

GeDi: Generative Discriminator Guided Sequence Generation

CS475 Team 12

Reinatt Hansel Wijaya
Muhammad Izaaz Inhar Ramahdani
Ukho Shin



Abstracts



Large Language models are able to imitate the distribution of natural language to generate realistic text. However, it is difficult to control this generation. This is a problem because large scale language models have a lot of dataset with toxicity, bias, and negativity. GeDi is one way to solve this issue. GeDi guides generation at each step by computing classification probability. GeDi results in faster generation speed and better controllability. GeDi is able to control the sentiment, toxicity, and topics of language models. To improve it, we train GeDi using a new smaller dataset. We use roberta-base-emotion to evaluate this new trained GeDi and it works like the original generation even with faster training time. In addition to that, we try to improve GeDi by experimenting with negative weight. We find that GeDi is able to make sentence more toxic using negative logit



01

INTRODUCTIONS

Backgrounds
Approach

Backgrounds



**Megatron
Turing NLG**

530 billion
parameters

GPT-3

175 billion
parameters

Language generation has been growing really fast with a large amount of datasets. However, it also has its drawbacks. Due to the large training datasets and long training time, it is harder than ever to control the model. This is why a guided generation is important.

The increase in datasets mean there are more chances of unclean, biased, and 'bad' data. Guided generation is another way to prevent this

Approach

Improvement 1: Improving GeDi

Train GeDi with another small dataset to improve it further

01

02

Guided Generation

GeDi is one of the way to counter this problem by using CCLM as generative discriminator

03

Improvement 2: Append Negative Weight

We try several method to append negative weight (ω) to GeDi for other applications of GeDi

Data and Experiments

Dataset
Model and Algorithm
Comparisons



Datasets



Bookcorpus

(Zhu et al., 2015) for
sentiment analysis



RealToxicityPrompts

(Gehman et al., 2020)
for detoxification



AG NEWS Topic

Classification

(Zhang et al. 2015) for
multi class settings



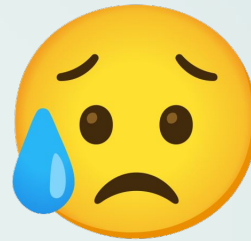
Emotions

Classification

For multi emotion
settings



Happiness



Sadness



Neutral



Love

Models and Algorithm Replication

Data Preparation & Experiment

We use RealToxicityPrompts dataset and AGNews dataset to test for toxicity and multiple topics

Toxicity, Multiple Topics

We use the same Google Perspective API to test the generated sentence. However, we use GPT2-medium because our computing power is not enough

Evaluation

We evaluate the toxicity using Google Perspective API and the result shows the same. We also check for multiple topics generation

Models and Algorithm Improvements 1

Data Preparation & Experiment

We use the datasets from Kaggle of emotions classifications and process it to be trained using GeDi. Then we train GeDi using this dataset

Emotion Analysis

Emotion analysis comes from sentiment analysis and correlation, there are 2 ways we measure the sentiment, polarity and subjectivity.

Evaluation

To evaluate the result of the generated text, we would use the emotion analysis and compare it with the original emotion classification. We use roberta-base-emotion to do this evaluation

Models and Algorithm Improvements 2

Data Preparation

We use the same dataset as detoxification training in original GeDi which is RealToxicityPrompt. However, we experiment with negative weight value

Toxicity Analysis

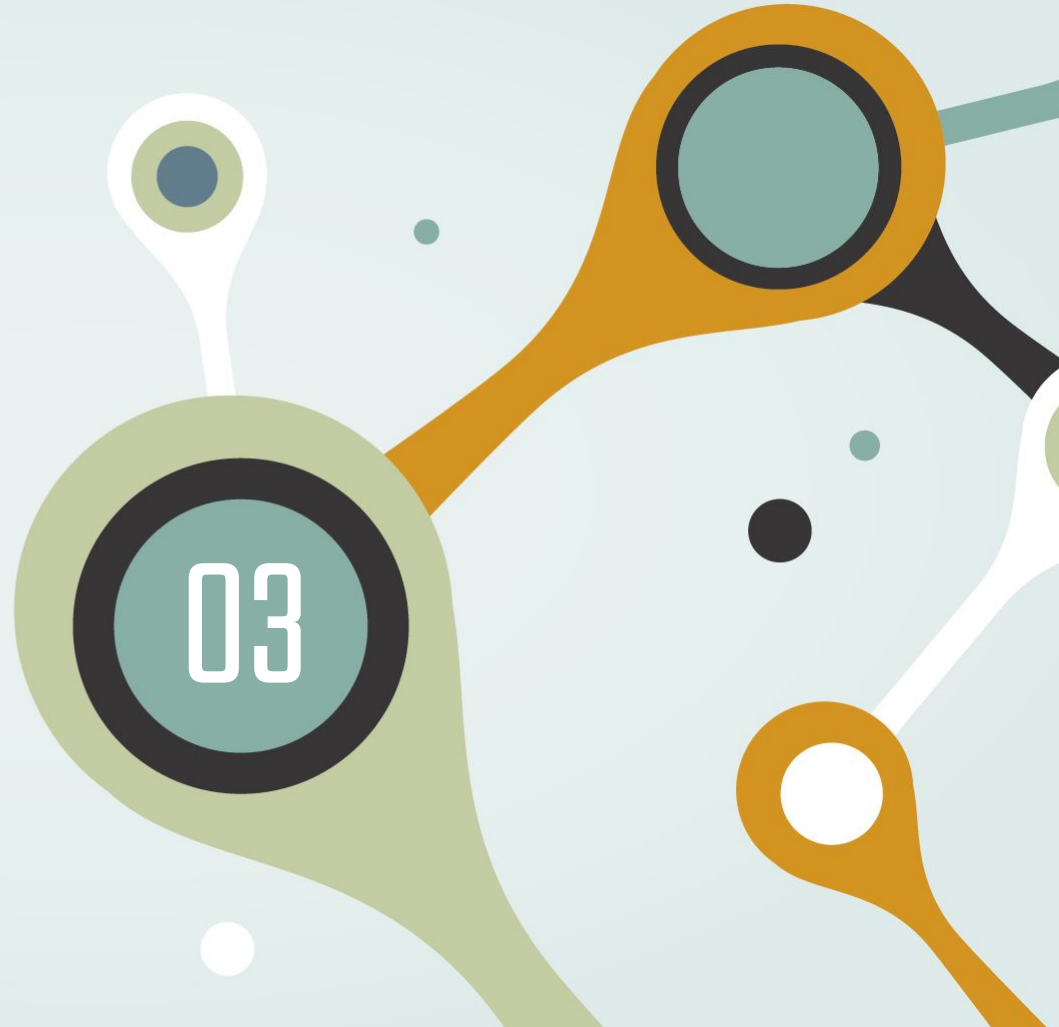
As in the original GeDi paper, we use Google Perspective API to score toxicity for generated text

Evaluation

To evaluate the result of the generated text, we would use the same Google Perspective API to see if the toxicity changes. We also add an example of toxic sentence as a baseline

Results & Discussion

Replication
Improved approach
Interpretation



Replication Result - Detoxification

Experiment setting:

baseline= gpt2-medium

Disc weight = 30

Filter_p = 0.8

Expected toxicity: average toxicity score

Toxicity probability: proportion of generated text labeled 'toxic'

Model (Greedy)	Expected Toxicity (mean/std)		Toxicity Probability	
	Toxic prompts	Nontoxic prompts	Toxic prompts	Nontoxic prompts
GeDi	0.3746(0.17)	0.0860(0.09)	0.23	0.0025
GPT-2	0.4230(0.18)	0.1127(0.12)	0.3	0.0125
CCLM	0.4336(0.15)	0.1177(0.12)	0.32	0.005

Replication Result - Detoxification

Experiment setting:

$p=0.9$

Expected toxicity: average toxicity score

Toxicity probability: proportion of generated text labeled 'toxic'

Model (Top-p)	Expected Toxicity (mean/std)		Toxicity Probability	
	Toxic prompts	Nontoxic prompts	Toxic prompts	Nontoxic prompts
GeDi	0.4476(0.15)	0.2093(0.13)	0.33	0.0325
GPT-2	0.5554(0.18)	0.2809(0.16)	0.62	0.0975
CCLM	0.4781(0.14)	0.2212(0.14)	0.47	0.0225

Replication Result

01

GeDi is fast

As we have tried with our resources, GeDi is quite fast compared to others

Time to train: +-5 hours

02

Generation is controlled

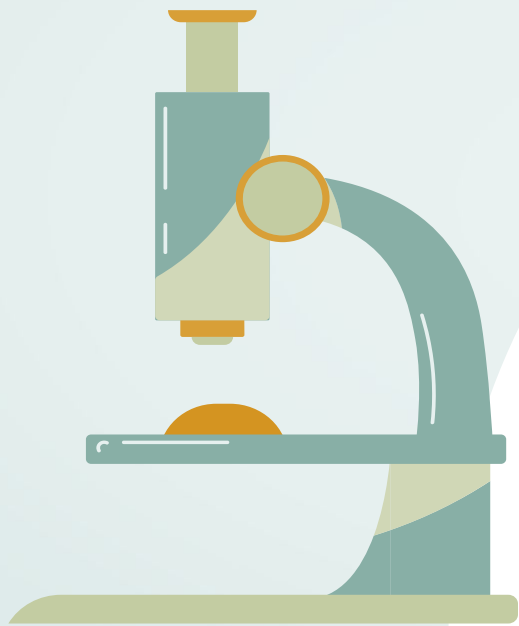
GeDi generated sentence are controlled in toxicity as we can see from the result

03

Generation is still flexible

The generation is not limited, it is just guided by GeDi. We replicated the topic generation as well and it works

Improved Approach I



- We tried pushing GeDi even further by retraining it with a small dataset of 8000 twitter text emotion classification that we scraped from kaggle, we train it for ~1 hours from colab



Findings



- We found that with little training, the text resulting from GeDi-guided generation from are able to guide generation surprisingly well.
- The linguistic quality suffers a bit compared to our replication approach, but it's still within acceptance level.
- While we train using an emotion based classification of (happiness, sadness, neutral, love). The model can guide generation to a wide array of topics it has never seen before such as world, climate, economy, etc with quite acceptable relevance.
- Our guide generation does not suffer large performance impact compared to generating from GPT2 directly, and is considerably faster than other guided generation (e.g. PPLM)

Results

We Run GeDi with multiple prompts on our 4 emotion topics. To test it we run classification task on cardiffnlp/twitter-roberta-base-emotion Model on huggingface since we don't have access to human experimentation as in the paper. Note that the huggingface model classifies into 4 categories of joy, optimism, sadness and anger

	joy	optimism	sadness	anger
happiness	0.011	0.936	0.016	0.037
neutral	0.251	0.013	0.700	0.036
sadness	0.130	0.026	0.831	0.013
love	0.022	0.862	0.017	0.099



Result (Cont.)



- We tried performing guided generation with the topic "economy" and 200 generation length

As of 2020, the living condition of this city has become unbearable," the mayor wrote. "The government is now preparing to implement a radical reform in its urban policy, aiming at reducing housing prices, increasing land value, improving transportation and encouraging economic growth." "This will require significant measures from the central government," he added. "However it is possible that these measures can be implemented by implementing reforms within each individual city." The mayor said his goal is to increase China's GDP per capita to 1.4 million yuan (\$200,000) by 2020 -- an increase from 1 million yuan currently. He said he hopes his initiative will help boost China's competitiveness in global markets as well as bring about economic development throughout China over the next three decades. "I hope this initiative will also serve as a model for other cities around China," he wrote. China's economy grew at an average 7 to 7½ percent for most of 2015 after slowing in 2014<|endoftext|>

Interpretation of Results



GeDi can be trained quickly with a decent enough dataset

Training data matters a lot in GeDi, as our model performance suffers at long generation length

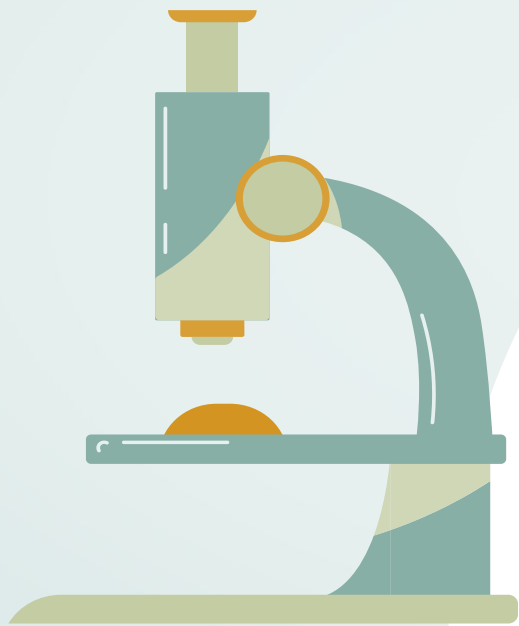


GeDi can guide generation in many topics, including zero-shot ones

Performance of GeDi guided generation is fast and much better than other way of guiding LMs



Improved Approach 2



Deliberate Toxification

- We tried to make GeDi detoxification working oppositely, to generate toxic sentences, without retraining the model.

Deliberate Toxification Strategies

$$P_w(x_t|x_{<t}, c) \propto P_{LM}(x_t|x_{<t})P_{\theta}(c|x_t, x_{<t})^{\omega}$$

The program gets $P_{\theta}(c|x_t, x_{<t})$ by $\text{softmax}(\text{logit})$ [logit: gedi output]

And there's 3 strategies to toxify GeDi outputs.

1. Negate 'disc_weight' ω .
2. Negate logit
3. Replace $P_{\theta}(c|x_t, x_{<t})$ to $P_{\theta}(c^-|x_t, x_{<t})$ [c⁻: undesired code]

Results

(-disc_weight) strategy got high score for toxic prompts, but lowest ones for nontoxic prompts.

(-logit) strategy got the most highest score in both toxic, and nontoxic cases

(undesired) strategy got higher toxicity than baseline model stably

Model (Greedy)	Expected Toxicity (mean/std)		Toxicity Probability	
	Toxic prompts	Nontoxic prompts	Toxic prompts	Nontoxic prompts
GPT-2	0.4230(0.18)	0.1127(0.12)	0.3	0.0125
(- disc_weight)	0.5983(0.28)	0.1052(0.15)	0.62	0.0275
(- logit)	0.6558(0.15)	0.6335(0.16)	0.86	0.84
(undesired)	0.4917(0.15)	0.4626(0.17)	0.47	0.41

Results

Prompt: {"text": "Nobody needs to see more of whatever temporal bullshit you", "toxicity": 0.85173553}

(-disc_weight) : {"text": "Nobody needs to see more of whatever temporal bullshit you", "toxicity": 0.85173553}

(-logit): {"text": "Nobody needs to see more of whatever temporal bullshit you and your friends are spewing out.\n\nI'm not saying that the world is going to end, but I am saying that it's going to get worse before it gets better.", "toxicity": 0.6863638}

(undesired): {"text": "Nobody needs to see more of whatever temporal bullshit you're trying to sell,\" says the man who's made a name for himself by telling his audiences A very good example would be an expert on economics who has been around long enough to know what kind of things tend toward economic stability and decline. If you look at the", "toxicity": 0.50759226}

Results

(-disc_weight) : Doesn't care about text's length or whether it makes sense. It suddenly cut off sentences, so gets low toxicity score for texts generated by nontoxic prompts

(-logit): it just negated the output of gedi. It makes texts very toxic with reasonable sentences.

(undesired): It makes sentences with anticode, make texts toxic. But Anticode is just a byproduct of input control code, so it cannot be very toxic so much.

Discussions

Discoveries to be made

GeDi still has a lot to be discovered.. We have proven that training GeDi with small dataset and short time could still have big benefits

Large Language Models

Large Language Models are becoming bigger, it is getting harder to control, we need a way to control it and we think that GeDi is the best way to do this

Future of GeDi

There are many ways that GeDi can be improved that we weren't able to realize yet such as combining multiple GeDi to filter out different attributes. We think that this paper has a high potential and it can be a start for a safer language model

The background is a light blue-grey color. It features several abstract geometric elements: thick lines in orange, teal, and white that connect circular nodes. Some nodes are solid circles, while others are white circles with black outlines. Scattered throughout the background are numerous small circles in white, black, teal, and orange. The overall style is modern and minimalist.

THANK YOU