

Sentiment Analysis of Reddit Comments in Pre & During COVID-19

Nitish Kumar Singh*

Khoury College of Computer Sciences, Northeastern University, Boston, MA 02115

Email: *singh.niti@northeastern.edu

Abstract—In the beginning of 2020, the world was introduced to a new disease called COVID-19 (*CO* for corona, *VI* for virus, *D* for disease and 19 for when the outbreak was first identified (31 December 2019)) [1]. In early March of that year, the World Health Organization (WHO) declared COVID-19 a pandemic — a disease outbreak occurring over a wide geographic area and affecting an exceptionally high proportion of the population [2]. Countries were put under strict lockdown, schools, offices, malls and other places of businesses were shutdown for months. As a result of these measures, Internet and social media usage has been observed in unprecedented magnitudes, when compared with the pre-pandemic period [3]. In this project, I have analyzed the change in sentiment and traffic of social media usage, especially the reddit comments, in pre and during pandemic (Jan'2019 - Mar'2021) for 5 specific subreddits - Boston, LosAngeles, Texas, Depression, and MentalHealth. Sentiment analysis has grown into one of the most important sub-areas in Natural Language Processing (NLP). One can effectively mine the implicit emotion in the text, which can help organizations to make an effective decision, and the explosive growth of data undoubtedly brings more opportunities and challenges to the sentiment analysis. At the same time, transfer learning has emerged as a new machine learning technique that uses the existing knowledge to solve different domain problems and produces state-of-the-art prediction results.

Index Terms—NLP, Sentiment Analysis, Transfer Learning, COVID-19, Reddit

I. INTRODUCTION

The ongoing pandemic of the COVID-19 virus, also referred to as the coronavirus pandemic, is the major global health event of 2020. The lengthy absence of a vaccine, alongside a rising death toll and the absence of a globally effective response to the outbreak, renders this pandemic a clear example of an existential risk [4]. In order to limit the propagation of the virus, quarantine measures have been adopted by many countries across the world. Qiu et al. [5] have explored the various forms of psychological distress associated with the strict quarantine measures applied in China, concluding that protocols have triggered problems such as panic disorder, anxiety, and depression. Brooks et al. [6] published a comprehensive and informative meta-analysis of the negative psychological impacts of quarantine, noting in particular the prevalence of post-traumatic stress symptoms, confusion, and anger.

During the outbreak, different forms and degrees of lockdown measures were deployed worldwide. As a result of these measures, Internet and social media usage has been observed

in unprecedented magnitudes, when compared with the pre-pandemic period [7]. Social media can play a crucial, positive role by providing a platform for people to share their opinions and to relay facts about the crisis, but it also provides an outlet for voicing fear about the pandemic. One of these frequently used social media platforms is Reddit.

Reddit is a social news aggregation, web content rating, and discussion website, and it claims to be "the front-page of the internet" as its moniker, recently including livestream content through Reddit Public Access Network. As of February 2021, Reddit ranks as the 18th most-visited website in the world and 7th most-visited website in the U.S., according to Alexa Internet. Subreddits are user-created areas of interest where discussions on Reddit are organized. There are about 138,000 active subreddits (among a total of 1.2 million) as of July 2018. [8]

Sentiment analysis refers to the process of extracting explicit or implicit polarity of opinions expressed in textual data. Sentiment analysis has been used for information seeking and demand addressing needs on the consumer side, whereas for business owners and other stakeholders for operational decision making (e.g., branding, preventive/reversal actions). Traditional sentiment analysis focus on extracting opinion polarities at a coarse level, which cannot fully satisfy aforementioned purposes. Sentiments are normally domain dependent (e.g. delicious indicates positive sentiment in the food domain, where it does not indicate any sentiment in the laptop domain). [9]

II. RELATED WORK

A. Transfer Learning

Transfer learning uses domain-specific data to fine tune the pre-trained deep learning models. The benefit of conducting transfer learning is twofold: The time spent in training is much less than the time used in training from scratch. In computer vision, it is a common practice to use transfer learning: Parameters of the fully connected layers of a pre-trained CNN are replaced with randomly initialized values. A fine-tuning process is then performed by updating the new values only, using backpropagation, while the parameters in the convolutional layers stay untouched [10], [11]. Transfer learning in NLP, however, has been shown as a somewhat difficult task. One early successful case involves fine-tuning the pre-trained word embeddings [12], has had a large impact

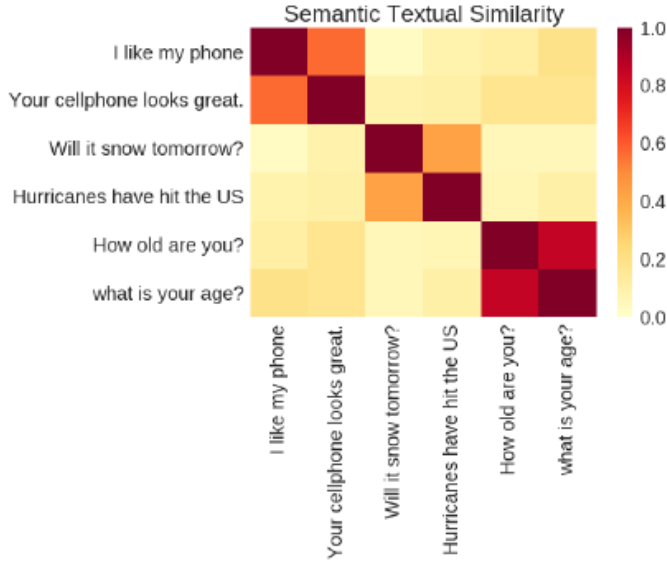


Figure 1. Sentence similarity scores using embeddings from the universal sentence encoder [17]

in practice. Howard and Ruder [13] proposed a transfer learning method that fine-tunes a three-layered LSTM language model [14] for text classification. In the first step of their approach, sentences of a training set are used to fine-tune the parameters in the LSTM layers. The labels of the training set are then used in step two to update the parameters of the fully connected layers. This approach reuses the embeddings of the original language model. Radford et al. [15] used the Generative Pre-trained Transformer (OpenAI GPT) to achieve state-of-the-art results on many sentence-level tasks from the GLUE benchmark [16].

B. Universal Sentence Encoder

In the paper[17], the authors present models for encoding sentences into embedding vectors that specifically target transfer learning to other NLP tasks. The models are efficient and result in accurate performance on diverse transfer tasks. Two variants of the encoding models allow for trade-offs between accuracy and compute resources. One makes use of the transformer [18] architecture, while the other is formulated as a deep averaging network (DAN) [19]. The models take as input English strings and produce as output a fixed dimensional embedding representation of the string. The sentence embeddings can be trivially used to compute sentence level semantic similarity scores that achieve excellent performance on the semantic textual similarity (STS) Benchmark as seen in fig 1. Both encoders have different design goals. One based on the transformer architecture targets high accuracy at the cost of greater model complexity and resource consumption. The other targets efficient inference with slightly reduced accuracy.

1) *Transformer*: The transformer based sentence encoding model constructs sentence embeddings using the encoding

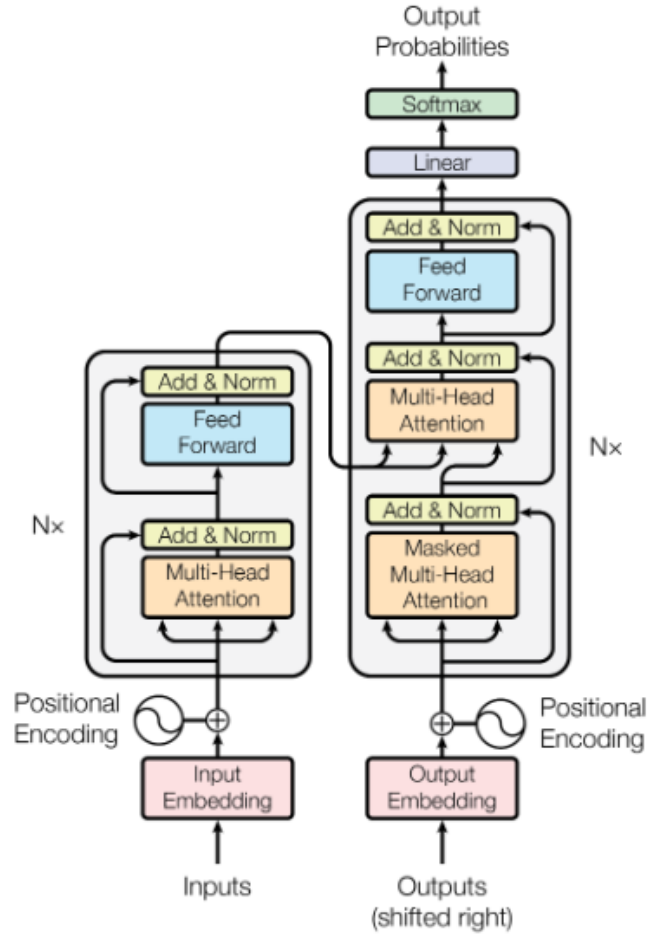


Figure 2. The Transformer - model architecture [18]

sub-graph of the transformer architecture [18]. This sub-graph uses attention to compute context aware representations of words in a sentence that take into account both the ordering and identity of all the other words. The context aware word representations are converted to a fixed length sentence encoding vector by computing the element-wise sum of the representations at each word position. The encoder takes as input a lowercased PTB tokenized string and outputs a 512 dimensional vector as the sentence embedding. The encoding model is designed to be as general purpose as possible. This is accomplished by using multi-task learning whereby a single encoding model is used to feed multiple downstream tasks. The supported tasks include: a Skip-Thought like task [20] for the unsupervised learning from arbitrary running text; a conversational input-response task for the inclusion of parsed conversational data [21]; and classification tasks for training on supervised data. The Skip-Thought task replaces the LSTM [22] used in the original formulation with a model based on the Transformer architecture.

2) *Deep Averaging Network*: The second encoding model makes use of a deep averaging network (DAN) [19] whereby input embeddings for words and bi-grams are first averaged

together and then passed through a feedforward deep neural network (DNN) to produce sentence embeddings. Similar to the Transformer encoder, the DAN encoder takes as input a lowercased PTB tokenized string and outputs a 512 dimensional sentence embedding. The DAN encoder is trained similarly to the Transformer based encoder. We make use of multitask learning whereby a single DAN encoder is used to supply sentence embeddings for multiple downstream tasks. The primary advantage of the DAN encoder is that compute time is linear in the length of the input sequence.

III. APPROACH

A. Pushshift API

The Pushshift Reddit API was designed and created by the /r/datasets mod team to help provide enhanced functionality and search capabilities for searching Reddit comments and submissions. The project lead, /u/stuck_in_the_matrix, is the maintainer of the Reddit comment and submissions archives located at <https://files.pushshift.io>.

This RESTful API gives full functionality for searching Reddit data and also includes the capability of creating powerful data aggregations. With this API, you can quickly find the data that you are interested in and find fascinating correlations.

I used Pushshift API to get all comments between *1st January 2019 and 31st March 2021* from the following subreddits:

- Boston - It has 268K members and was created on 25th January 2008.
- LosAngeles - It has 320K members and was created on 14th April 2008.
- Texas - It has 297K members and was created on 27th March 2008.
- Depression - It has 753K members and was created on 1st January 2009.
- MentalHealth - It has 293K members and was created on 12th June 2008.

B. Sentiment140 dataset

It contains 1,600,000 tweets extracted using the twitter api. The tweets have been annotated (0 = negative, 2 = neutral, 4 = positive) and they can be used to detect sentiment. It contains the following 6 fields:

- target: the polarity of the tweet (0 = negative, 2 = neutral, 4 = positive)
- ids: The id of the tweet (2087)
- date: the date of the tweet (Sat May 16 23:58:44 UTC 2009)
- flag: The query (lyx). If there is no query, then this value is NO_QUERY.
- user: the user that tweeted (robotickilldozr)
- text: the text of the tweet (Lyx is cool)

According to the authors [23], their approach was unique because their training data was automatically created, as opposed to having humans manually annotate tweets. They assumed that any tweet with positive emoticons, like :), were positive, and tweets with negative emoticons, like :(, were negative. They

Layer (type)	Output Shape	Param #
keras_layer (KerasLayer)	(None, 512)	147354880
dense (Dense)	(None, 256)	131328
dense_1 (Dense)	(None, 64)	16448
dense_2 (Dense)	(None, 1)	65
Total params: 147,502,721		
Trainable params: 147,841		
Non-trainable params: 147,354,880		

Figure 3. Transformer based sentiment prediction model architecture

Layer (type)	Output Shape	Param #
keras_layer_1 (KerasLayer)	(None, 512)	256797824
dense_3 (Dense)	(None, 256)	131328
dense_4 (Dense)	(None, 64)	16448
dense_5 (Dense)	(None, 1)	65
Total params: 256,945,665		
Trainable params: 147,841		
Non-trainable params: 256,797,824		

Figure 4. DAN based sentiment prediction model architecture

used the Twitter Search API to collect these tweets by using keyword search.

C. Sentiment Prediction Model

For the sentiment prediction model as seen in Fig 3, I have used the transformer based Universal Sentence Encoder which takes as input a string and outputs a 512 dimensional vector as the sentence embedding. This output is then forwarded as input to two dense layers combined with Relu activation where the final layer uses a sigmoid activation to output a value between 0 to 1. Only the last three layers are trainable.

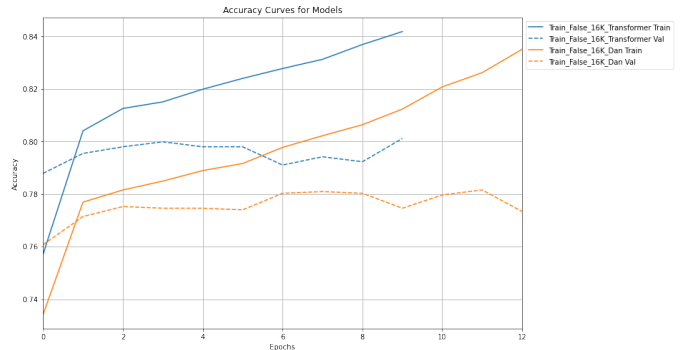


Figure 5. Accuracy of models based on Transformer and DAN during training

This model is then trained on the Sentiment140 dataset. I used 16000 items as the training set and 1584 items as the validation set, followed by 1584 items as the test set.

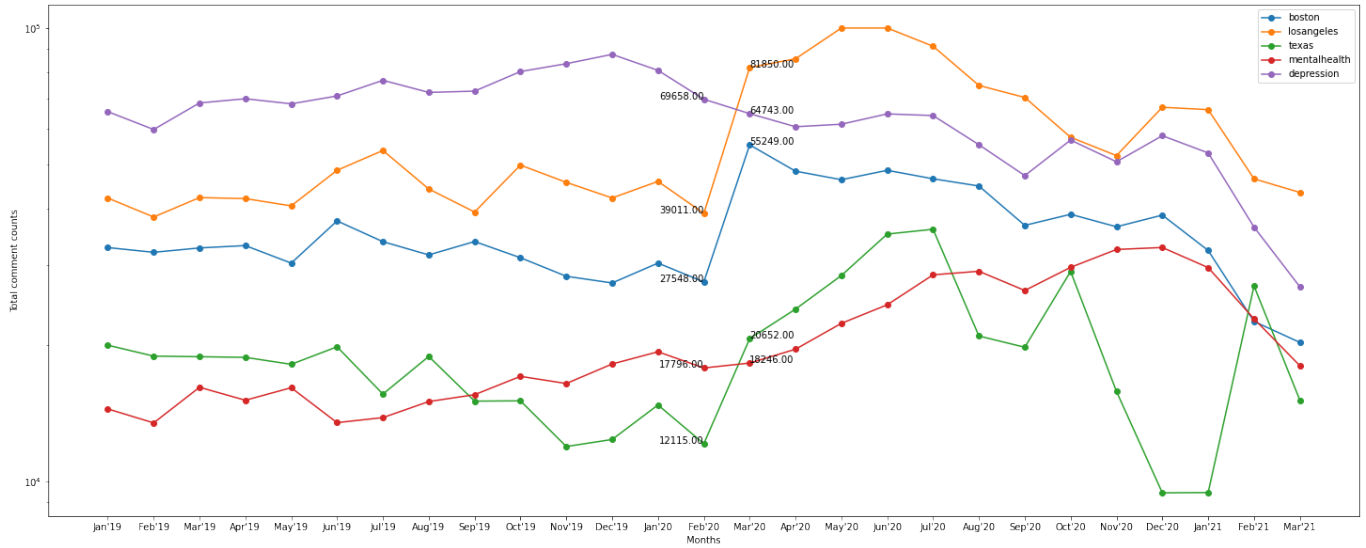


Figure 6. Total comment counts of subreddits per month

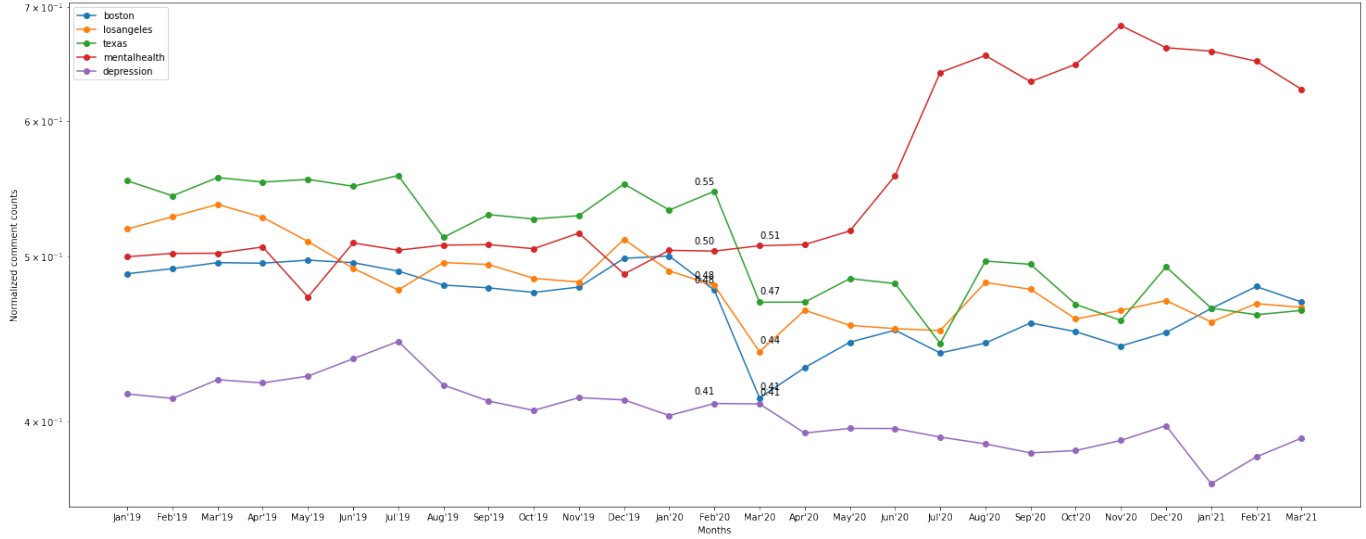


Figure 7. Percentage of positive sentiment comments of subreddits per month

To select the final sentiment prediction model, I compared a transformer based model vs DAN based model as seen in Fig 5. As explained earlier, transformer based model performed better than the DAN based model at sentiment prediction. The test accuracy for Transformer based model is 0.813 vs 0.798 of DAN based model.

IV. RESULTS

The Sentiment Prediction model, Fig 3, takes a subreddit comment as input and gives a value between 0 and 1 as output. We assign a class between 0 to 4 to the comments based on their predicted value. So, if the predicted value is less than or equal to 0.2 then I assign a class value of 0, if the predicted value is greater than 0.2 and less than or equal to 0.4 then I assign a class of 1, and similarly 2, 3, and 4. I consider any

comment with class less than or equal to 2 as comments with a negative sentiment while the rest are considered as comments with positive sentiment.

I did two main comparisons based on these predicted classes. Firstly, I analyzed the change in reddit traffic based on the number of comments posted in pre and during pandemic. The change in total count of comments per subreddit per month is visible in Fig 6. Secondly, I analyzed the change in percentage of positive sentiment comments in pre and during pandemic. The change in percentage of positive sentiment comments is visible in Fig 7.

Based on Fig 6, we can notice a sudden increase in the total number of subreddit comments for all subreddits under consideration as we move from Feb'2020 to Mar'2020. The only exception is the /r/mentalhealth, which increased

Subreddit	Feb 2020	Mar 2020
Boston	27548	55249
LosAngeles	39011	81850
Texas	12115	20652
Depression	69658	64743
MentalHealth	17796	18246

Table I

CHANGE IN TOTAL NUMBER OF COMMENTS PER SUBREDDIT

Subreddit	Feb 2020	Mar 2020
Boston	0.48	0.41
LosAngeles	0.48	0.44
Texas	0.55	0.47
Depression	0.41	0.41
MentalHealth	0.50	0.51

Table II

CHANGE IN PERCENTAGE OF POSITIVE SENTIMENT COMMENTS PER SUBREDDIT.

gradually.

Based on Fig 7, we can notice a drastic decrease in the percentage of sentiment counts for */r/boston*, */r/losangeles*, and */r/texas*. */r/depression* had almost negligible change, while we can notice a unique behavior in */r/mentalhealth*. There was no decrease in positive sentiment comments in case of */r/mentalhealth*, infact it increased a lot in the following months.

V. CONCLUSION

From the results, we can see that as the spread of Covid-19 gained speed in USA, the percentage of positive sentiment comments decreased a lot in the impacted cities like Boston, Los Angeles, and Texas. We can also notice that as various kinds of quarantine measures were applied in these cities, the number of comments increased a lot, signifying the starting of lockdowns and increase in social media usage. Also, the unique behavior of */r/mentalhealth* can be attributed to the fact that more people sought for help and provided support as the lockdowns continued and Covid-19 spread.

I also tried to understand the topic of comments in */r/mentalhealth* using a topic model based on Latent Dirichlet Allocation. The words like 'people', 'please' 'thank' and 'help' had the highest weight which points to people seeking and showing support to each other. This still needs further analysis and will be a part of future work.

APPENDIX

The code is available at <https://github.com/Rathore25/Reddit-Sentiment-Pre-Post-Covid>

REFERENCES

- [1] W. contributors, "Covid-19 — wikipedia, the free encyclopedia," 2021, accessed: 04/28/2021. [Online]. Available: <https://en.wikipedia.org/w/index.php?title=COVID-19&oldid=1020229208>
- [2] K. Katella, "5 things everyone should know about the coronavirus outbreak," 2020, accessed: 04/28/2021. [Online]. Available: <https://www.yalemedicine.org/news/2019-novel-coronavirus>
- [3] E. Koeze and N. Popper, "The virus changed the way we internet," 2020, accessed: 04/28/2021. [Online]. Available: <https://www.nytimes.com/interactive/2020/04/07/technology/coronavirus-internet-use.html>
- [4] A. F. . V. G. Marzouki, Y., "Understanding the buffering effect of social media use on anxiety during the covid-19 pandemic lockdown," 2021, accessed: 04/28/2021. [Online]. Available: <https://www.nature.com/articles/s41599-021-00724-x>
- [5] J. Qiu, B. Shen, M. Zhao, Z. Wang, B. Xie, and Y. Xu, "A nationwide survey of psychological distress among chinese people in the covid-19 epidemic: implications and policy recommendations," *General Psychiatry*, vol. 33, no. 2, 2020. [Online]. Available: <https://gpsych.bmj.com/content/33/2/e100213>
- [6] L. E. S. L. W. S. W. N. G. G. J. R. Samantha K Brooks, Rebecca K Webster, "The psychological impact of quarantine and how to reduce it: a rapid review of the evidence," 2020, accessed: 04/28/2021. [Online]. Available: [https://doi.org/10.1016/S0140-6736\(20\)30460-8](https://doi.org/10.1016/S0140-6736(20)30460-8)
- [7] J. I. S. G. M. H. T. P. P. Maria Effenberger, Andreas Kronbichler, "Association of the covid-19 pandemic with internet search volumes: a google trends analysis," 2020, accessed: 04/28/2021. [Online]. Available: <https://doi.org/10.1016/j.ijid.2020.04.033>
- [8] Wikipedia contributors, "Reddit — Wikipedia, the free encyclopedia," 2021, [Online; accessed 28-April-2021]. [Online]. Available: <https://en.wikipedia.org/w/index.php?title=Reddit&oldid=1020261853>
- [9] X. F. Jie Tao, "Toward multi-label sentiment analysis: a transfer learning based approach," 2020, accessed: 04/28/2021. [Online]. Available: <https://journalofbigdata.springeropen.com/articles/10.1186/s40537-019-0278-0>
- [10] G. Huang, Z. Liu, and K. Q. Weinberger, "Densely connected convolutional networks," *CoRR*, vol. abs/1608.06993, 2016. [Online]. Available: <http://arxiv.org/abs/1608.06993>
- [11] M. Long and J. Wang, "Learning transferable features with deep adaptation networks," *CoRR*, vol. abs/1502.02791, 2015. [Online]. Available: <http://arxiv.org/abs/1502.02791>
- [12] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," *CoRR*, vol. abs/1310.4546, 2013. [Online]. Available: <http://arxiv.org/abs/1310.4546>
- [13] J. Howard and S. Ruder, "Fine-tuned language models for text classification," *CoRR*, vol. abs/1801.06146, 2018. [Online]. Available: <http://arxiv.org/abs/1801.06146>
- [14] S. Merity, C. Xiong, J. Bradbury, and R. Socher, "Pointer sentinel mixture models," *CoRR*, vol. abs/1609.07843, 2016. [Online]. Available: <http://arxiv.org/abs/1609.07843>
- [15] A. Radford and K. Narasimhan, "Improving language understanding by generative pre-training," 2018, accessed: 04/28/2021. [Online]. Available: https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/language_understanding_paper.pdf
- [16] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, "GLUE: A multi-task benchmark and analysis platform for natural language understanding," *CoRR*, vol. abs/1804.07461, 2018. [Online]. Available: <http://arxiv.org/abs/1804.07461>
- [17] D. Cer, Y. Yang, S. Kong, N. Hua, N. Limtiaco, R. S. John, N. Constant, M. Guajardo-Cespedes, S. Yuan, C. Tar, Y. Sung, B. Strope, and R. Kurzweil, "Universal sentence encoder," *CoRR*, vol. abs/1803.11175, 2018. [Online]. Available: <http://arxiv.org/abs/1803.11175>
- [18] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *CoRR*, vol. abs/1706.03762, 2017. [Online]. Available: <http://arxiv.org/abs/1706.03762>
- [19] M. Iyyer, V. Manjunatha, J. Boyd-Graber, and H. Daumé III, "Deep unordered composition rivals syntactic methods for text classification," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Beijing, China: Association for Computational Linguistics, Jul. 2015, pp. 1681–1691. [Online]. Available: <https://www.aclweb.org/anthology/P15-1162>
- [20] R. Kiros, Y. Zhu, R. Salakhutdinov, R. S. Zemel, A. Torralba, R. Urtasun, and S. Fidler, "Skip-thought vectors," *CoRR*, vol. abs/1506.06726, 2015. [Online]. Available: <http://arxiv.org/abs/1506.06726>
- [21] M. L. Henderson, R. Al-Rfou, B. Strope, Y. Sung, L. Lukács, R. Guo, S. Kumar, B. Miklos, and R. Kurzweil, "Efficient natural language response suggestion for smart reply," *CoRR*, vol. abs/1705.00652, 2017. [Online]. Available: <http://arxiv.org/abs/1705.00652>
- [22] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, pp. 1735–80, 12 1997.

- [23] B. R. Go, A. and L. Huang, "Twitter sentiment classification using distant supervision," 2009, accessed: 04/28/2021. [Online]. Available: <https://www-cs.stanford.edu/people/alecmgo/papers/TwitterDistantSupervision09.pdf>