# Submission and Formatting Instructions for International Conference on Machine Learning (ICML 2024)

**Firstname1 Lastname1** [*][1]  **Firstname2 Lastname2** [*][1][2]  **Firstname3 Lastname3** [2]  **Firstname4 Lastname4** [3]
**Firstname5 Lastname5** [1]  **Firstname6 Lastname6** [3][1][2]  **Firstname7 Lastname7** [2]  **Firstname8 Lastname8** [3]
**Firstname8 Lastname8** [1][2]

## 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.

## 1. Related Literature

Guo et al. (2023) 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. NB: 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 (because partial gradients with respect to the classification loss component and the counterfactual validity component point in opposite directions) and suffers from poor adversarial robustness. 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).

Ross et al. (2024) 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). NB: 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).

### 1.1. Data Augmentation for Generalization

Sauer & Geiger (2021) 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*. The demonstrate that using these generated counterfactuals during classifier training improves model robustness. Similarly, Abbasnejad et al. (2020) 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. Teney et al. (2020) 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).

## References

Abbasnejad, E., Teney, D., Parvaneh, A., Shi, J., and van den Hengel, A. Counterfactual vision and language learning.

---

[*]Equal contribution [1]Department of XXX, University of YYY, Location, Country [2]Company Name, Location, Country [3]School of ZZZ, Institute of WWW, Location, Country. Correspondence to: Firstname1 Lastname1 <first1.last1@xxx.edu>, Firstname2 Lastname2 <first2.last2@www.uk>.

In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10041–10051, 2020. doi: 10.1109/CVPR42600.2020.01006.

Guo, H., Nguyen, T. H., and 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*, KDD '23, pp. 577–589, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599290. URL https://doi.org/10.1145/3580305.3599290.

Ross, A., Lakkaraju, H., and Bastani, O. Learning models for actionable recourse. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, NIPS '21, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9781713845393.

Sauer, A. and Geiger, A. Counterfactual generative networks, 2021. URL https://arxiv.org/abs/2101.06046.

Teney, D., Abbasnedjad, E., and 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*, pp. 580–599, Berlin, Heidelberg, 2020. Springer-Verlag. ISBN 978-3-030-58606-5. doi: 10.1007/978-3-030-58607-2_34. URL https://doi.org/10.1007/978-3-030-58607-2_34.

## A. You *can* have an appendix here.

You can have as much text here as you want. The main body must be at most 8 pages long. For the final version, one more page can be added. If you want, you can use an appendix like this one.

The \onecolumn command above can be kept in place if you prefer a one-column appendix, or can be removed if you prefer a two-column appendix. Apart from this possible change, the style (font size, spacing, margins, page numbering, etc.) should be kept the same as the main body.