



# Marketing Email Subject Evaluation



# Problem

Having gathered data on marketing emails, Element 451 wants to find whether there is a relationship between marketing email subject line and email open rate.



# Aim

Find out the **features** of the **Email Subject** that result in a **higher opening rate**

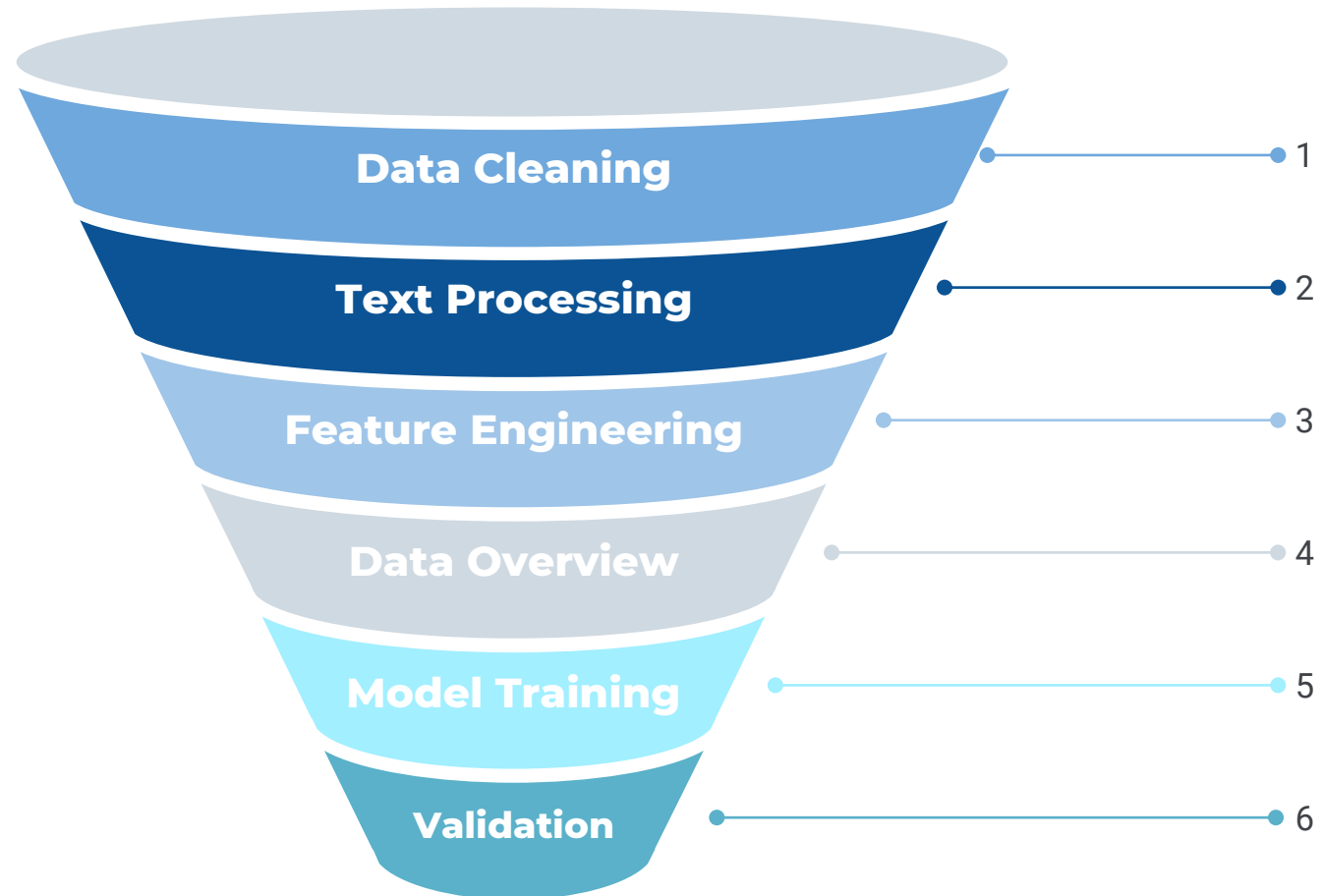
# Brief

***“The team created a pipeline which can clean, update useful features and predict using autotuned ML models.***

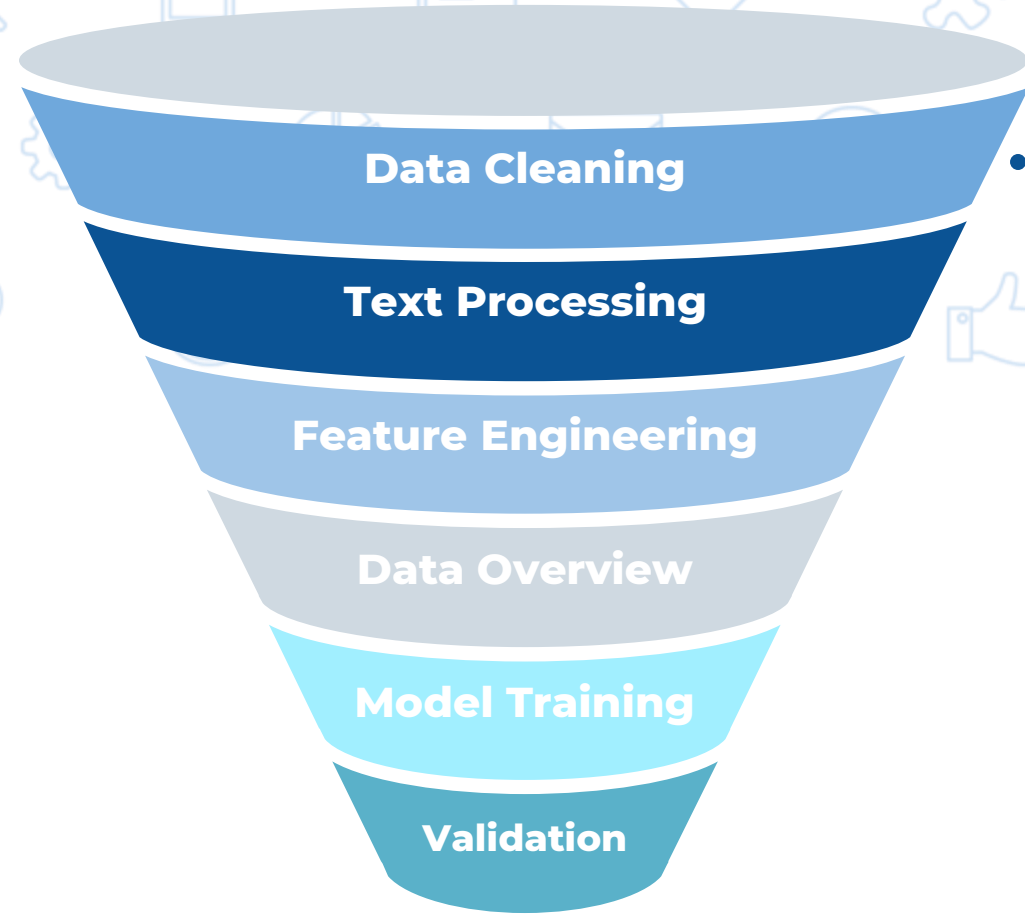
***The final model can provide an approximate impact of the email, after you input email subject, number of recipients and the sending date”***

# 0

## Our Pipeline



1



Data Cleaning

Text Processing

Feature Engineering

Data Overview

Model Training

Validation



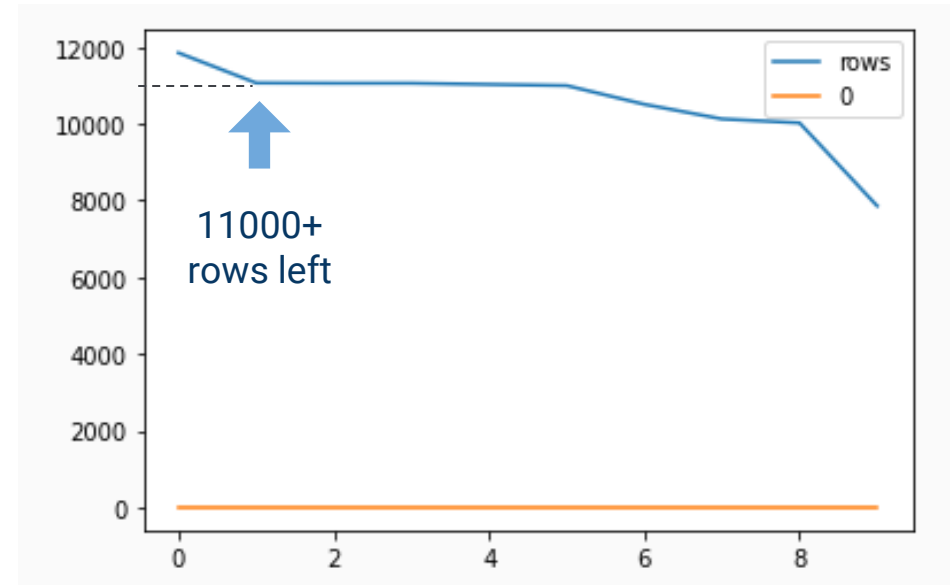
# Data Cleaning

# 1

## Data Cleaning

Step 1 / 4

**Remove NAs**



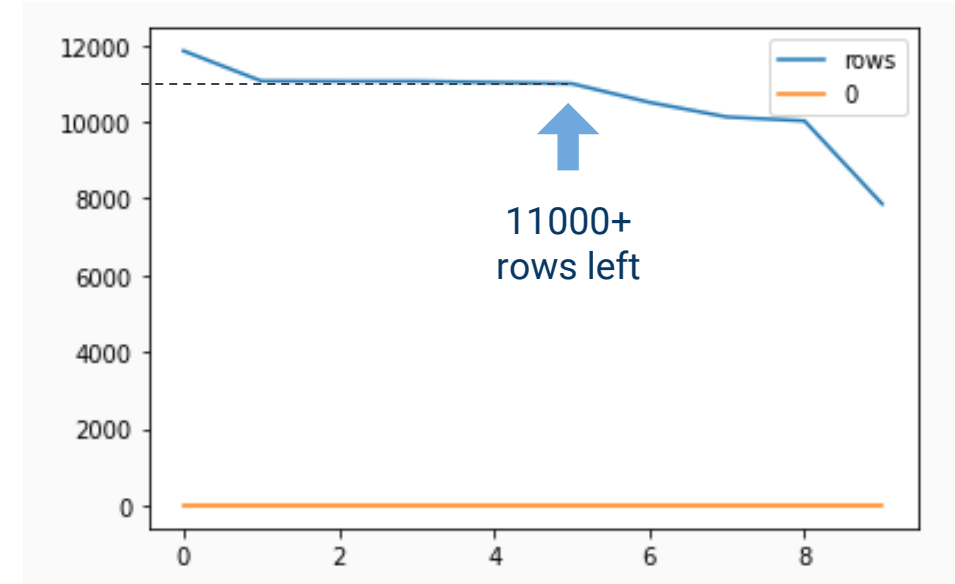
# 1

## Data Cleaning

Step 2 / 4

### Logic Check

- ▶ “opens” should be **greater** than “unsubscribed”
- ▶ “total” count should be **greater** than “unique” count
- ▶ “unique rate” should be **less** than 1





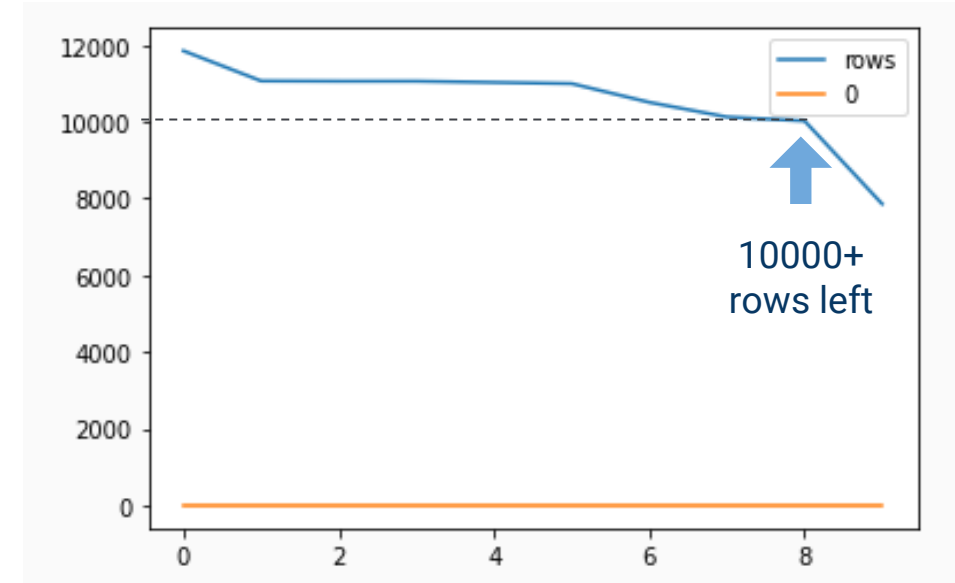
# 1

## Data Cleaning

### Step 3 / 4

### Removing Precision Error

- ▶  $CTV = \frac{\text{unique clicked}}{\text{unique opens}}$
- ▶  $\text{unique open rate} = \frac{\text{unique opens}}{\text{total delivered}}$
- ▶  $\text{open rate} = \frac{\text{total opens}}{\text{total delivered}}$



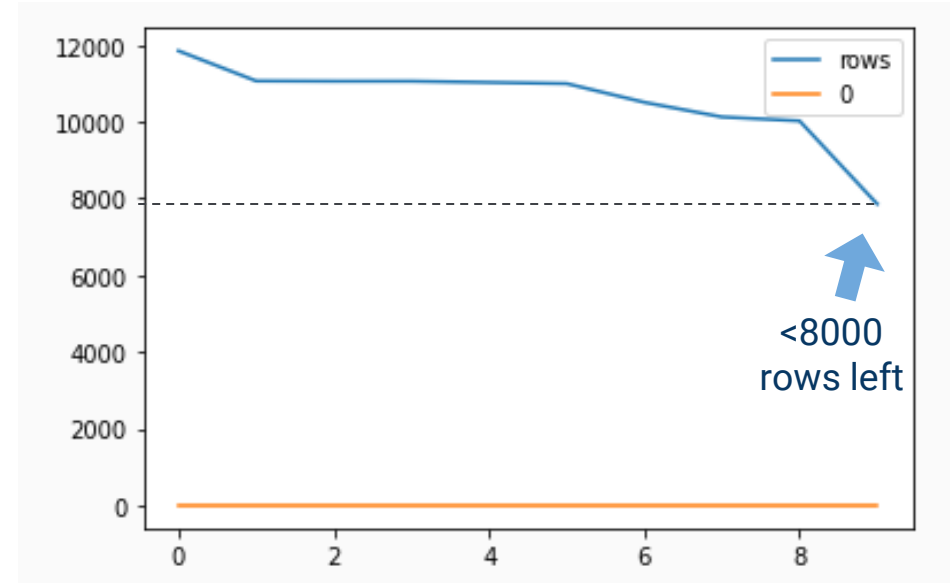
# 1

## Data Cleaning

Step 4 / 4

### Remove Test Case

Deleting rows that total delivered number  $\leq 100$

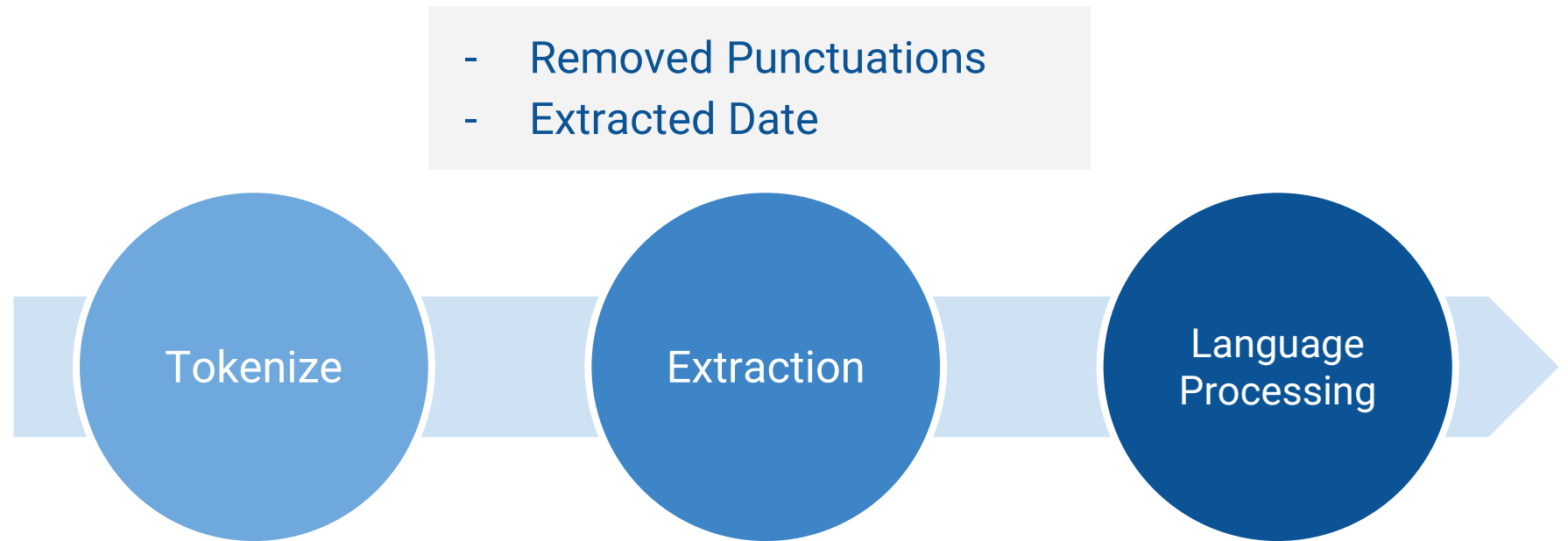




# 2

## TEXT PROCESSING

Focusing on the Email Subject to understand the meaning



- Removed Punctuations
- Extracted Date

Tokenize

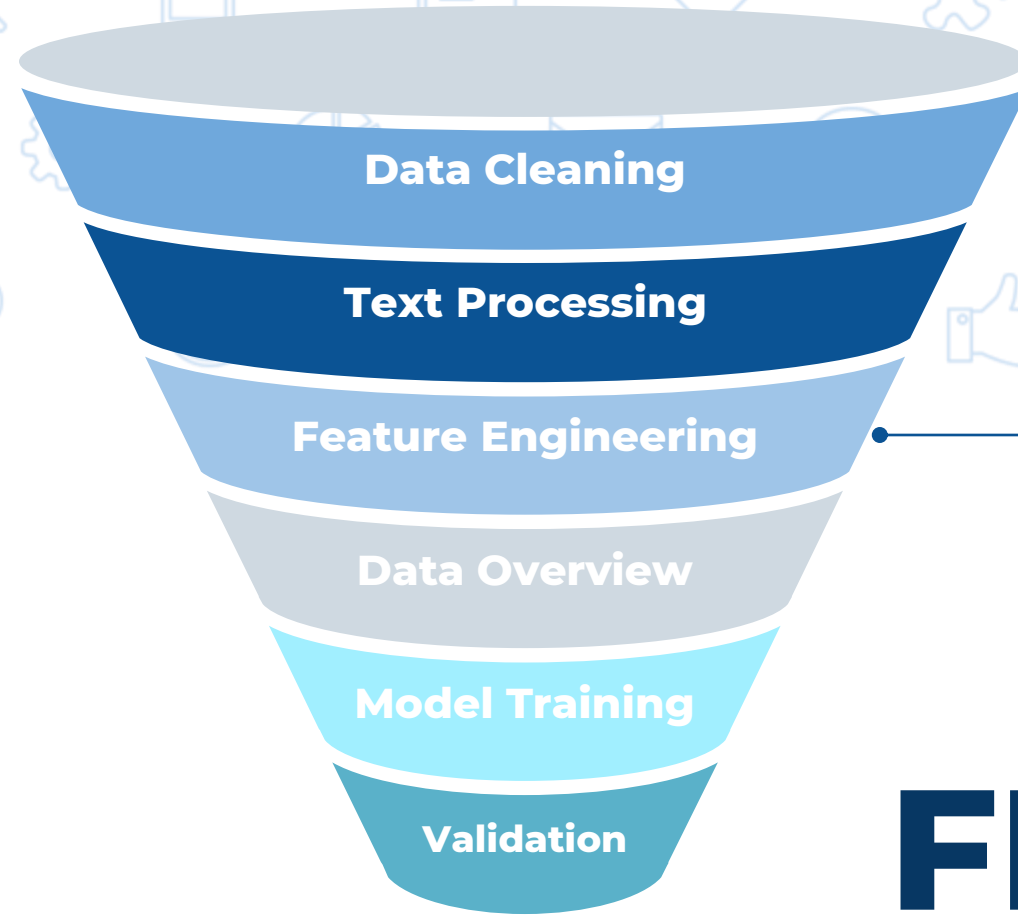
Extraction

Language Processing

- Extracted Username
- Extracted emojis
- Turned them into separate tokens

- Stemming
- Lemmatizing
  - **Clubbing words with similar meaning**

3



# FEATURE ENGINEERING

## Generating Features from Email Subject and Parameters

# 3

### Feature from Text

#### **Count**

Word / Char/ Emoji  
Emoji Type  
Punctuation  
User Token

#### **NER**

Total Count  
Organization  
Location  
Person  
Facility

#### **Date**

Weekday  
Month  
Day

#### **WORDS**

Stemming  
Lemmatizing

# 3

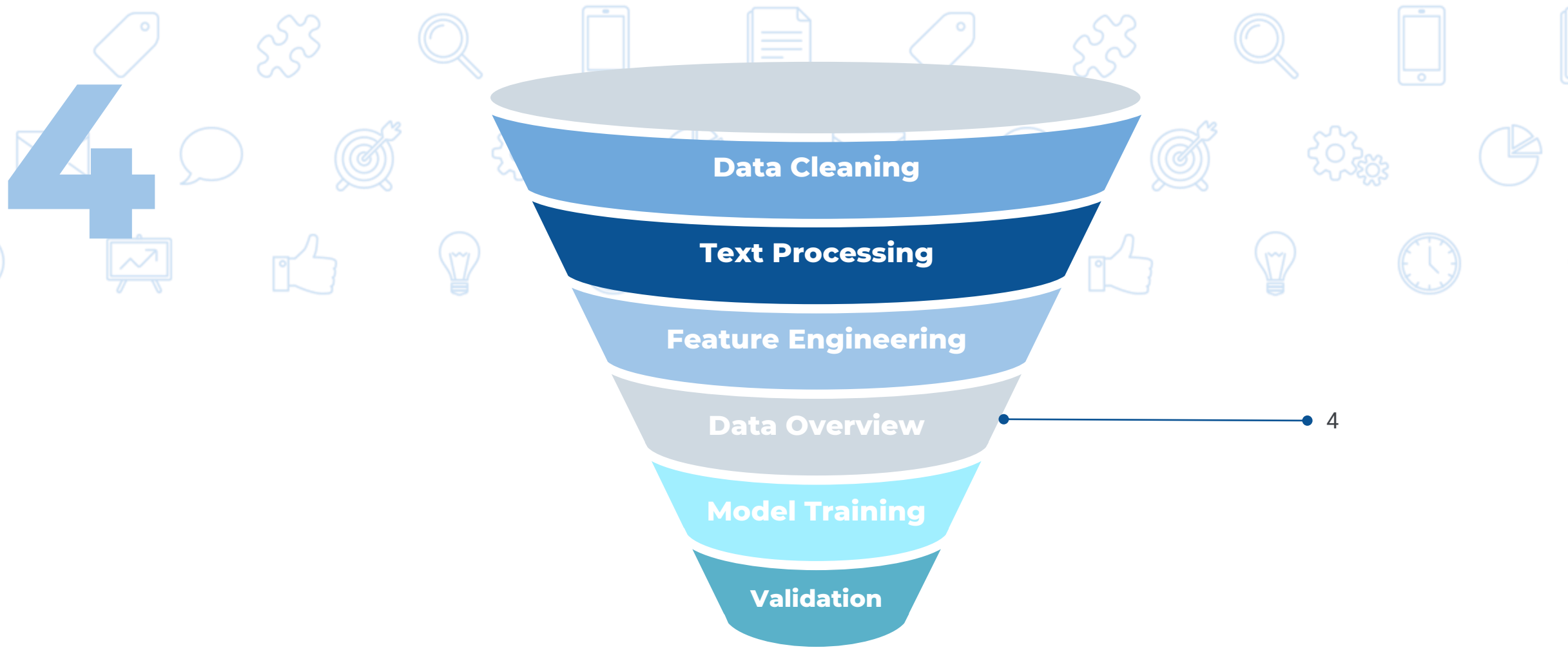
## Feature from Numbers

### Feature Transformed

unsubscribe  
click  
delivered  
open

### Fixed feature

Successful boolean  
unique open rate without  
subscribe

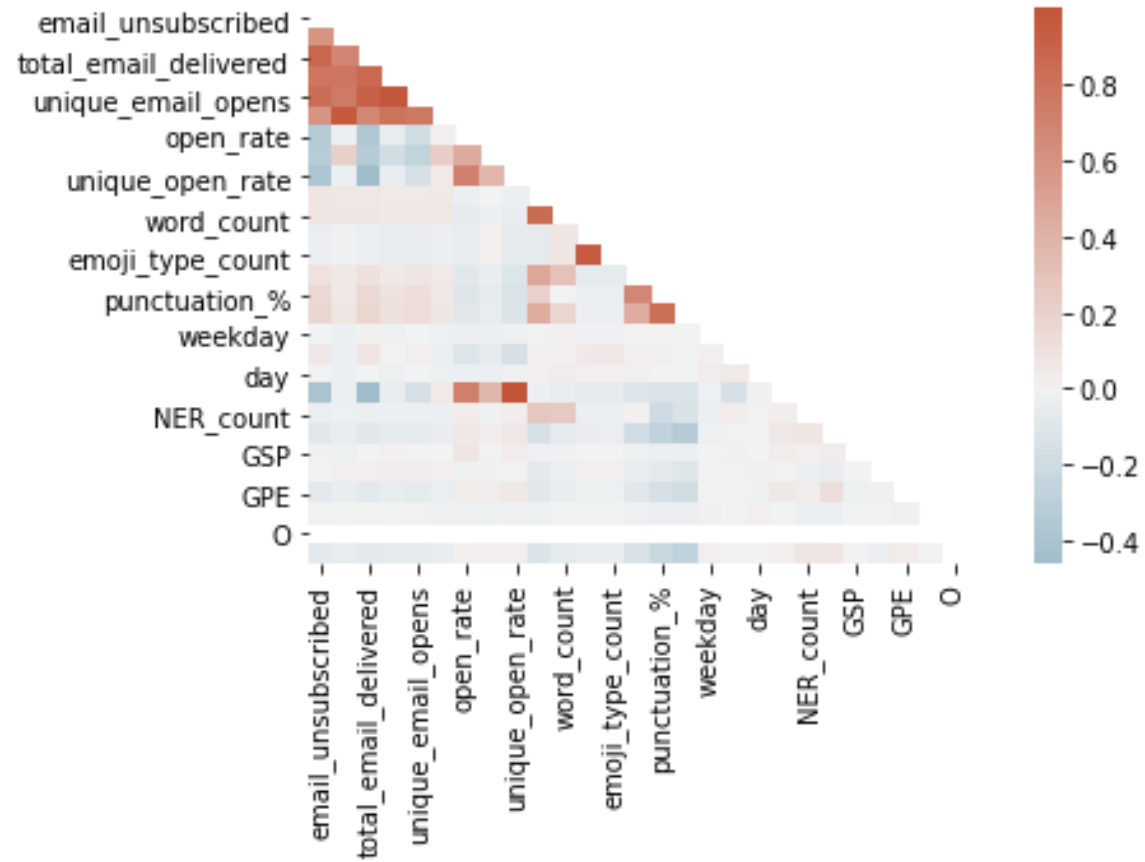


# DATA OVERVIEW



# 4

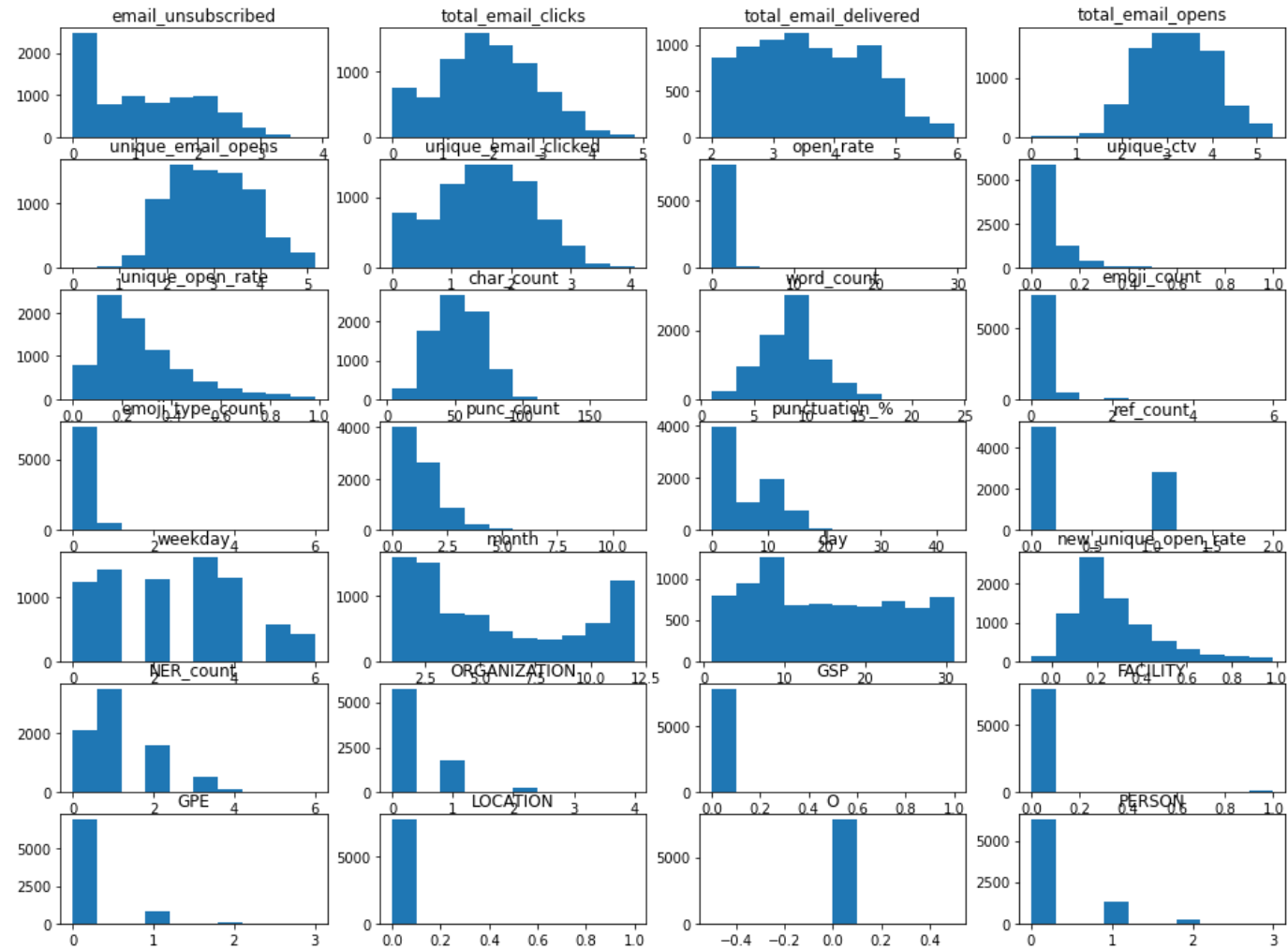
## Correlation plot



Day is highly correlated with Unique open rate

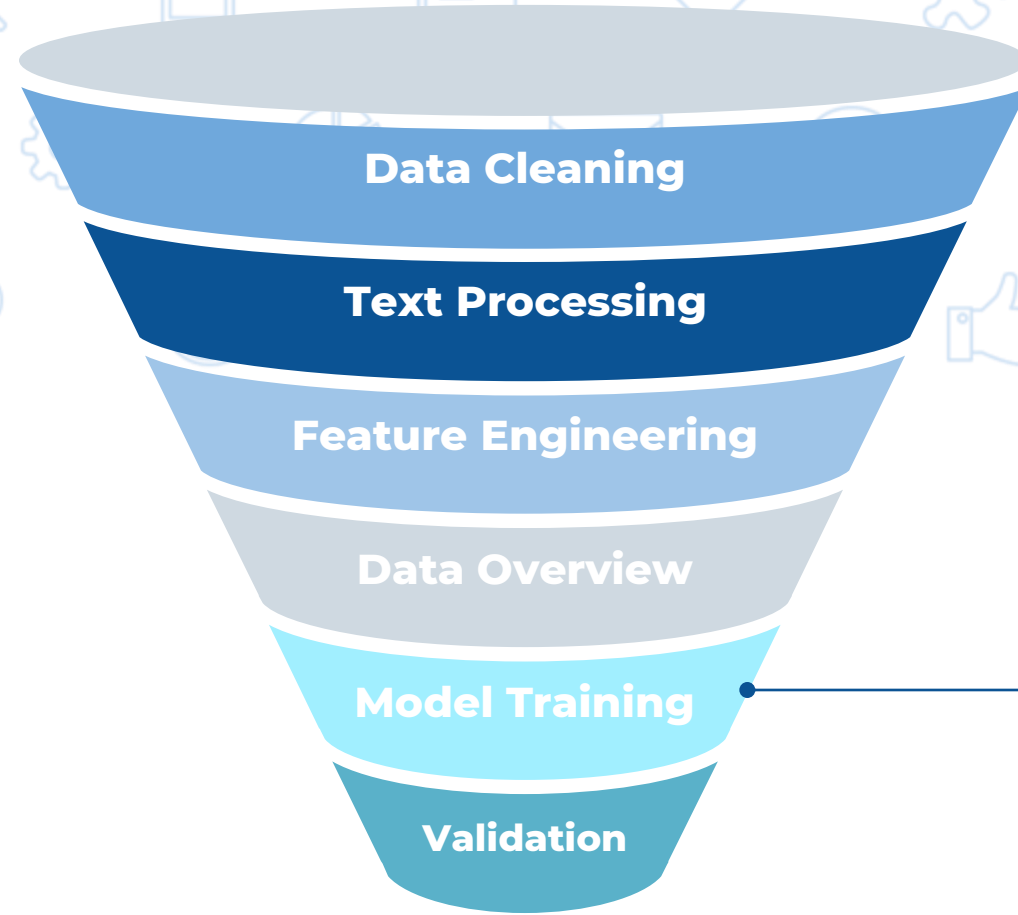
# 4

## Histogram



After transformation, no feature is highly skewed but most are not normally distributed

5



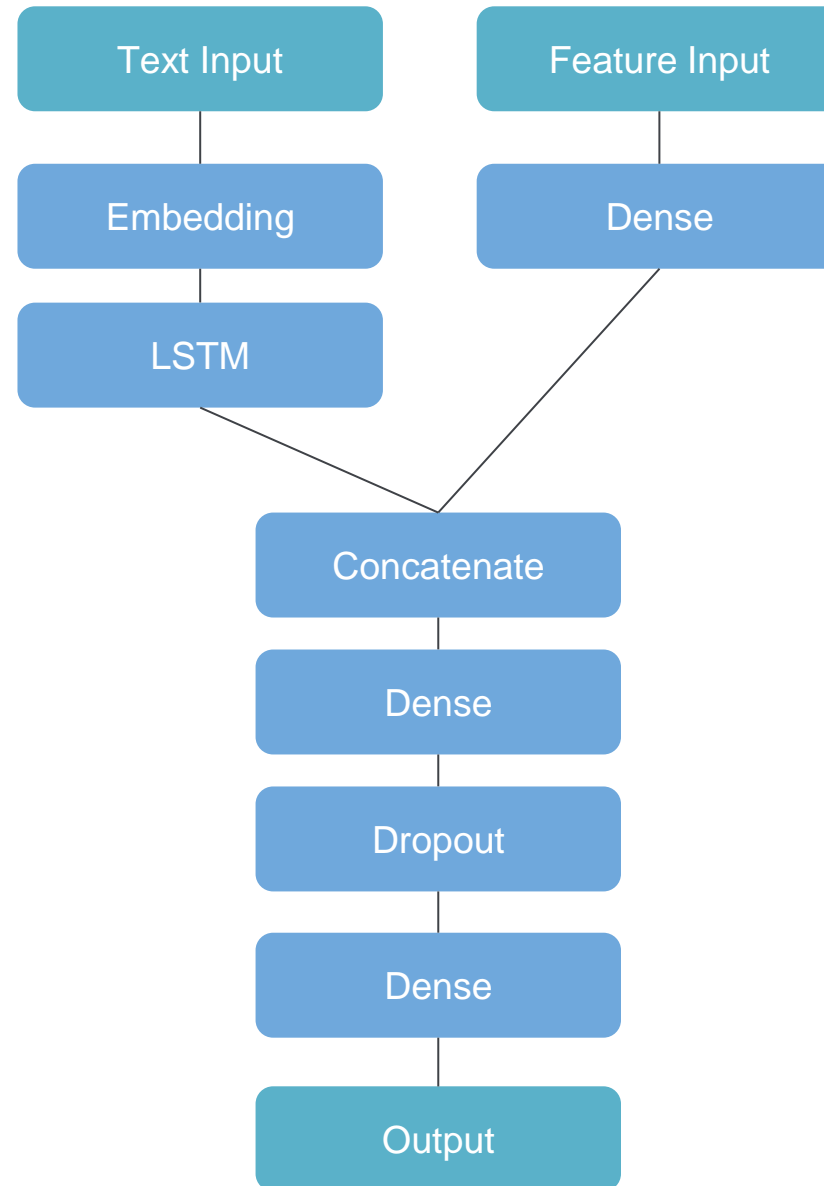
5

# MODEL TRAINING

# 5

## Model setup

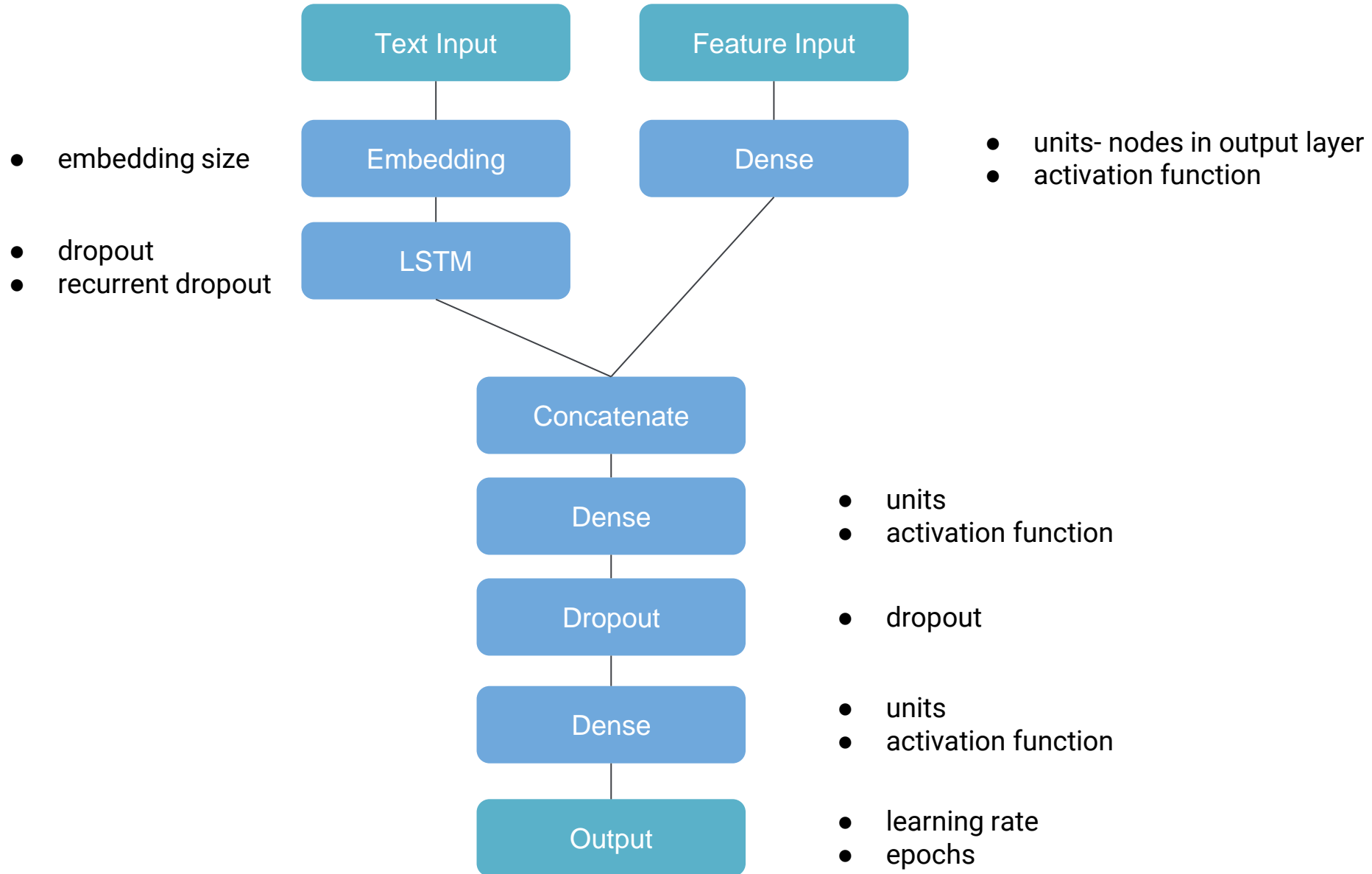
### Layer Overview of the ML Model



# 5

## Auto Tuning

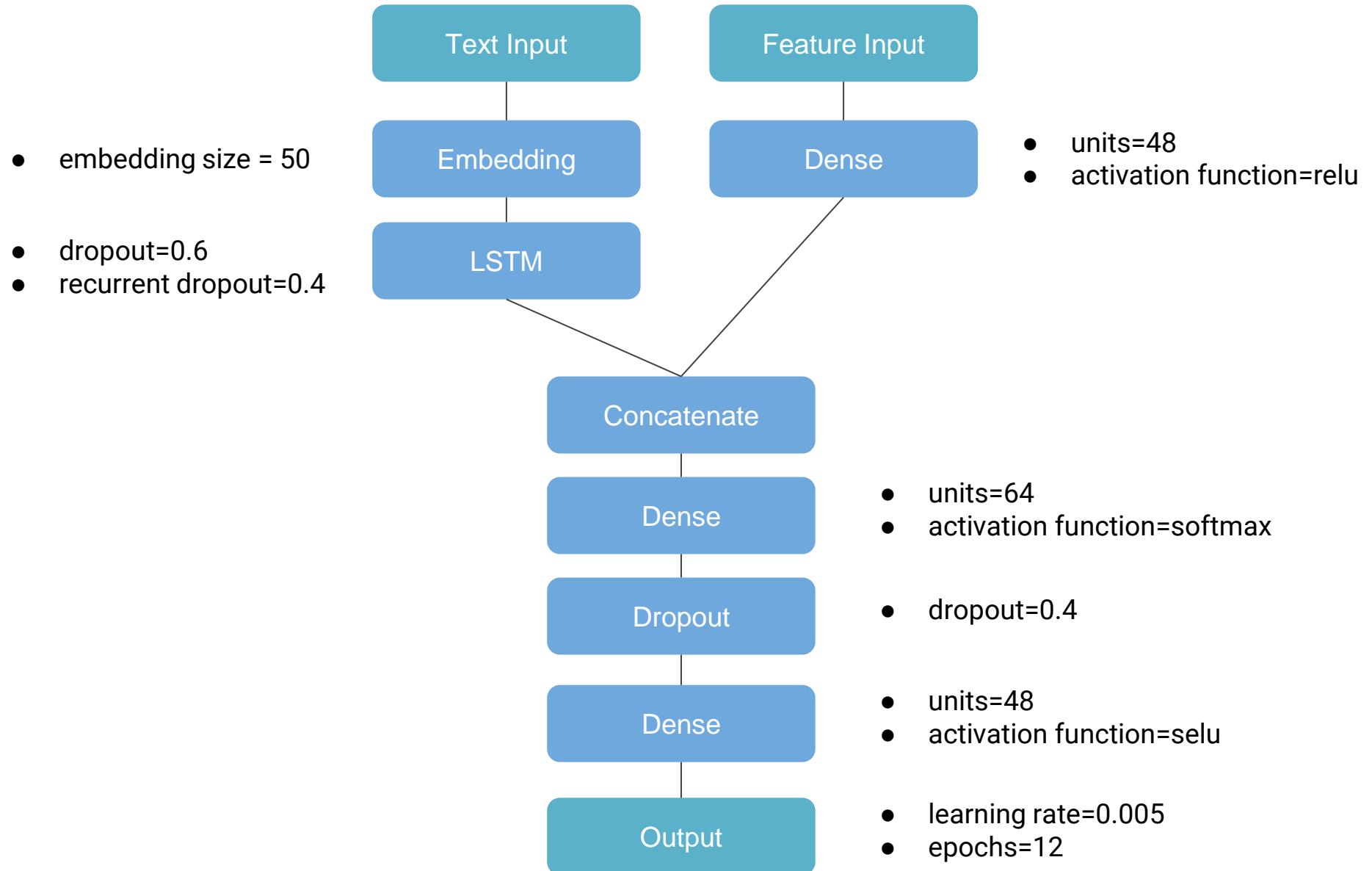
### HyperParameters for Auto Tuning Model



# 5

## Optimized Parameter Auto Tuning

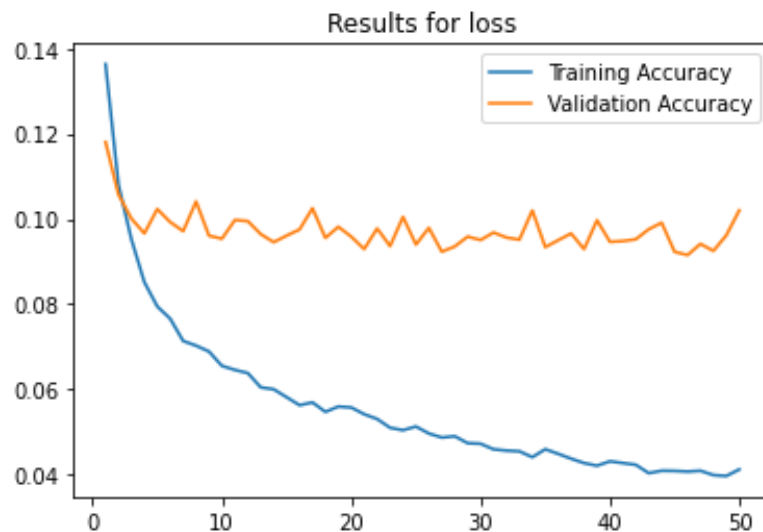
### Optimized Parameters for open rate prediction



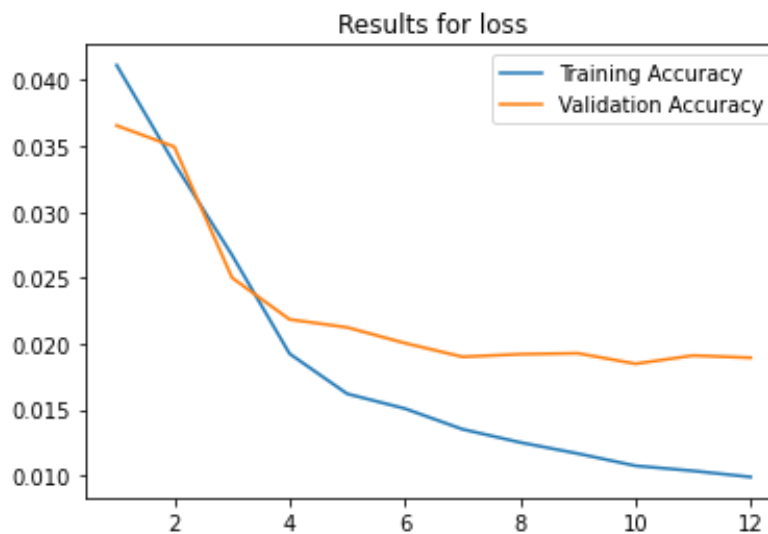
# 5

## Text Processing

### Before Autotuning



### After Autotuning

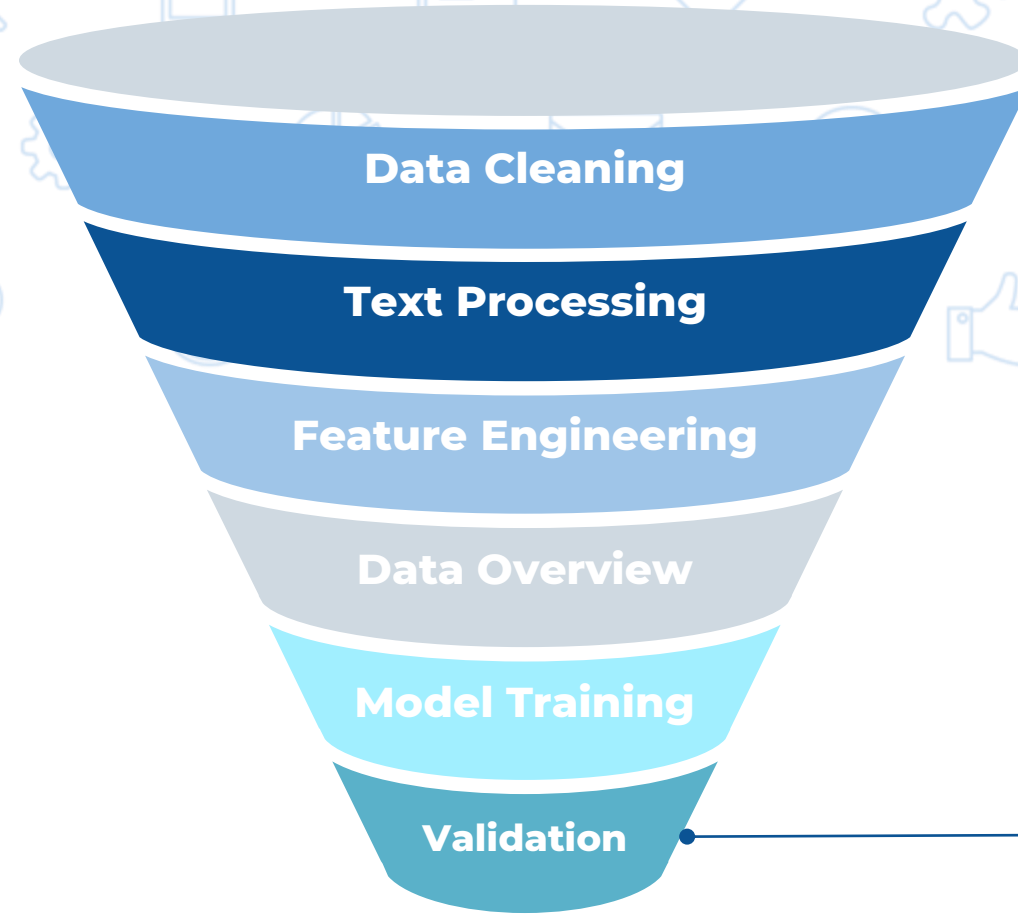


Loss :

0.1 -> 0.02

(-80%)

6



Data Cleaning

Text Processing

Feature Engineering

Data Overview

Model Training

Validation

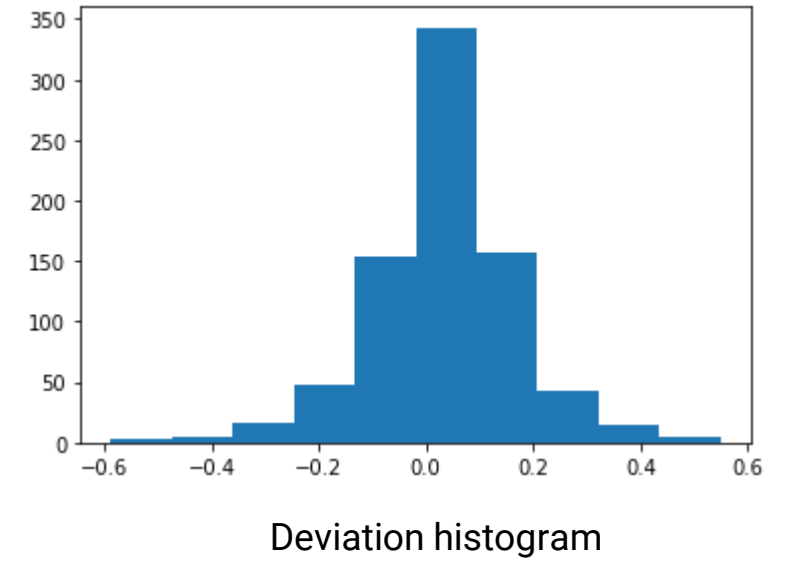
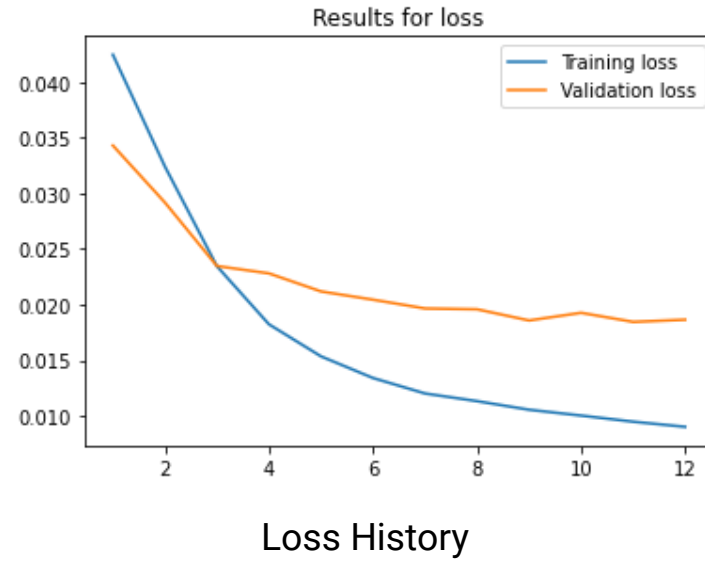
6

**VALIDATION**



# 6

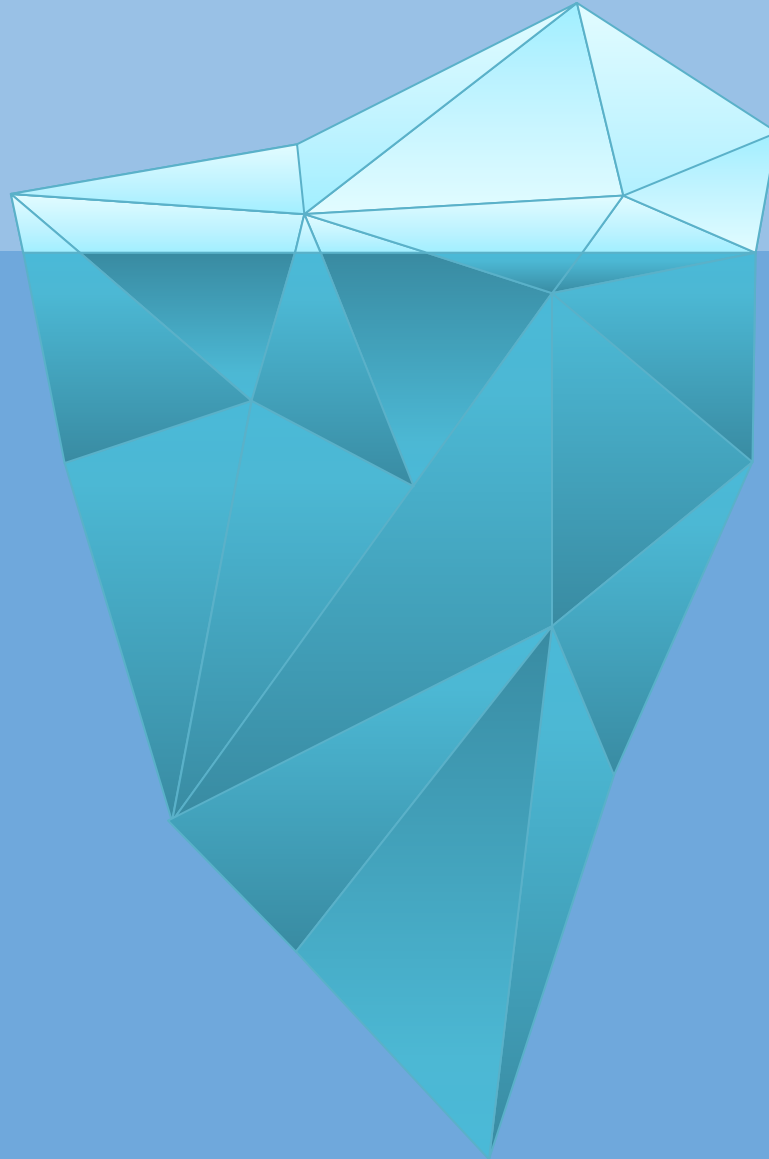
## Text Processing



Validation loss: 0.189  
Deviation : 0.0185

# 7

**From the deviation plot we can conclude 75% prediction is fall within  $\pm 0.2$  of the real situation**



## Conclusion

Our model measured by error. Since we are unaware of the standard of the loss, we do not know how good our model is.

## To Explore

Bigger dataset can get us a better accuracy as there will be bigger data to test and train

Defining a successful email will allow us to create a binary column which can help us run classifiers