

Factor Detection Task of Cyberbullying Using the Deep Learning Model

Yu-Hsuan Wu¹, Sheng-Wei Huang¹, Wei-Yi Chung², Chen-Chia Yu³, Jheng-Long Wu¹

¹Dept. of Data Science, Soochow University,

²School of Big Data Management, Soochow University,

³Dept. of Language Science, The University of Edinburgh



Introduction



- The large amount of data posted on social media offers much flexibility and opportunities to be discovered in exploring various social issues.
- Cyberbullying is highly lethal to individuals and is difficult to effectively punish or prevent due to its vague definition.
- In this study, textual cyberbullying is the main research object. The main purpose of this study is to **provide a set of analysis methods for cyberbullying**.



Main contributions:

1. Multiple cyberbullying factor corpus can be created by the study because there are no benchmark dataset about the cyberbullying factors to a classification task.
2. The factor detection of cyberbullying classification can be solved by the powerful learning models such as the BERT classifier.
3. The proposed cyberbullying detection model can train from the proposed corpus, and can apply to the big data analysis on the Internet for cyberbullying.

Related work – Definitions and Factors of Cyberbullying



- Bullying can be defined as the bully's intention to cause harm to the victim, repeated bullying behaviors, and power imbalance between the bully and the victim. (Olweus 1994)
- However, cyberbullying is more specific than traditional bullying, and the repetition and power imbalance are more difficult to judge, due to the anonymity mechanism of online discourse.
- Menin et al (2021) made important findings on the definition of bullying based on adolescents' perceptions. They found that "intent" and "perceived harm" are important factors that constitute cyberbullying.
- Define cyberbullying as the commenters have intent to cause harm via words or other media online, and the victim can feel themselves got harm by the comments.

Related work – Extracting Cyberbullying Factor from Definition



- The concept of "intent to harm" can be categorized as aggressive or offensive language and other directly harmful language such as **insults**, **intimidate**, **harassment**. (Rosenthal et al 2021 ; Kalraa et al 2021)
- Zhong et al (2022) found that the language used in cyberbullying contains direct and indirect hurtful sentences. Based on the example of indirectly harmful sentences in their study, this type of sentence can be specifically defined as **sarcasm**.
- **Arousal** refers to the degree of intensity or aggressiveness of the words used and whether the commenter has strong emotions.
- **Valence** is based on the emotion of reading the sentences while being the victim's perspective.

Related work – Existing Corpus for Bullying Detection



- Most of existing corpus for bullying detection are only classified articles or sentences as bullying/non-bullying or only add toxic labels. (Ask.fm, Formspring, Myspace, Twitter)
- Some studies have labeled some factors of bullying such as inflammatory words, offensive language, toxicity, etc., but not enough to be comprehensive. (Maskat et al 2020; Zhang et al 2020; Tang et al 2020)

Related work – Bullying Detection



- Classification method such as SVM, Randomforest, etc., are used to detect bullying articles or sentences. (Kumar et al 2019; Kanan et al 2020; Abarna et al 2022)
- However recent studies have shown that deep neural network-based approaches are more effective than traditional machine learning techniques on detecting bullying. (Emon et al 2022; Ahmed et al 2022; Raj et al 2022)
- Neural network methods have advantages in classifying a wide range of cyberbullying features. Therefore, this study will use a neural network approach in combination with a pre-trained model to build the model.

Model – Data



- Social media platform : PTT
- Post : 125
- Comment : 5926
- Time : Board opening date – January 2022

批踢踢實業坊				聯絡資訊	關於我們
熱門看板	分類看板				
Gossiping	14469	綜合	◎[八卦] 不要政問了捏		
C_Chat	3667	閒談	◎[希洽] 2022GOTY樂透進行中		
Stock	3271	學術	◎[股票] 新板規上路 新聞格式嚴審		
NBA	1964	NBA.	◎[NBA] Dwight Howard 來台打球		
Lifeismoney	1880	省錢	◎[省錢] 省錢板 發文要有省錢點		
WorldCup	1814	足球	◎[世足] Richarlison 六星級射門		
HatePolitics	1549	Hate	◎[政黑]發文前請先詳閱板規!		
Baseball	1499	棒球	◎[棒球] Baseball is life		
basketballTW	796	籃球	◎[台籃] 例行賽 - 嗚呼!板主徵選中		
LoL	730	遊戲	◎[LoL] 帳號轉移開放		
car	665	車車	◎[汽車] 大家要注意板規1-8喔		

← PTT board list

PTT post list →

批踢踢實業坊

> 看板 Gossiping

聯絡資訊

關於我們

看板

精華區

最舊

< 上頁

下頁 >

最新

搜尋文章...

Re: [新聞] 高虹安掃街幫推館長產品 網友質疑:利益交換 uavan

11/25

...

2 [問卦] 南向政策是不是先知 eric112

11/25

...

[新聞] 看我可愛想要我? 通緝犯騎車違規遭小女警 taiwan08

11/25

...

Re: [爆卦] 朱學恆FB:數位中介法還在推?

Model – Data



Author
Post title
Post time

Post content

Comment content

批踢踢實業坊 > 看板 Gossiping

聯絡資訊 關於我們

看板 Gossiping

作者 goodday5566 (好天五六)
標題 [問卦] 小智現在在想什麼?
時間 Sat Nov 26 20:12:51 2022

雖然小智拿到八大師賽冠軍
已經是實質上的世界第一了
不過到了帕底亞後
皮神又要被路邊的新葉喵喵虐了
身為世界第一卻又要在西班牙被菜雞嘲諷
小智現在在想什麼?

--
※ 發信站: 批踢踢實業坊(ptt.cc), 來自: 114.136.239.5 (臺灣)
※ 文章網址: <https://www.ptt.cc/bbs/Gossiping/M.1669464774.A.DE3.html>

推 johnwu: 有人在搞他	123.195.122.52	11/26	20:13
推 kuma525566: 宵夜要吃什麼	42.79.219.105	11/26	20:13
推 akko76815: 晚點吃麥當勞	1.200.113.87	11/26	20:13
推 milkyway168: 看光光了	106.64.130.9	11/26	20:14
→ wison4451: 洗洗睡	111.243.9.241	11/26	20:14
→ sck3612575: 安慰母狗	114.36.218.69	11/26	20:14
推 za755029: 可以寫一篇世界冠軍為例的論文了	118.166.52.11	11/26	20:14
推 onejune: 可以跟比鵬回家了	114.42.12.131	11/26	20:14
推 jkduke: 可以回桃園 接 比雕了	36.236.148.46	11/26	20:15
推 a97586266: 弊案查起來!!!!!!!!!!!!	49.159.166.249	11/26	20:15
推 nkibond: 小智會去帕底亞嗎? 要去也小豪去吧。	101.12.89.183	11/26	20:15

Board

Model – Data Annotation

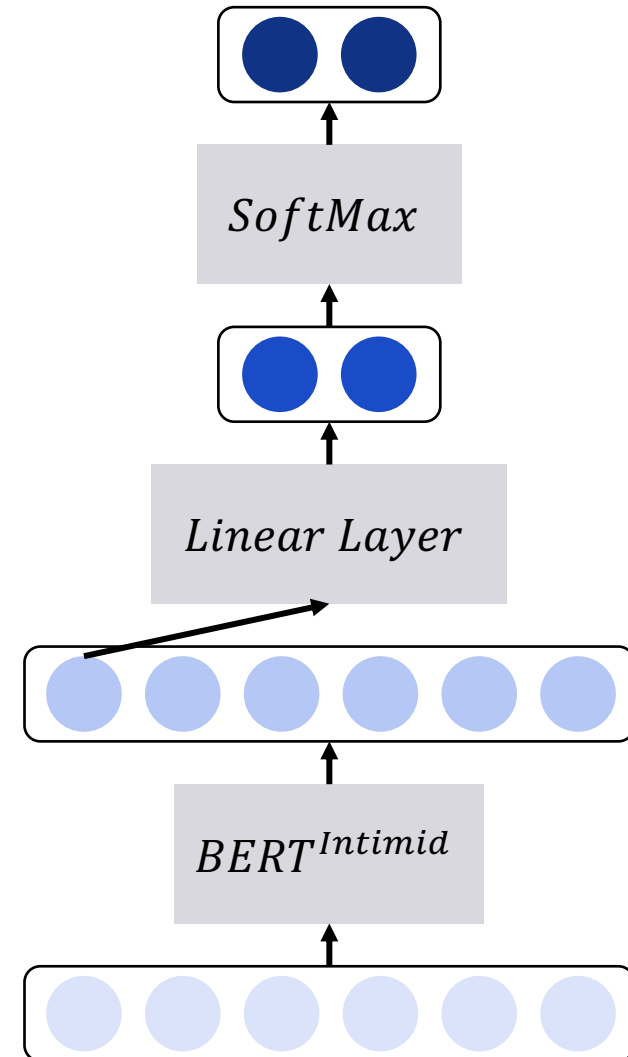


Factors	Labeled as	Define
Intimidation	intimidation/non-intimidation	The use of words to intimidate the main character of the article.
Harassment	harassment/non-harassment	A variety of offensive behavior that is repetitive and unwilling.
Insults	insult/non-insult	Words that contain insulting or demeaning words.
Sarcasm	sarcasm/non-sarcasm	Sarcasm intended to expose the contradiction or shortcomings of the article's main character.
Arousal	arousal/non-arousal	The commenter has clear emotional fluctuation from the sentence.
Valence	positive, neutral, and negative	The emotion that will result from reading a statement from the perspective of the person being bullied.

Model – Classifier



- The input of model is a comment.
- Use **BERT** to extract hidden features vector of input comment.
- Using the first hidden feature vector to estimate the classification probability of cyberbullying factor through **Linear Layer** and **SoftMax**.



Experiment – Annotation Consistency



- Each comment were annotated by **three annotators**. Since each annotator's judgment of the sentences could be biased, the **gold standard was determined by the results of the three annotators' votes**.
- Cronbach's alpha was used to evaluate the quality and consistency of the annotations. Each score has been maintained at above 0.5.

TABLE I. DATA DISTRIBUTION ON THREE ANNOTATORS

Factor	Yes	No	
Intimidation	48	5878	
Harassment	2175	3751	
Insult	1603	4323	
Sarcasm	3641	2285	
Arousal	1618	4308	
	Positive	Neutral	Negative
Valence	37	2598	2349

TABLE II. ANNOTATION CONSISTENCY ON THREE ANNOTATORS

Factor	Cronbach's Alpha
Intimidation	0.56
Harassment	0.65
Insult	0.69
Sarcasm	0.61
Arousal	0.53
Valence	0.63

Experiment – model



The classifiers compared in the experiments:

- BERT_ α : The pre-trained model [distilbert-base-multilingual-cased](#) is used.
- BERT_ β : The pre-trained model [bert-base-multilingual-uncased](#) is used in this BERT model.
- RFC: Considered as a baseline. The TF-IDF method is used to obtain the feature vectors and the random forest classification algorithm is used.
- SVM: Considered as the baseline. The TF-IDF method is used to obtain the feature vectors.

Experiment – All models for experiments



- The results of BERT _{β} is better than another model. The possible reason for this result is that the task of learning whether a comment sentence has cyberbullying factors are relatively complex and difficult, so more hidden features need to be learned to obtain better results.

TABLE III. MACRO AVERAGE PERFORMANCE COMPARISON OF MODELS

	<i>Arousal</i>			<i>Harassment</i>			<i>Insult</i>			<i>Intimidate</i>			<i>Sarcasm</i>			<i>Valence</i>		
	<i>P</i>	<i>R</i>	<i>F</i>	<i>P</i>	<i>R</i>	<i>F</i>	<i>P</i>	<i>R</i>	<i>F</i>	<i>P</i>	<i>R</i>	<i>F</i>	<i>P</i>	<i>R</i>	<i>F</i>	<i>P</i>	<i>R</i>	<i>F</i>
BERT _{α}	0.60	0.60	0.60	0.63	0.63	0.63	0.66	0.66	0.66	0.70	0.70	0.69	0.58	0.58	0.58	0.45	0.42	0.43
BERT _{β}	0.61	0.61	0.61	0.65	0.62	0.63	0.66	0.66	0.66	0.80	0.78	0.77	0.59	0.59	0.59	0.49	0.46	0.47
SVM	0.61	0.54	0.40	0.68	0.52	0.43	0.79	0.52	0.47	0.60	0.51	0.52	0.71	0.51	0.44	0.43	0.26	0.20
RFC	0.62	0.54	0.38	0.67	0.52	0.43	0.77	0.52	0.46	0.60	0.51	0.52	0.67	0.51	0.44	0.41	0.26	0.20

Case study



- Using big data analysis of **two cases** to describes the cyberbullying pattern of PTT.
- The time period was from the beginning of the case until January 2022.
- A total of 224 posts and 39,541 comments related to two events on PTT

TABLE IV. DATA DISTRIBUTION ON EACH EVENT

Event	Post	Comment
Johny Depp	17	2,499
Alisasa	207	37,042

Two kind of graph to analyze the the cyberbullying factor:

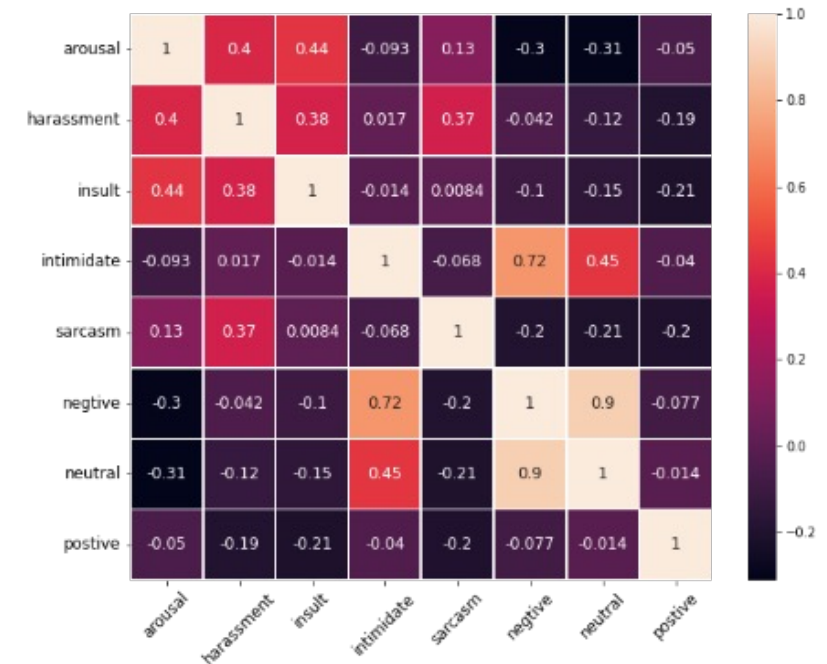
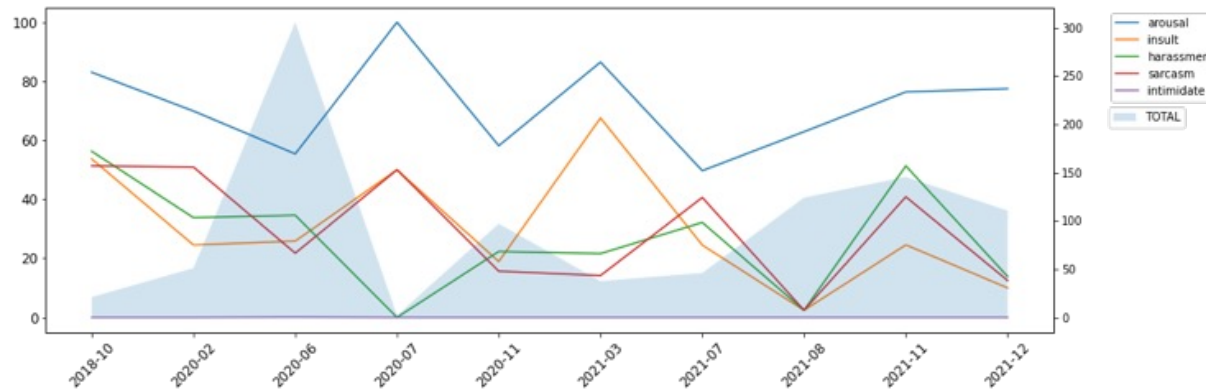
1. **Line Graph**: with the X-axis being the time and the Y-axis being the percentage of comments that match a particular cyberbullying factor to the total number of comments for that time period.
2. **Heatmap**: show the relevance of each factor in a specific case.

Case study – (Celebrity)

The Johnny Depp and Amber Heard Event



- Time: October 2018 to December 2021.
- Each month has a high percentage of comments with arousal and most of the comment are with emotions.
- A similar trend for insults, sarcasm and harassment in the case, with a weak to moderate correlation between these three factors.

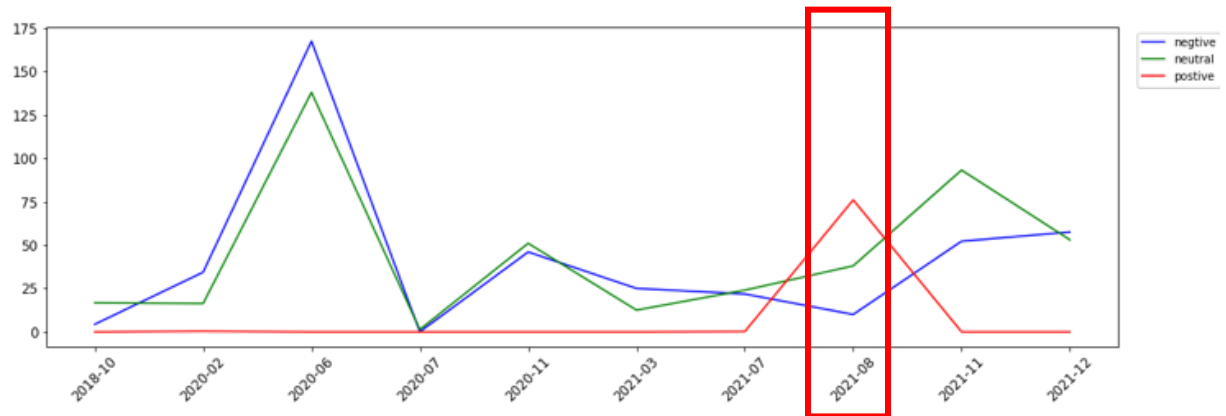
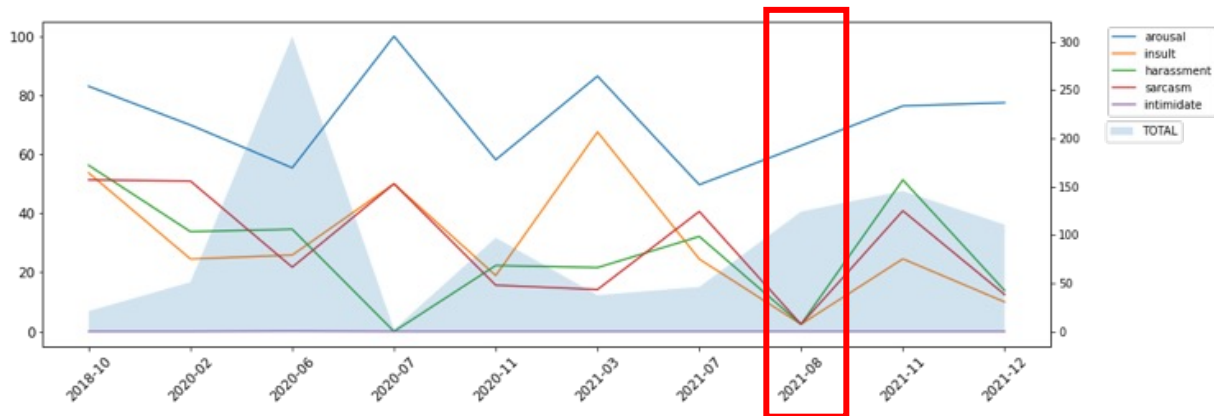


Case study – (Celebrity)

The Johnny Depp and Amber Heard Event



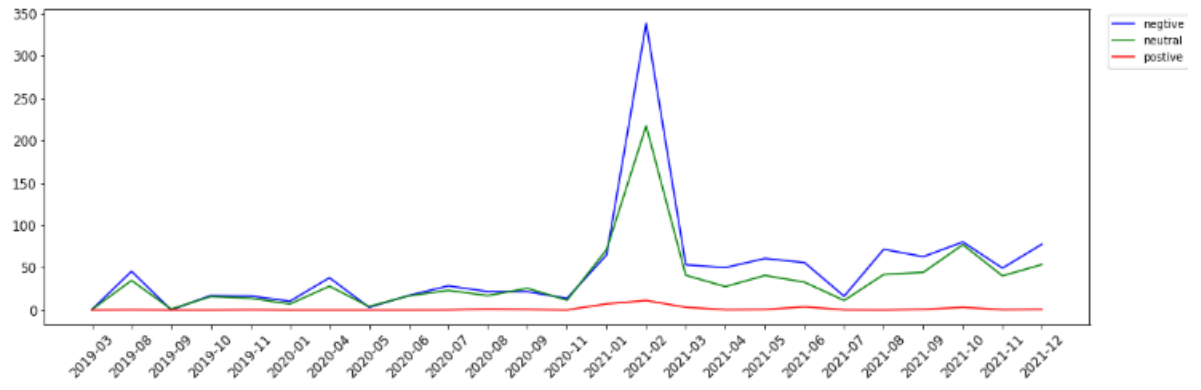
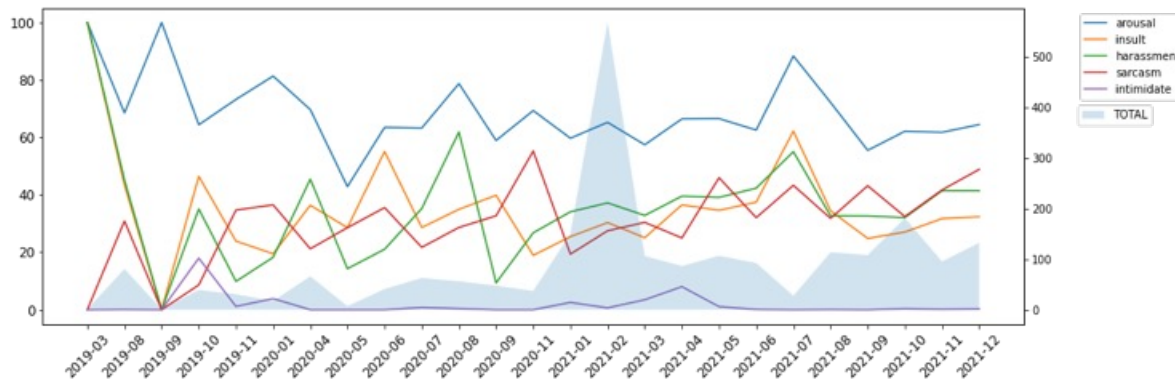
- The condemnation of Johnny in the comments at the outbreak of the case was one-sided due to Amber's accusations against him.
- In August 2021, more neutral tone of the discussion in the forum later in the case, when Johnny fought back against Amber's false allegations.
- Valence reflects the victim's sentiment when reading the comments, and a significant increase in the total number of positive sentiments can be observed, which reflects the beginning of solidarity with Johnny.



Case study – (Internet Celebrity) Alisasa Event



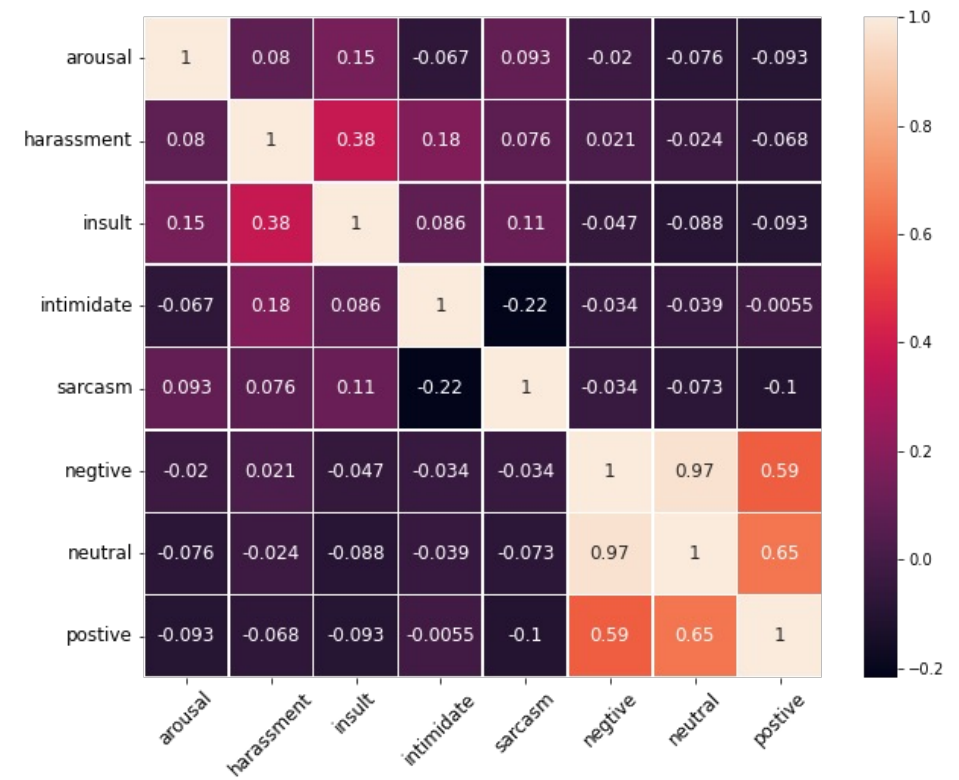
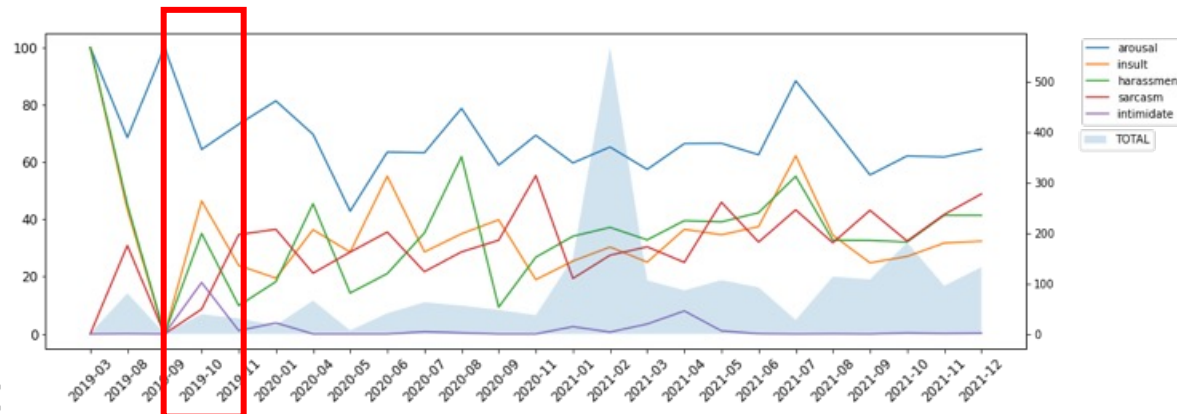
- Time: 2019 to the end of 2021.
- Comments discussing her controversial statements had significant mood waves, which is reflected in the high percentage of arousal.
- Valence reflects the fact that most of the comments made her feel negative emotions after reading them, which means that most of the comments directed at her controversial statements were offensive.



Case study – (Internet Celebrity) Alisasa Event



- The two factors of **harassment** and **insult** show moderate correlations, while sarcasm is less correlated with them, but the trends are partially similar.
- While there are few comments that match intimidation, they are often accompanied by explicit harassment and insults, as seen in the October 2019 comments.





- Based on the case studies, the proposed cyberbullying detection model constructed in this research can detect cyberbullying patterns in real cases.
- The weak prediction performance of arousal and sarcasm by the proposed model reflected the useful hidden features which are difficult to learn.
- The most difficult aspect to the model was the valence factor. Since the definition of valence in this study was **judged from the perspective of the victim**.
- Most of the valence factors in the dataset were negative sentences, which may have made the model difficult for the model to judge positive sentences.

Conclusion



- This study successfully created the Chinese cyberbullying factors corpus, built the BERT classifier, and analyzed the big data of two real cyberbullying cases.
- The predictions of the model can be consistent with the temporal development.
- Victims of cyberbullying comments will be considered, as well as the role each responder takes in cyberbullying.

Thank you for listening.

Thank you for listening.