# Cross-Dataset Generalization for COVID-Related Hate Speech Detection

Student: Xuanyu Su
Supervisor: Isar Nejadgholi

**Text Analytics team**

National Research Council Canada    Conseil national de recherches Canada

Canada

# Table of Content

National Research Council Canada

# Motivation

❖ **Online hate:** Hateful language occurs on the Internet and leads to severe harm for the targeted group and the society in large. NLP is often used to detect this type of language.

❖ **Online hate linked with Covid-19**: Detecting hateful comments with more specific themes; Anti-Asian racism during COVID-19.



Neighboring keywords for the target word 'China' on Twitter. (Rabiul, et al.,2020).

# Research Question

❖ Data annotation is laborious, costly and time consuming.

❖ **Cross-Dataset Generalization:**

➢ Can we use the datasets annotated for other types of abusive behaviour to detect COVID-related hate speech?

❖ **Train Datasets:**

➢ 3 Previously annotated datasets for other types of abusive behaviour

➢ 1 dataset annotated for COVID-related hate speech

❖ **Test Set :Covid Racism dataset**

➢ too small for splitting into train, valid and test sets

# Challenges

Datasets are created with  different task formulations and definitions.

Examples of previously used definitions:

- ➤ **Attack:**  Violent, harmful, or destructive comment against other people.
- ➤ **Toxic**: Makes people leave the platform.
- ➤ **Hate Speech**: Hateful words based on race, ethnicity, nationality, religion and etc..
- ➤ **East Asian Prejudice**: Prejudice against East Asians.

# Train Datasets

## Attack Comment (binary)

- ❖ <u>Source</u>: English Wikipedia
- ❖ <u>Ratio of positive class</u>: 0.10
- ❖ 15362(pos): 159686(neg)
- ❖ <u>Labels</u>: Normal / Personal attack

## Toxicity (binary)

- ❖ <u>Source</u>: English Wikipedia (larger than AC)
- ❖ <u>Ratio of positive class</u>: 0.15
- ❖ 232055(pos): 1366234(neg)
- ❖ <u>Labels</u>: Normal /Toxic

## GAB

- ❖ <u>Source</u>: from social website: <u>gab.ai</u>
- ❖ <u>Ratio of positive class</u>: 0.13
- ❖ 11249(pos): 75280(neg)
- ❖ <u>Labels:</u> Normal / Hate /(13 subcategories)

## East Asian (covid-related)

- ❖ <u>Source:</u> collected from Twitter
- ❖ <u>Ratio of positive class</u>: 0.27
- ❖ 5331(pos): 14669(neg)
- ❖ <u>Labels</u>: Normal / hostile/critical/ counterhate/discussion

National Research Council Canada

# Test Dataset

## Covid Racism

- ❖ <u>Source</u>: Twitter
- ❖ <u>Ratio of positive class</u>: 0.29
- ❖ 678(pos): 1641(neg)
- ❖ <u>Labels</u>: Hate/ Counter-hate/ Neutral/ others

# Binarizing the labels

| Database | Positive | Negative |
|---|---|---|
| **Toxicity** | Toxic | Normal |
| **Attack Comment** | Personal attack | Normal |
| **Gab** | Vulgarity or Offensive language<br>Hate based on race or ethnicity<br>Hate based on nationality/regionalism | Others |
| **Esat Asian** | Entity_directed_hostility<br>Entity_directed_critism | None_of_above,<br>Counter_speech,<br>Discussion_of_eastasian_prejudice |
| **Covid Racism**<br>(test dataset) | Hate | Counter-hate<br>Neutral<br>Others |

# Classifiers

## BERT + Linear

**Pretraining Data:**

- ❖ Wikipedia
- ❖ Books Corpus

## RoBERTa + Linear

**Pretraining Data:**

- ❖ Wikipedia
- ❖ Books Corpus
- ❖ News
- ❖ Other texts from Web

# Classifiers

## BERT + Linear

**Hyperparameters:**
- Epochs: 5
- Learning Rate: Ie-5
- Weight-decay: 0.01
- Batch size: 16
- Optimizer: AdamW
- Criterion: CrossEntropy()

**Evaluations:**
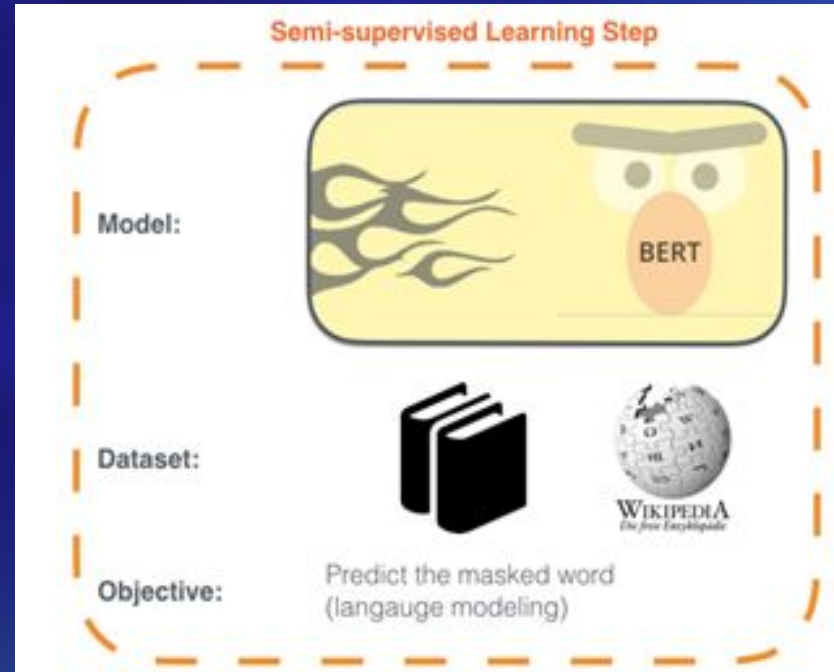- ROC-AUC score
- F1 scores
- Macro Avg

National Research Council Canada



Figure 2: The pretrained datasets and objective of BERT model

# Classifiers



Figure 3: The pretrained datasets and objective of RoBERTa model
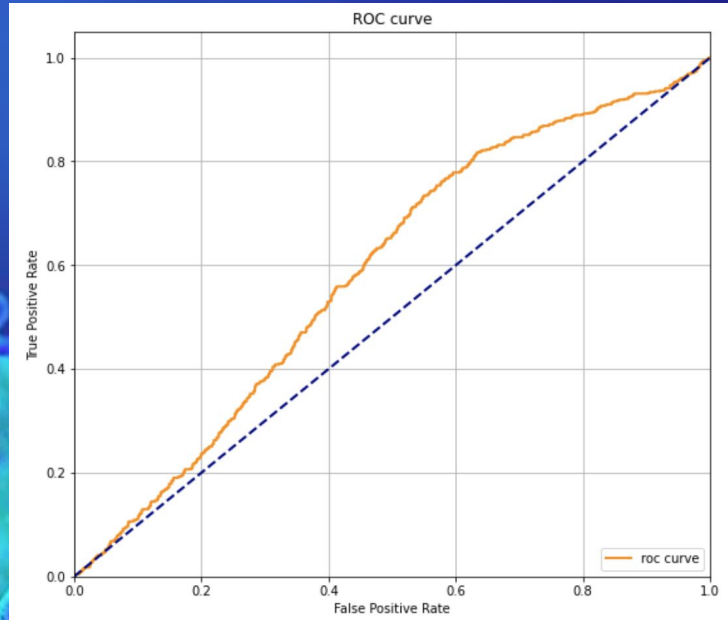
## RoBERTa + Linear

**Hyperparameters:**
- Epochs: 3
- Learning Rate: le-5
- Weight-decay: 0.01
- Batch size: 16
- Optimizer: AdamW
- Criterion: CrossEntropy()

**Evaluations:**
- ROC-AUC score
- F1 scores
- Macro Avg

# Example

E.g.: RoBERTa model trained with EA data, tested on Covid Racism data:



The ROC curve of model

```
Classification Report:
              precision    recall  f1-score   support

           1     0.6142    0.7419    0.6720       678
           0     0.8833    0.8074    0.8437      1641

    accuracy                         0.7883      2319
   macro avg     0.7487    0.7747    0.7578      2319
weighted avg     0.8046    0.7883    0.7935      2319
```
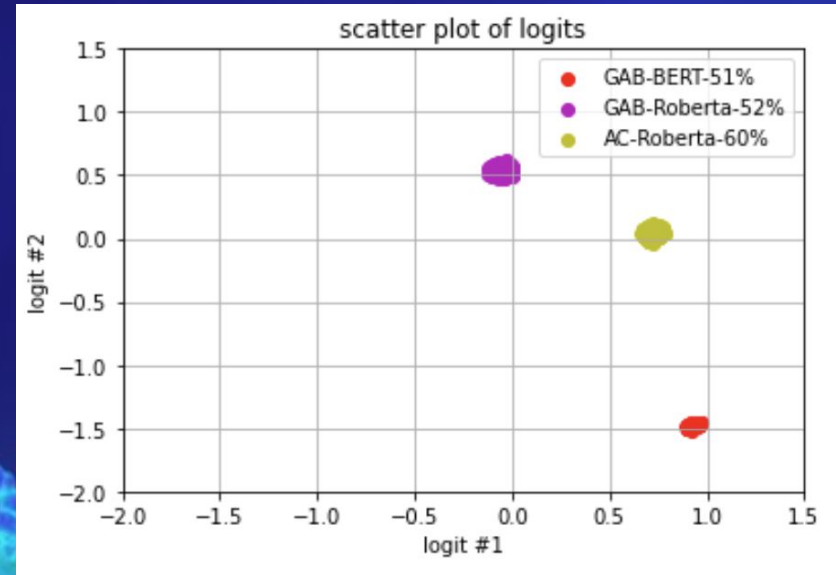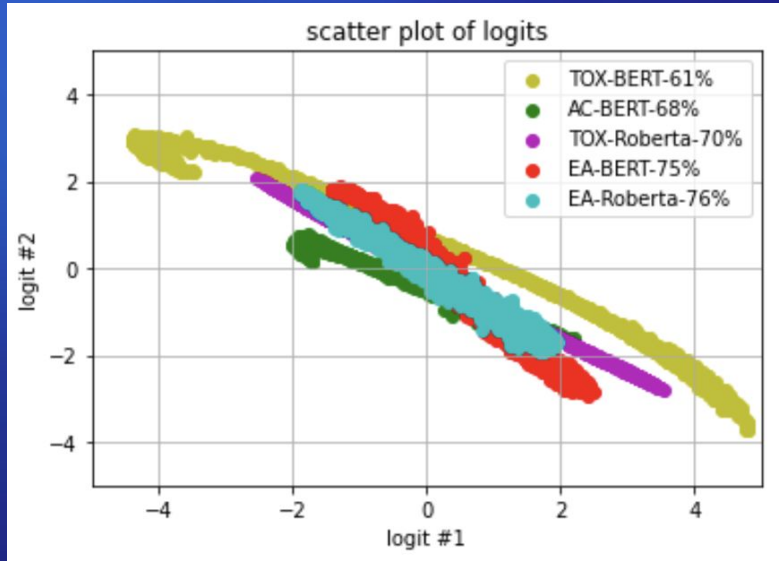
The classification report of model

# Results

Accuracy per class and macro-averaged F1-score on the covid-racism dataset for different models and training datasets

| Training dataset<br><br>classifier | Attack Comment | Gab | EA | Toxicity |
|---|---|---|---|---|
| **BERT+Linear** | 0.6460(pos)<br>0.7228(neg)<br>0.6844(MA) | 0.2439(pos)<br>0.7876(neg)<br>0.5157(MA) | **0.6683(pos)**<br>**0.8242(neg)**<br>**0.7462(MA)** | 0.5240(pos)<br>0.7060(neg)<br>0.6150(MA) |
| **RoBERTa+Linear** | 0.5278(pos)<br>0.6804(neg)<br>0.6041(MA) | 0.2676(pos)<br>0.8023(neg)<br>0.5294(MA) | **0.6720(pos)**<br>**0.8437(neg)**<br>**0.7578(MA)** | 0.6173(pos)<br>0.7974(neg)<br>0.7060(MA) |

# Future Direction

We visualized the distribution of classifier output for each dataset.

Is the distribution of classifier output an indicator of generalizability?

# Conclusion

❖ The EA data set generalized best to the Covid Racism data (both detect attack languages towards Asians during COVID-19).

❖ For the Toxicity, Gab, EA data sets, the RoBERTa model performs better than Bert.

❖ Overall, the Roberta+Linear model trained with EA dataset is the best performing one.

❖ More adequate preprocessing of data (selection and classification of hashtags) could help the model achieve better results in the prediction phase.
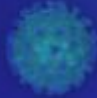
National Research Council Canada

# Knowledge learned during this internship

- ❖ **<u>Academic:</u>**
  - ➢ The structure and use of Bert and RoBERTa model
  - ➢ Steps and operations of data preprocessing
  - ➢ Practical methods to prevent overfitting: weight-decay and dropout
  - ➢ How to evaluate the imbalanced class, etc..
- ❖ **<u>Work and Daily:</u>**
  - ➢ Professional vocabularies in text writing and oral expression
  - ➢ Regular work summary could help more efficient learning and task completion

National Research Council Canada

Special thanks to my supervisor:
Isar Nejadgholi
and other colleagues for their
help during this internship!

National Research Council Canada

# Thank you!

Xuanyu Su
Master student from University of Ottawa
xsu072@uottawa.ca

**Questions?**

National Research Council Canada · Conseil national de recherches Canada

Canada