

Team project #2

The summary of your team project consists of the following parts:

1. The title of your project

**Using BERT and Ruled Based Sentiment Analysis to Predict PTT Public Opinion Responding Level (透過BERT與Ruled Based進行多標籤極性分析以預測輿論影響程度)**

2. Your team member names

110078501 李念紘 Joanne Lee

110078506 邱可人 Karen Chiu

110078509 施宇軒 Leo Shr

3. The source of text documents, how do you get them.

First, we searched for some famous Taiwanese athletes on PTT and choose five athletes which represent five different sports. We used BeautifulSoup and a for loop to scrape and set the number of pages(in our final project we set five pages for each athlete) on the PTT Gossiping and got the titles of each article. There are twenty articles on each page, which means we got one hundred articles for each athlete. Then, we used BeautifulSoup and a for loop to capture the content, comments, and comment number of each article through the title stored them in the list and then wrote them into a .csv file separately.

4. The sequence of processing the text

1. Data cleaning

An article of an athlete is a document. We collected a lot of documents and did some data preprocessing on them. First, we added new terms for tokenization and tokenized those words through CKIP, and we also defined new stop words and removed those stop words. We applied NER to tag the part of speech and only kept the adjectives, nouns and proper nouns.

2. BERT

We adopted BERT as our unsupervised learning method for the training dataset. We used articles about athletes as the input and want to get the result of positive or negative sentiment. But the difficulty we face was the data we have were too little. Thus, we used transfer learning which was the trained model from other professionals to help us finish training our own data and then got the result.

3. Keymoji

We used Keymoji to help us process Rule-Based optimization of our documents. KeyMoji estimates are based on the eight major paired polar core emotions proposed by American psychologist Robert Plutchik, namely Anger, Anticipation, Disgust, Fear, Joy, Sadness, Sur

prise, and Trust. First, we only use the titles of each article as our input, but the result came out with a lot of 0. We guessed that the sentiment score could not be effectively obtained due to the small vocabulary of the title. Therefore, we add content to our input to make our document longer and enrich sentiment scores in our results.

#### 4. Decision Tree

First, we use the Decision Tree as our training model to determine the relationship between the eight Keymoji sentiment scores and the “negative” or “positive” category.

#### 5. Random Forest

Bagging Decision Tree are too stable, we applied Random Forest (include Bagged Decision Tree & Drop out) as our another prediction model to predict how many comments an article has.

#### 6. Predictive Model Selection

Which model or skill to apply is critical for accuracy, these are called the tuning skill. To accomplish this goal, we built customized functions to execute the techniques. We embraced K-fold Validation, bagging, and boosting to check the models’ RMSE (Root-Mean-Square Error) which became the basis for our selection of models.

#### 7.

### 5. The objectives of the text processing; that is, how would you like to use the outputs for?

First, we want to understand how the Taiwanese think of those famous athletes, and how they would describe them (Text Cloud). Go one step further, we want to use the sentiment score that appear in the titles and the content of the article to predict PTT public opinion responding level, which means how many comments an article has. Also, we try to find out which sentiments affect the number of comments the most. The result can give some suggestions to those who want to attract people’s attention like Youtuber or journalists, etc.

### 6. The output of the results.

#### ● Decision Tree + Kimoji

According to our result in Decision Tree, we find out the relationship between eight Keymoji sentiment indicators and the “negative” or “positive” category:

1. Most of the Gossiping version articles have less relationship with sadness and are in line with common sense.
2. Joy has the highest average value among the sentiment indicators, which is in line with the characteristics of the Gossiping version.
3. Joy and Anger represent pride occupying most of the right half of the picture. We speculate that it may be related to the training materials we use, which were articles of well-known athletes, and the villager culture of PTT.
4. Anger and Disgusting represent the emotion of contempt, occupying most of the articles with Anger emotion.
5. Because the Gini coefficient is used for segmentation, the trust score is the most evenly distributed among the sentiment indicators in each segment. Therefore, 4 of the

e 8 layers are classified using Trust.

- **Feature Importance**

According to our result in Random Forest, we find out that the most important sentiment feature to affect the number of comments are **Anticipation**, and **Surprise**, which then is **Fear**.

- **Prediction**

Through our prediction, we can now use sentiment scores to predict the number of comments on a post.

- 

## **7. The self-evaluation of the work, including the contributions of the results, the lessons learned, and the potential application domains.**

We have used BERT and Keymoji to get the sentiment score of each article, then tried the supervised learning approach to predict the number of comments. First, BERT has some restrictions like a word limit so we have to do more data cleaning in the beginning. Secondly, we used the free version of Keymoji which have a usage limit, so we have to spend more time and cut our document into a smaller dataset as our input. In the future, we should find some way to deal with this problem or write code to increase our efficiency.

Also, the problem we met is our data were too little which lead to our prediction accuracy not being high enough. So in the future, we can extend the number of documents, for example, we can grasp 100 pages of articles on the Gossip page in PTT, which is about 2000 articles.

Finally, in the future, we can further enhance our model. For example, we can add mediators to the model, or have a discussion about interaction.

## **8. Your link to your codes that TA and I can access.**

