# Sentiment Analysis on LGU Kayuan

Yingyao Liu, Ying Xue, Yunfei Ke, and Shijie Shao School of Data Science The Chinese University of Hong Kong, Shenzhen {120090628, 120090452, 120090277, 120090491}@link.cuhk.edu.cn

## Abstract

Aiming to improve the mental health of CUHKSZ students, we trained a sentiment analysis model based on LGU Kayuan, a popular anonymous social website targeting at CUHKSZ students. Based on 15000 texts scratched from Kayuan, we trained a BERT-based model with several techniques solving data imbalance, and finally our model achieved an accuracy of 85.33% and an F1 score of 74.37% on the test set. In the future, we will collaborate with student counselling departments to monitor students' emotional levels and provide them with timely mental service support according to the predictions of our model.

# 1 Motivation

In the recent two years, affected by the pandemic and increasingly intense peer pressure, the mental health state of CUHKSZ students is becoming more and more concerning. Aiming for an efficient approach to monitor students' mental health states, we performed sentiment analysis on LGU Kayuan, an anonymous social website targeted at CUHKSZ students. Due to its anonymity, more and more students prefer to express their real feelings on it, which made it the best public platform to observe CUHKSZ students' emotional levels. Therefore, we managed to train a model based on LGU Kayuan corpus to reveal the underlying emotion under every post on Kayuan.

# 2 Data Preparation

## 2.1 Data Collection

Based on Python web scraping and the software Charles, we extract 15,000 posts from LGU Kayuan from March 17th to April 7th and stored the data into JSON files.

<u>ଜ</u>	cuhk	••• - ••	100	
			420	node: {
_			421	"id": "UG9zdDo10TAzOA",
Q198 最新	〒 飛口		422	"content": "提个问,如果感染了甲流医务室会开药吗,还是直接转诊,总觉得自己的症状怪怪的",
-			423	"images": [],
● 字前的卡卡		38 5999900	424	"isAnonymous": true,
TA长得好可爱像小熊			425	"creator": {
听他上课心情都很好 🔁			426	"id": "VXNlcjoyNTQxMw",
好评好评			427	"username": null,
用个卡卡	(会校祖	🐵 🔿 👘	428	"createdAt": null,
			429	"profile": {
□ 那个课, uu是几号 □ 他也想看着可爱小照		430	"id": "UHJvZmlsZToyNTQxMw",	
-			431	"nickname": "罗小黑",
🖶 道扬的卡卡		39 分钟前	432	"avatarUrl": "https://cdn-v4.szlikeyou.com/img/avatar/25413/f355b880-dbd1-42b
			433	"bio": "你猜我是什么猫",
基金不知道能不能做得出;	来、是不是没致了		434	"gender": "FEMALE",
			435	"grade": 2019,
学业求助室	<i>(</i> २ स.स.	⊕r ()	436	"institute": "经管学院",
🚗 不会是我的队友吧			437	" typename": "Profile"
•			438	

Figure 1: LGU Kayuan

Figure 2: Posts obtained from LGU Kayuan in JSON file

#### 2.2 Data Preprocess

To balance the text length, we cut off fractions of texts exceeding 128 words. Then, we manually labeled each post according to their underlying emotions into four categories: neutral, positive, upset, and angry. With examples shown in Figure 4, the detailed logistics of labeling are as follows:

- **Neutral**: Texts without obvious emotions are labeled as "neutral", including suggestionseeking posts, opinion expressions, etc.
- **Positive**: It refers to posts conveying appreciation and gratitude towards others, optimism, or other positive emotions.
- **Upset**: To be upset is to be disturbed or very unhappy. In Kayuan, these posts are usually correlated with academic pressure and interpersonal relationships.
- Angry: Compared with upset, it is an intense emotion against unsatisfactory outer space. In Kayuan, students usually express angry emotions towards some courses or peers.

#### 2.3 Data Analysis

As shown in Figure 3, neutral posts (77%) occupy the largest proportion of Kayuan corpus, which leads to data imbalance. We managed to deal with this unbalance issue in the follow up sections.





Figure 3: Percentage of emotions

Figure 4: Examples of posts with different emotions

# 3 Model

To achieve a good performance, we chose BERT (Bidirectional Encoder Representations from Transformers) [2], a popular model for NLP tasks nowadays. BERT is a pre-trained model that can be fine-tuned for a wide range of NLP tasks, including text classification, named entity recognition, and question answering. Moreover, Bert uses the Transformer architecture. According to previous studies, it has achieved state-of-the-art results on a wide range of NLP benchmarks.

## 3.1 Training Process

After data processing, the input text is first tokenized into subwords and then mapped to embedding vectors. A classification layer is added on top of the pre-trained BERT model. Then on top of the BERT model, there is a classification layer. The output of the final encoder layer is fed into the classification layer and makes a prediction on the sentiment label, which is shown in Figure 6.

## 3.2 Architecture of Bert

The transformer architecture of BERT greatly contributes to its power. Bert consists of a series of encoder layers, each consisting of two sub-layers: a self-attention layer and a feed-forward neural network layer.

As shown in Figure 5 the self-attention layer is the key component of the Transformer model. It allows the model to capture dependencies between different words in a sentence by computing attention weights for each word. These attention weights indicate how important each word is for predicting the output of the model.

The feed-forward neural network layer applies a non-linear transformation to the output of the self-attention layer. This allows the model to learn complex non-linear relationships between words.

BERT has two parameter intensive settings [4]:

- **BERT-base**: 12 layers, 768 hidden dimensions and 12 attention heads (in transformer) with the total number of parameters, 110M;
- **BERT-large**: 24 layers, 1024 hidden dimensions and 16 attention heads (in transformer) with the total number of parameters, 340M

We chose a popular BERT model for Chinese corpus, RoBERTa-wwm-ext-large (Chinese), which is pre-trained with the whole word masking (WWM) task in Chinese [1].

	Masking	Туре	Data Source	# Tokens	# Steps
BERT (Google)	WordPiece	BERT-base	wiki	0.4B	?
RoBERTa-wwm-ext-large	WWM	BERT-large	wiki+ext	5.4B	2M <sup>MAX512</sup>

Table 1: Comparison of BERT (Google) and RoBERTa-wwm-ext-large

It is a large bert model which contains 24 encoder layers and is connected with 3 fully connected layers for classification. The output of our model is a 4-d vector representing the possibility of 4 emotions: positive, neutral, angry and upset.

#### 3.3 Techniques for solving unbalanced data

As shown in Figure 3, the labels of our training sample is unbalanced. Therefore, it is vital to observe the F1 score of the predicted labels, which takes precision and recall into account:

$$F1 = \frac{2 * (precision * recall)}{precision + recall}$$

In our chosen base models, we found that some of them have high accuracy, but much lower F1 scores(e.g., word2vec+KNN); In order to solve this issue, we used the following two techniques:

- **Resample**: we repeat sampling from a given sample to augment data with rare labels. After our attempts, we find that by resampling texts with labels "positive" and "angry" to 1800, and texts labeled as "upset" to 2500, the model performance improves.
- Focal Loss [3]: we used focal loss ( $\gamma = 2, \alpha_t = 1$  for all t) to substitute cross entropy loss as follows:

$$FL(p_t) = -\alpha_t * (1 - p_t)^{\gamma} * log(p_t)$$

, where  $p_t$  represents the probability that the model predicts the true label t,  $\alpha_t$  is the weight of class t, and  $\gamma$  is a tunable parameter used to adjust the ratio of difficult and easy samples. When  $\gamma = 0$ , Focal Loss is equivalent to cross entropy loss; When  $\gamma$  gets larger, Focal Loss pays more attention to the rare labels.

#### **4** Experiments and Discussion

#### 4.1 Model performance

As shown in Table 2, we choose 2 baseline methods: word-to-vector plus KNN, and Chinese Bert finetuned by the Weibo dataset. For the former, it is a classical machine learning method for sentiment analysis. For the latter one, Chinese Bert is an established Chinese sentiment analysis pre-trained model. While Weibo, the most popular social media in China, is usually used for sentiment analysis. Meanwhile, the data of Weibo is similar to data of LGU Kayuan; both consist of posts and comments and are relative to our daily life.

Method	Training Data	Val		Test		
Method		Acc	F1	Acc	F1	
word2vec+KNN	LGU Kayuan	0.7718	0.4021	0.7561	0.3882	
Chinese Bert	Weibo <sup>1</sup>	0.1593	0.2275	0.1615	0.2387	
Chinese Bert without techniques	LGU Kayuan	0.8123	0.6529	0.8223	0.6911	
Chinese Bert(Ours)	LGU Kayuan	0.8472	0.7745	0.8533	0.7437	

Table 2: Model Performance on LGU Kayuan

However, from Table 2, we see that word-to-vector plus KNN shows low F1 score, and the model finetuned by Weibo data has both low accuracy and F1 score. Therefore, it is necessary to train a new model targeting Kayuan. We tried several techniques to improve the performance, which is shown in Table 3. The result of our final model is much more accurate than the other two.

ID	Weibo Corpus	Resampling	Weighted Cross	Focal Loss	Acc	F1
			Entropy Loss			
i					0.8223	0.6911
ii	$\checkmark$				0.7729	0.6466
iii		$\checkmark$			0.8467	0.7115
iv			$\checkmark$		0.8127	0.6812
v				$\checkmark$	0.8434	0.7242
vi		$\checkmark$		$\checkmark$	0.8533	0.7437

Table 3: Chinese Bert ablation experiment

#### 4.2 Discussion

Table 3 shows the performance of the Chinese Bert model under different data and loss functions, using accuracy (Acc) and F1 score as evaluation metrics.

- (i) The default setting (LGU Kayuan corpus & cross entropy loss) gives 0.8223 accuracy and 0.6911 F1 score
- (ii) Compared to not using Weibo corpus, the accuracy drops by 0.0494, and the F1 score drops by 0.0445. This may be because the Weibo corpus has a large difference from the target dataset, leading to the model's inability to adapt to the target domain.
- (iii) Using the resampling method will improve the performance of the model. Compared to not using resampling, the accuracy increases by 0.0244, and the F1 score increases by 0.0204. This may be because the resampling method can balance the distribution of different categories in the dataset, reducing the impact of class imbalance.
- (iv) Using the weighted cross entropy loss function will reduce the performance of the model. Compared to not using the weighted cross entropy loss function, the accuracy drops by 0.0096, and the F1 score drops by 0.0099. This may be because the weighted cross entropy loss function gives too much penalty to misclassified samples, leading to overfitting of the model.
- (v) Using the focal loss function will improve the performance of the model. Compared to not using the focal loss function, the accuracy increases by 0.0211, and the F1 score increases by 0.0331. This is because the focal loss function can transfer model's attention from easy-to-classify samples to hard-to-classify samples, thereby improving the predictive strength across all classes.
- (vi) Using resampling and focal loss function together can achieve the best effect, with an accuracy of 0.8533 and an F1 score of 0.7437. This shows that these two methods can complement each other and jointly solve the problems of class imbalance and the difference between easy and hard samples in the dataset.

## 5 Conclusion

Although limited by insufficient data due to the time-consuming labeling work, we managed to train a well-performing BERT-based model which can accurately detect the emotion under every Kayuan post. We are planning to share our contributions with student counseling departments so that we can monitor students' mental state and provide them with timely support accordingly.

<sup>&</sup>lt;sup>1</sup>Weibo (a wideused social media) corpus is attained from https://smp2020ewect.github.io/

## References

- [1] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. Revisiting pretrained models for Chinese natural language processing. In *Proceedings of the 2020 Conference* on Empirical Methods in Natural Language Processing: Findings, pages 657–668, Online, November 2020. Association for Computational Linguistics.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [3] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [4] Hu Xu, Bing Liu, Lei Shu, and Philip S Yu. Bert post-training for review reading comprehension and aspect-based sentiment analysis. *arXiv preprint arXiv:1904.02232*, 2019.

# A Appendix

Here are some figures about the detailed architecture of Bert model:



Figure 5: Attention layer



Figure 6: BERT for sentiment analysis