

SparCAssist

A Model Risk Assessment Assistant Based on Sparse Generated Counterfactuals

Zhang, Zijian; Setty, Vinay; Anand, Avishek

DOI

[10.1145/3477495.3531677](https://doi.org/10.1145/3477495.3531677)

Publication date

2022

Document Version

Final published version

Published in

SIGIR 2022 - Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval

Citation (APA)

Zhang, Z., Setty, V., & Anand, A. (2022). SparCAssist: A Model Risk Assessment Assistant Based on Sparse Generated Counterfactuals. In *SIGIR 2022 - Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (pp. 3219-3223). (SIGIR 2022 - Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval). Association for Computing Machinery (ACM). <https://doi.org/10.1145/3477495.3531677>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Green Open Access added to TU Delft Institutional Repository

'You share, we take care!' - Taverne project

<https://www.openaccess.nl/en/you-share-we-take-care>

Otherwise as indicated in the copyright section: the publisher is the copyright holder of this work and the author uses the Dutch legislation to make this work public.

SparCAssist: A Model Risk Assessment Assistant Based on Sparse Generated Counterfactuals

Zijian Zhang

zzhang@l3s.de

L3S Research Center

Hannover, Lower Saxony, Germany

Vinay Setty

vsetty@acm.org

L3S Research Center

Germany

University of Stavanger

Stavanger, Norway

Avishek Anand

avishek.anand@tudelft.nl

L3S Research Center

Germany

Delft University of technology

Delft, Netherlands

ABSTRACT

We introduce SparCAssist, a general-purpose risk assessment tool for the machine learning models trained for language tasks. It evaluates models' risk by inspecting their behavior on counterfactuals, namely out-of-distribution instances generated based on the given data instance. The counterfactuals are generated by replacing tokens in rational subsequences identified by ExPred, while the replacements are retrieved using HotFlip or Masked Language Model based algorithms. The main purpose of our system is to help the human annotators to assess the model's risk on deployment. The counterfactual instances generated during the assessment are the by-product and can be used to train more robust NLP models in the future.

CCS CONCEPTS

• **Information systems** → *Online analytical processing*; • **Human-centered computing** → *User interface toolkits*; • **Computing methodologies** → *Reasoning about belief and knowledge*.

KEYWORDS

Interpretable Machine Learning, Counterfactual Interpretation, Data-annotation tools, Human-in-the-loop Machine Learning

ACM Reference Format:

Zijian Zhang, Vinay Setty, and Avishek Anand. 2022. SparCAssist: A Model Risk Assessment Assistant Based on Sparse Generated Counterfactuals. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*, July 11–15, 2022, Madrid, Spain. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3477495.3531677>

1 INTRODUCTION

In this paper we present a demonstration of SparCAssist: a tool for human-assisted debugging of machine learning (ML) models for language tasks. Debugging machine learning models in terms of estimating if a trained model is fit for deployment is considered more of an art than a science with a lot depending on the heuristics

employed by the model developer [1]. Standard deployment checks often involve heavy reliance on performance metrics on validation and test sets to evaluate the generalization and robustness of systems. That means there is heavy focus on inspecting incorrectly predicted instances. We argue that ensuring good performance is not enough. Deep neural networks for language tasks are known to be particularly fragile on out-of-distribution instances even if they deliver impressive validation set accuracy [24].

Some of the known risks of deploying such over-parameterized models such as transformers is that they tend to memorize undesirable shortcuts in the training data that lead to reasonable validation performance but leads to failures in real-world scenarios [37]. Unlike common practice of ML metrics, we propose to provide a system that uses counterfactuals as examples to help the user, in this case the model developer, to understand the risks of the model under consideration. Specifically, given a trained ML model, our intention is to allow the user to identify instances that have potential risk during deployment. An instance is considered *risky* if its minor human-understandable modifications result in undesirable predictions.

Consider a *document classification task* of classifying review text of movies into *positive* and *negative* sentiments. Say a trained sentiment classifier classifies the input “the movie was awesome” as a positive sentiment. However, by slightly perturbing it to a counterfactual sentence “the movie ~~was~~ *is* awesome” results in a negative sentiment. This counterfactual sentence exposes two potential risks in the model, firstly, the model focuses on non-essential words towards making a prediction and secondly, the model is not robust under non-decisive perturbation. Other approaches that provide insights into the reasons behind predictions of complex machine learning models, a large variety of post-hoc explanation methods have been proposed [17, 22, 32]. However, all these explanation approaches operate on existing instances on the training or the test set. Therefore they cannot effectively assess the potential risk when the models are deployed in the field.

The aim of SparCAssist is to modify existing counterfactual generation approaches to help humans identify instances that might be particularly risky. We allow users enough control to choose input instances that are of interest and pieces of text out of them within that might be important for the task under consideration. In allowing a fine-grained control to the inputs choosing and feedback providing, the users are able to both get a sense of the functioning of the model as well as be educated about the potential risks it might possess. In SparCAssist, we work on the task of sentiment classification as an example but in principle SparCAssist can be

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '22, July 11–15, 2022, Madrid, Spain

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-8732-3/22/07...\$15.00

<https://doi.org/10.1145/3477495.3531677>

extended to a large set of language tasks, for example question answering, fact checking, textual entailment etc.

The system then provides multiple “*perturbations*” of the input sentence that result in a change in the prediction. The perturbations should be small enough so that we can effectively assess the model’s local risk, specifically in the close vicinity surrounding a given instance. Therefore, we are only interested in *sparse* perturbations or we only allow a small number of perturbations to the original input. The users’ role is to assess whether the generated counterfactual is *plausible*, *natural* and if the model’s prediction is *correct* on counterfactuals. Plausible here means that the input text is believed to be written or conceived by a human. In terms of approach we start from a trained model and identify sentences that need to be scrutinized. We choose this subset of sentences using ExPred[35], a select-then-predict system that is interpretable by design. We then use adversarial approaches like HotFlip [5] and Masked Language Model approaches to generate counterfactuals. These generated counterfactuals are then shown to humans who in turn interact with the suggestions to establish its plausibility and validity. Finally, we compute the *risk* of the model under consideration from the trail of user interactions.

The primary goal is to assess the model’s deficiencies using the perturbations, and then our secondary goal is to generate a dataset of counterfactual explanations, which are grammatically plausible and meaningful from human perspective, could be the by-product of the system.

1.1 Related Work

We divide the related work into two parts – counterfactual generation, and automated approaches for model debugging. The counterfactual generation approaches can be categorized into two classes: automatic approaches and human-annotation approaches. Automatic approaches for generating counterfactuals are predominantly based on gradient-based optimization to generate adversarial examples like HotFlip [5], AutoPrompt [25], and Imitation Attacks [33]. Another set of approaches to automatically generate counterfactuals is to use prompt-based approaches like Polyjuice [34]. However, the main disadvantage of automatic approaches is that the generated counterfactuals are not natural, i.e., likely to be not plausible from human perspective. Prompt-based approaches are slightly more natural than gradient-based approaches, but they rely on heuristic prompts that have to be manually provided.

On the other hand, some approaches conduct user-study with crowd workers to obtain human-understandable counterfactual datasets [10, 18]. These approaches are understandably non-scalable and exhibit high variance due to subjectivity and the well-known Rashomon effect [3]. We use both gradient-based and prompt-based approaches for counterfactual generation but instead use humans-in-the-loop methods for quality assessment. More precisely, we take a hybrid approach for producing counterfactuals with utility of the counterfactual as the most desirable outcome.

The second piece of related work is regarding debugging ML models that have been used in text classification [21], search [27, 29, 31], and many language tasks in general [6, 9] by using explanations or interpretable machine learning approaches called explanation-based human debugging (EBHD). Lertvittayakumjorn and Toni

[12] recently review EBHD approaches that exploit explanations to enable humans to give feedback and debug NLP models. Unlike adversarial algorithms [20, 33] the goals of these systems is to explicitly expose model’s deployment risk by ensuring that models are *right for the right reason*. In addition to academic efforts, model debugging tools like [28, 36], SageMaker¹ provided by Amazon are also available. Other human-in-the loop systems include [8, 30] to update models and knowledge bases. However, they mostly focus on training process and do not provide a post-hoc model risk assessment.

2 SOLUTION APPROACH

The primary goal of our SparCAssist tool is to assist users to assess the deployment risk of the NLP classification model trained on specific datasets using counterfactual instances (counterfactuals). The input to SparCAssist is a model together with a dataset using which the risk assessment of the model is to be performed. The SparCAssist generates counterfactuals based on a selected data instance and tries to flip the model’s prediction for that instance. The tool assesses the risks of the model based on the user’s feedback on the plausibility and meaningfulness of the counterfactuals. The secondary output SparCAssist tool is plausible and meaningful counterfactual dataset for future training. In this section, we state our problem by first defining the counterfactual formally with its three essential properties and then formalize the model’s risk.

2.1 Interpretation Using Counterfactuals

Interpretation-by-Counterfactuals, also known as interpretation using minimal contrastives [24] is a model-interpretation technique that allows humans to understand the reason behind prediction by investigating models’ behavior under counterfactual inputs [14, 19]. According to [7, 24], in order to assess the risk of the models, the counterfactuals need to satisfy following rules:

- 1 The number of replaced tokens should be as few as possible.
- 2 The counterfactual instance should be plausible and meaningful [2].

Like in HotFlip [5] and Polyjuice [34], we mainly focus on the counterfactual generated by replacing the tokens. These two approaches, however, do not constrain the number of replaceable tokens while generating the counterfactuals. This will result in an overwhelming number of alternative counterfactuals for human annotators. In SparCAssist, we address this issue by firstly focusing on sentences containing rationale subsequences identified by the ExPred model [35]. We then only allow *rationale* tokens identified by ExPred in the selected sentence to be replaced, while the other tokens are kept fixed. Alternatively to the ExPred, users can manually select the tokens allowed to be replaced in the rationale sentences.

More formally, we use $\mathbf{x} = [x_0, x_1, \dots, x_l]$ to represent the input sequence (rationale sentence) to our tool and the model to be assessed is represented as \mathcal{M} , while it’s predicted label for a given instance x is denoted as $\hat{y} = \mathcal{M}(\mathbf{x} \odot \mathbf{m})$, where $\mathbf{m} \in \{0, 1\}^l$ is the rationale mask, \odot stands for element-wise multiplication. We then define a counterfactual generation function $C_{x_i \rightarrow v}$ which replaces

¹<https://docs.aws.amazon.com/sagemaker/latest/dg/train-debugger.html>

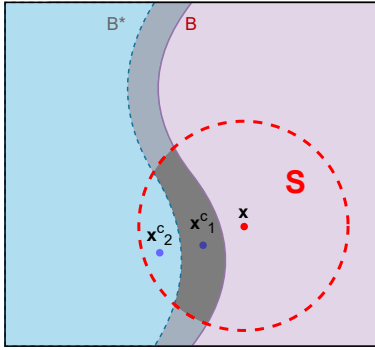


Figure 1: Risk and unexpected behavior in the vicinity of an inside-distribution-instance x . The dark gray area is the risky area.

the rationale token x_i with a counterfactual token v from the vocabulary V . A counterfactual with only one token replaced is called *1-counterfactual* and represented as $\mathbf{x}_x^c = C_{x_i \rightarrow v | i \in \mathcal{R}(m), v \in V}(\mathbf{m} \odot \mathbf{x})$. The condition $i \in \mathcal{R}(m)$ restricts the tokens to be replaced to the rationale mask \mathbf{m} . We repeat the above processes until the prediction on the counterfactual sentence changes.

2.2 Model Risk Assessment

We consider a model as “risky”, when its behavior on the counterfactuals do not match domain experts’ (users’) decision, despite its high accuracy on the given, or inside-distribution, instances, including instances from the train, dev-test, and test set [11]. Since the global risk is usually intractable [23], we consider its local counterpart instead, namely how wrong does the model predict on the counterfactuals sampled from the vicinity of an instance in the dataset.

As illustrated in Fig. 1, for a counterfactual \mathbf{x}^c sampled from the vicinity of an inside-distribution instance $\mathbf{x}^c \in S(\mathbf{x})$, if the decision of the model $\mathcal{M}(\mathbf{x}^c)$ is dissonant with the decision of domain experts’, we call this an “risky behavior” and the area resulting in risky behavior “risky volume”. The local risk score of the model is the ratio between risky volume and the volume of S . In more detail, we assume an “optimal” decision boundary B^* as a hyperplane resulting in decisions aligning with domain experts’ prediction, while B is the decision boundary of the model. The model’s local risk in S is then the ratio between the volume of V and S , where V is the interior enclosed by S and both decision boundaries as V , namely $\partial V = (B^* \cup B) \cap S$. We can estimate it by dividing the number of samples resulting in risky behaviors with the number of samples drawn universally from S using the Monte-Carlo strategy. However, it is difficult to draw samples from such vicinity due to the discreteness of human language sentences. Therefore, instead of drawing counterfactual samples universally, we focus on such counterfactuals resulting in the model’s prediction flipping. Such examples locate right across the model’s decision boundary and will most likely expose the model’s risk.

For example, the sentiment classifier should predict \mathbf{x} = “I love this movie.” as positive and the counterfactuals \mathbf{x}_2^c = “I love hate this movie.” to negative. Because this is in accord with users’ decisions,

the model’s risk does not increase. But if the model predicts \mathbf{x}_1^c = “I like this movie film.” to negative, the risk score increases, as users will not agree with this prediction. We ask the annotator to rate the prediction ($f \in 1 - 5$, f for faithfulness) by asking to how much extent does the model predicts in concord with her expectation. We do not ask for a binary rating because the optimal decision boundaries are vague and subjective. Then for each annotator, we calculate her risk score around the instance \mathbf{x} by: $\text{risk} = \sum_{\mathbf{x}^c \in \mathcal{X}^c} (5 - f(\mathbf{x}^c)) / |\mathcal{X}^c|$, where \mathcal{X}^c represents the set of counterfactual samples provided by this annotator. The final risk of the model in S is estimated by the weighted average of risk scores from all annotators, where the weight is the number of the counterfactuals generated in her/his session. It is worth mentioning that it requires further research to decide the exact number of counterfactuals needed to effectively evaluate the model’s risk.

3 SYSTEM DESIGN

As illustrated in Figure 2, we decompose the SparCAssist system into three components: **Rationale Selection**, **Position Attribution** and **Word Replacement**. The SparCAssist generates counterfactuals are as follows:

- (1) It first randomly selects an instance \mathbf{x} from the given dataset which may or may not be labeled. For the demo, we use the “Movies” datasets for the sentiment analysis task [4].
- (2) It then leverages the ExPred model [35] which is a two-phase select-then-predict model pre-trained on the movies dataset to select the **rationale tokens**.
- (3) For all those rationale tokens, we **attribute their positions** and select candidate position(s) based on the attribution. In our SparCAssist, we are only interested in such counterfactuals changes the model’s prediction ($\mathcal{M}(c) \neq \hat{y}$), therefore:
- (4) It **replaces the tokens** on each selected position so that the estimated loss regarding the alternative label other than the original prediction is minimized.
- (5) It then selects the final single counterfactual with only one token replaced based on their actual prediction loss.
- (6) We repeat steps 3 to 5 until the prediction changes or the number of word replacements exceeds a pre-defined upper limit. Based on our empirical observation, the set the upper limit to five replacements.

We integrate two counterfactual generation approaches to SparCAssist, the Taylor-expansion based HotFlip and a Masked Language Model (MLM)-based algorithm used in Polyjuice. These two approaches are decomposed into two phases, namely **position attribution** and **token selection** as step 3 and 4 introduced above. We explain these two approaches below:

HotFlip approach: For this approach, we follow [5] to search for positions and words using the beam search. More formally, we attribute positions with the dot-product of corresponding token embedding with the partial derivative of the loss regarding the alternated prediction label with respect to the token embedding, i.e. $\frac{\partial}{\partial \mathbf{x}_i} l(\hat{y}, \hat{y} \neq \hat{y}) \cdot \mathbf{e}_i$, where \mathbf{e}_i represents the embedding of the i -th token. We select the top- p , $p \leq r$ of the position, where r is the number of rationale tokens. Then in token selection, for

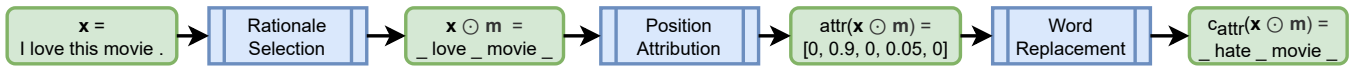


Figure 2: The system diagram of the counterfactual generation assistant

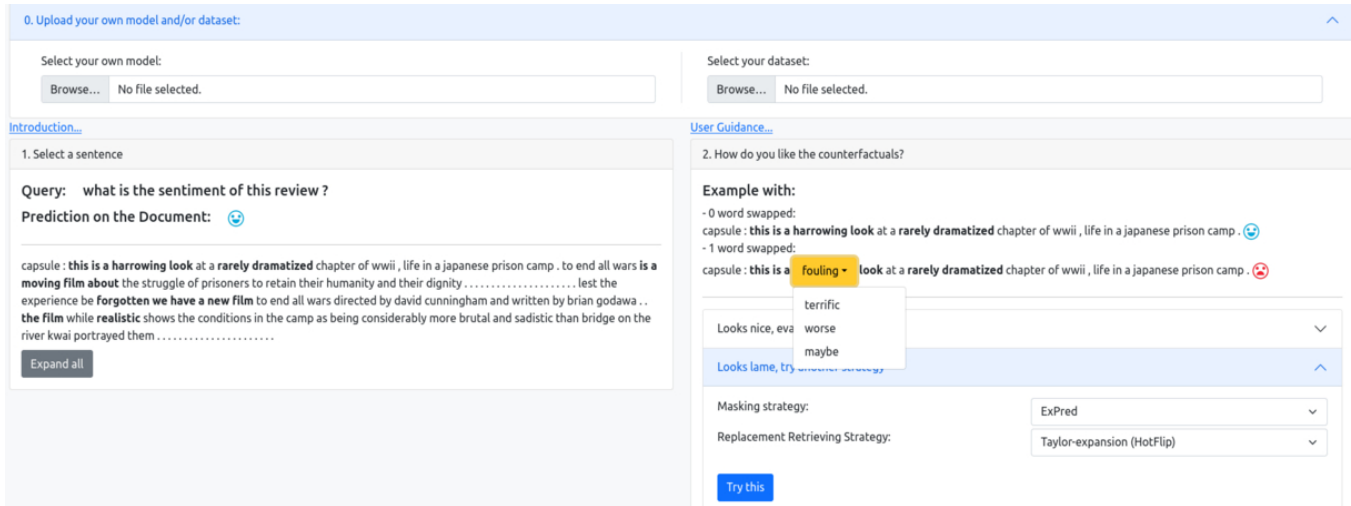


Figure 3: The screenshot of the SparCAssist user interface

each select position we choose the token that maximizes the dot-product of the new token’s embedding with the partial derivation, i.e. $v_{replace} = \operatorname{argmax}_{v \in V} \frac{\partial}{\partial x_i} l(\hat{y}, \tilde{y} \neq \hat{y}) e^v$.

Masked Language Modelling-based approach: For this method, we first construct the MLM model’s input with the prompt template as suggested in [13, 15, 26] as follows: “<masked sequence>” the sentiment of this review is “<alternative label>”. For example, “this movie is [mask]” the sentiment of this review is “negative”.

Unlike in the hotflip, where the positions are attributed with the product of gradient and the embedding, there is no such close-form position attribution in the MLM-based method. Therefore we trivially attribute all tokens with the same score. For the token replacement, we mask each rationale word and apply the prompt template, then predict the masked word using a pre-trained blank-filling model based on RoBERTa [16] (with pre-trained weights from the Huggingface framework²). For each word position, we keep the token with the highest score, then we select the final counterfactual which minimizes the loss with respect to the alternative label.

A screenshot of our system is illustrated in Fig. 3. The top part shows widgets to upload the model and the corresponding dataset. The lower part is split into two columns:

1. The left part (“1. Select a sentence”) presents a document selected from the dataset, where rationale tokens are highlighted with boldface font. By default, all sentences containing at least one rationale token are displayed. The users can expand the full review by clicking the “Expand all” button. Even after restricting the sentences using the rationale masks, the length of the whole review document is still large, users can therefore select a single sentence at a time to inspect its counterfactuals.

²<https://huggingface.co/roberta-large>

2. In the right hand side, (“2. How do you like the counterfactuals?”) the counterfactuals are displayed following the order of word replacement, with a replaced word highlighted. As introduced in the previous section, the users are asked to evaluate the *plausibility* and *meaningfulness*, together with the *faithfulness* of the prediction, if the user is satisfied with the current counterfactual generation. The users may also choose to use the default ExPred rationale masks or provide a custom mask by selecting the alternative rationale tokens according to their preference.

4 CONCLUSION AND FUTURE WORK

In this paper, we propose SparCAssist, an ML model debugging tool. SparCAssist allows users to construct word-replacement-based counterfactual instances of rationale masked sentences. These counterfactuals are helpful for better understanding of model’s behavior and identifying deployment risks of the underlying model. The SparCAssist prunes the choices of counterfactuals by restricting the range of variable tokens to rationale tokens only. This minimizes the number of counterfactual choices presented to the users. Optionally, the users can also specify custom masks for the variable tokens. We implement HotFlip and Masked Language Model-based approaches to retrieve out-of-distribution (OOD) tokens and enable the users to evaluate generation results. It is also possible to deploy the SparCAssist system and gather counterfactual annotations with evaluations from large crowd-workers through web browsers. Another advantage of our tool is that it is flexible and generic. It allows debugging any ML model with a user chosen dataset. It easily supports replacing the rationale-masking, position attribution and OOD tokens selection components with other methods.

As for the future work, as mentioned in Sect. 2, we plan to conduct the user study to figure out the exact number of counterfactual samples needed to evaluate models. Furthermore, a visualization of all counterfactual generated by current user, nominated with their softmax score, would be helpful to get a better understanding of the model's local behavior. Finally, we plan to collect and release a large-scale human plausible counterfactual dataset using the crowd-sourcing deployment of our tool.

5 ACKNOWLEDGMENT

This work is partially funded by project MIRROR under grant agreement No. 832921 (project MIRROR from the European Commission: Migration-Related Risks caused by misconceptions of Opportunities and Requirement) and project ROXANNE, the European Union's Horizon 2020 research and innovation program under grant agreement No. 833635.

REFERENCES

- [1] Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Yunhan Jia, Joydeep Ghosh, Ruchir Puri, José M. F. Moura, and Peter Eckersley. 2020. Explainable Machine Learning in Deployment. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (FAT* '20). Association for Computing Machinery, New York, NY, USA, 648–657. <https://doi.org/10.1145/3351095.3375624>
- [2] Noam Chomsky. 2009. *Syntactic structures*. De Gruyter Mouton.
- [3] Blair Davis, Robert Anderson, and Jan Walls. 2015. *Rashomon effects: Kurosawa, Rashomon and their legacies*. Routledge.
- [4] Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C. Wallace. 2020. ERASER: A Benchmark to Evaluate Rationalized NLP Models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 4443–4458. <https://doi.org/10.18653/v1/2020.acl-main.408>
- [5] Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. HotFlip: White-Box Adversarial Examples for Text Classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Melbourne, Australia, 31–36. <https://doi.org/10.18653/v1/P18-2006>
- [6] Thorben Funke, Megha Khosla, and Avishek Anand. 2021. Zorro: Valid, Sparse, and Stable Explanations in Graph Neural Networks. *arXiv preprint arXiv:2105.08621* (2021).
- [7] Denis Hilton. 2017. Social attribution and explanation. (2017).
- [8] Helge Holzmann and Avishek Anand. 2016. Tempas: Temporal archive search based on tags. In *Proceedings of the 25th International Conference Companion on World Wide Web*. 207–210.
- [9] Maximilian Idahl, Lijun Lyu, Ujwal Gadiraju, and Avishek Anand. 2021. Towards benchmarking the utility of explanations for model debugging. *arXiv preprint arXiv:2105.04505* (2021).
- [10] Divyansh Kaushik, Eduard Hovy, and Zachary Lipton. 2020. Learning The Difference That Makes A Difference With Counterfactually-Augmented Data. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=SkIgs0NFvr>
- [11] Masanori Koyama and Shoichiro Yamaguchi. 2020. When is invariance useful in an Out-of-Distribution Generalization problem? *arXiv: 2008.01883* (2020).
- [12] Piyawat Lertvittayakumjorn and Francesca Toni. 2021. Explanation-Based Human Debugging of NLP Models: A Survey. *Transactions of the Association for Computational Linguistics* 9 (12 2021), 1508–1528. https://doi.org/10.1162/tacl_a_00440
- [13] Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 3045–3059. <https://doi.org/10.18653/v1/2021.emnlp-main.243>
- [14] Peter Lipton. 1990. Contrastive Explanation. *Royal Institute of Philosophy Supplement* 27 (1990), 247–266. <https://doi.org/10.1017/S1358246100005130>
- [15] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021. GPT Understands, Too. *CoRR* abs/2103.10385 (2021). [arXiv:2103.10385](https://arxiv.org/abs/2103.10385)
- [16] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [17] Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Proceedings of the 31st international conference on neural information processing systems*. 4768–4777.
- [18] Tyler McDonnell, Matthew Lease, Mucahid Kutlu, and Tamer Elsayed. 2016. Why is that relevant? collecting annotator rationales for relevance judgments. In *Proceedings of the AAAI Conf. on Human Computation and Crowdsourcing*, Vol. 4.
- [19] Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence* 267 (2019), 1–38. <https://doi.org/10.1016/j.artint.2018.07.007>
- [20] Vedant Misra. 2019. Black box attacks on transformer language models. In *ICLR 2019 Debugging Machine Learning Models Workshop*.
- [21] Robert Munro Monarch. 2021. *Human-in-the-Loop Machine Learning: Active learning and annotation for human-centered AI*. Simon and Schuster.
- [22] Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. *CoRR* abs/1602.04938 (2016). [arXiv:1602.04938](https://arxiv.org/abs/1602.04938) <http://arxiv.org/abs/1602.04938>
- [23] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [24] Alexis Ross, Ana Marasović, and Matthew Peters. 2021. Explaining NLP Models via Minimal Contrastive Editing (MiCE). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 3840–3852. <https://doi.org/10.18653/v1/2021.findings-acl.336>
- [25] Taylor Shin, Yasaman Razeghi, Robert L. Logan IV au2, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. [arXiv:2010.15980](https://arxiv.org/abs/2010.15980) [cs.CL]
- [26] Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 4222–4235. <https://doi.org/10.18653/v1/2020.emnlp-main.346>
- [27] Jaspreet Singh and Avishek Anand. 2018. Posthoc Interpretability of Learning to Rank Models using Secondary Training Data. *arXiv preprint arXiv:1806.11330* (2018).
- [28] Jaspreet Singh and Avishek Anand. 2019. EXS: Explainable Search Using Local Model Agnostic Interpretability. In *ACM WSDM*.
- [29] Jaspreet Singh and Avishek Anand. 2020. Model agnostic interpretability of rankers via intent modelling. In *Conf. on Fairness, Accountability, and Transparency*.
- [30] Jaspreet Singh, Johannes Hoffart, and Avishek Anand. 2016. Discovering entities with just a little help from you. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*. 1331–1340.
- [31] Jaspreet Singh, Megha Khosla, Zhenye Wang, and Avishek Anand. 2021. Extracting per Query Valid Explanations for Blackbox Learning-to-Rank Models. In *ICTIR '21: The 2021 ACM SIGIR International Conference on the Theory of Information Retrieval, Virtual Event, Canada, July 11, 2021*, Faegheh Hasibi, Yi Fang, and Akiko Aizawa (Eds.). ACM, 203–210. <https://doi.org/10.1145/3471158.3472241>
- [32] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic Attribution for Deep Networks. *CoRR* abs/1703.01365 (2017). [arXiv:1703.01365](https://arxiv.org/abs/1703.01365) <http://arxiv.org/abs/1703.01365>
- [33] Eric Wallace, Mitchell Stern, and Dawn Song. 2020. Imitation Attacks and Defenses for Black-box Machine Translation Systems. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 5531–5546. <https://doi.org/10.18653/v1/2020.emnlp-main.446>
- [34] Tongshuang Wu, Marco Túlio Ribeiro, Jeffrey Heer, and Daniel S. Weld. 2021. Polyjuice: Automated, General-purpose Counterfactual Generation. *CoRR* abs/2101.00288 (2021). [arXiv:2101.00288](https://arxiv.org/abs/2101.00288) <https://arxiv.org/abs/2101.00288>
- [35] Zijian Zhang, Koustav Rudra, and Avishek Anand. 2021. Explain and Predict, and Then Predict Again. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining (Virtual Event, Israel) (WSDM '21)*. Association for Computing Machinery, New York, NY, USA, 418–426. <https://doi.org/10.1145/3437963.3441758>
- [36] Zijian Zhang, Koustav Rudra, and Avishek Anand. 2021. FaxPlainAC: A Fact-Checking Tool Based on EXPLAINable Models with HumAn Correction in the Loop. In *CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*, Gianluca Demartini, Guido Zuccon, J. Shane Culpepper, Zi Huang, and Hanghang Tong (Eds.). ACM, 4823–4827. <https://doi.org/10.1145/3459637.3481985>
- [37] Yiming Zheng, Serena Booth, Julie Shah, and Yilun Zhou. 2021. The Irrationality of Neural Rationale Models. *arXiv preprint arXiv:2110.07550* (2021).