# Analyzing Sentiment of Reddit Comments for Content Creators And Users

Team member:
Ta-Wei Wang, Ting-Hsuan Chen,
Ting-Ann Lu, Luke Hong

# Table of contents

1

Dataset

3

Conclusion

Introduction

2

Model

4

# 01

## Introduction

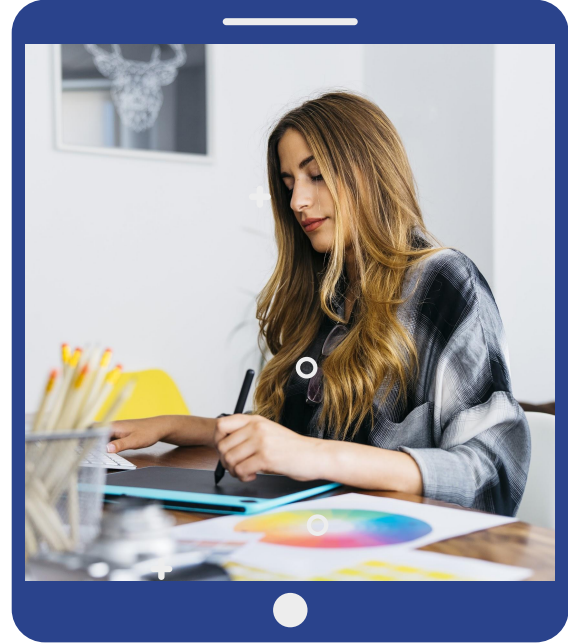Background And Motivation

# About Reddit

- one of the largest forum community social networks on the internet

- over 430 million monthly active users and over 100,000 active communities

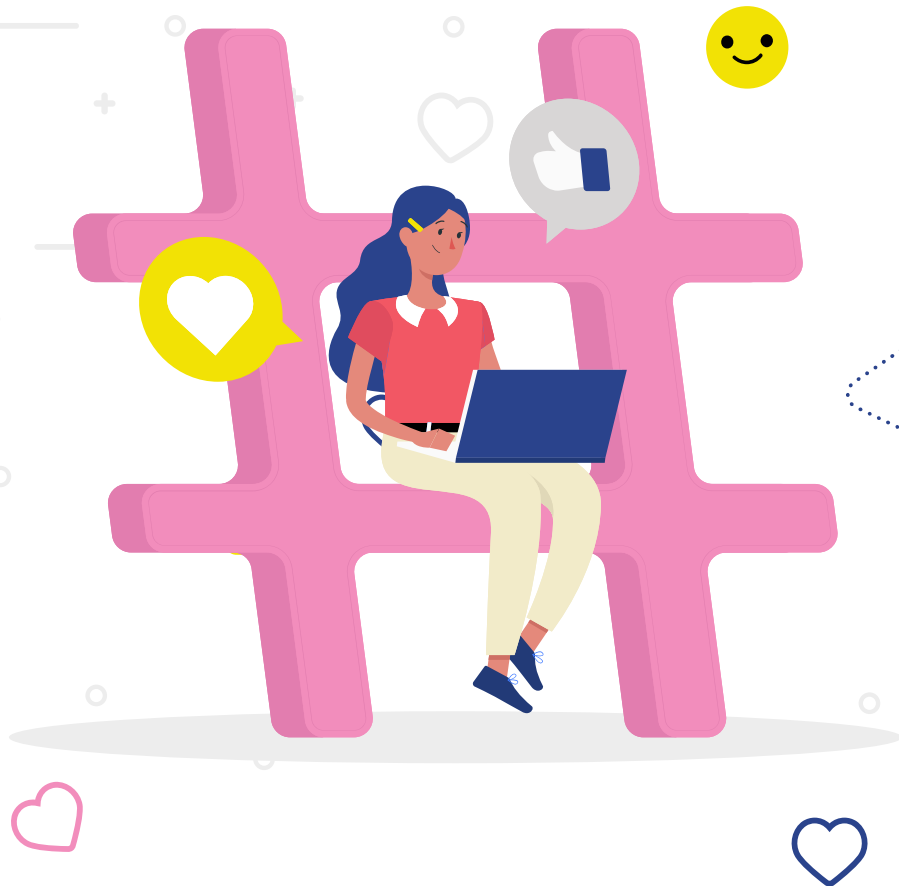- users can comment on and upvote/downvote posts and comments

# Motivation

From the perspective of content creators

understand the sentiment of Reddit comments on their posts to assess engagement and reach

# Motivation

From the perspective of Reddit users

understand the sentiment of Reddit comments by classifying the emotions of the posts' comments

# 02

........................

# Dataset

Data Collection And Data Cleaning

Data source: GoEmotions Dataset

3 datasets:
     1st: 70,000 rows / 37 columns
     2nd: 70,000 rows / 37 columns
     3rd: 71,225 rows / 37 columns

| text | id | author |
|------|-----|--------|
| subreddit | link_id | parent_id |
| created_utc | rater_id | example_very_unclear |

| Positive | | Negative | | Ambiguous |
|----------|----------|----------|----------|-----------|
| admiration 👏 | joy 😃 | anger 😡 | grief 😢 | confusion 😕 |
| amusement 😂 | love ❤️ | annoyance 😒 | nervousness 😬 | curiosity 🤔 |
| approval 👍 | optimism 🤞 | disappointment | remorse 😔 | realization 💡 |
| caring 🤗 | pride 😌 | disapproval 👎 | sadness 😞 | surprise 😮 |
| desire 😍 | relief 😅 | disgust 🤮 | | |
| excitement 🤩 | | embarrassment 😳 | | |
| gratitude 🙏 | | fear 😨 | | |

# Data Cleaning and preprocessing

| | text | id | author | subreddit | created_utc | rater_id | anger | disgust | fear |
|---|---|---|---|---|---|---|---|---|---|
| 18277 | She's horrible | ee1r7x8 | HoldenCaulfield7 | EDAnonymous | 1.547476e+09 | 16 | 0 | 0 | 0 |
| 38338 | She's horrible | ee1r7x8 | HoldenCaulfield7 | EDAnonymous | 1.547476e+09 | 1 | 0 | 0 | 1 |
| 50778 | She's horrible | ee1r7x8 | HoldenCaulfield7 | EDAnonymous | 1.547476e+09 | 23 | 0 | 0 | 1 |
| 55497 | She's horrible | ee1r7x8 | HoldenCaulfield7 | EDAnonymous | 1.547476e+09 | 11 | 0 | 0 | 0 |
| 63965 | She's horrible | ee1r7x8 | HoldenCaulfield7 | EDAnonymous | 1.547476e+09 | 74 | 1 | 0 | 0 |

```
len(clean['text'].unique())
```

57732

# Data Cleaning and preprocessing

## Punctuation

## Contraction

## Possible misspelling

### A surprise, to be sure, but a welcome one

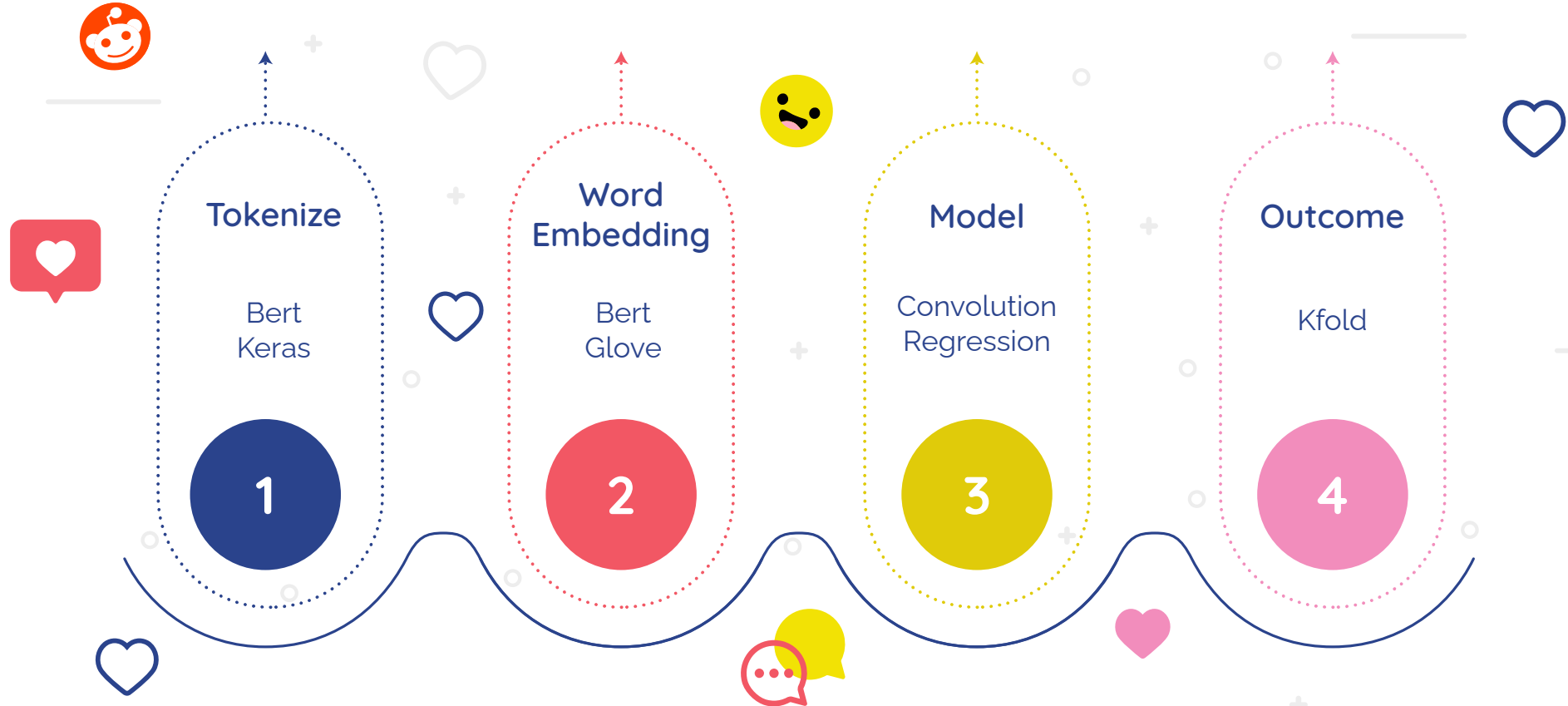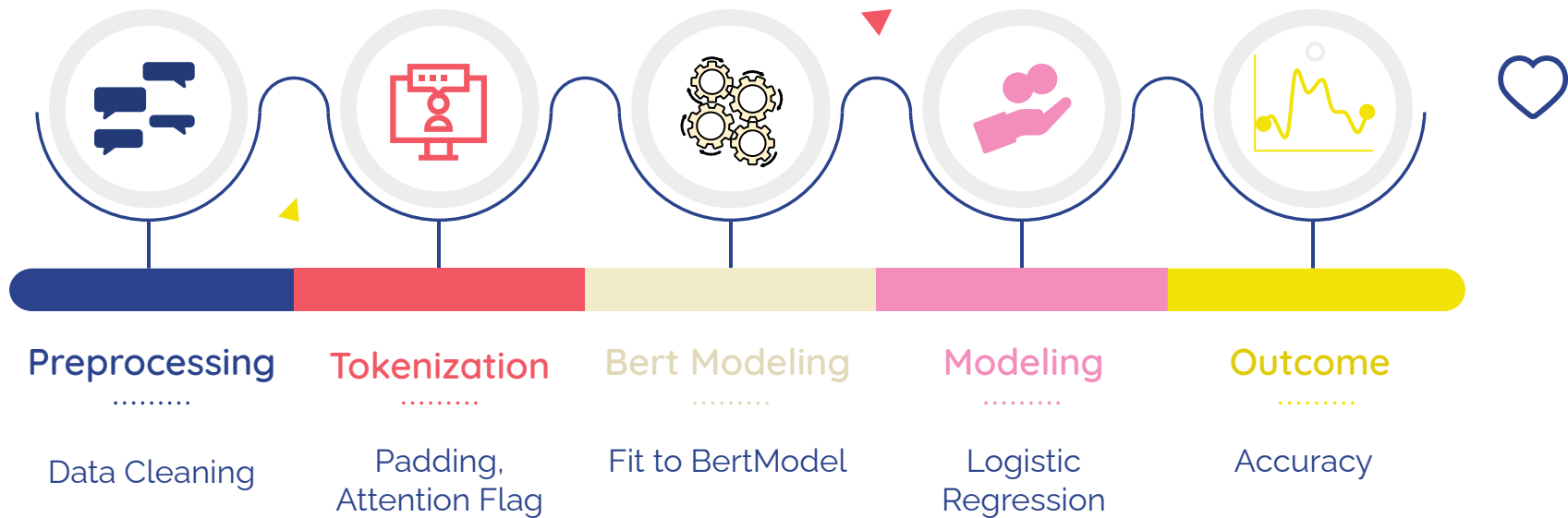→ a surprise to be sure but a welcome one

Lol! &#x200B;

→ lol

# 03

Modeling

Own model And Pre-trained Model

# Process

## Tokenize

Bert
Keras

**1**

## Word Embedding

Bert
Glove

**2**

## Model

Convolution
Regression

**3**

## Outcome

Kfold

**4**

# First Model



**Preprocessing**
..........

Data Cleaning

**Tokenization**
..........

Padding,
Attention Flag

**Bert Modeling**
..........

Fit to BertModel

**Modeling**
..........

Logistic
Regression

**Outcome**
..........

Accuracy

# Tokenization

- We used the 'BertTokenizer' function to tokenize the text
- Align the results with padded 0's and added an attention flag for BERT model
- Take the first 1500 rows as example:
  - Total number of words: 18740
  - Max length of a sentence: 41
  - 3705 unique words

**Sentence**

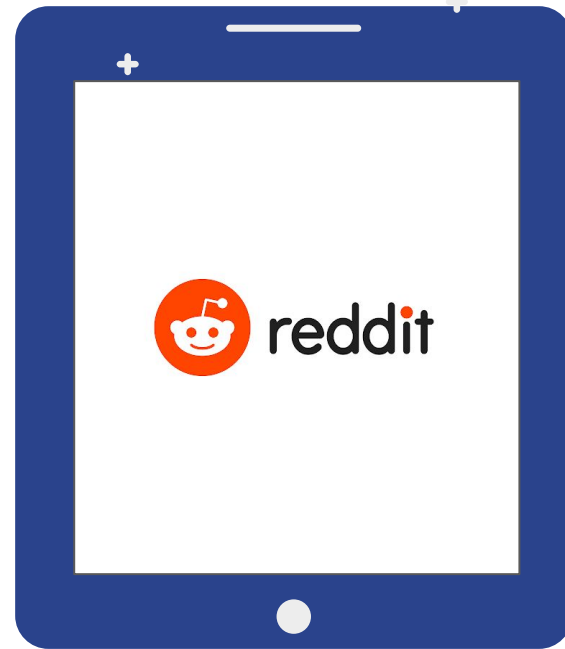| | |
|---|---|
| 0 | if you do not wear brown and orange you do not matter we need a tshirt with that on it asap |
| 1 | what do scottish people look like how i would love to have been there to take a swing at that softball |
| ... | ... |
| 1499 | they look like such fucking ocs designed them i do not know what you expected he phones in video games all the time |

**Tokenization & Padding**

**Sentence Tokenized**

| | | | | |
|---|---|---|---|---|
| 0 | 1045 | 2097 | ... | 0 | 0 |
| 1 | 1998 | 2024 | ... | 0 | 0 |
| ... | ... | ... | ... | ... | ... |
| 1499 | 1057 | 3501 | ... | 0 | 0 |

# BERT Model

- Run our sentences through BERT
  - BERT adds a token for classification at the beginning of every sentence.
  - The output corresponding to that token can be thought of as an embedding for the entire sentence.
- The results are returned as 'last_hidden_states' for calculations

Limitation:
- BERT requires huge quantity of CUDA memory.
- Though we managed to run our calculations on dedicated GPU, we are only able to import proportions of the dataset without exceeding the RAM limit.

## Sentence Tokenized

| | | | | |
|---|---|---|---|---|
| 0 | 1045 | 2097 | ... | 0 | 0 |
| 1 | 1998 | 2024 | ... | 0 | 0 |
| ... | ... | ... | ... | ... | ... |
| 1499 | 1057 | 3501 | ... | 0 | 0 |

**Bert**

## Sentence embeddings

| | | | |
|---|---|---|---|
| 0 | 6.2360e-02 | 3.9872e-01 | ... | 4.4866e-01 |
| 1 | 2.1780e-01 | 2.3444e-01 | ... | 3.0756e-01 |
| ... | ... | ... | ... | ... |
| 1499 | 1.0603e-01 | 1.2828e-01 | ... | 5.3088e-01 |

# Testbench Evaluation

- Fit the embedded results to a logistic regression
- Evaluate the accuracy of each label.

## 90.45%

| | | | |
|---|---|---|---|
| admiration | 0.869 | fear | 0.974 |
| amusement | 0.935 | gratitude | 0.951 |
| anger | 0.882 | grief | 0.993 |
| annoyance | 0.794 | joy | 0.915 |
| approval | 0.752 | love | 0.951 |
| caring | 0.915 | nervousness | 0.984 |
| confusion | 0.902 | optimism | 0.895 |
| curiosity | 0.846 | pride | 0.98 |
| desire | 0.977 | realization | 0.873 |
| disappointment | 0.843 | relief | 0.99 |
| disapproval | 0.82 | remorse | 0.99 |
| disgust | 0.931 | sadness | 0.928 |
| embarrassment | 0.964 | surprise | 0.925 |
| excitement | 0.922 | neutral | 0.611 |

# Second Model

**Tokenization**
..........
Tensorflow.
Tokenize

**Word
Embedding**
..........
GloVe.6B
100 x 100

**Define Model**
..........
Convolutional
Layers

**Outcome**
..........
Kfold
K folds
N epochs

**Overfitting**
..........
Plots

So we have another model…

We ran the model with every single labels, so we will have binary classification for every emotion instead of 28 multiple classifications.

We will have 28 denses in this layer

# Performance

With 3 K-folds and 150 epochs, we have all the average accuracy for every emotion in every fold.

| admiration | amusement | anger | annoyance | approval | caring | confusion | curiosity | desire | disappointment | disapproval | disgust | embarrassment | excitement |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.814617872 | 0.904684007 | 0.882320881 | 0.78298223 | 0.712211072 | 0.913683772 | 0.884434462 | 0.87748 | 0.941569507 | 0.857094169 | 0.820344985 | 0.923229039 | 0.966046214 | 0.919956386 |
| 0.849117041 | 0.923569918 | 0.896843255 | 0.808072567 | 0.744733095 | 0.913956523 | 0.892070651 | 0.894525 | 0.941092253 | 0.864798546 | 0.83357197 | 0.925615311 | 0.962841749 | 0.915797353 |
| 0.861652792 | 0.918723583 | 0.902495563 | 0.825855732 | 0.765239358 | 0.917087138 | 0.902972877 | 0.90495 | 0.948247671 | 0.879994571 | 0.854902506 | 0.923087418 | 0.965157509 | 0.921110034 |

| fear | gratitude | grief | joy | love | nervousness | optimism | pride | realization | relief | remorse | sadness | surprise | neutral |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.95459193 | 0.928615272 | 0.990727484 | 0.885934412 | 0.935637832 | 0.970818818 | 0.870525658 | 0.981046 | 0.849662483 | 0.978591383 | 0.974227846 | 0.902502239 | 0.927728891 | 0.623099446 |
| 0.954796493 | 0.934819639 | 0.989363849 | 0.893911481 | 0.946342111 | 0.967478037 | 0.889411628 | 0.97975 | 0.856821418 | 0.979750454 | 0.974296033 | 0.91647917 | 0.932501554 | 0.688961625 |
| 0.958270848 | 0.943065584 | 0.989976823 | 0.902904689 | 0.951247811 | 0.96904403 | 0.905018389 | 0.980567 | 0.860152721 | 0.983567417 | 0.979476333 | 0.923769236 | 0.936792552 | 0.717168987 |

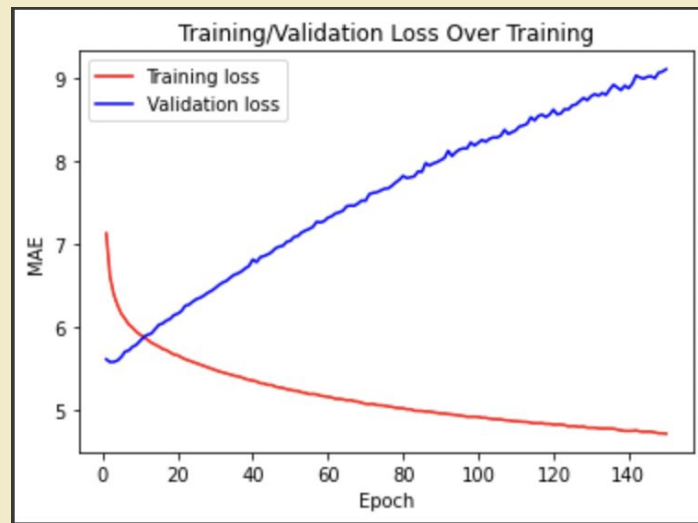Also, we have the overall performance...

**93.42%**
Training

**89.59%**
Validation

# Performance

## 93.42%

We have a great performance on accuracy!

**However...**

# 04

................

## Conclusion

```
Sum of each label
neutral              31449
approval             13235
annoyance            10024
admiration            9912
disapproval           8399
realization           7248
disappointment        6656
curiosity             6203
optimism              6200
joy                   5688
anger                 5644
confusion             5311
gratitude             5298
amusement             5180
sadness               4667
love                  4349
excitement            4336
caring                4330
disgust               4053
surprise              3823
desire                2838
example_very_unclear  2731
fear                  2136
embarrassment         2003
remorse               1663
nervousness           1556
pride                 1127
relief                1085
grief                  558
```
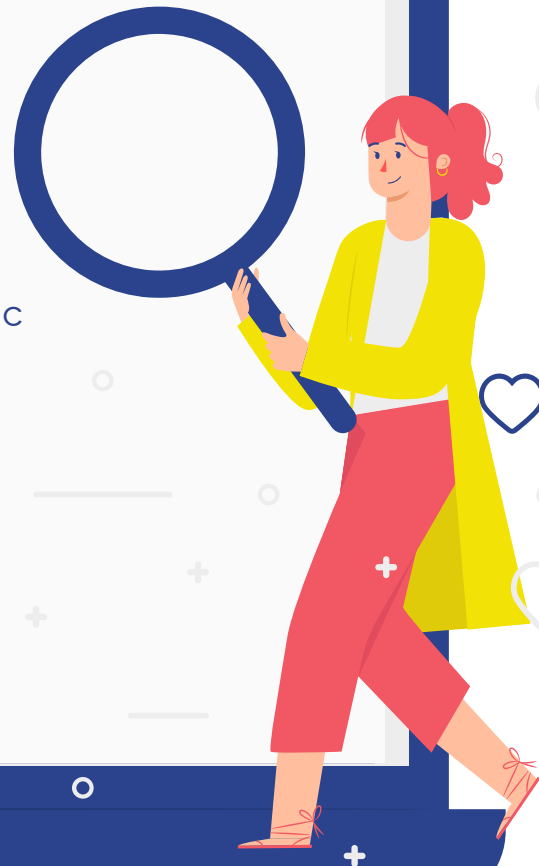
# Conclusion

1.  Own Accuracy(CNN) > Bert Logistic Regression Accuracy
    **93.42% vs 90.45%**

2.  Overfitting Problem

3.  Imbalanced Issue

# Thanks

**Credits:** This presentation template was created by Slidesgo, including icons by Flaticon, and infographics & images by Freepik and illustrations by Storyset