
On Generating Plausible Counterfactual and Semi-Factual Explanations for Deep Learning (Supplementary Material)

Anonymous Author(s)

Affiliation

Address

email

1 S1 Model details

2 Here we detail all model architectures, their accuracy when appropriate, and the data pre-processing
3 steps. The Adam optimizer was used throughout. All URLs were verified to be available on the 3rd
4 of June 2020. For the reader's convenience, all code and pre-trained models can be found at our
5 repository https://github.com/after_anon_review.

6 S1.1 MNIST CNN

7 The CNN architecture used to train a single model, which in turn was used for both Expt. 1 and
8 Expt. 2, is detailed in Table 1. This model achieved an accuracy of 99.59% on the 10,000 test
9 images for MNIST. The model was trained for 10 epochs at learning rate (lr) $1r=1e-1$, 10 epochs
10 at $1r=1e-2$, and 10 more epochs at $1r=1e-3$. The data was normalized between -1 and 1, and the
11 batch size was 8.

12 S1.2 CIFAR-10 CNN

13 The CNN used in Expt. 1 for CIFAR-10 was a modification of the popular ResNet18 architecture. The
14 code was found at [https://github.com/kuangliu/pytorch-cifar/blob/master/models/](https://github.com/kuangliu/pytorch-cifar/blob/master/models/resnet.py)
15 [resnet.py](https://github.com/kuangliu/pytorch-cifar/blob/master/models/resnet.py). The fourth block of the architecture was removed from our model. In addition, dropout
16 (with $p=0.1$) was added after every convolutional layer to enable Monte Carlo Dropout, echoing [3].
17 The accuracy of the network on the 10,000 test images was 92.2%. The model was trained for
18 30 epochs at $1r=1e-1$, 30 epochs at $1r=1e-2$, and 30 more epochs at $1r=1e-3$. The data was
19 normalized between -1 and 1, and the batch size was 32. Random crops (padding=4), and random
20 horizontal flips were also used.

21 S1.3 Autoencoders for IM1 and IM2

22 To train autoencoders (AEs) for these evaluation metrics (IM1 and IM2) to use on CIFAR-10,
23 the architecture found at [https://github.com/jellycsc/PyTorch-CIFAR-10-autoencoder/](https://github.com/jellycsc/PyTorch-CIFAR-10-autoencoder/blob/master/main.py)
24 [blob/master/main.py](https://github.com/jellycsc/PyTorch-CIFAR-10-autoencoder/blob/master/main.py) was used, alongside all its hyperparameter and training choices. To train
25 the AEs for MNIST, the design in Table 2 was used; they were trained for 50 epochs at $1r=1e-2$, and
26 batch size 16.

27 S1.4 Autoencoder for CEM and Proto-CF

28 To train the AEs necessary for these techniques, the examples at the documentation <https://docs.seldon.io/projects/alibi/en/stable/> were used.

30 S1.5 GANs

31 The GANs used for MNIST and CIFAR-10 were found pre-trained at <https://github.com/csinva/gan-vae-pretrained-pytorch>, the reader is referred to this repository for the architecture and training details.

Table 1: The CNN architecture used to train a model on MNIST, this trained model was then used in both experiments.

MNIST CNN	
Layer	Layer Parameters
Conv2d	8 filters, (5x5), (1x1), padding=2
Dropout2d	p=0.1
BatchNorm2d	
ReLU	
Conv2d	8 filters, (5x5), (1x1), padding=2
Dropout2d	p=0.1
BatchNorm2d	
ReLU	
Conv2d	16 filters, (5x5), (2x2), padding=2
Dropout2d	p=0.1
BatchNorm2d	
ReLU	
Conv2d	32 filters, (5x5), (1x1), padding=2
Dropout2d	p=0.1
BatchNorm2d	
ReLU	
Conv2d	64 filters, (5x5), (2x2), padding=2
Dropout2d	p=0.2
BatchNorm2d	
ReLU	
Conv2d	128 filters, (3x3), (1x1), padding=1
GAP	
Linear	128, 10
SoftMax	

34 S2 Method hyperparameters

35 Here the details of all hyperparameter choices for the different methods are detailed for both experi-
 36 ments. Namely, the first experiment involving counterfactual explanations, and the second experiment
 37 involving semi-factual explanations.

38 S2.1 Counterfactual experiment

39 In this experiment five different methods (including our own) were compared.

40 **PIECE** The number of epochs for each explanation was 300 in both datasets. The learning rate
 41 used was $1r=1e-2$ in both datasets. The alpha threshold α was 0.05 in all tests. The Adam optimizer
 42 was used.

43 **Min-Edit** On MNIST $1r=1e-3$ was used for incorrect and close-correct instances, and $1r=1e-2$
 44 for correct instances. On CIFAR-10 $1r=1e-1$ was used on correct instances, and $1r=1e-2$ on
 45 incorrect instances. The Adam optimizer was used.

Table 2: The autoencoder architecture used to train the models for IM1 and IM2 on MNIST, subsequently used in Expt. 1 and Expt. 2.

MNIST IM1 and IM2 autoencoder	
Layer	Layer Parameters
Encoder	
Conv2d	16 filters, (3x3), (1x1), padding=1
ReLU	
Conv2d	16 filters, (2x2), (2x2), padding=0
ReLU	
Conv2d	16 filters, (3x3), (1x1), padding=1
ReLU	
Conv2d	16 filters, (2x2), (2x2), padding=0
Decoder	
Upsample	Scale Factor 2 - Bilinear
ReLU	
Conv2d	16 filters, (3x3), (1x1), padding=1
ReLU	
Upsample	Scale Factor 2 - Bilinear
ReLU	
Conv2d	1 filters, (3x3), (1x1), padding=1
Sigmoid	

46 **C-Min-Edit** In MNIST $lr=1e-3$ was utilized for incorrect and close-correct instances, and
 47 $lr=1e-2$ for correct instances. For CIFAR-10 $lr=1e-1$ was chosen on correct instances, and
 48 $lr=1e-2$ on incorrect ones. The lambda parameter λ was set to 0.1 and incrementally increased by
 49 0.1 each epoch in all tests. The distance function $d(\cdot)$ used the L_2 norm, and the Adam optimizer
 50 was used.

51 **CEM** The default parameters suggested in the official documentation at <https://docs.seldon.io/projects/alibi/en/stable/> were used for both MNIST and CIFAR-10. However, `c_init`
 52 was set to 0, and `c_steps` set to 1 on CIFAR-10.

54 **Proto-CF** The default parameters suggested in the official documentation at <https://docs.seldon.io/projects/alibi/en/stable/> were used for both MNIST and CIFAR-10. However,
 55 `c_init` was again set to 0, and `c_steps` set to 1 on CIFAR-10.

57 S2.2 Semi-factual experiment

58 In this experiment three different methods (including our own method PIECE) were compared.

59 **PIECE** The number of epochs for each explanation was 300. The learning rate was $lr=1e-2$. The
 60 alpha threshold α was 0.05, and the Adam optimizer was used.

61 **Min-Edit** The learning rate used was $lr=1e-2$ alongside the Adam optimizer.

62 **C-Min-Edit** The learning rate used was $lr=1e-2$. The lambda parameter λ was set to 0.1 and
 63 incrementally increased by 0.1 each epoch. The distance function $d(\cdot)$ used the L_2 norm, and lastly
 64 the Adam optimizer was used.

65 S3 Machine specifications

66 Both experiments used the same machine to generate all explanations. The machine was a: MacBook
 67 Pro; Processor 2.9 GHz Intel Core i5; Memory 16 GB 2133 MHz LPDDR3; 500GB storage. (n.b., no
 68 GPU was used, only a CPU).

69 However, to train the CNNs and AE models, Google Colab <https://colab.research.google.com/> was utilized alongside its built-in GPU capability.

71 S4 Additional plots

72 To quote Section 3 of the main paper:

73 ‘Interestingly, for all results on MNIST, a plot of the NN-Dist measure against
74 the MC-Mean/MC-STD scores show a significant linear relationship $r = -0.8/0.82$.
75 So, the more a generated counterfactual is grounded in the training data, the more
76 likely it is to be plausible...’

77 Here, the plots for these correlations are shown, which demonstrate that there is a strong correlation
78 between plausibility (measured with MC Dropout), and the distance of an explanation’s latent
79 representation from the training data, as many have previously eluded to (e.g., see [4, 2, 9, 10]).
80 This discovery lends much credibility to NN-Dist as a valid evaluation metric (when generating
81 explanations with GANs, as is becoming popular [8, 5, 1, 6, 7]) to measure the plausibility of
82 generated counterfactual explanations.

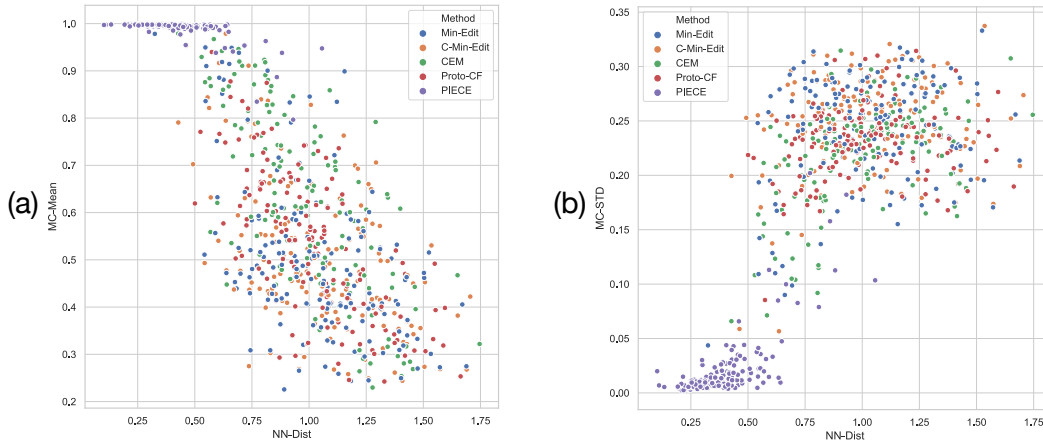


Figure 1: The latent representation x of every explanation image for every method in MNIST is shown; specifically, its NN-Dist plotted against its MC-Mean and MC-STD scores. The results show strong linear correlations between how close an explanation (n.b., explanations generated with a GAN) is to a training instance, and how plausible it is.

83 References

- 84 [1] A. Feghahati, C. R. Shelton, M. J. Pazzani, and K. Tang. Cdeepex: Contrastive deep explanations.
85 2018.
- 86 [2] A.-H. Karimi, G. Barthe, B. Belle, and I. Valera. Model-agnostic counterfactual explanations
87 for consequential decisions. *arXiv preprint arXiv:1905.11190*, 2019.
- 88 [3] A. Kendall and Y. Gal. What uncertainties do we need in bayesian deep learning for computer
89 vision? In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and
90 R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5574–5584.
91 Curran Associates, Inc., 2017.
- 92 [4] T. Laugel, M.-J. Lesot, C. Marsala, X. Renard, and M. Detryniecki. The dangers of post-hoc
93 interpretability: Unjustified counterfactual explanations. *arXiv preprint arXiv:1907.09294*,
94 2019.
- 95 [5] S. Liu, B. Kailkhura, D. Loveland, and H. Yong. Generative counterfactual introspection
96 forexplainable deep learning. Technical report, Lawrence Livermore National Lab.(LLNL),
97 Livermore, CA (United States), 2019.

- 98 [6] P. Samangouei, A. Saeedi, L. Nakagawa, and N. Silberman. Explaingan: Model explanation via
99 decision boundary crossing transformations. In *Proceedings of the European Conference on*
100 *Computer Vision (ECCV)*, pages 666–681, 2018.
- 101 [7] J. Seah, J. Tang, A. Kitchen, and J. Seah. Thinking like a machine—generating visual rationales
102 through latent space optimization. 2018.
- 103 [8] S. Singla, B. Pollack, J. Chen, and K. Batmanghelich. Explanation by progressive exaggeration.
104 In *International Conference on Learning Representations*, 2020.
- 105 [9] A. Van Looveren and J. Klaise. Interpretable counterfactual explanations guided by prototypes.
106 *arXiv preprint arXiv:1907.02584*, 2019.
- 107 [10] S. Wachter, B. Mittelstadt, and C. Russell. Counterfactual explanations without opening the
108 black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.