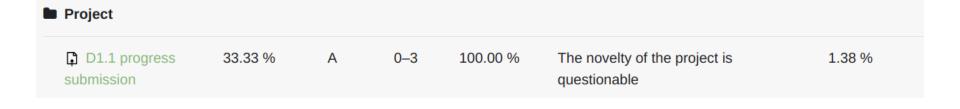


Real or Not? NLP with Disaster Tweets

[F23] PracticalMLandDL Project 2023

Roman Makarov & Adela Krylova





Real or Not? NLP with Disaster Tweets

[F23] PracticalMLandDL Project 2023

Roman Makarov & Adela Krylova



Disaster tweet Generation and Detection

[F23] PracticalMLandDL Project 2023

Roman Makarov & Adela Krylova

PROBLEM

19.4% Of active Twitter accounts are fake or spam

38.2% Americans have accidentaly shared fake news

Of Americans use Twitter as a main platform to read news

11/1

Fake tweets detection on Kaggle

- Competition on Kaggle "Natural Language Processing with Disaster Tweets"
- 10K of real and fake tweets dataset
- Best solution has 96% accuracy
- However, the dataset is quite limited and it contains low quality fake tweets.
- Correct labels were leaked and utilized to mess up leaderboard



Our solution

Generator

GPT2 pre-trained on a few real tweets and fine-tuned along with Discriminator model



Discriminator

Small Bert fine-tuned on GPT-generated data and 300 real tweets.



Examples of generated Tweets









Examples of tweets which tricked the detector





Tweet1

Iranian nuclear weapon found:
Iranian government officials
said Thursday that they
recovered two nuclear-capable
ICBMs from a remote area of
southwest Iran.
http://t.co/rH6J3Wy1MtLpic







Tweet2

Emergency responders
evacuate hostages from a
Saudi Arabia-led coalition
warring...
https://t.co/IXoR1VXyKm8
#bbc0My heart goes out to all
the families affected by this









Tweet3

I understand everyone is upset about the MH17 disaster but why not take a trip down memory lane and send the message that computers can be trusted? #TrustFund







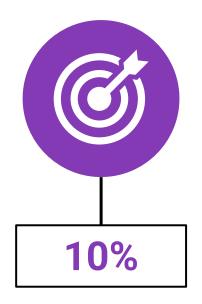
How good is our solution?



Accuracy Increase

Was achieved after fine-tuning Discriminator and Generator

How good is our solution?



Accuracy Increase

Was achieved after fine-tuning Discriminator and Generator

DistilBERT trained on 10000

samples: 0.84 accuracy

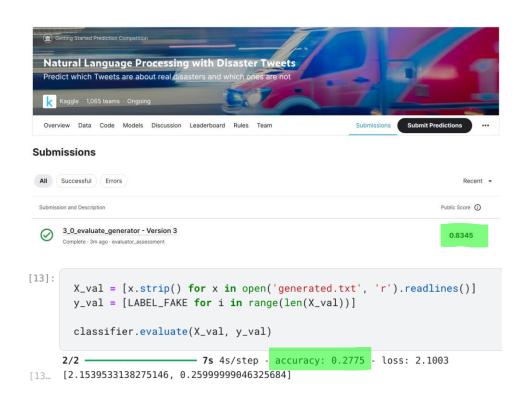
Discriminator trained on 300 real

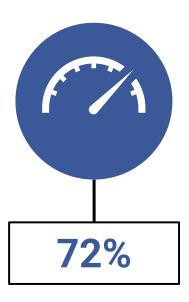
and 100 generated samples:

0.68 accuracy

With only 3% of data!

How good is our solution?





Generations

Were able to fool one of the top Kaggle competition model

Limitations

Link generation

All generated links are not working; hence, the fakeness could be check by the links

Merged Tweets

Some generated Tweets comprise information of several disasters at the same time. This tricks the discriminator



Unstable convergence

Time and memory

The solution consume quite a large number of memory and time

Thank you for your attention



"Times change! Before, no one believed the weather forecast.

Today in the news only the weather forecast can be trusted."

— Ljupka Cvetanova, Yet Another New Land

