

# Team members:

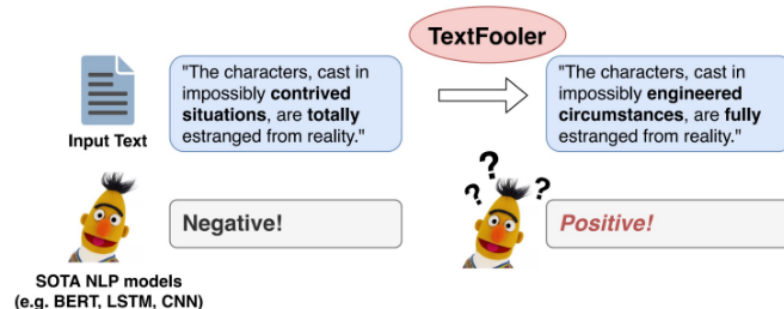
- Gia Trong Nguyen
- Nikita Sergeev

# Current progress:

Implementing attack algorithm from scratch

- Algorithm was taken from this [article](#). In the previous week we used an implemented version of this algorithm by *textattack* python package, but it's not good to use the already implemented code, because it gives zero understanding. So, we decided to implement this algorithm from scratch for better understanding of nlp adversarial attack.
- Some information about this algorithm:

Goal	Constraints	Transformation	Search Method
Untargeted Classification	1. Word embedding distance 2. USE sentence similarity 3. POS consistency	Word substitution by counter-fitted GloVe embedding space	Greedy search with word importance ranking



- Algorithm look like this:

## Evasion Attacks: TextFooler

### • Algorithm

#### Algorithm 1 Adversarial Attack by TEXTFOOLER

**Input:** Sentence example  $X = \{w_1, w_2, \dots, w_n\}$ , the corresponding ground truth label  $Y$ , target model  $F$ , sentence similarity function  $\text{Sim}(\cdot)$ , sentence similarity threshold  $\epsilon$ , word embeddings  $\text{Emb}$  over the vocabulary  $\text{Vocab}$ .

**Output:** Adversarial example  $X_{\text{adv}}$

```

1: Initialization:  $X_{\text{adv}} \leftarrow X$ 
2: for each word  $w_i$  in  $X$  do
3:   Compute the importance score  $I_{w_i}$  via Eq. (2)
4: end for
5:
6: Create a set  $W$  of all words  $w_i \in X$  sorted by the descending
7:   order of their importance score  $I_{w_i}$ .
8: Filter out the stop words in  $W$ .
9: for each word  $w_j$  in  $W$  do
10:   Initiate the set of candidates  $\text{CANDIDATES}$  by extracting
11:     the top  $N$  synonyms using  $\text{CosSim}(\text{Emb}_{w_j}, \text{Emb}_{\text{word}})$  for
12:     each word in  $\text{Vocab}$ .
13:    $\text{CANDIDATES} \leftarrow \text{POSFilter}(\text{CANDIDATES})$ 
14:    $\text{FINCANDIDATES} \leftarrow \{\}$ 
15:   for  $c_k$  in  $\text{CANDIDATES}$  do
16:      $X' \leftarrow \text{Replace } w_j \text{ with } c_k \text{ in } X_{\text{adv}}$ 
17:     if  $\text{Sim}(X', X_{\text{adv}}) > \epsilon$  then
18:       Add  $c_k$  to the set  $\text{FINCANDIDATES}$ 
19:        $Y_k \leftarrow F(X')$ 
20:        $P_k \leftarrow F_{Y_k}(X')$ 
21:     end if
22:   end for
23:   if there exists  $c_k$  whose prediction result  $Y_k \neq Y$  then
24:     In  $\text{FINCANDIDATES}$ , only keep the candidates  $c_k$  whose
25:       prediction result  $Y_k \neq Y$ 
26:      $c^* \leftarrow \underset{c_k \in \text{FINCANDIDATES}}{\text{argmax}} \text{Sim}(X, X'_{w_j \rightarrow c})$ 
27:      $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$ 
28:     return  $X_{\text{adv}}$ 
29:   else if  $P_{Y_k}(X_{\text{adv}}) > \min_{c_k \in \text{FINCANDIDATES}} P_k$  then
30:      $c^* \leftarrow \underset{c_k \in \text{FINCANDIDATES}}{\text{argmin}} P_k$ 
31:      $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$ 
32:   end if
33: end for
34: return None
  
```

- Let's talk a little bit about constraints. As you can see, one of the steps of implementing is to implement constraints with cosine similarity. We need to define a corpus of synonym words. Our choice is the dataset from this [article](#). Then we compute for every word top 50 synonyms. It was a hard part of our work, because we did not have a lot of resources. POS filter we do using *nltk* package. Stop words we also take from *nltk* package. To compute similarity of sentences we use [Universal Sentence Encoder](#).
- Implementation of algorithm you can find in **TextFooler FROM SCRATCH** block in our notebook.
- **We attack only toxic labels.**

## Get the result:

- To understand if our attack works or not we try to attack 5% of train and valid datasets. Then make the pandas dataframe to visualize the results. Here is an example:

Original Label	Attacked Label	Original	Attacked	Original target probability	Attacked target probability	Similarity of Sentences	Is Attack succeed?	
25677	1	1	<p>while i would dearly <b>love</b> to have my way with her in a wild and <b>playfulness</b> means , join needs to <b>re-exploring</b> she is chinese and thus not <b>that</b> <b>chou</b> american . <b>theology</b> <b>fucked</b> that woman is confused . and <b>observes</b> her with luo zhiqiang in that <b>video</b> <b>shagging</b> to dance like some <b>lips</b> <b>breaks</b> slut is just embarrassing . <b>diskindofsexist</b> <b>mntokhere</b> ??? ? ? ? ? ? ? ? wmf dapits ((( ( 202 . 86 . 135 . 18</p>	<p>while i would dearly <b>love</b> to have my way with her in a wild and <b>playfulness</b> means , join needs to <b>re-exploring</b> she is chinese and thus not <b>that</b> <b>chou</b> american . <b>theology</b> <b>fucked</b> that woman is confused . and <b>observes</b> her with luo zhiqiang in that <b>unreliable</b> video <b>shagging</b> to dance like some <b>lips</b> <b>breaks</b> slut is just embarrassing . <b>diskindofsexist</b> <b>mntokhere</b> ??? ? ? ? ? ? ? ? wmf dapits ((( ( 202 . 86 . 135 . 18</p>	0.977109	0.887330	0.500560	No
71182	1	0	<p>look at your own <b>absurd</b> behavior first . by the way , i'm a different user than the one you seem to be writing to about , even though i agree with this user more than i agree with you .</p>	<p>look at your own <b>absurd</b> behavior first . by the way , i'm a different user than the one you seem to be writing to about , even though i agree with this user more than i agree with you .</p>	0.693431	0.228694	0.813941	Yes
135376	1	1	<p>she 's probaly a gurgurshite .</p>	<p>she 's probaly a gurgurshite .</p>	0.716491	0.716491	0.999655	No
125802	0	0	<p> , bulgarians , have nothing in common with any mongols or turks . this is an evil campaign started by our enemies about a hundred years ago to divide the bulgarians . only in the russian federation there are at least 20 million people of direct bulgarian origin , and the russian authorities simply do n't want them united because this is a threat to their empire , no wonder this article is within the scope of wikiproject russia , they support this anti bulgarian turkic crappy propaganda that has nothing to do with reality and is extremely abusive against us . bulgarians . we know very well who we are and where do we come from , we need no stupid shitopedia to brainwash us with this pathetic neo turkic crap .</p>	<p> , we , bulgarians , have nothing in common with any mongols or turks . this is an evil campaign started by our enemies about a hundred years ago to divide the bulgarians . only in the russian federation there are at least 20 million people of direct bulgarian origin , and the russian authorities simply do n't want them united because this is a threat to their empire , no wonder this article is within the scope of wikiproject russia , they support this anti bulgarian turkic crappy propaganda that has nothing to do with reality and is extremely abusive against us . bulgarians . we know very well who we are and where do we come from , we</p>	0.888947	0.888947	1.000000	Skipped
86723	1	1	<p>yo you are a <b>big</b> mother <b>idiot</b> <b>losses</b> .</p>	<p>yo you are a <b>grasly</b> big mother <b>damned</b> idiot <b>losses</b> .</p>	0.981596	0.955354	0.505714	No

Meaning of columns:

\* *Original Label*\* - Label, which model predict

\* *Attacked Label*\* - Label, after attacking

\* *Original*\* - Original text

\* *Attacked*\* - Attacked text

\* *Original target probability*\* - Probability of Label, which model predict

\* *Attacked target probability*\* - Probability of original Label, after attacking

\* *Similarity of Sentences* \* - Cosine similarity, which show us semantic similarity of sentences

\* *Is Attack succeed?* \* - Attack succeeded or not

- Some statistic:  
5% of Train dataset and 5% of Val dataset

```
Yes      279
No       170
Skipped   51
Name: Is Attack succeed?, dtype: int64
```

```
No       22
Yes      21
Skipped   7
Name: Is Attack succeed?, dtype: int64
```

As you can see, our attack worked.

- \* Yes - Attack is success
- \* No - Attack is failed
- \* Skipped - Model predict toxic label as non-toxic

## Train robust model (adversarial training)

- To train a robust model we add success attack samples to our original dataset and then we retrain the model.

- So, after the previous step we compare the *performance* of two models (**Compute accuracy only for toxic label, exclude non-toxic samples**).

```
Comparing performance of original and adversarial trained models on valid dataset with adversarial examples.
```

```
[ ] 1 from IPython.display import clear_output
    2
    3 true, pred = evaluate(model, prepared_toxic_val_iterator, verbose=False)
    4 clear_output()
    5 print(f"Accuracy score of the Original model: {accuracy_score(true, pred)}")
```

```
Accuracy score of the Original model: 0.8349328214971209
```

```
[ ] 1 true, pred = evaluate(adversarial_model, prepared_toxic_val_iterator, verbose=False)
    2 clear_output()
    3 print(f"Accuracy score of the Adversarial trained model: {accuracy_score(true, pred)}")
```

```
Accuracy score of the Adversarial trained model: 0.9385796545105566
```

**Conclusion:** Here, as you can see, model, which not trained on adversarial examples give bad performance comparing with model trained on adversarial attack.

```
Comparing performance of original and adversarial trained models on valid original dataset.
```

```
[ ] 1 true, pred = evaluate(model, original_toxic_val_iterator, verbose=False)
    2 clear_output()
    3 print(f"Accuracy score of the Original model: {accuracy_score(true, pred)}")
```

```
Accuracy score of the Original model: 0.87
```

```
[ ] 1 true, pred = evaluate(adversarial_model, original_toxic_val_iterator, verbose=False)
    2 clear_output()
    3 print(f"Accuracy score of the Original model: {accuracy_score(true, pred)}")
```

```
Accuracy score of the Original model: 0.938
```

- To evaluate *robustness* of two models we use this formula:

$$\text{attack success rate} = \frac{\text{\# of successful attacks}}{\text{\# of total attacks}}$$

```
Original model attack success rate: 0.621380846325167
Adversarial trained model attack success rate: 0.4863157894736842
```

As we can see, adversarial trained model become more robust than the original one.

## Each member contribution:

Gia Trong - Implementing main algorithm of adversarial attack, Train models, Evaluate the models, Write Report

Nikita - Compute top 50 synonyms of the corpus, Prepare data for training and evaluating processes.

## Future work:

- Fixing some little bugs. Bring beauty.
- Try to use other methods of attack from *textattack* package
- Prepare presentation for defending project