

Recent Advances in Underwater Basket Weaving Under the Extreme Pressure of the Mariana Trench

André Lauren Benjamin¹, Calvin Cordozar Broadus Jr.^{2,3} (✉), and Antwan
André Patton¹[0000–1111–2222–3333]

¹ Fictional Southern University, Savannah GA 31404, USA
`{a.l.benjamin,a.a.patton}@fsu.fake`

² Fictional West Coast University, Long Beach CA 90840, USA `ccb@fwcu.fake`

³ Secondary European Affiliation, Tiergartenstr. 17, 69121 Heidelberg, Germany
`lncs@springer.com`

Abstract. This document provides a basic paper template and submission guidelines. Abstracts must be a single paragraph, ideally between 4–6 sentences long. Gross violations will trigger corrections at the camera-ready phase.

Keywords: First keyword · Second keyword · Another keyword.

1 Related Literature

1.1 Background on Counterfactual Explanations

[11, 5, 2]

1.2 Learning Representations

For example, joint-energy models

1.3 Generalization and Robustness

[9] generate counterfactual images for MNIST and ImageNet through independent mechanisms (IM): each IM learns class-conditional input distributions over a specific lower-dimensional, semantically meaningful factor, such as *texture*, *shape* and *background*. They demonstrate that using these generated counterfactuals during classifier training improves model robustness. Similarly, [1] argue that counterfactuals represent potentially useful training data in machine learning, especially in supervised settings where inputs may be reasonably mapped to multiple outputs. They, too, demonstrate the augmenting the training data of image classifiers can improve generalization.

[10] propose an approach using counterfactuals in training that does not rely on data augmentation: they argue that counterfactual pairs typically already

exist in training datasets. Specifically, their approach relies on, firstly, identifying similar input samples with different annotations and, secondly, ensuring that the gradient of the classifier aligns with the vector between pairs of counterfactual inputs using the cosine distance as a loss function (referred to as *gradient supervision*)

This might be useful for our task as well...

In the natural language processing (NLP) domain, counterfactuals have similarly been used to improve models through data augmentation: [12], propose POLYJUICE, a general-purpose counterfactual generator for language models. They demonstrate empirically that augmenting training data through POLYJUICE counterfactuals improves robustness in a number of NLP tasks.

1.4 Link to Adversarial Training

[3] propose two definitional differences between Adversarial Examples (AE) and Counterfactual Explanations (CE): firstly, and more importantly according to the authors, the term AE implies missclassification, which is not the case for CE;

this might be a useful notion for use to distinguish between adversarials and explanations during training

secondly, they argue that closeness plays a more critical role in the context of CE but confess that even counterfactuals that are not close might be relevant explanations. [7] show that CE and AE are equivalent under certain conditions and derive upper bounds on the distances between them.

1.5 Closely Related

[4] are the first to propose end-to-end training pipeline that includes counterfactual explanations as part of the training prodeduce. In particular, they propose a specific network architecture that includes a predictor and CE generator network,

akin a GAN?

where the parameters of the CE generator network are learnable. Counterfactuals are generated during each training iteration and fed back to the predictor network (**here we are aligned**). In contrast, we impose no restrictions on the neural network architecture at all.

to ensure the one-hot encoding of categorical features is maintained, they simple use softmax (might be interesting for CE.jl)

Interestingly, the authors find that their approach is sensitive to the choice of the loss function: only MSE seems to lead to good performance. They also demonstrate theoretically, that the objective function is difficult to optimize due to divergent gradients and suffers from poor adversarial robustness.

because partial gradients with respect to the classification loss component and the counterfactual validity component point in opposite directions

To mitigate these issues, the authors use block-wise gradient descent: they first update with respect to classification loss and then use a second update with respect to the other loss components

this might be useful for our task as well

[8] propose a way to train models that are guaranteed to provide recourse for individuals with high probability. The approach builds on adversarial training,

here we are aligned

where in this context adversarial examples are actively encouraged to exist, but only target attacks with respect to the positive class. The proposed method allows for imposing a set of actionable recourse ex-ante: for example, users can impose mutability constraints for features.

here we are aligned

To solve their objective function more efficiently, they use a first-order Taylor approximation to approximate the recourse loss component (might be applicable in our case

[6] introduce Counterfactual Adversarial Training (CAT) with intention of improving generalization and robustness of language models. Specifically, they propose to proceed as follows: firstly, identify training samples that are subject to high predictive uncertainty (entropy); secondly, generate counterfactual explanations for those samples; and, finally, finetune the model on the augmented dataset that includes the generated counterfactuals.

2 Introduction

2.1 Main Contributions

Please note that the first paragraph of a section or subsection is not indented. The first paragraph that follows a table, figure, equation etc. does not need an indent, either.

Subsequent paragraphs, however, are indented.

Table 1. Table captions should be placed above the tables.

Heading level	Example	Font size and style
Title (centered)	Lecture Notes	14 point, bold
1st-level heading	1 Introduction	12 point, bold
2nd-level heading	2.1 Printing Area	10 point, bold
3rd-level heading	Run-in Heading in Bold. Text follows	10 point, bold
4th-level heading	<i>Lowest Level Heading.</i> Text follows	10 point, italic

3 Related Work

3.1 Basket Weaving

Underwater Basket Weaving Only two levels of headings should be numbered. Lower level headings remain unnumbered; they are formatted as run-in headings.

Underwater Basket Weaving Under Difficult Circumstances The contribution should contain no more than four levels of headings. Table 1 gives a summary of all heading levels.

4 Recent Advances from the Mariana Trench

Displayed equations are centered and set on a separate line.

$$x + y = z \tag{1}$$

Please try to avoid rasterized images for line-art diagrams and schemas. Whenever possible, use vector graphics instead (see Fig. 1).

Theorem 1. *This is a sample theorem. The run-in heading is set in bold, while the following text appears in italics. Definitions, lemmas, propositions, and corollaries are styled the same way.*

Proof. Proofs, examples, and remarks have the initial word in italics, while the following text appears in normal font.

For citations of references, we prefer the use of square brackets and consecutive numbers. Citations using labels or the author/year convention are also acceptable.

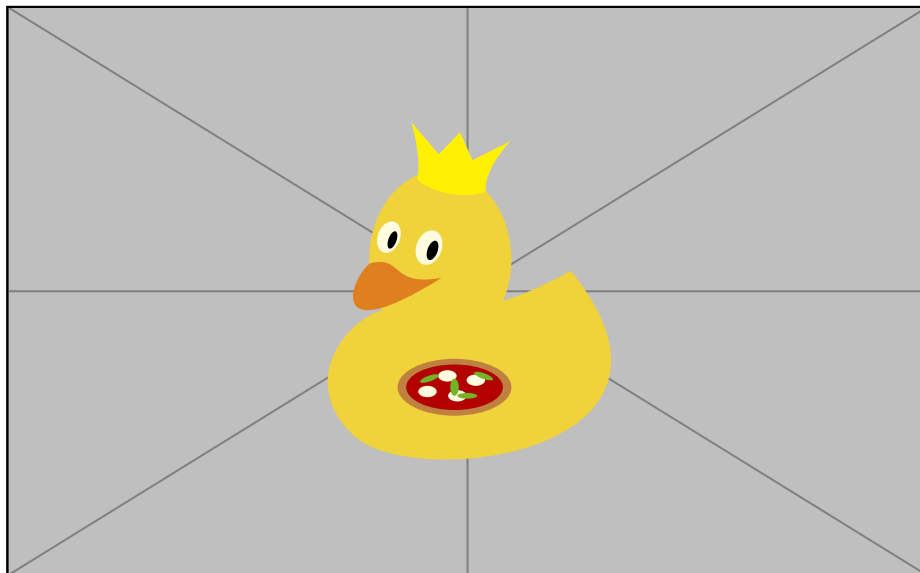


Fig. 1. A figure caption is always placed below the illustration. Please note that short captions are centered, while long ones are justified by the macro package automatically.

5 Experiments

5.1 Experimental Setup

5.2 Experimental Results

6 Discussion

7 Conclusion

Of course, authors have complete freedom on how they choose to structure their paper. Section headers from Introduction up to and including Conclusions are completely up to the discretion of the authors; use whichever structure you see fit. Title, Abstract, the credits environment, and References, however, are mandatory.

Acknowledgments. A bold run-in heading in small font size at the end of the paper is used for general acknowledgments, for example: This study was funded by X (grant number Y).

Disclosure of Interests. It is now necessary to declare any competing interests or to specifically state that the authors have no competing interests. Please place the statement with a bold run-in heading in small font size beneath the (optional) acknowledgments, for example: The authors have no competing interests to declare

that are relevant to the content of this article. Or: Author A has received research grants from Company W. Author B has received a speaker honorarium from Company X and owns stock in Company Y. Author C is a member of committee Z.

Bibliography

- [1] Abbasnejad, E., Teney, D., Parvaneh, A., Shi, J., van den Hengel, A.: Counterfactual vision and language learning. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10041–10051 (2020). <https://doi.org/10.1109/CVPR42600.2020.01006>
- [2] Altmeyer, P., Farmanbar, M., van Deursen, A., Liem, C.C.: Faithful model explanations through energy-constrained conformal counterfactuals. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 38, pp. 10829–10837 (2024)
- [3] Freiesleben, T.: The intriguing relation between counterfactual explanations and adversarial examples. *Minds and Machines* **32**(1), 77–109 (2022)
- [4] Guo, H., Nguyen, T.H., Yadav, A.: Counternet: End-to-end training of prediction aware counterfactual explanations. In: Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. p. 577–589. KDD '23, Association for Computing Machinery, New York, NY, USA (2023). <https://doi.org/10.1145/3580305.3599290>, <https://doi.org/10.1145/3580305.3599290>
- [5] Joshi, S., Koyejo, O., Vijitbenjaronk, W., Kim, B., Ghosh, J.: Towards realistic individual recourse and actionable explanations in black-box decision making systems (2019)
- [6] Luu, H.L., Inoue, N.: Counterfactual adversarial training for improving robustness of pre-trained language models. In: Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation. pp. 881–888 (2023)
- [7] Pawelczyk, M., Agarwal, C., Joshi, S., Upadhyay, S., Lakkaraju, H.: Exploring counterfactual explanations through the lens of adversarial examples: A theoretical and empirical analysis. In: Camps-Valls, G., Ruiz, F.J.R., Valera, I. (eds.) Proceedings of The 25th International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research, vol. 151, pp. 4574–4594. PMLR (28–30 Mar 2022), <https://proceedings.mlr.press/v151/pawelczyk22a.html>
- [8] Ross, A., Lakkaraju, H., Bastani, O.: Learning models for actionable recourse. In: Proceedings of the 35th International Conference on Neural Information Processing Systems. NIPS '21, Curran Associates Inc., Red Hook, NY, USA (2024)
- [9] Sauer, A., Geiger, A.: Counterfactual generative networks (2021), <https://arxiv.org/abs/2101.06046>
- [10] Teney, D., Abbasnejad, E., van den Hengel, A.: Learning what makes a difference from counterfactual examples and gradient supervision. In: Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16. pp. 580–599. Springer (2020)

- [11] Wachter, S., Mittelstadt, B., Russell, C.: Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. JL & Tech.* **31**, 841 (2017). <https://doi.org/10.2139/ssrn.3063289>
- [12] Wu, T., Ribeiro, M.T., Heer, J., Weld, D.: Polyjuice: Generating counterfactuals for explaining, evaluating, and improving models. In: Zong, C., Xia, F., Li, W., Navigli, R. (eds.) *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. pp. 6707–6723. Association for Computational Linguistics, Online (Aug 2021). <https://doi.org/10.18653/v1/2021.acl-long.523>, <https://aclanthology.org/2021.acl-long.523>