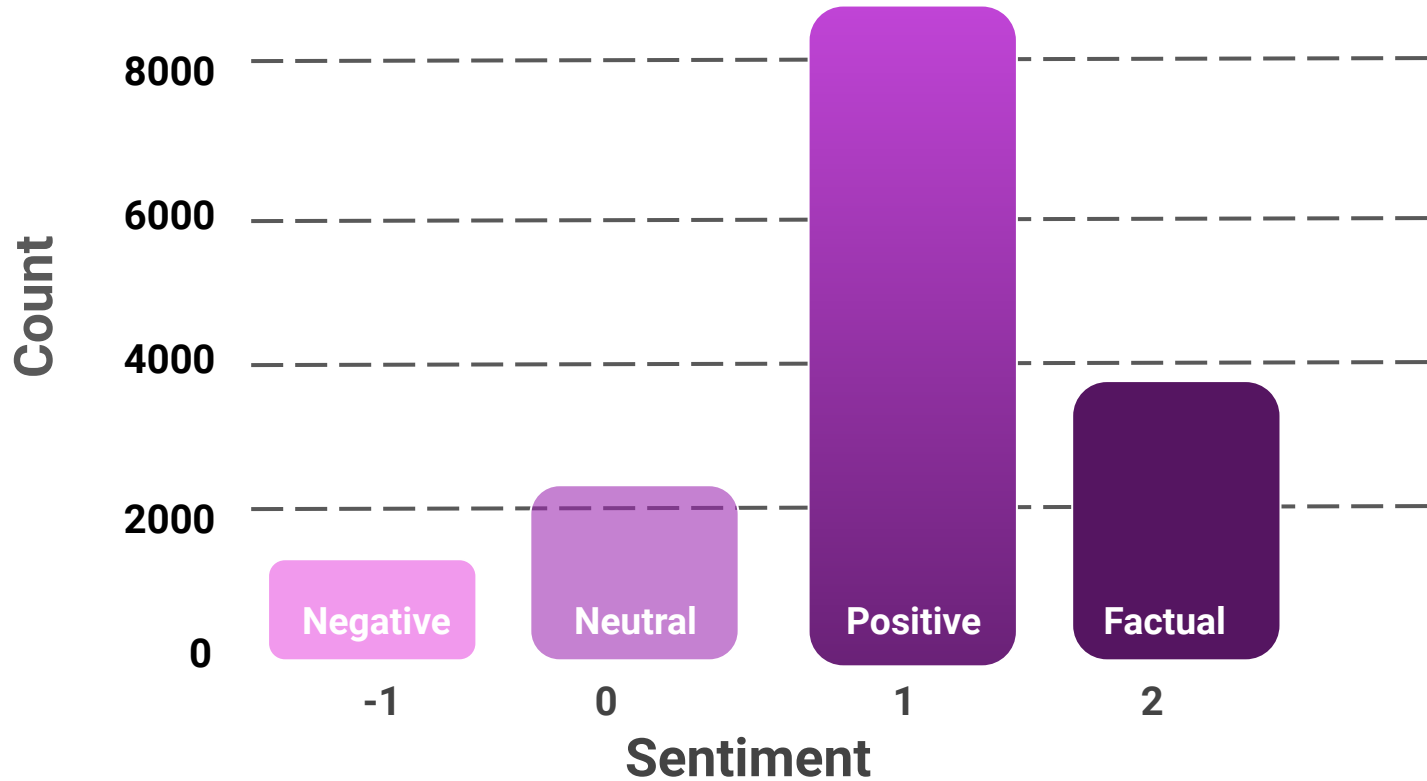LESS NOISE **VS** CLEAN SENTENCES
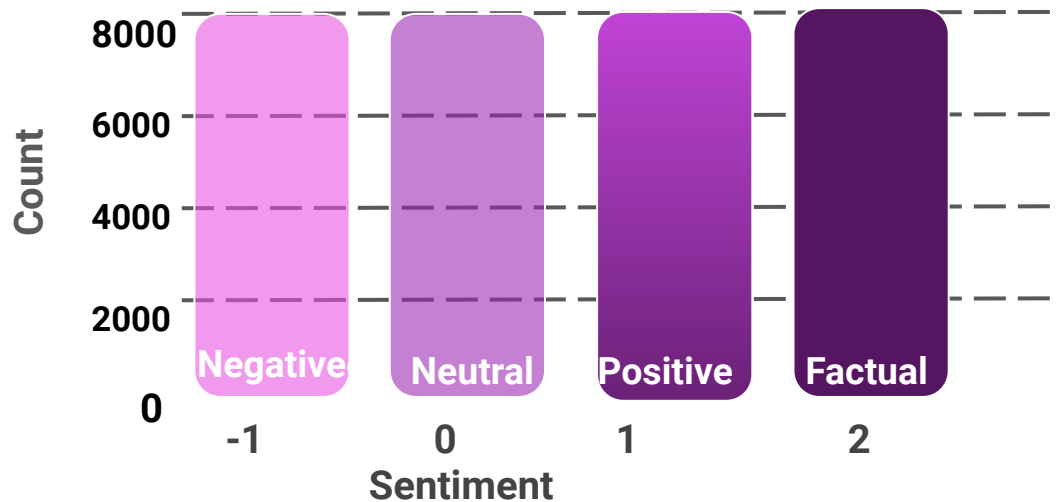
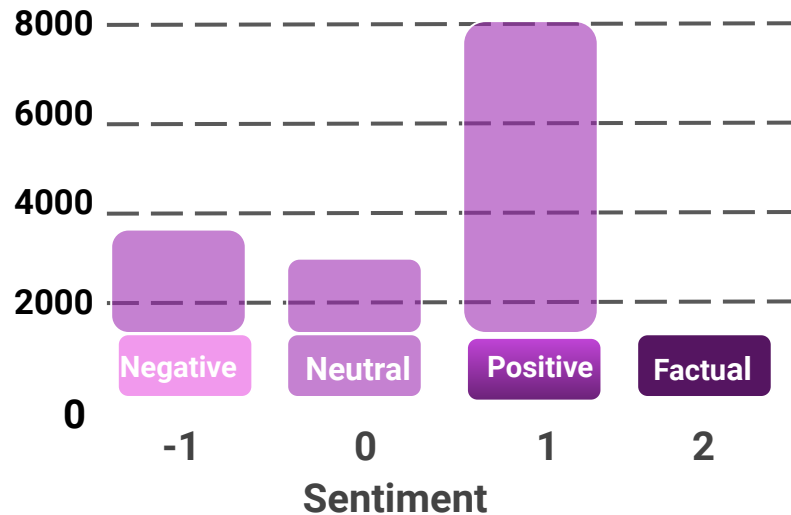CLEAN SENTENCES **VS** STEMMED CLEAN SENTENCES **VS** LEMMATIZED CLEAN SENTENCES

# Text Imbalance

- Models become better at predicting one class over the others.
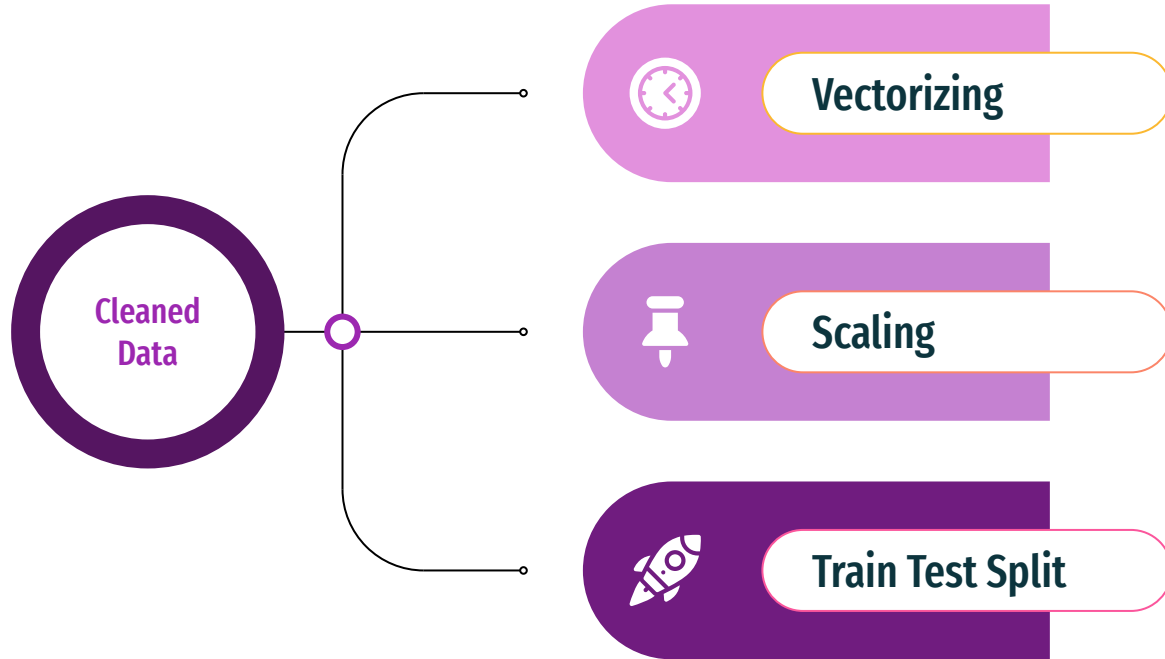- Inherent predictive bias.

# Upsampled Data

Count

8000
6000
4000
2000
0

Negative   Neutral   Positive   Factual

-1   0   1   2

Sentiment

# Downsampled Data

8000
6000
4000
2000
0

Negative   Neutral   Positive   Factual

-1   0   1   2

Sentiment

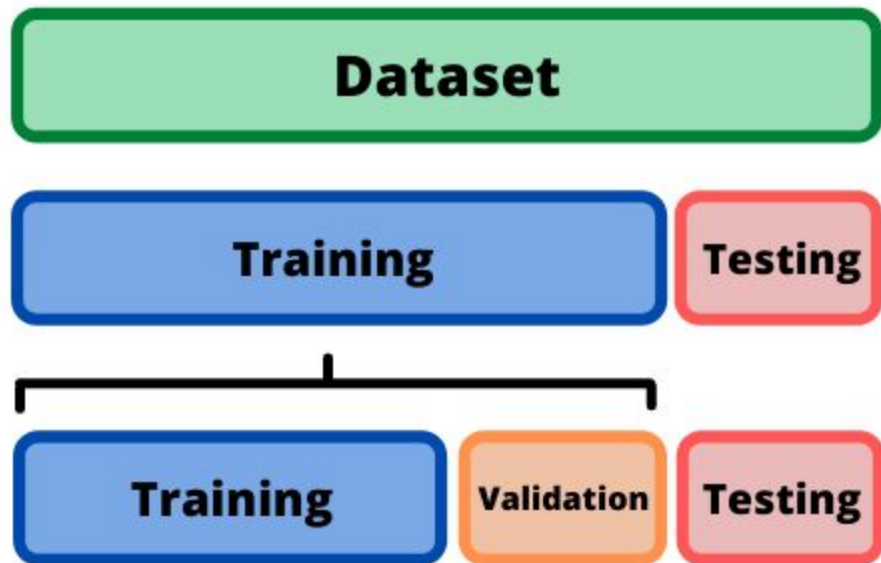- Risk overfitting, can check against Test data.
- Generally better than downsampling

- Risk losing valuable information.
- Reduces dataset to a more manageable size.

# Feature Text Engineering

**Cleaned Data**

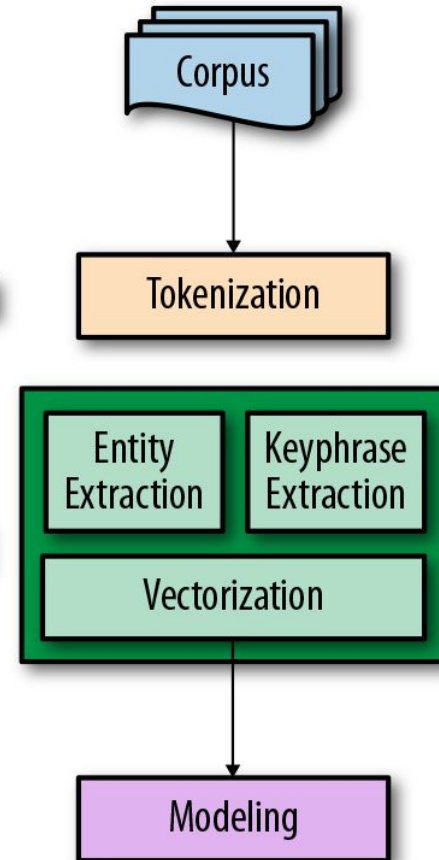Vectorizing

Scaling
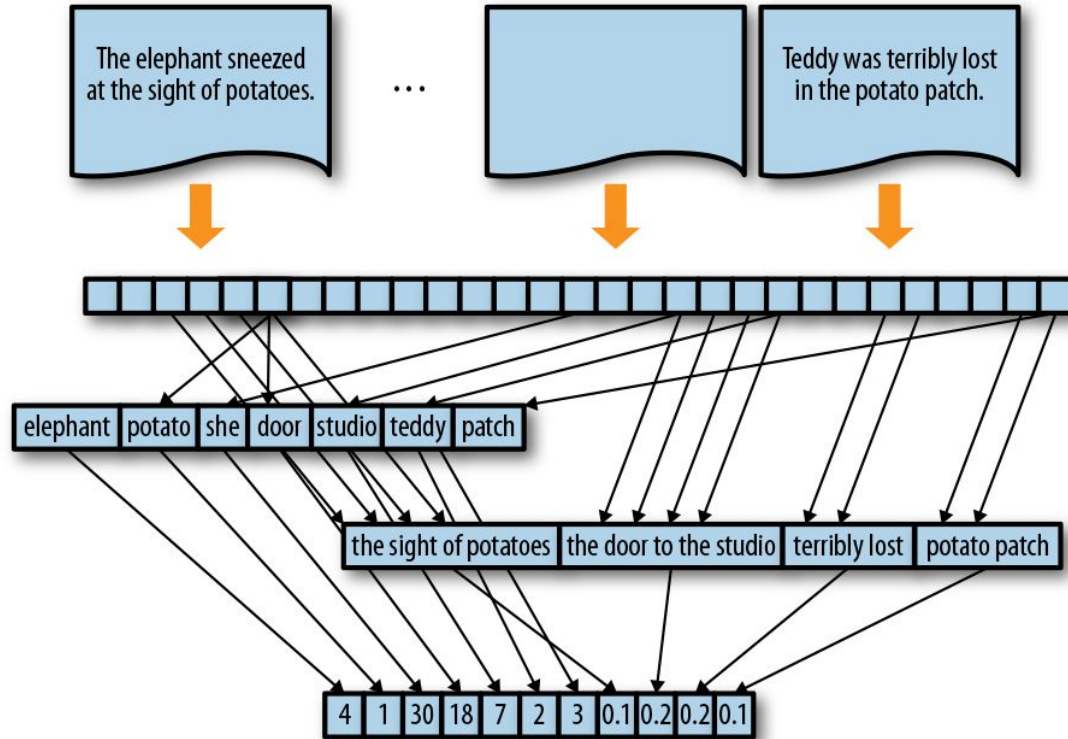
Train Test Split

# Train Test Split

# Text Vectorisation

- A document term matrix is generated and each column represents an individual unique word.
- Each cell contains a weight value that signifies how important a word is for an individual text message or document.
- Different from the count vectorization in the sense that it takes into considerations not just the occurrence of a word in a single document but in the entire corpus.
- TF-IDF gives more weight to less frequently occurring events and less weight to expected events. So, it penalizes frequently occurring words that appear frequently in a document such as "the", "is" but assigns greater weight to less frequent or rare words.
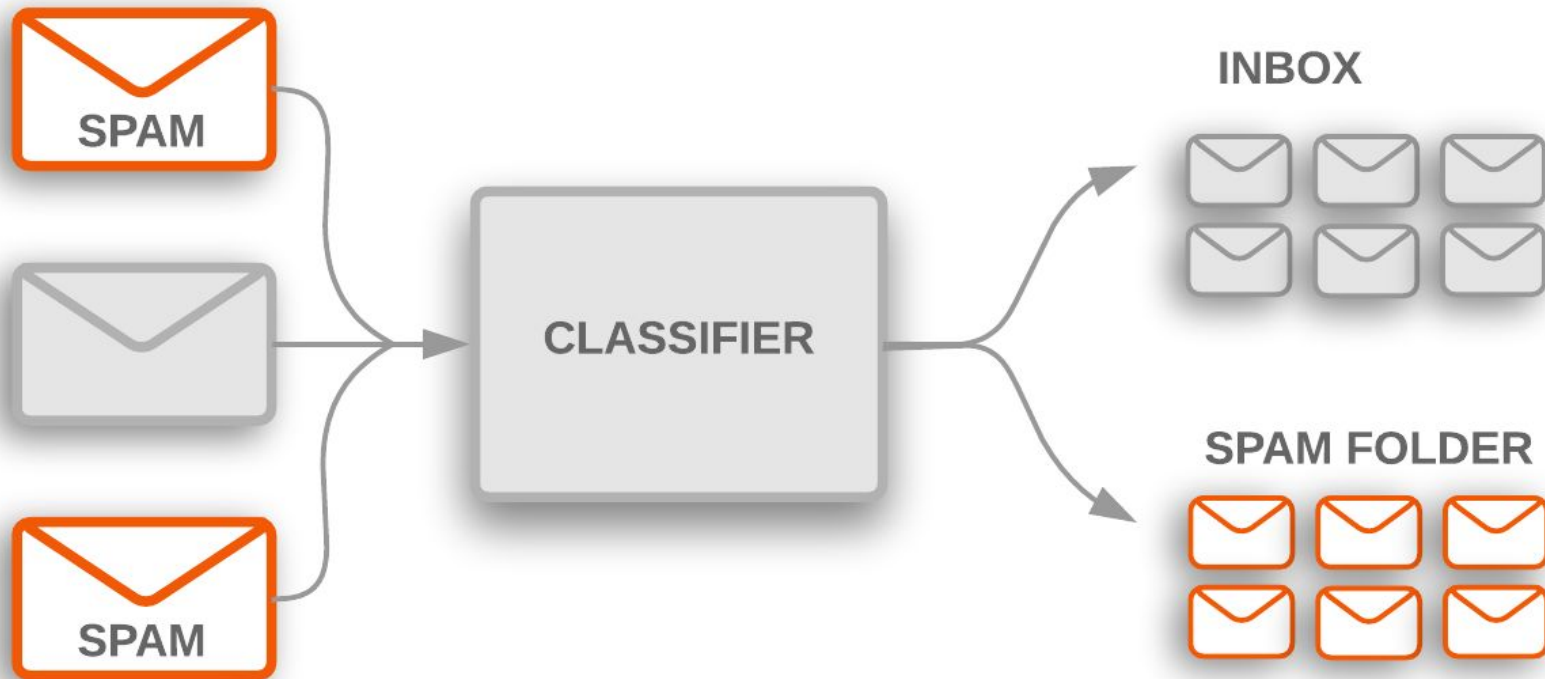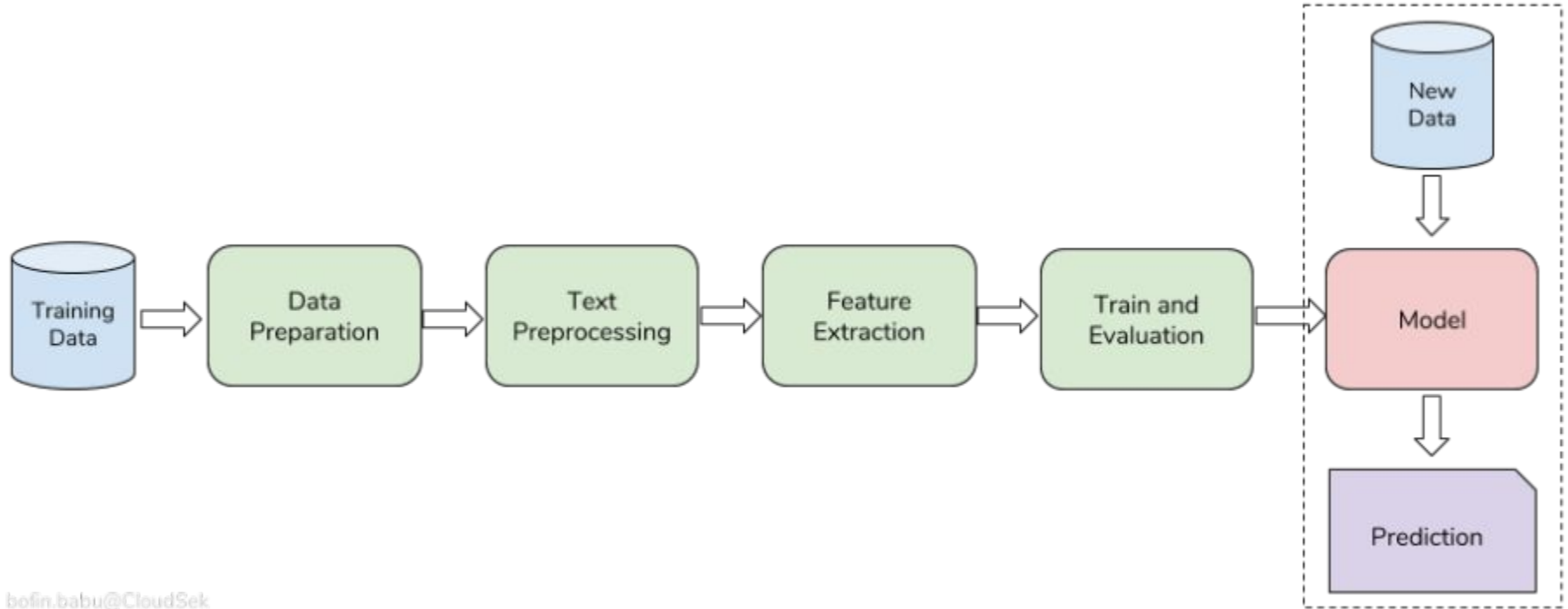
# Text Vectorisation

# Data Scaling

## Min Absolute Scaler

- Scale each feature by its maximum absolute value.
- This estimator scales and translates each feature individually such that the maximal absolute value of each feature in the training set will be 1.0. It does not shift/center the data, and thus does not destroy any sparsity.
- This scaler can also be applied to sparse matrices.

# Modelling

# Modelling

# Classification Algorithms

Logistic
Regression

Naive Bayes

Support Vector

K-Nearest
Neighbor