

Men Have Feelings Too: Debiasing Sentiment Analyzers using Sequence-to-Sequence Generative Adversarial Networks

Athiya Deviyani

adeviyan

Mehak Malik

mehakm

Haris Widjaja

iwidjaja

Dan Hoskins

dhoskins

Abstract

Natural Language Processing (NLP) models often magnify the bias with respect to race, gender and age present in datasets that they are trained on. Furthermore, it is becoming increasingly challenging to collect an unbiased dataset given that sexist and racist content are ubiquitous in common sources of data such as social media. In this work, we propose a Generative Adversarial Network (GAN) based approach to augment a sentiment analysis dataset with unbiased samples, mitigating the gender bias present in the original dataset. We show that our method successfully reduces the disparity displayed by a downstream model trained on the augmented dataset, as measured by various fairness criteria compared across genders.

1 Motivation

As Natural Language Processing (NLP) systems rapidly improve across many subdomains (e.g. translation, sentiment analysis, topic classification), they earn more trust among the general public. It's easy to see why: when a system has a very high accuracy, it seems to be functioning very well. As a result of this perception, more and more of society believes the positives outweigh the risks of these systems, and their adoption increases. While these systems are unquestionably useful in many scenarios, a number of pernicious effects of these NLP-based systems exist, including biased predictions. These biased predictions result in unfair outcomes for marginalized groups. As the adoption of NLP models increases, the deleterious effects of these biased predictions will commensurately be propagated.

For example, human resources groups in private corporations are adopting sentiment analysis tools for gauging employee feedback. These systems are used to better inform decisions about the future directions of the companies that use them.

Many sentiment analysis systems exhibit gender bias. One review by Maurer (2019) found that, of all classifiers that took part in a particular SemEval task (Affect in Tweets) (Mohammad et al., 2018), most systems output “higher sentiment intensity predictions” for one gender than another. Since the sentiment associated with any given piece of feedback can determine the extent to which it is considered in decisions, these biased predictions can result in one gender’s feedback being considered more than the other. This will, in turn, impact company decisions, which affects employees. Clearly, NLP-based systems’ biases can impact decisions, resulting in downstream effects on humans. It’s therefore critical to develop methods to combat these biases, enabling the machine learning community to build fairer systems.

Of the existing bias-mitigation techniques, many focus on modifying the dataset that the systems are trained on. These techniques have serious limitations. For example, counterfactual data augmentation (i.e. adding sentences with swapped pronouns of existing data points to the dataset) requires manual labeling. This must be performed for every new task that the technique is applied to. This means that developers of these systems have to dedicate significant effort to modify their datasets to mitigate the biases.

Clearly, it would be valuable to have a bias-mitigation technique that generalizes well to arbitrary contexts, eliminating the need for significant rework. One technique uses generative adversarial learning to generate data points with which to augment the training dataset. In this technique, the relationships necessary for bias mitigation are learned, rather than specified through manual labels, suggesting it might generalize to arbitrary contexts. However, this generalizability hasn’t been thoroughly tested.

2 Problem definition

2.1 Problem

The high level goal of our project is to neutralize bias in datasets through data augmentation using the adversarial learning objective. We believe that our fairness metrics can be improved for twitter datasets by augmenting the dataset with artificially generated non-sexist and non-offensive tweets.

We will analyze our fairness metrics using sentiment analysis as the downstream task. Given a tweet (or any piece of text sequence), a sentiment analyzer will decide whether the tweet expresses negative or positive sentiment. This is essentially a binary text classification task. Following from the observation made by [Kiritchenko and Mohammad \(2018\)](#), we hypothesize that tweets containing female-related terms (such as mother, woman, girl), amplifies the sentiment of the tweet. For example, the sentence “my sister is sad” will have a higher negative sentiment score when compared to the sentence “my brother is sad”, given that the former contains a female-related term (“sister”).

We believe that the amplification of sentiment can have adverse effects on people identifying as female. In addition to enforcing dangerous stereotypes, there are many cases where sentiment analysis is used to automatically detect depression from social media or blog post entries to advance psychological research and mental health aid, similar to what [Husseini Orabi et al. \(2018\)](#) and [Deshpande and Rao \(2017\)](#) has done in their work. The misclassification of sentiment for women (and potentially other minority groups) may lead to problems such as having one gender group being flagged as more depressed than the other.

Natural Language Processing applications are also commonplace in professional settings. As mentioned previously, some human resource departments within organizations increasingly use automated sentiment analyzers to assess employee feedback ([Maurer, 2019](#)). This entails that biased sentiment analysis results can lead to the silencing of female voices, wherein their feedback might not be taken into account. Additionally, the over-amplification of sentiment on text sequences containing instances of female-related words can also mean that negative words are scored more severely, therefore a female worker who received negative feedback might be scored lower than a male worker who received a similar kind of feedback. Conversely, if positive words are also scored more

highly, then it would be unfair for male workers who receive similar positive feedback but may receive a much lower positive score.

Therefore, our main objective in this project is to use the adversarial training objective to artificially generate text that simulates real tweets with the aim of neutralizing the effect of data disparity to the over-amplification of sentiment in the sentiment analysis task on text containing female-related terms. Given that our method works in maximizing the fairness of the downstream task, our method can potentially be adapted to various tasks and validates the generalizability of a novel, reproducible debiasing technique.

2.2 Dataset

We decided to use the Sentiment140 dataset as our baseline collected by [Go et al. \(2009\)](#). The dataset contains over 1.6 million tweets, alongside other relevant information such as the tweetID, tweet date, query used to obtain the tweet, and the tweet author. Each tweet comes with the corresponding sentiment/polarity label of NEGATIVE and POSITIVE. Below are examples of tweets reflecting the aforementioned sentiment labels:

NEGATIVE: “I’ve just spent 1 hour to enter all the bureaucratic nonsense for March. What a waste of my time.”

POSITIVE: “I’m meeting up with one of my besties tonight! Can’t wait!! - GIRL TALK!!”

Contrary to most datasets collected for the task of sentiment analysis, the Sentiment140 dataset’s labels are not hand-annotated. In fact, they have automated the collection process fully. The authors collected a handful of tweets from twitter and automatically labeled tweets containing positive emoticons such as :) as POSITIVE and tweets containing emoticons that signify negative emotions such as :(as NEGATIVE. We have performed Twitter-specific preprocessing on top of normal text processing (which involves lowercasing, tokenizing, and removing stop words), such as obfuscating and removing usernames and links for privacy purposes, as well as separating hashtags into individual words (e.g. #BlackLivesMatter → Black Lives Matter) using the help of the `ekphrasis` Python library ([Baziotis et al., 2017](#)). This step is highly crucial as we would like to minimize the amount of noise in our dataset (such as the presence of stop

3.2 GANs for natural language generation

To train the generator to avoid generating sequences with strong indication of ethnicity, Agrawal (2022) uses an ethnicity prediction model as a feedback mechanism in a GAN-like framework. One challenge in adopting the original GAN framework (Goodfellow et al., 2014) to natural language generation is the inability to backpropagate through the sampling operation which language models use to sample words from logits.

To address this difficulty, the authors of Agrawal (2022) adopt a reinforcement learning-based approach, SeqGAN (Yu et al., 2017), which bypasses the need to backpropagate through the sampling operation. Motivated by shortcomings of maximum likelihood training (Welleck et al., 2019), the original SeqGAN uses a real-fake discriminator to train a generator to generate realistic language. Since this allows for an arbitrary form and number of discriminators, this makes SeqGAN an attractive framework for controlling the characteristics of generated text.

In this work, we aim to investigate the SeqGAN framework’s ability to train generators to generate data useful for specific downstream tasks, while ensuring that the generated samples adhere to a non-sexism constraint, leading to an overall less biased augmented dataset for downstream models.

4 Baseline

4.1 Text classification model

To determine the effectiveness of our proposed debiasing approach, we train a deep learning model on the sentiment analysis task with no further intervention. Our baseline model consists of two Bidirectional Gated Recurrent Unit (GRU) layers followed by a convolutional layer, a Recurrent CNN architecture shown to be effective for text classification by Zhang et al. (2018). The GRU preserves historical information in long text sequences, while the final convolution layer extracts local features, leading to a better representation for the sentiment classification task. We use a hidden size of 64 for both the GRU cells and the convolution layer, with a dropout layer after each Bidirectional GRU layer to regularize our model and prevent overfitting. The network architecture is shown in Figure 2.

The model takes in special GloVe word embeddings with 100 dimensions which were trained on tweets as input (Pennington et al., 2014). We believe that this is the most suitable representation for

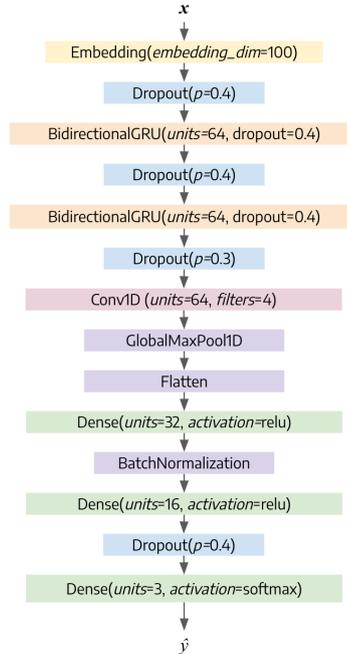


Figure 2: Recurrent CNN architecture for sentiment classification

our training data as it consists of highly colloquial text from tweets, and that this feature representation will capture the relationships between the words most appropriately, retaining as much of its original meaning as possible in a vectorized form.

4.2 Results and analysis

We measured our baseline results using the performance and fairness metrics mentioned in section 2.3. The classification performance is tabulated in Table 1.

	Precision	Recall	F1-score	Support
Negative	0.81	0.83	0.82	3625
Positive	0.73	0.69	0.71	2375

Table 1: Baseline sentiment analysis performance

		Mean	Minimum	Maximum
Negative	Male	0.27	0.03	0.92
	Female	0.08	0.01	0.82
Positive	Male	0.73	0.08	0.96
	Female	0.92	0.18	0.99

Table 2: Baseline fairness metrics

As we can see, the model performs with an average test accuracy of 78% on the test dataset. It is also interesting to see that the model performs better in predicting negative sentiment than positive

sentiment. The F-1 score for predicting negative sentiment is 0.82 and 0.81 for positive sentiment.

Upon examining the fairness metric of our baseline in Table 2, we see that there is a significant difference between the mean sentiment scores for different genders. We can see that for statements containing nouns or pronouns associated with females, the mean positive sentiment is 0.92 which is much higher than 0.73 associated with males. Similarly for negative sentiment, the mean negative sentiment is much lower for the female class. It can be seen that usually sentiments for females are more extreme as observed by the minimum and maximum values shown in Table 2. The mean positive sentiment score being higher and the mean negative score being lower for females can be attributed to the stereotype that women are more emotional than men and it is possible that the dataset contains more data for women that is extreme in terms of sentiment. This is the disparity that our work aims to decrease using a more balanced dataset.

5 Methodology

Our approach focuses on using Generative Adversarial learning to generate an unbiased dataset which can be used to augment a biased dataset to be used in a downstream Natural Language Processing (NLP) task. The downstream NLP task to be used here is sentiment analysis and the objective will be to generate a gender-unbiased dataset to augment with the Sentiment140 dataset.

5.1 Overview

The overall architecture of our approach is described in Figure 3. The approach involves using a generative adversarial network (GAN) to generate synthetic tweets that are similar to our original dataset, the Sentiment140 dataset. We implement the SeqGAN as our generative adversarial network (GAN) as described by Yu et al. (2017). It is further explained in the following section. The generator used in the GAN is DistilGPT2 model trained on tweets (Dayma, 2021). The generated synthetic tweets are then subject to a selection process by a sexism detector. For this, we used the classifier described by Safi Samghabadi et al. (2020). It is explained in depth in the upcoming sections of this paper. We use this classifier to remove tweets that might increase the gender bias in our dataset. Finally, we use a sentiment analyzer from Heitmann et al. (2020) as our ‘gold standard’ sentiment anal-

ysis tool to label our generated data. The process of generating data using the GAN and purging the gender biased data to produce unbiased data can be repeated many times to increase the size of our augmented dataset. We then retrain our baseline model using the augmented dataset and evaluate our results.

5.2 SeqGAN: Sequence Generative Adversarial Networks with Policy Gradient

Our proposed approach utilizes a GAN approach on natural language generation, which faces the difficulty of backpropagating through the sampling operation in our language model-based generator. In particular, the gradient of the loss with respect to our generator’s parameters θ is

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{y \sim G_{\theta}} [Q_{\phi}(y)]$$

In the equation above, G_{θ} is the generator, Q_{ϕ} is the discriminator, and y is the sequence sampled by the generator G_{θ} . Since the expectation is taken over the distribution G_{θ} , in which the sequence y undergoes a sampling step, this gradient cannot be computed analytically.

To address this challenge, we adopt a reinforcement learning-based approach called SeqGAN Yu et al. (2017), which bypasses the need to back-propagate through the sampling operation. In Yu et al. (2017), the gradient is approximated by the REINFORCE gradient Williams (1992):

$$\nabla_{\theta} J(\theta) \approx \sum_y \nabla_{\theta} G_{\theta}(y) \cdot Q_{\phi}(y)$$

Observe that, in this formulation, all quantities can be computed analytically. In particular, we only need the gradient of the probability distribution induced by the generator G_{θ} , instead of the gradient with respect to the sampled sequences y . For a full discussion of this gradient estimation, the reader is referred to Yu et al. (2017).

An additional complication for adapting the GAN framework for discrete token generation, as is the case in natural language generation, is that it is non-trivial to specify rewards for partly-generated sequences. In particular, the discriminator only assigns rewards for fully completed sentences, and not partially-generated sentences. This makes the reward signals ‘sparse’, making it harder to train the generator due to not receiving intermediate rewards immediately after generating a token.

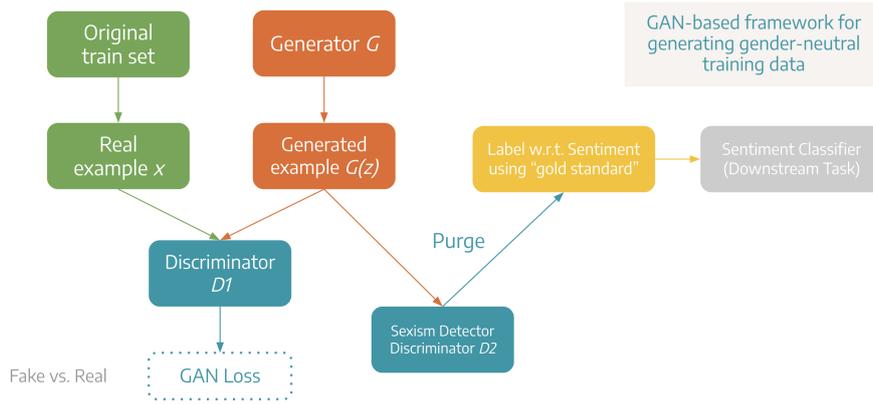


Figure 3: GAN-based framework for generating gender-neutral training data

To address the challenge of sparse rewards, Yu et al. (2017) estimates the intermediate rewards using Monte Carlo Tree Search, a technique shown to be successful in estimating intermediate rewards in the sequential nature of turn-based games (Silver et al. (2016)). For full details of the Monte Carlo Tree Search policy used, we refer the reader to Yu et al. (2017).

In order to adapt the SeqGAN framework to the task of generating synthetic and unbiased data useful for a downstream task, we make the following modifications to the original SeqGAN framework:

- We retool the real-fake discriminator to distinguish between data coming from the original downstream dataset and synthetic data generated by the generator. We do this to ensure that the generator generates data useful for specific downstream tasks. We further make use of the Monte Carlo Tree Search rollouts as described in Yu et al. (2017) to convert the sparse rewards to dense ones.
- We attach a sexism detector as an additional discriminator, and use it to provide feedback to the generator that encourages it to generate samples which are not sexist. This helps to reduce the overall bias in the downstream dataset, which reduces the bias in models trained on that dataset. We also convert this sparse reward to a dense one using Monte Carlo Tree Search rollouts. We describe our sexism detector in a subsequent section.

However, due to the instability of adversarial training, we decided to omit the sexism detector-discriminator from the final version of our approach. In the analysis section, we justify this

decision, and propose an alternative which results in a reasonable amount of debiasing.

5.3 Sexism detection

The sexism detector is based on the system presented by Safi Samghabadi et al. (2020). It uses a BERT based model to detect aggression and misogyny as two separate tasks. The BERT based layers are used to extract contextual information. The output of this layer is fed to an attention layer followed by a fully connected layer. Finally, the output is fed to two different classification layers: one for detection of aggression and the other misogyny. We use the outputs from the misogyny classification as our sexism detector to 'purge' gender-biased examples generated by the SeqGAN.

5.4 Results and analysis

	Precision	Recall	F1-score	Support
Negative	0.77	0.74	0.76	1662
Positive	0.59	0.63	0.61	976

Table 3: Sentiment analysis performance on debiased dataset

		Mean	Minimum	Maximum
Negative	Male	0.61	0.00	1.00
	Female	0.59	0.00	1.00
Positive	Male	0.39	0.00	1.00
	Female	0.41	0.00	1.00

Table 4: Fairness metrics on debiased dataset

We have tabulated the performance metrics of the sentiment analyzer trained on our augmented dataset in Table 3. It can be seen that our performance does fall from the baseline with a test accu-

Sentence	Sentiment
"I can't think of a better portrayal of a lonely lonely lone wolf than the guy who plays the guitar. He has no real life."	Negative
"I never thought that a man could be so mad as me."	Negative
"A girl laying on the floor, crying, and all the people are just warmly kissing her while I convey my sad, broken heart."	Negative
"There are beautiful women and beautiful men ."	Positive

Table 5: Sample sentences generated by SeqGAN. Words that denote gender are emboldened.

racy of 70%. This drop is seen across all metrics. It is also interesting to note that even here our model performs better in predicting negative sentiment than positive. However, we can see from Table 4 that our augmented dataset does indeed improve the considered fairness metrics when compared to our baseline. We see that the difference between the mean negative and positive sentiment for different gender classes drastically reduces with mean negative sentiment being 0.59 and 0.61 respectively for females and males respectively.

The tweet-based SeqGAN is successfully able to generate synthetic tweet data that was close to the source corpus. This could be attributed to the fact that GPT-2 models have been trained on huge corpora and are able to generate data that is coherent, logical and believable. The sexism detector is also able to select tweets that would lower the gender bias in our dataset and thus reduce bias in our downstream task, in this case sentiment analysis. We present some of the sentences generated by our system in Table 5.

Additionally, we decided to further evaluate the sentences that were generated by our SeqGAN. We obtained a set of words related to emotion as defined by Shaver et al. (1987), where they defined six primary emotions, 25 secondary emotions and 135 tertiary emotions. We identified the generated sentences that contain 'emotion' words, and found out that over 2000 samples containing 'emotion' words are associated with male-related words while only around 500 samples are associated with female-related words. We also observed that there are plenty of sentences which contain 'emotion' words and both male and female-related words, for example, "there are beautiful women and beautiful men".

However, it was a non-trivial task to train the GAN. We had initially planned to use the sexism detector as an additional discriminator for our SeqGAN, but we could not achieve convergence in the joint training loss. Instead, we train the GAN only using the real-fake discriminator, and instead use the sexism detector to filter away generated

examples that exhibit sexism.

We found that a majority of the generated samples (around 90,000 out of 100,000) were not sexist. However, it is desirable to obtain a generator which inherently knows how to generate non-offensive examples. Integrating the sexism detector as an additional discriminator during the GAN training could be a valuable future extension to our work.

6 Ethical implications

There are clear benefits of a technique that could remove bias from datasets in a fully generalizable fashion. However, there are some potential drawbacks. First, alleviating bias along one axis (e.g. gender) can impact the bias exhibited along other axes. The bias along any given axis can only be tested by explicitly partitioning the dataset by that axis. This requires labeling each data-point according to that axis. It is nontrivial to get these labels. Therefore, it's difficult to analyze the bias of a system along all axes. Because of this, it may be difficult to identify if and when this bias-mitigation approach exacerbates bias along a different axis.

Our GAN-based approach also makes it easy to generate extremely biased synthetic samples by a simple modification to the objective function: instead of penalizing the generator for generating sexist samples, malicious actors can instead reward it.

There are also a few problems associated with synthetic data, in general. Synthetic data can be used maliciously, even if its unbiased. For example, the availability of high-quality, unbiased data would make it easier to impersonate a member of a given group. Additionally, synthetic data generation sometimes produces nonsensical results. Depending on the downstream application, this can impact the performance on the downstream task in unpredictable ways, in turn resulting in an impact to the people impacted by that model. Finally, synthetic data doesn't always capture outliers' characteristics. In situations where detection of outliers is especially important, this would be a significant problem.

References

- Adarsh Agrawal. 2022. [Mitigating bias in ai using debias-gan](#).
- Christos Baziotis, Nikos Pelekis, and Christos Douk-eridis. 2017. Datastories at semeval-2017 task 4: Deep lstm with attention for message-level and topic-based sentiment analysis. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 747–754, Vancouver, Canada. Association for Computational Linguistics.
- Boris Dayma. 2021. [Distilgpt2 tweet bot](#).
- Mandar Deshpande and Vignesh Rao. 2017. [Depression detection using emotion artificial intelligence](#). In *2017 International Conference on Intelligent Sustainable Systems (ICISS)*, pages 858–862.
- Alec Go, Richa Bhayani, and Lei Huang. 2009. Twitter sentiment classification using distant supervision. *Processing*, 150.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. *Advances in neural information processing systems*, 27.
- Mark Heitmann, Christian Siebert, Jochen Hartmann, and Christina Schamp. 2020. More than a feeling: Benchmarks for sentiment analysis accuracy. *Available at SSRN 3489963*.
- Ahmed Husseini Orabi, Prasadith Buddhitha, Mahmoud Husseini Orabi, and Diana Inkpen. 2018. [Deep learning for depression detection of Twitter users](#). In *Proceedings of the Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic*, pages 88–97, New Orleans, LA. Association for Computational Linguistics.
- Svetlana Kiritchenko and Saif M. Mohammad. 2018. [Examining gender and race bias in two hundred sentiment analysis systems](#).
- Roy Maurer. 2019. [Employee sentiment analysis shows hr all the feels](#).
- Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. [SemEval-2018 task 1: Affect in tweets](#). In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 1–17, New Orleans, Louisiana. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Niloofer Safi Samghabadi, Parth Patwa, Srinivas PYKL, Prerana Mukherjee, Amitava Das, and Tamar Solorio. 2020. [Aggression and misogyny detection using BERT: A multi-task approach](#). In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*, pages 126–131, Marseille, France. European Language Resources Association (ELRA).
- Phillip Shaver, Judith Schwartz, Donald Kirson, and Cary O’connor. 1987. Emotion knowledge: further exploration of a prototype approach. *Journal of personality and social psychology*, 52(6):1061.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. 2016. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489.
- Sean Welleck, Ilya Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2019. Neural text generation with unlikelihood training. *arXiv preprint arXiv:1908.04319*.
- Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256.
- Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31.
- Jiarui Zhang, Yingxiang Li, Juan Tian, and Tongyan Li. 2018. Lstm-cnn hybrid model for text classification. In *2018 IEEE 3rd Advanced Information Technology, Electronic and Automation Control Conference (IAEAC)*, pages 1675–1680. IEEE.