
Our Evaluation Metric Needs an Update to Encourage Generalization

Swaroop Mishra¹ Anjana Arunkumar¹ Chris Bryan¹ Chitta Baral¹

Abstract

Models that surpass human performance on several popular benchmarks display significant degradation in performance on exposure to Out of Distribution (OOD) data. Recent research has shown that models overfit to spurious biases and ‘hack’ datasets, in lieu of learning generalizable features like humans. In order to stop the inflation in model performance – and thus overestimation in AI systems’ capabilities – we propose a simple and novel evaluation metric, *WOOD Score*, that encourages generalization during evaluation.

1. Introduction

Training and evaluation in Machine Learning have been based on IID data. Encountering OOD data while testing is however, inevitable. This is due to several practical reasons such as: (i) difficulty in finding an ‘ideal’ evaluation set that captures the entire distribution, and (ii) an ‘ideal’ evaluation set identified during benchmark creation may not hold up against evolving benchmarks (Torralla & Efron, 2011; Quionero-Candela et al., 2009; Hendrycks et al., 2020).

Recent language models have surpassed human performance across several popular AI benchmarks such as Imagenet (Russakovsky et al., 2015), SNLI (Bowman et al., 2015) and SQUAD (Rajpurkar et al., 2016). However, a substantial performance drop is seen in these models on exposure to OOD data (Bras et al., 2020; Eykholt et al., 2018; Jia & Liang, 2017). A growing number of recent works (Gururangan et al., 2018; Poliak et al., 2018; Kaushik & Lipton, 2018; Tsuchiya, 2018; Tan et al., 2019; Schwartz et al., 2017) have shown that models are not truly learning the underlying task; their performance in topping leaderboards can rather be attributed to their exploitation of spurious biases. This fundamental issue in learning leads to the overestimation of AI system performance (Bras et al., 2020; Sakaguchi et al., 2019; Hendrycks et al., 2019), hence restricting AI

deployment in several critical domains such as medicine (Hendrycks & Gimpel, 2016a).

Several approaches have been proposed to address this issue at various levels: (i) *Data* – filtering of biases (Bras et al., 2020; Li & Vasconcelos, 2019; Li et al., 2018; Wang et al., 2018), quantifying data quality, controlling data quality, using active learning, and avoiding the creation of low quality data (Mishra et al., 2020; Nie et al., 2019; Gardner et al., 2020; Kaushik et al., 2019), and (ii) *Model* – utilizing prior knowledge of biases to train a naive model exploiting biases, and then subsequently training an ensemble of the naive model and a ‘robust’ model to learn rich features (Clark et al., 2019; He et al., 2019; Mahabadi & Henderson, 2019). However, addressing this issue via an *evaluation metric perspective* remains underexplored.

In pursuit of this metric, we first analyze recent works (Hendrycks et al., 2020; Bras et al., 2020; Wang et al., 2020; Hendrycks & Dietterich, 2019; Talmor & Berant, 2019) involving datasets that have been paired with an OOD counterpart, to isolate factors that distinguish OOD samples from IID samples. In this process, we identify four major problems that the Machine Learning community needs to address: (q_1) evaluating a model’s generalization capability involves the *overhead of OOD dataset identification*, (q_2) *absence of a clear boundary separating IID from OOD* makes OOD identification even harder, (q_3) generalization is merely evaluated, not *encouraged during evaluation*, and (q_4) the issue of *inflated model performance* (and thus AI systems’ overestimation) is yet to be addressed effectively.

Our analysis yields *Semantic Textual Similarity (STS) of test data with respect to the training data as a distinguishing factor between IID and OOD*. We show the same across the prediction probabilities of ten models in two datasets. We then divide the datasets into several hierarchies, based on STS (and thus the degree of OOD characteristics). We can therefore reasonably identify samples within the dataset which have higher OOD characteristic levels. This allows the same dataset to be used to evaluate OOD (q_1). STS can not only be used to draw a boundary between IID and OOD, but also to control the degree of OOD characteristics in a dataset (q_2). We further propose a metric, *Weighting Out of Distribution Score (WOOD Score)*, by weighting each test sample in proportion to its degree of OOD characteristics;

¹Department of Computer Science, Arizona State University. Correspondence to: Swaroop Mishra <srnmishr1@asu.edu>.

higher levels of OOD characteristics imply higher weightage. Our intuition is simple: *models must solve data with higher weightage to have higher accuracy*. This compels a model to generalize in order to dominate leaderboards (q_3). Our results show a decrease in model performance when evaluated using WOOD Score, resulting in lower benchmark accuracy (q_4). Our work inspires several potential solutions to handle OOD, using an evaluation metric perspective.

2. Differentiating IID and OOD using STS

We use two movie review datasets: SST-2 (Socher et al., 2013) and IMDb (Maas et al., 2011), which contain succinct expert reviews and full length general reviews respectively. We utilize SST-2 as the IID dataset and IMDb as the OOD dataset, and evaluate them using ten models: Bag-of-words (BoW) model (Harris, 1954), word embedding - word2vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014) encoded with three models - word averages (Wieting et al., 2015), LSTM (Hochreiter & Schmidhuber, 1997) and CNN (LeCun et al., 1995), and pretrained transformer models -BERT Base and Large (Devlin et al., 2018) along with GELU (Hendrycks & Gimpel, 2016b) and RoBERTA (Liu et al., 2019), following a recent work on OOD Robustness (Hendrycks et al., 2020).

Experimental Setup: We perform a complete analysis by finding the STS between every pair of training set and test set samples. We sort samples of the test set in a descending order based on the average STS value with varying percentages of the train set samples. We consider the top 1% – 100% of the training data (obtained by sorting train set samples in descending order of STS against each test set sample) with nine total steps, as similarity between the train and test sets is a task dependent hyperparameter, that trades off between inductive bias and spurious bias (Mishra et al., 2020; Gorman & Bedrick, 2019). We train models on the IID data (SST-2) and evaluate on both the IID test set (SST-2) and the OOD test set (IMDb). We compare model predictions with the average STS value for each sample.

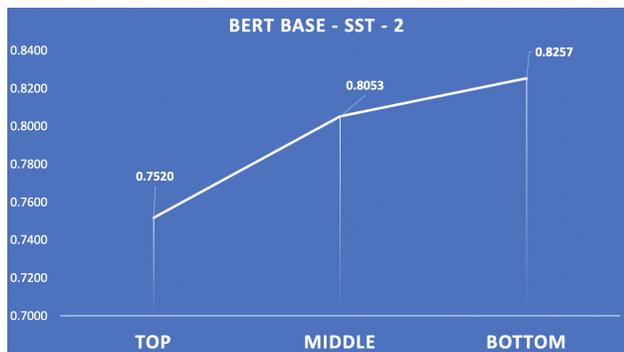


Figure 1. Softmax probabilities of incorrect classifications using BERT-Base model across test samples of SST-2 and IMDb in decreasing order train (SST-2)-test similarity.

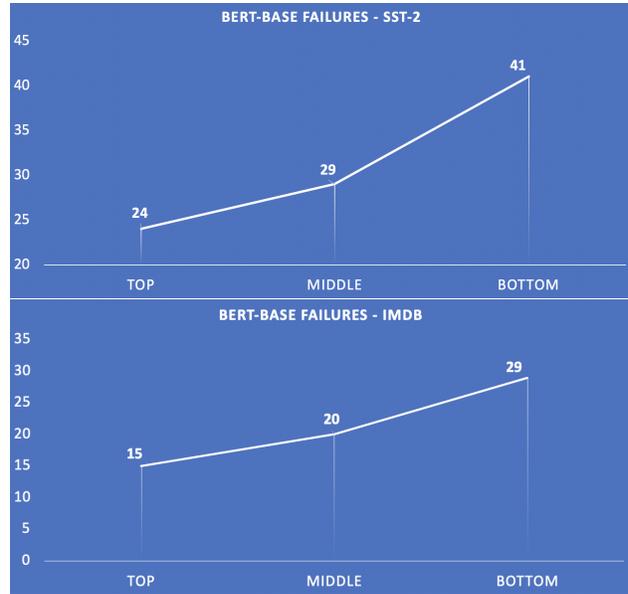


Figure 2. Number of incorrect classifications using BERT-Base model across test samples of SST-2 and IMDb in decreasing order of train (SST-2)-test similarity.

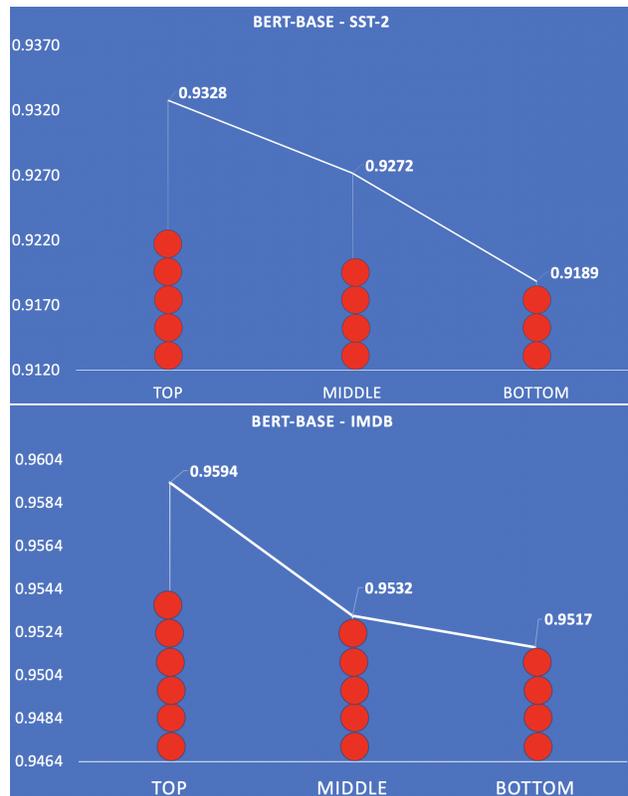


Figure 3. Softmax probabilities of correct classifications using BERT-Base model across test samples of SST-2 and IMDb in decreasing order of train (SST-2)-test similarity. The circles help visualize the impact of uncertainty in classification by representing the potential number of incorrectly classified samples per 20 samples for each split of data.

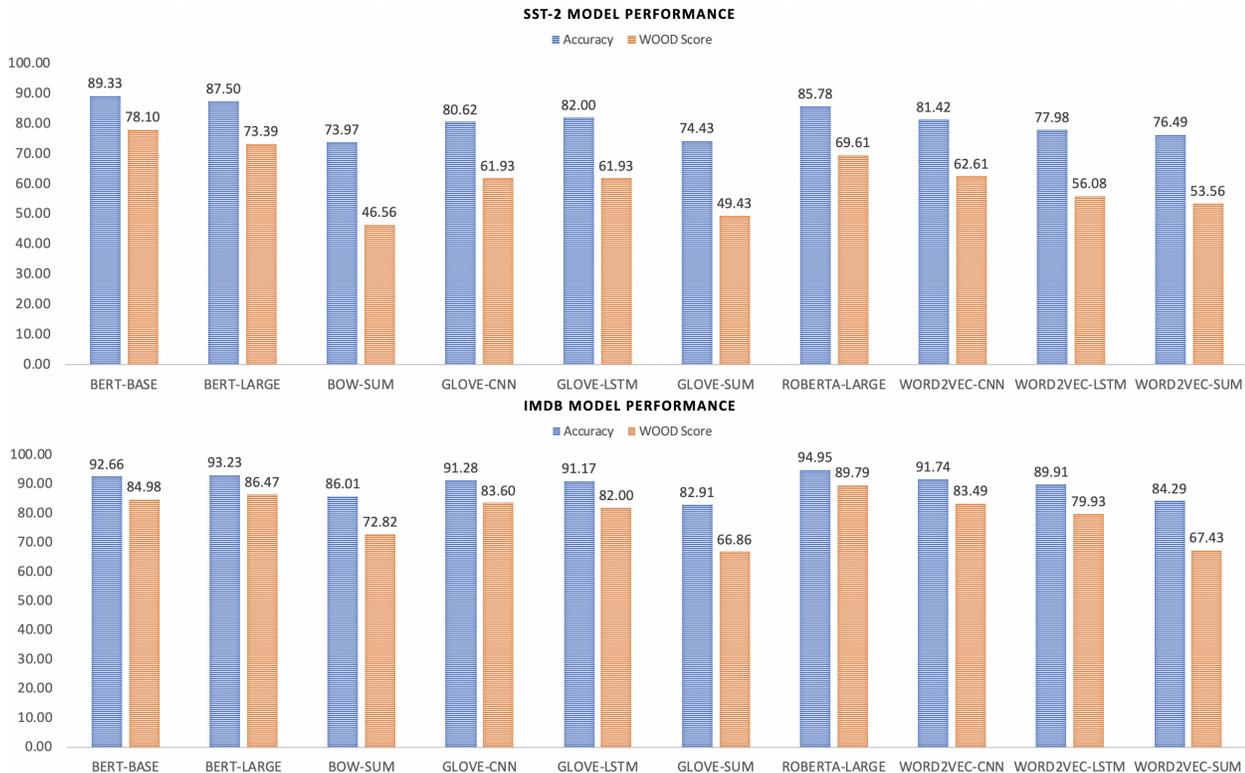


Figure 4. Accuracy and WOOD Score of SST-2 and IMDB across models.

Results: We find three broad patterns: (i) Samples with higher average STS value are classified correctly with higher confidence (Figure 3), (ii) Incorrect classification frequency increases as STS value decreases (Figure 2), and (iii) Confidence value (softmax probability) of an incorrect classification increases as we move towards the region having lower STS (Figure 1).

3. Metric for Equitable Evaluation of Data

Benchmarks and leaderboards guide model development. We propose *equitable data evaluation* using a weighted metric, in lieu of the conventional metric that does uniform evaluation. Our intuition is to weight the evaluation score of a sample in proportion to its ‘hardness’, i.e. the level of OOD characteristics it displays. In this paper, we define a metric, *WOOD Score*, using STS (with the training set) as an inverse indicator of ‘hardness’, based on the findings of Section 2.

Formalization: Let X represent a dataset where X_{Test} is the test set spanned by i and X_{Train} is the train set. E represents the evaluation metric (which depends on the application), p is the degree of OOD characteristics a sample has, S represents STS, a allows for the control of p based on S , b is the number of train samples considered that have higher similarity values than rest of the dataset. W_{opt} rep-

resents our proposed metric in generic form, and W_{acc} is the proposed accuracy metric in this paper. A represents accuracy. A_3 , A_2 , and A_1 represent the accuracies of the categories of samples having the highest, moderate, and lowest degrees of OOD characteristics respectively.

$$W_{opt} = \sum_{X_{Test}} E_i p_i \quad (1)$$

$$p = \frac{a}{\sum_{X_{Train}} \max_b S} \quad (2)$$

$$W_{acc} = A_1 + 2A_2 + 3A_3 \quad (3)$$

Controlling Benchmark Accuracy Using Hyperparameters: Benchmark accuracy can be controlled using appropriate values of a and b . Using W_{acc} for both datasets across ten models has resulted in a significant reduction in accuracy, as illustrated in Figure 4.

4. Discussion

Extension to Other Metrics Beyond Accuracy: Our focus in this paper is on accuracy, but the same issue of ‘uniform evaluation score’ prevails in all metrics such as Pearson’s correlation score, BLEU, and F_1 score. A potential future work is to extend our metric beyond accuracy.

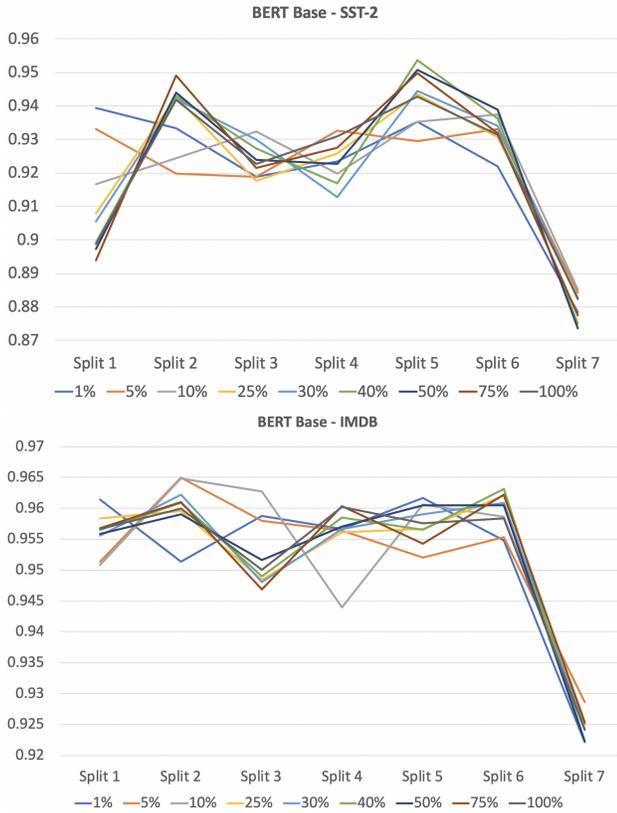


Figure 5. The top $b\%$ of training samples is obtained by sorting in descending order of STS with each test set sample. The test set samples are then divided into seven splits, based on decreasing STS values averaged over the top $b\%$ of training samples considered.

Augmenting STS: Detailed analysis shows that STS’s distinguishing capacity may not follow monotonic behavior for certain cases, such as those illustrated in Figure 5 and 6. Similarity across several granularities – such as word, bigram, and trigram – can be used to augment STS and increase the robustness of ‘hardness’ evaluation.

OOD Drop is More Than IID with the Updated Metric: Figure 5 illustrates that Benchmark accuracy drops more in the IMDB dataset (OOD) than SST-2 (IID). IMDB has a different writing style and sample size in contrast to SST-2. This may indicate that, our metric is more effective for applications involving diverse data considered as OOD.

Leveraging the Metric to Learn a Task Better: Hardness of an anticipated OOD task can be reduced by adding similar data to the training set. Our metric can guide the data augmentation process to learn a task better.

Strengthening ‘in-house’ IID (acting OOD): We further observe that, IID data, even with STS calibration, may not represent many properties of an OOD data sample – such as variations in writing style, topic, vocabulary, sentence length, and number of sentences. Contrast sets (Gardner et al., 2020) can be utilized to partially address this issue by strengthening the test set. We recommend that dataset

creators go beyond the common patterns found in a dataset, and draw patterns from other

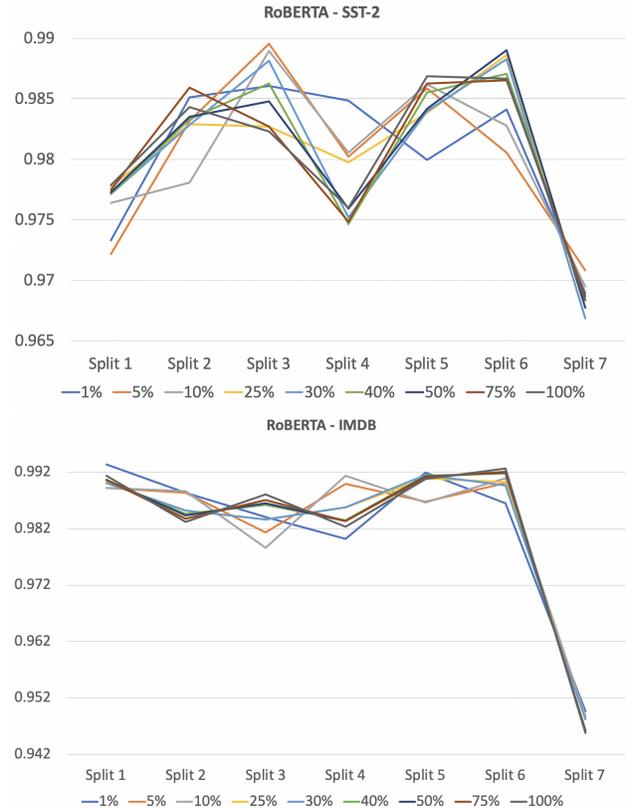


Figure 6. Detailed analysis where the test set is divided in to seven splits. Figure shows variation across a range of b values for the RoBERTA model.

datasets intended for the same task, while creating contrast sets, in order to fully address this issue. Robustness of a dataset can also be improved by adding more application-centric variations – such as anticipating various possibilities in question format, structure, language, reasoning skills required, syntax, and numerical reasoning requirements.

5. Conclusion

We propose a STS based approach to address certain important issues in robustness research: finding a suitable OOD dataset, and drawing a valid boundary between IID and OOD. Our approach also helps in controlling the degree of OOD characteristics. We also propose *WOOD Score*, a metric that implicitly encourages generalization by weighting the evaluation scores of samples in proportion to their hardness. This results in reduced model performance and benchmark accuracy, addressing the issue of model performance inflation and overestimation of progress in AI. We show the efficacy of our work on ten popular models across two NLP datasets. Finally, we provide insights into several future works on encouraging generalization and improving robustness from a metric perspective.

Acknowledgements

We thank the anonymous reviewers, Dan Hendrycks (UC Berkeley) and Xiaoyuan Liu (Shanghai Jiao Tong University) for their thoughtful feedback. We also thank Jason Yalim and ASU HPC for their consistent support. The support of DARPA SAIL-ON program (W911NF2020006) is gratefully acknowledged.

References

- Bowman, S. R., Angeli, G., Potts, C., and Manning, C. D. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*, 2015.
- Bras, R. L., Swayamdipta, S., Bhagavatula, C., Zellers, R., Peters, M. E., Sabharwal, A., and Choi, Y. Adversarial filters of dataset biases. *arXiv preprint arXiv:2002.04108*, 2020.
- Clark, C., Yatskar, M., and Zettlemoyer, L. Don't take the easy way out: Ensemble based methods for avoiding known dataset biases. *arXiv preprint arXiv:1909.03683*, 2019.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., and Song, D. Robust physical-world attacks on deep learning visual classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1625–1634, 2018.
- Gardner, M., Artzi, Y., Basmova, V., Berant, J., Bogin, B., Chen, S., Dasigi, P., Dua, D., Elazar, Y., Gottumukkala, A., et al. Evaluating nlp models via contrast sets. *arXiv preprint arXiv:2004.02709*, 2020.
- Gorman, K. and Bedrick, S. We need to talk about standard splits. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pp. 2786–2791, 2019.
- Gururangan, S., Swayamdipta, S., Levy, O., Schwartz, R., Bowman, S. R., and Smith, N. A. Annotation artifacts in natural language inference data. *arXiv preprint arXiv:1803.02324*, 2018.
- Harris, Z. S. Distributional structure. *Word*, 10(2-3):146–162, 1954.
- He, H., Zha, S., and Wang, H. Unlearn dataset bias in natural language inference by fitting the residual. *arXiv preprint arXiv:1908.10763*, 2019.
- Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016a.
- Hendrycks, D. and Gimpel, K. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016b.
- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., and Song, D. Natural adversarial examples. *arXiv preprint arXiv:1907.07174*, 2019.
- Hendrycks, D., Liu, X., Wallace, E., Dziedzic, A., Krishnan, R., and Song, D. Pretrained transformers improve out-of-distribution robustness. *arXiv preprint arXiv:2004.06100*, 2020.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Jia, R. and Liang, P. Adversarial examples for evaluating reading comprehension systems. *arXiv preprint arXiv:1707.07328*, 2017.
- Kaushik, D. and Lipton, Z. C. How much reading does reading comprehension require? a critical investigation of popular benchmarks. *arXiv preprint arXiv:1808.04926*, 2018.
- Kaushik, D., Hovy, E., and Lipton, Z. C. Learning the difference that makes a difference with counterfactually-augmented data. *arXiv preprint arXiv:1909.12434*, 2019.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- LeCun, Y., Bengio, Y., et al. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- Li, Y. and Vasconcelos, N. Repair: Removing representation bias by dataset resampling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9572–9581, 2019.
- Li, Y., Li, Y., and Vasconcelos, N. Resound: Towards action recognition without representation bias. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 513–528, 2018.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

- Maas, A. L., Daly, R. E., Pham, P. T., Huang, D., Ng, A. Y., and Potts, C. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies-volume 1*, pp. 142–150. Association for Computational Linguistics, 2011.
- Mahabadi, R. K. and Henderson, J. simple but effective techniques to reduce biases. *arXiv preprint arXiv:1909.06321*, 2019.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119, 2013.
- Mishra, S., Arunkumar, A., Sachdeva, B., Bryan, C., and Baral, C. Dqi: Measuring data quality in nlp. *arXiv*, pp. arXiv–2005, 2020.
- Nie, Y., Williams, A., Dinan, E., Bansal, M., Weston, J., and Kiela, D. Adversarial nli: A new benchmark for natural language understanding. *arXiv preprint arXiv:1910.14599*, 2019.
- Pennington, J., Socher, R., and Manning, C. D. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543, 2014.
- Poliak, A., Naradowsky, J., Haldar, A., Rudinger, R., and Van Durme, B. Hypothesis only baselines in natural language inference. *arXiv preprint arXiv:1805.01042*, 2018.
- Quionero-Candela, J., Sugiyama, M., Schwaighofer, A., and Lawrence, N. D. *Dataset shift in machine learning*. The MIT Press, 2009.
- Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3): 211–252, 2015.
- Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. Winogrande: An adversarial winograd schema challenge at scale. *arXiv preprint arXiv:1907.10641*, 2019.
- Schwartz, R., Sap, M., Konstas, I., Zilles, L., Choi, Y., and Smith, N. A. The effect of different writing tasks on linguistic style: A case study of the roc story cloze task. *arXiv preprint arXiv:1702.01841*, 2017.
- Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A. Y., and Potts, C. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1631–1642, 2013.
- Talmor, A. and Berant, J. Multiqa: An empirical investigation of generalization and transfer in reading comprehension. *arXiv preprint arXiv:1905.13453*, 2019.
- Tan, S., Shen, Y., Huang, C.-w., and Courville, A. Investigating biases in textual entailment datasets. *arXiv preprint arXiv:1906.09635*, 2019.
- Torralba, A. and Efros, A. A. Unbiased look at dataset bias. In *CVPR 2011*, pp. 1521–1528. IEEE, 2011.
- Tsuchiya, M. Performance impact caused by hidden bias of training data for recognizing textual entailment. *arXiv preprint arXiv:1804.08117*, 2018.
- Wang, T., Zhu, J.-Y., Torralba, A., and Efros, A. A. Dataset distillation. *arXiv preprint arXiv:1811.10959*, 2018.
- Wang, Y., Chen, X., You, Y., Li, L. E., Hariharan, B., Campbell, M., Weinberger, K. Q., and Chao, W.-L. Train in germany, test in the usa: Making 3d object detectors generalize. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11713–11723, 2020.
- Wieting, J., Bansal, M., Gimpel, K., and Livescu, K. Towards universal paraphrastic sentence embeddings. *arXiv preprint arXiv:1511.08198*, 2015.