

# The Dark Side of Language: Using NLP to Combat Hate Speech

Shalini Chawla, shalini\_chawla@berkeley.edu

Adam Hyman, adamhyman@berkeley.edu

Sreeram Ravinoothala, sreeram@berkeley.edu

## Abstract

With the advent of social media, people have had access to propagate hate anonymously, making it more dangerous than ever as the impact is not contained by a social group or geographic area any more. It is very important for social media companies to accurately identify hateful content and a lot of work has been already done to classify online content as toxic vs non-toxic for content moderation. We would like to extend this classification from binary to a multi label classification and identify the specific hate categories toxic content belongs to. This information can be used by socio-political studies exploring the relationship between specific social-political events and their impact on triggering hateful content in online media. In this paper, we fine-tune three pre-trained large language models for a multi-label classification task to categorize text comments as toxic, severe\_toxic, obscene, threat, insult and identity\_hate. To tackle the imbalance of available training data, we use back translation using multiple intermediate to augment training data for the minority classes.

## 1 Introduction

Social media brings people from all areas of society to shared spaces. While a wide-variety of view points is usually a good thing, there is a small minority of speech, hate speech, that social media platforms consider to be hurtful to the platform and the community they're trying to build. Hate speech can be targeted against an individual or group based on inherent characteristics, like race or gender. The offensive speech disrupts social peace on these platforms and upsets the users. The social media platforms use considerable resources to remove hate speech, by allowing users to flag offensive content and with automated algorithms, so the negative effects can be mitigated as soon as possible.

With the advances in machine learning, many of the social media platforms have done a lot of great work in controlling the proliferation of hate speech by making use of classification algorithms but there is much more to be done to reduce the exposure to a great extent. They are successful in many ways but classifying the content to check if it indeed is hate is very complex as it depends on the context as well as whether its used loosely. There is lot of research done or going on classifying the text as well as augmenting such text.

The ability in building accurate models to correctly classify hateful content is a challenging task due to the limited availability of labeled data required for training. The training data needs to be labeled by human annotators. In addition to being a tedious task, it also exposes the annotators to the disturbing content in the data they are required to label. The public data sets available today are highly imbalanced and have very limited samples for some of the minority categories including threat, obscene and severe\_toxic compared to the large amount of non-toxic samples. We follow two approaches to deal with the limitations of data availability.

1) Pre-trained large language models have made it feasible to achieve better accuracy in many NLP tasks where the amount of training data is limited, by providing word embeddings that have already learned characteristics of natural language. Using these pre-learned embeddings as the input to a classification model enables the model to learn to classify the content using much lesser amount of training data than it would have needed to learn the same relationship from scratch.

We use three large language models: BERT, T5 and XLNet and fine-tune them for a multi label classification task using our data set that has been

labeled with the 6 toxicity classes.

2) Our source data set is highly imbalanced and has very limited data for three out of the six categories. We use a combination of approaches to balance our data set to help the model balance learning across all categories. We have three minority categories: threat, severe-toxic and identity\_hate that have very low representation in our training data set. We use the TOXIGEN (Hartvigsen et al., 2022) data set to augment data for the identity\_hate class. To augment data for the other classes, we use back translation to generate additional samples from the existing data.

## 2 Background

There has been a lot of research and previous work done with focus on toxic content classification and a few different pathways have been explored to improve the accuracy of classification. The unavailability of enough quality labeled data has been widely accepted as a limitation in the ability to create accurate models for toxic content classification. (Rastogi et al., 2020) has explored generating synthetic data using EDA and back translation to augment the training data and reported improved recall and F1 scores. (Hartvigsen et al., 2022) have explored the complexity of identifying implicit identity\_hate targeted at minority groups and used GPT-2 to generate additional data to augment the current available human labeled data sets. Another issue has been the inconsistency in labeling across the publicly available data sets limiting the generalization of trained models (Fortuna et al., 2020).

## 3 Data

### 3.1 Source Data Set

We used the dataset from the Jigsaw Toxic Comment Classification Challenge held by Kaggle, that contains Wikipedia comments, that were labelled by human raters for various types of toxic content. The data set consist of 3 files that includes the labeled training data, test data and the test labels (that was made available after the competition concluded). We combine the test data and labels file by joining them on the unique comment identifiers. The training data set has 160,000 comments while the test data set has 64,000 comments. We converted the comments into lowercase and cleaned the contents by removing punctuation and special

characters. Once processed both the training and test data consisted of the text comments and binary labels for the 6 categories:

- toxic
- severe\_toxic
- obscene
- threat
- insult
- identity\_hate

### 3.2 Data Balancing

Our source data is highly imbalanced. Out of the 6 labeled classes, 3 are majority classes with and 3 are minority classes.

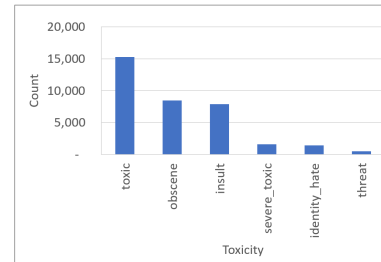


Figure 1: Toxicity Label Count

### 3.3 Augmenting Minority Class Data

We augmented the training data for minority classes to bring it closer in count to the majority classes.

For *severe\_toxic* and *threat*, we used back-translation with Google Translate API to supplement the training data. For the smallest class, *threat*, we back-translated all 476 comments using 3 languages - Chinese, Hindi and French - which along with the original comment quadrupled the number of comments in our training data. For *severe\_toxic*, we back translated 500 comments using Chinese, to supplement our training data.

We supplemented *identity\_hate* with 640 records created by Microsoft’s Toxigen repository, which uses a large scale pre-trained language model and a classifier to generate synthetic hate speech. The repository includes two methods for generating hate speech: Demonstration-Based Prompting and an adversarial set up between a classifier and a generator. We included 320 records from each of the methods.

### 3.4 Undersampling Majority Class Data

To further balance the data set for the 6 identified classes, we selected all comments labelled with the minority classes, including augmented data. Since many of these comments are also labelled with the majority classes, we only needed to select a small amount of additional comments, that were labelled with the majority classes, but not the minority classes, in order to get a data set that was approximately balanced.

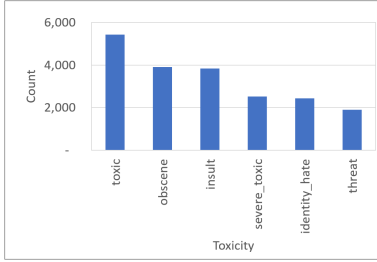


Figure 2: Label Count for Balanced Data

## 4 Methods

### 4.1 Baseline Model - BERT fine-tuned on source data

Our baseline model is BERT base model(cased). We fine-tune the model on the full training data set that has been split into training and validation using 80:20 split and stratified on the 6 label classes. We add a hidden layer with 200 neurons and dropout layer with 0.3 before our final classification layer with 6 outputs to perform a multi label classification. We train the model for 2 epochs using learning rate of .00005 and batch size of 32.

The model performs really poorly for our 3 minority classes, completely failing to identify any of the severe\_toxic, threat and identity\_hate comments.

### 4.2 BERT fine-tuned with balanced data

Our baseline model’s poor performance is due to the much smaller amount of training data for the minority classes. To model is favoring the majority classes and is learning to just return the label as one of the majority classes to maximize accuracy since that is what most of the data is. To help our model learn the features of the minority classes, we create a new balanced data set by undersampling our majority classes and augmenting the minority classes. We run the BERT model with the same hyperparameters as in our baseline with the balanced

class	prec	recall	f1-score
toxic	0.62	0.40	0.49
severe_toxic	0.00	0.00	0.00
obscene	0.66	0.32	0.43
threat	0.00	0.00	0.00
insult	0.65	0.25	0.36
identity_hate	0.00	0.00	0.00
micro avg	0.64	0.31	0.42
macro avg	0.32	0.16	0.21
weighted avf	0.58	0.31	0.40
samples avg	0.04	0.03	0.03

Table 1: Baseline Model - Classification Report

data set. Balancing the data set did not improve the classification score for minority classes at all it still stayed flat at 0 with none of the 3 minority classes being identified. Our data set size was reduced drastically after undersampling but we still had 13,000 records which is an acceptable amount of data for fine-tuning BERT. Since BERT was originally trained on wikipedia data and has not previously seen a majority of the vocabulary that is in the toxic data, it is not able to learn enough from the provided data to identify the minority classes.

### 4.3 Retraining BERT layers

Our next step was to systematically unfreeze the pre-trained BERT layers 2 at a time and retrain them using our data set. We used the original unbalanced data set and the balanced dataset we created and re-trained the top 2 BERT layers, followed by 4, 6 and all the layers. The first classifier with retraining only the top 2 BERT layers gave as a huge improvement in score for all three of the minority classes. As we retrained deeper layers, the total score kept moving upwards. For the full unbalanced dataset, the improvement stagnates after 4 layers. This has been observed in previous work done by (Singh, 2020). But we observed that with the smaller balanced dataset, the improvement continues as we retrain more layers with full retraining of all 12 layers giving us the best scores.

### 4.4 T5

We used T5-small (Raffel et al., 2020) on the unbalanced training dataset similar to BERT. The model was adopted from (Jaunjale, 2020). As part of the multilabel classification a language modeling head is used on top of the decoder. The model was run 2

epochs with 8 and 16 batch size. The lesser epochs was because of the time the model was taking to train and predict. The model was very biased towards the higher dataset and also failed to predict the lower classes (insult, identity\_hate) which can be seen from the lower ROC-AUC scores.

The same experiment was repeated with balanced dataset for 2 epochs and 16 batch size. The results were similar to what was seen with unbalanced dataset. It was still not predicting the classes insult and identity\_hate properly.

With above results, it was not prudent for us to pursue with T5 due to no improvement in the metrics.

#### 4.5 XLNet

Once T5 didn't give an improvement over Bert, we examined models that use a different approach to pre-training and implemented XLNet ((Yang et al., 2019), which is an auto-regressive based model that does a better job of learning dependency pairs.

We ran XLNet on both the unbalanced and balanced data, for 6 epochs on each.

We found that XLNet gave us a modest improvement over Bert, both with balanced and unbalanced data.

### 5 Results and Discussion

Metrics were captured across multiple models that were tested as shown in Table 2 and Table 3. Each model was run for multiple epochs to ensure proper outcome of metrics. We captured macro AUC-ROC scores, macro F1 score and AUC-ROC scores for different labels across different runs.

The baseline BERT model (with no re training) performed very poorly on the minority classes for both Unbalanced and Augmented Balanced dataset. We fine tuned the model by unfreezing multiple layers and training BERT to learn the toxic data. As can be seen in the tables, there was a significant improvement in the F1, AUC-ROC scores. We also saw the minority classes performing well in this situation.

T5 was applied on unbalanced data and balanced data. The observations showed that T5 does worse on all of our metrics compared to BERT and XLNet with the exception of classifying se-

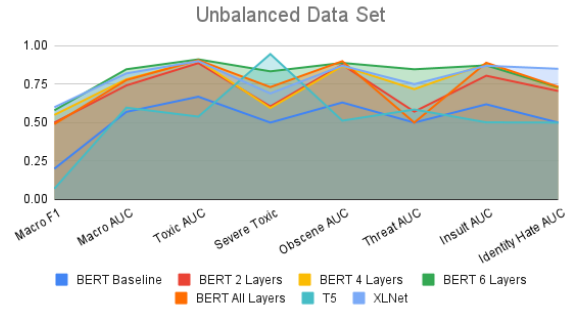


Figure 3: Unbalanced Data set

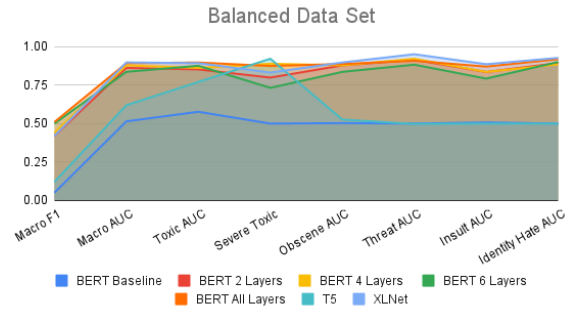


Figure 4: Balanced Data set

vere\_toxic data where it outperforms all the other runs with huge margin.

Similarly XLNet model was applied on both unbalanced as well as balanced dataset. We observed that XLNet outperformed all of our runs with other models for both the datasets.

Overall, the metrics across all the labels are more uniformly distributed when the models are trained on the balanced dataset. Referring to Figures 3 and 4 there is a clear indication of how each model shines given a balanced dataset.

Model	Macro F1	Macro AUC	Toxic AUC	Severe Toxic AUC	Obscene AUC	Threat AUC	Insult AUC	Id Hate AUC
BERT Baseline	0.20	0.56	0.66	0.50	0.63	0.50	0.61	0.50
BERT 2-layers	0.50	0.74	0.88	0.60	0.87	0.57	0.80	0.71
BERT 4-layers	0.55	0.78	0.90	0.59	0.86	0.71	0.88	0.72
BERT 6-layers	0.58	0.85	0.91	0.83	0.89	0.85	0.87	0.73
BERT 12-layers	0.49	0.77	0.90	0.73	0.89	0.49	0.88	0.73
T5	0.07	0.59	0.53	0.94	0.51	0.58	0.50	0.50
XLNet	0.60	0.82	0.90	0.69	0.87	0.75	0.87	0.85

Table 2: Metrics - All Models trained with original unbalanced data

Model	Macro F1	Macro AUC	Toxic AUC	Severe Toxic AUC	Obscene AUC	Threat AUC	Insult AUC	Id Hate AUC
Bert Baseline	0.05	0.51	0.57	0.5	0.50	0.49	0.50	0.50
Bert 2 layers	0.42	0.86	0.85	0.79	0.87	0.91	0.83	0.89
Bert 4 layers	0.44	0.87	0.60	0.88	0.87	0.92	0.83	0.88
Bert 6 layers	0.50	0.83	0.87	0.73	0.83	0.88	0.79	0.90
Bert 12 layers	0.51	0.89	0.89	0.87	0.88	0.90	0.87	0.91
T5	0.12	0.61	0.77	0.92	0.52	0.49	0.50	0.50
XLNet	0.42	0.90	0.89	0.83	0.90	0.95	0.89	0.93

Table 3: Metrics - All Models trained with augmented balanced data

## 6 Conclusion

In conclusion we found that the benefits of a pre-trained large language model are limited by the vocabulary that it was trained on. With Bert being trained only on wikipedia data it hasn't seen the same level of toxic content as in our training dataset which makes it harder to classify toxic comments correctly. T5 was pre-trained on a big corpus of data (C4) where it would have picked up some toxic labels. The metrics for severe\_toxic in Table 3 as well as 2 clearly shows that it is biased towards that label due to some pre-learning whereas few other classes were left out. With XLNet outshining the other two models in the metrics, we can clearly say that XLNet can be used for multi label classification on toxic data

## 7 Future work

Our experiments with different models were based on the toxic dataset and some augmentation that was done to balance it. One area we did not focus on is unbiased training or providing unbiased dataset. So, there may be some biases that are exhibited with the training we did for the dataset we obtained. This is one area we can look at in the future.

For future work we would like to evaluate and optimize our model for fairness across sensitive groups inspired by (Baldini et al., 2021).

We have seen that each of the model has its own area of strength. So, we would like to create and ensemble model derived from the strengths shown by BERT, T5 and XLNet.

## References

- Baldini, I., Wei, D., Ramamurthy, K. N., Yurochkin, M., and Singh, M. (2021). Your fairness may vary: Pretrained language model fairness in toxic text classification. *arXiv preprint arXiv:2108.01250*.
- Fortuna, P., Soler, J., and Wanner, L. (2020). Toxic, hateful, offensive or abusive? what are we really classifying? an empirical analysis of hate speech datasets. In *Proceedings of the 12th language resources and evaluation conference*, pages 6786–6794.
- Hartvigsen, T., Gabriel, S., Palangi, H., Sap, M., Ray, D., and Kamar, E. (2022). Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. *arXiv preprint arXiv:2203.09509*.
- Jaunjale, P. (2020). T5 - Multi Label Classification. <https://www.kaggle.com/code/prithvijjaunjale/>

[t5-multi-label-classification/notebook](#). [Online; accessed Apr-2023].

- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Rastogi, C., Mofid, N., and Hsiao, F.-I. (2020). Can we achieve more with less? exploring data augmentation for toxic comment classification. *arXiv preprint arXiv:2007.00875*.
- Singh, T. (2020). How much does pre-trained information help? partially re-initializing bert during fine-tuning to analyze the contribution of layers.
- Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. R., and Le, Q. V. (2019). Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.