

# Fake-EmoReact

## Team - November

# Outline

1. Data Analyze
2. Data Preprocessing
3. Experiment
4. Result (leaderboard)

# Data Analyze

	real	fake
the number of samples (A)	31799	136722
the numbers of unique text (B)	23134	3225
the average number of replies ( $=A/B$ )	1.37	42.39
number of unique texts which are repeated posted	2495	2
the average of the length of the text	113	199

# Data Preprocessing

- feature: “text” and “reply”
- word embedding: TF-IDF、GloVe
- data set: split train.json (training : validation : testing) = (8 : 1 : 1)

# Data set

	real	fake
training	25470	109513
validation	3149	13536
testing	3180	13673

# Hardware

CPU : Intel(R) Xeon(R) CPU E5-2676 v3 @ 2.40GHz

RAM : 32G

GPU : GeForce GTX 1080 11G

OS : Ubuntu

# Experiment 1

Different pre-process in BERT-Base (L-12\_H-768\_A-12)

1. text only
2. text + reply
3. text + reply with replacement Tag and URL:  
e.g. "... leading these terrorists! @GenFlynn @BarbaraRedgate ...  
<https://t.co/VK5hGfZmPf> @JosephJFlynn1 @GenFlynn"  
  
=> "... leading these terrorists! TAGUSER TAGUSER ... URL  
TAGUSER TAGUSER"

# Experiment 1 result

method	accuracy	precision	recall
text	0.9738	0.990	0.716
text + reply	0.9949	0.9980	0.9965
text + reply with replacement	0.9923	0.9959	0.9941



# Experiment 2

Compare different models

1. TF-IDF + Xgboost
2. GloVe(twitter) + LSTM
3. GloVe(twitter) + Bi-LSTM
4. BERT-Mini (L-4\_H-256\_A-4) + dropout + classifier
5. BERT-Base (L-12\_H-768\_A-12) + dropout + classifier

## Experiment 2

model	accuracy	precision	recall
TF-IDF + Xgboost	0.966	0.971	0.82
GloVe + LSTM	0.971	0.976	0.989
GloVe + Bi-LSTM	0.971	0.971	0.989
BERT-Mini + dropout + classifier	0.971	0.976	0.989
BERT-Base + dropout + classifier	0.993	0.997	0.998

## Result (leaderboard)

BERT-Base (L-12\_H-768\_A-12) + dropout + classifier

Precision score	Recall score	F1 score
0.7394(18)	0.6664(13)	0.6353(14)

# Reference

- Transformers: <https://huggingface.co/transformers/>
- Classify text with BERT:  
[https://www.tensorflow.org/text/tutorials/classify\\_text\\_with\\_bert](https://www.tensorflow.org/text/tutorials/classify_text_with_bert)
- bert-as-service: <https://github.com/hanxiao/bert-as-service>
- predict fake news:  
<https://towardsdatascience.com/predicting-fake-news-using-nlp-and-machine-learning-scikit-learn-glove-keras-lstm-7bbd557c3443>

team November 陳菀渝 309555007 請問遇到表情符號，有作什麼特殊處理嗎

報告組別: team november 孫維佑 309511065問題: 考慮過完全平衡data? 像是讓dev裡的real/fake 1:1這樣

team November 羅禾洲 0816166 請問你們覺得為什麼replacement沒有增加precision嗎

報告組別: team November, 林哲宇 0616018 請問訓練時間大概多久呢, 我自己訓練BERT 開 GPU 只訓練 10% 測資就花了兩小時左右

請問取代掉以後效果更差的可能原因? 或是說有可能有其他更好的取代方式? 或是直接刪除效果可能更好?報告組別 :team november, 問題:0686011。

報告組別: Team November 問題: 309551179 你們如果不是用replacement而是直接刪掉網址, 會不會訓練效果比較好

team november 黃靖 309706022 想請問為什麼沒有考慮用Categories內容去做分析?有甚麼考量嗎?

有沒有想過可能有甚麼原因造成eval和practice phase的表現有落差? 報告組別: Team November 問題: 0713309