

---

# COUNTERFACTUAL TRAINING: TEACHING MODELS PLAUSIBLE AND ACTIONABLE EXPLANATIONS

---

A PREPRINT

**Patrick Altmeyer** 

Faculty of Electrical Engineering, Mathematics and Computer Science  
Delft University of Technology

[p.altmeyer@tudelft.nl](mailto:p.altmeyer@tudelft.nl)

**Arie van Deursen**

Faculty of Electrical Engineering, Mathematics and Computer Science  
Delft University of Technology

**Cynthia C. S. Liem**

Faculty of Electrical Engineering, Mathematics and Computer Science  
Delft University of Technology

January 29, 2025

## ABSTRACT

Counterfactual Explanations (CE) have emerged as a popular method to explain the predictions made by opaque machine learning models in a post-hoc fashion. We propose a novel approach that leverages counterfactuals during the training phase of models.

**Keywords** Counterfactual Explanations • Explainable AI

**1 1 Introduction**

**2 2 Related Literature**

**3 2.1 Background on Counterfactual Explanations**

(Wachter, Mittelstadt, and Russell 2017; Joshi et al. 2019; Altmeyer et al. 2024)

**4 2.2 Learning Representations**

For example, joint-energy models

**5 2.3 Generalization and Robustness**

Sauer and 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*. They 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 that augmenting the training data of image classifiers can improve generalization.

19 Tenev, Abbasnejad, and Hengel (2020) propose an approach using counterfactuals in training that does not rely on  
 20 data augmentation: they argue that counterfactual pairs typically already exist in training datasets. Specifically, their  
 21 approach relies on, firstly, identifying similar input samples with different annotations and, secondly, ensuring that the  
 22 gradient of the classifier aligns with the vector between pairs of counterfactual inputs using the cosine distance as a loss  
 23 function (referred to as *gradient supervision*) (*this might be useful for our task as well*). In the natural language pro-  
 24 cessing (NLP) domain, counterfactuals have similarly been used to improve models through data augmentation: Wu et  
 25 al. (2021), propose POLYJUICE, a general-purpose counterfactual generator for language models. They demonstrate  
 26 empirically that augmenting training data through POLYJUICE counterfactuals improves robustness in a number of  
 27 NLP tasks.

## 28 **2.4 Link to Adversarial Training**

29 Freiesleben (2022) propose two definitional differences between Adversarial Examples (AE) and Counterfactual Ex-  
 30 planations (CE): firstly, and more importantly according to the authors, the term AE implies missclassification, which  
 31 is not the case for CE (*this might be a useful notion for use to distinguish between adversarials and explanations*  
 32 *during training*); secondly, they argue that closeness plays a more critical role in the context of CE but confess that  
 33 even counterfactuals that are not close might be relevant explanations. Pawelczyk et al. (2022) show that CE and AE  
 34 are equivalent under certain conditions and derive upper bounds on the distances between them.

## 35 **2.5 Closely Related**

36 Guo, Nguyen, and Yadav (2023) are the first to propose end-to-end training pipeline that includes counterfactual ex-  
 37 planations as part of the training procedure. In particular, they propose a specific network architecture that includes  
 38 a predictor and CE generator network (*akin a GAN?*), where the parameters of the CE generator network are learn-  
 39 able. Counterfactuals are generated during each training iteration and fed back to the predictor network (*here we are*  
 40 *aligned*). In contrast, we impose no restrictions on the neural network architecture at all. (*to ensure the one-hot en-*  
 41 *coding of categorical features is maintained, they simple use softmax (might be interesting for CE.jl)*) Interestingly,  
 42 the authors find that their approach is sensitive to the choice of the loss function: only MSE seems to lead to good  
 43 performance. They also demonstrate theoretically, that the objective function is difficult to optimize due to divergent  
 44 gradients and suffers from poor adversarial robustness. (*because partial gradients with respect to the classification*  
 45 *loss component and the counterfactual validity component point in opposite directions*). To mitigate these issues,  
 46 the authors use block-wise gradient descent: they first update with respect to classification loss and then use a second  
 47 update with respect to the other loss components (*this might be useful for our task as well*). Ross, Lakkaraju, and  
 48 Bastani (2024) propose a way to train models that are guaranteed to provide recourse for individuals with high proba-  
 49 bility. The approach builds on adversarial training (*here we are aligned*), where in this context adversarial examples  
 50 are actively encouraged to exist, but only target attacks with respect to the positive class. The proposed method allows  
 51 for imposing a set of actionable recourse ex-ante: for example, users can impose mutability constraints for features  
 52 (*here we are aligned*). (*To solve their objective function more efficiently, they use a first-order Taylor approximation*  
 53 *to approximate the recourse loss component (might be applicable in our case)*)

54 Luu and Inoue (2023) introduce Counterfactual Adversarial Training (CAT) with intention of improving generalization  
 55 and robustness of language models. Specifically, they propose to proceed as follows: firstly, identify training samples  
 56 that are subject to high predictive uncertainty (entropy); secondly, generate counterfactual explanations for those  
 57 samples; and, finally, finetune the model on the augmented dataset that includes the generated counterfactuals.

## 58 **3 Counterfactual Training**

## 59 **4 Experiments**

### 60 **4.1 Experimental Setup**

### 61 **4.2 Experimental Results**

## 62 **5 Discussion**

## 63 **6 Conclusion**

## 64 **7 Appendix**

### 65 **7.1 Training Details**

### 66 **7.2 Initial Grid Search**

67 For the initial round of experiments we

68 **7.2.1 Generator Parameters**

69 The hyperparameter choices are shown in Parameters 7.1:

70 **Parameters 7.1** (Parameters).

- 71 • **generator\_params**:
- 72   – **lambda\_cost**: 0.0, 0.001, 0.1  
 73   – **lambda\_energy**: 0.01, 0.05, 0.1, 0.5, 1.0, 5.0, 10.0, 15.0  
 74   – **lr**: 1.0  
 75   – **maxiter**: 20, 50, 100  
 76   – **opt**: sgd
- 77 • **generator\_type**: ecco, generic, omni, revise  
 78 • **training\_params**:
- 79   – **objective**: full, vanilla

80 **7.2.1.1 Linearly Separable**

- 81 • **Energy Penalty** (Table 1): *ECCo* generally does yield better results than *Vanilla* for higher choices of the  
 82 energy penalty (10,15) during training. *Generic* performs poorly across the board. *Omni* seems to have an  
 83 anchoring effect, in that it never performs terribly but also never as good as the best *ECCo* results. *REVISE*  
 84 performs poorly across the board.
- 85 • **Cost (distance penalty)**: Results for all generators (except *Omni*) are quite bad, which can likely be attributed  
 86 to extremely bad results for some choices of the **Energy Penalty** (results here are averaged). For *ECCo* and  
 87 *Generic*, higher cost values generally lead to worse results.
- 88 • **Maximum Iterations**: No clear patterns recognizable, so it seems that smaller choices are ok.
- 89 • **Validity**: *ECCo* almost always valid except for very low values during training and high values at evaluation  
 90 time. *Generic* often has poor validity.
- 91 • **Accuracy**: Seems largely unaffected.

Table 1: Results for Linearly Separable data by energy penalty.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
full	0.01	<i>ECCo</i>	$-9.91 \cdot 10^{11}$	$2.25 \cdot 10^{12}$
full	0.01	<i>Generic</i>	$-5.71 \cdot 10^{17}$	$1.3 \cdot 10^{18}$
<b>full</b>	<b>0.01</b>	<b>Omniscient</b>	<b>-2.54</b>	<b>0.116</b>
full	0.01	<i>REVISE</i>	-15.6	13.2
vanilla	0.01	<i>ECCo</i>	-4.28	3.52
vanilla	0.01	<i>Generic</i>	-4.45	3.47
vanilla	0.01	<i>Omniscient</i>	-5.12	4.46
vanilla	0.01	<i>REVISE</i>	-4.91	4.24
full	0.05	<i>ECCo</i>	$-5.63 \cdot 10^5$	$1.28 \cdot 10^6$
full	0.05	<i>Generic</i>	$-8.35 \cdot 10^{17}$	$1.9 \cdot 10^{18}$
<b>full</b>	<b>0.05</b>	<b>Omniscient</b>	<b>-2.53</b>	<b>0.114</b>
full	0.05	<i>REVISE</i>	-15	12.6
vanilla	0.05	<i>ECCo</i>	-4.4	3.66
vanilla	0.05	<i>Generic</i>	-4.38	3.48
vanilla	0.05	<i>Omniscient</i>	-5.25	4.62
vanilla	0.05	<i>REVISE</i>	-4.94	4.22
full	0.1	<i>ECCo</i>	$-6.74 \cdot 10^5$	$1.53 \cdot 10^6$
full	0.1	<i>Generic</i>	$-1.72 \cdot 10^{11}$	$3.9 \cdot 10^{11}$
<b>full</b>	<b>0.1</b>	<b>Omniscient</b>	<b>-2.56</b>	<b>0.124</b>
full	0.1	<i>REVISE</i>	-15.6	13.2
vanilla	0.1	<i>ECCo</i>	-4.28	3.52
vanilla	0.1	<i>Generic</i>	-4.45	3.48
vanilla	0.1	<i>Omniscient</i>	-5.12	4.46
vanilla	0.1	<i>REVISE</i>	-4.91	4.25

Continuing table below.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
full	0.5	<i>ECCo</i>	-11.8	9.83
full	0.5	<i>Generic</i>	$-1.06 \cdot 10^{18}$	$2.42 \cdot 10^{18}$
<b>full</b>	<b>0.5</b>	<b>Omniscient</b>	<b>-2.54</b>	<b>0.123</b>
full	0.5	<i>REVISE</i>	-15	12.6
vanilla	0.5	<i>ECCo</i>	-4.4	3.65
vanilla	0.5	<i>Generic</i>	-4.38	3.48
vanilla	0.5	<i>Omniscient</i>	-5.25	4.61
vanilla	0.5	<i>REVISE</i>	-4.95	4.22
full	1	<i>ECCo</i>	-11.5	11.1
full	1	<i>Generic</i>	$-1.71 \cdot 10^{11}$	$3.88 \cdot 10^{11}$
<b>full</b>	<b>1</b>	<b>Omniscient</b>	<b>-2.59</b>	<b>0.117</b>
full	1	<i>REVISE</i>	-15.7	13.3
vanilla	1	<i>ECCo</i>	-4.28	3.51
vanilla	1	<i>Generic</i>	-4.44	3.47
vanilla	1	<i>Omniscient</i>	-5.11	4.46
vanilla	1	<i>REVISE</i>	-4.91	4.25
full	5	<i>ECCo</i>	-3.99	3.12
full	5	<i>Generic</i>	$-4.88 \cdot 10^{17}$	$1.11 \cdot 10^{18}$
<b>full</b>	<b>5</b>	<b>Omniscient</b>	<b>-2.53</b>	<b>0.117</b>
full	5	<i>REVISE</i>	-14.6	12.1
vanilla	5	<i>ECCo</i>	-4.4	3.65
vanilla	5	<i>Generic</i>	-4.38	3.48
vanilla	5	<i>Omniscient</i>	-5.25	4.61
vanilla	5	<i>REVISE</i>	-4.95	4.22
<b>full</b>	<b>10</b>	<b>ECCo</b>	<b>-2.31</b>	<b>0.735</b>
full	10	<i>Generic</i>	$-1.7 \cdot 10^{11}$	$3.86 \cdot 10^{11}$
full	10	<i>Omniscient</i>	-2.53	0.117
full	10	<i>REVISE</i>	-15.5	13
vanilla	10	<i>ECCo</i>	-4.28	3.51
vanilla	10	<i>Generic</i>	-4.44	3.47
vanilla	10	<i>Omniscient</i>	-5.12	4.46
vanilla	10	<i>REVISE</i>	-4.91	4.24
<b>full</b>	<b>15</b>	<b>ECCo</b>	<b>-2.01</b>	<b>0.488</b>
full	15	<i>Generic</i>	$-4.91 \cdot 10^{17}$	$1.12 \cdot 10^{18}$
full	15	<i>Omniscient</i>	-2.53	0.116
full	15	<i>REVISE</i>	-14.4	11.7
vanilla	15	<i>ECCo</i>	-4.4	3.65
vanilla	15	<i>Generic</i>	-4.38	3.48
vanilla	15	<i>Omniscient</i>	-5.25	4.6
vanilla	15	<i>REVISE</i>	-4.95	4.23

### 92 7.2.1.2 Moons

- 93 • **Energy Penalty** (Table 2): *ECCo* consistently yields better results than *Vanilla*, except for very low choices  
94 of the energy penalty during training for which it performs abysmal. *Generic* performs quite badly across  
95 the board for high enough choices of the energy penalty at evaluation time. *Omni* has small positive effect.  
96 *REVISE* performs poorly across the board.
- 97 • **Cost (distance penalty)**: *Generic* generally does better for higher values, while *ECCo* does better for lower  
98 values.
- 99 • **Maximum Iterations**: No clear patterns recognizable, so it seems that smaller choices are ok.
- 100 • **Validity**: *ECCo* generally achieves full validity except for very low choices the energy penalty during training  
101 and high choices at evaluation time. *Generic* performs poorly for high choices of the energy penalty during  
102 evaluation.
- 103 • **Accuracy**: Largely unaffected although *ECCo* suffers a bit for very low choices the energy penalty during  
104 training. *REVISE* suffers a lot in general (around 10 percentage points).

Table 2: Results for Moons data by energy penalty.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
full	0.01	<i>ECCo</i>	$-2.8 \cdot 10^{22}$	$6.39 \cdot 10^{22}$
full	0.01	<i>Generic</i>	$-4.89 \cdot 10^{30}$	$1.11 \cdot 10^{31}$
<b>full</b>	<b>0.01</b>	<b>Omniscient</b>	<b>-4.74</b>	<b>5.08</b>
full	0.01	<i>REVISE</i>	-572	$1.25 \cdot 10^3$
vanilla	0.01	<i>ECCo</i>	-15.5	17.3
vanilla	0.01	<i>Generic</i>	-10.9	11.9
vanilla	0.01	<i>Omniscient</i>	-12.7	14.4
vanilla	0.01	<i>REVISE</i>	-11.2	13
full	0.05	<i>ECCo</i>	$-1.55 \cdot 10^{16}$	$3.52 \cdot 10^{16}$
full	0.05	<i>Generic</i>	$-2.22 \cdot 10^{20}$	$5 \cdot 10^{20}$
<b>full</b>	<b>0.05</b>	<b>Omniscient</b>	<b>-4.41</b>	<b>4.48</b>
full	0.05	<i>REVISE</i>	$-1.04 \cdot 10^3$	$2.3 \cdot 10^3$
vanilla	0.05	<i>ECCo</i>	-15.5	17.2
vanilla	0.05	<i>Generic</i>	-11.7	12.8
vanilla	0.05	<i>Omniscient</i>	-12.4	14.1
vanilla	0.05	<i>REVISE</i>	-11.3	13.1
full	0.1	<i>ECCo</i>	$-3.41 \cdot 10^3$	$7.73 \cdot 10^3$
full	0.1	<i>Generic</i>	$-5.22 \cdot 10^{30}$	$1.19 \cdot 10^{31}$
<b>full</b>	<b>0.1</b>	<b>Omniscient</b>	<b>-4.78</b>	<b>5.12</b>
full	0.1	<i>REVISE</i>	-288	594
vanilla	0.1	<i>ECCo</i>	-15.5	17.2
vanilla	0.1	<i>Generic</i>	-10.9	11.9
vanilla	0.1	<i>Omniscient</i>	-12.7	14.4
vanilla	0.1	<i>REVISE</i>	-11.3	13.1
full	0.5	<i>ECCo</i>	-7.09	7.51
full	0.5	<i>Generic</i>	$-1.11 \cdot 10^{31}$	$2.53 \cdot 10^{31}$
<b>full</b>	<b>0.5</b>	<b>Omniscient</b>	<b>-4.58</b>	<b>4.83</b>
full	0.5	<i>REVISE</i>	$-1.19 \cdot 10^3$	$2.64 \cdot 10^3$
vanilla	0.5	<i>ECCo</i>	-15.5	17.2
vanilla	0.5	<i>Generic</i>	-11.7	12.8
vanilla	0.5	<i>Omniscient</i>	-12.4	14.1
vanilla	0.5	<i>REVISE</i>	-11.3	13.1
full	1	<i>ECCo</i>	-6.06	6.33
full	1	<i>Generic</i>	$-1.58 \cdot 10^{33}$	$3.59 \cdot 10^{33}$
<b>full</b>	<b>1</b>	<b>Omniscient</b>	<b>-4.66</b>	<b>4.89</b>
full	1	<i>REVISE</i>	$-1.16 \cdot 10^3$	$2.59 \cdot 10^3$
vanilla	1	<i>ECCo</i>	-15.5	17.3
vanilla	1	<i>Generic</i>	-10.9	11.9
vanilla	1	<i>Omniscient</i>	-12.7	14.4
vanilla	1	<i>REVISE</i>	-11.3	13.1
<b>full</b>	<b>5</b>	<b>ECCo</b>	<b>-2.57</b>	<b>2.07</b>
full	5	<i>Generic</i>	$-1.17 \cdot 10^{28}$	$2.66 \cdot 10^{28}$
full	5	<i>Omniscient</i>	-4.29	4.31
full	5	<i>REVISE</i>	-530	$1.16 \cdot 10^3$
vanilla	5	<i>ECCo</i>	-15.5	17.2
vanilla	5	<i>Generic</i>	-11.7	12.7
vanilla	5	<i>Omniscient</i>	-12.4	14.1
vanilla	5	<i>REVISE</i>	-11.3	13.1
<b>full</b>	<b>10</b>	<b>ECCo</b>	<b>-1.76</b>	<b>0.974</b>
full	10	<i>Generic</i>	$-1.54 \cdot 10^{33}$	$3.51 \cdot 10^{33}$
full	10	<i>Omniscient</i>	-4.44	4.56
full	10	<i>REVISE</i>	$-1.52 \cdot 10^3$	$3.4 \cdot 10^3$
vanilla	10	<i>ECCo</i>	-15.5	17.3

Continuing table below.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
vanilla	10	<i>Generic</i>	-10.9	11.9
vanilla	10	<i>Omniscient</i>	-12.7	14.4
vanilla	10	<i>REVISE</i>	-11.3	13.1
<b>full</b>	<b>15</b>	<b>ECCo</b>	<b>-1.37</b>	<b>0.365</b>
full	15	<i>Generic</i>	$-5.32 \cdot 10^{28}$	$1.21 \cdot 10^{29}$
full	15	<i>Omniscient</i>	-4.34	4.38
full	15	<i>REVISE</i>	-473	$1.03 \cdot 10^3$
vanilla	15	<i>ECCo</i>	-15.5	17.2
vanilla	15	<i>Generic</i>	-11.7	12.8
vanilla	15	<i>Omniscient</i>	-12.4	14.1
vanilla	15	<i>REVISE</i>	-11.3	13.1

105 **7.2.1.3 Circles**

- 106 • **Energy Penalty** (Table 3): *ECCo* consistently yields better results than *Vanilla*, though primarily for low to  
 107 medium choices of the energy penalty ( $<=5$ ) during training. The same goes for *Generic*, which sometimes  
 108 outperforms *ECCo* (for small energy penalty at evaluation time). *Omni* does alright for lower energy penalty  
 109 at evaluation time, but loses out for higher choices. *REVISE* performs poorly across the board (except very  
 110 low choices at evaluation time).
- 111 • **Cost (distance penalty)**: *ECCo* and *Generic* generally achieve the best results when no cost penalty is used  
 112 during training. Both *Omni* and *REVISE* are largely unaffected.
- 113 • **Maximum Iterations**: *ECCo* consistently yields better results for higher numbers of iterations. *Generic*  
 114 generally does best for a medium number (50). *Omni* is sometimes invalid (???).
- 115 • **Validity**: *ECCo* tends to outperform its *Vanilla* counterpart, though primarily for low to medium choices of  
 116 the energy penalty ( $<=5$ ) during training and evaluation. *Vanilla* typically worse across the board.
- 117 • **Accuracy**: Mostly unaffected, but *REVISE* again consistently some deterioration and *ECCo* deteriorates for  
 118 high choices of energy penalty during training, reflecting other outcomes above.

Table 3: Results for Circles data by energy penalty.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
<b>full</b>	<b>0.01</b>	<b>ECCo</b>	<b>-1.26</b>	<b>0.423</b>
full	0.01	<i>Generic</i>	-1.49	0.71
full	0.01	<i>Omniscient</i>	-5.21	5.25
full	0.01	<i>REVISE</i>	$-2.71 \cdot 10^{26}$	$6.37 \cdot 10^{26}$
vanilla	0.01	<i>ECCo</i>	-9.33	7.34
vanilla	0.01	<i>Generic</i>	-8.89	6.88
vanilla	0.01	<i>Omniscient</i>	-8.67	6.87
vanilla	0.01	<i>REVISE</i>	-8.65	6.8
full	0.05	<i>ECCo</i>	-1.29	0.397
<b>full</b>	<b>0.05</b>	<b>Generic</b>	<b>-1.21</b>	<b>0.356</b>
full	0.05	<i>Omniscient</i>	-5.08	5.09
full	0.05	<i>REVISE</i>	$-5.91 \cdot 10^{27}$	$1.36 \cdot 10^{28}$
vanilla	0.05	<i>ECCo</i>	-9.35	7.32
vanilla	0.05	<i>Generic</i>	-8.85	6.87
vanilla	0.05	<i>Omniscient</i>	-8.7	6.96
vanilla	0.05	<i>REVISE</i>	-8.52	6.76
<b>full</b>	<b>0.1</b>	<b>ECCo</b>	<b>-1.2</b>	<b>0.383</b>
full	0.1	<i>Generic</i>	-1.5	0.735
full	0.1	<i>Omniscient</i>	-5.17	5.23
full	0.1	<i>REVISE</i>	$-3.06 \cdot 10^{26}$	$7.7 \cdot 10^{26}$
vanilla	0.1	<i>ECCo</i>	-9.33	7.32
vanilla	0.1	<i>Generic</i>	-8.88	6.86
vanilla	0.1	<i>Omniscient</i>	-8.69	6.9

Continuing table below.

Objective	$\lambda_{\text{div}}(\text{train})$	Generator	Value	Std
vanilla	0.1	<i>REVISE</i>	-8.68	6.81
<b>full</b>	<b>0.5</b>	<b>ECCo</b>	<b>-1.12</b>	<b>0.217</b>
full	0.5	<i>Generic</i>	-1.21	0.352
full	0.5	<i>Omniscient</i>	-5.09	5.12
full	0.5	<i>REVISE</i>	$-5.97 \cdot 10^{27}$	$1.37 \cdot 10^{28}$
vanilla	0.5	<i>ECCo</i>	-9.35	7.3
vanilla	0.5	<i>Generic</i>	-8.89	6.92
vanilla	0.5	<i>Omniscient</i>	-8.68	6.93
vanilla	0.5	<i>REVISE</i>	-8.53	6.75
<b>full</b>	<b>1</b>	<b>ECCo</b>	<b>-1.1</b>	<b>0.163</b>
full	1	<i>Generic</i>	-1.49	0.726
full	1	<i>Omniscient</i>	-5.16	5.2
full	1	<i>REVISE</i>	$-3.09 \cdot 10^{26}$	$7.22 \cdot 10^{26}$
vanilla	1	<i>ECCo</i>	-9.34	7.36
vanilla	1	<i>Generic</i>	-8.86	6.85
vanilla	1	<i>Omniscient</i>	-8.7	6.9
vanilla	1	<i>REVISE</i>	-8.69	6.85
full	5	<i>ECCo</i>	-1.75	0.154
<b>full</b>	<b>5</b>	<b>Generic</b>	<b>-1.21</b>	<b>0.363</b>
full	5	<i>Omniscient</i>	-5.14	5.16
full	5	<i>REVISE</i>	$-1.1 \cdot 10^{28}$	$2.5 \cdot 10^{28}$
vanilla	5	<i>ECCo</i>	-9.36	7.32
vanilla	5	<i>Generic</i>	-8.88	6.91
vanilla	5	<i>Omniscient</i>	-8.7	6.93
vanilla	5	<i>REVISE</i>	-8.52	6.73
full	10	<i>ECCo</i>	$-1.02 \cdot 10^6$	$2.32 \cdot 10^6$
<b>full</b>	<b>10</b>	<b>Generic</b>	<b>-1.49</b>	<b>0.702</b>
full	10	<i>Omniscient</i>	-5.13	5.16
full	10	<i>REVISE</i>	$-3.74 \cdot 10^{26}$	$9.09 \cdot 10^{26}$
vanilla	10	<i>ECCo</i>	-9.31	7.33
vanilla	10	<i>Generic</i>	-8.87	6.86
vanilla	10	<i>Omniscient</i>	-8.7	6.89
vanilla	10	<i>REVISE</i>	-8.69	6.83
full	15	<i>ECCo</i>	$-3.31 \cdot 10^{13}$	$7.54 \cdot 10^{13}$
<b>full</b>	<b>15</b>	<b>Generic</b>	<b>-1.22</b>	<b>0.37</b>
full	15	<i>Omniscient</i>	-5.2	5.23
full	15	<i>REVISE</i>	$-9.01 \cdot 10^{27}$	$2.06 \cdot 10^{28}$
vanilla	15	<i>ECCo</i>	-9.38	7.34
vanilla	15	<i>Generic</i>	-8.86	6.87
vanilla	15	<i>Omniscient</i>	-8.69	6.96
vanilla	15	<i>REVISE</i>	-8.51	6.73

## 119 References

- 120 Abbasnejad, Ehsan, Damien Teney, Amin Parvaneh, Javen Shi, and Anton van den Hengel. 2020. “Counterfactual  
121 Vision and Language Learning.” In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition  
122 (CVPR)*, 10041–51. <https://doi.org/10.1109/CVPR42600.2020.01006>.
- 123 Altmeyer, Patrick, Mojtaba Farmanbar, Arie van Deursen, and Cynthia CS Liem. 2024. “Faithful Model Explanations  
124 Through Energy-Constrained Conformal Counterfactuals.” In *Proceedings of the AAAI Conference on Artificial  
125 Intelligence*, 38:10829–37. 10.
- 126 Friesleben, Timo. 2022. “The Intriguing Relation Between Counterfactual Explanations and Adversarial Examples.”  
127 *Minds and Machines* 32 (1): 77–109.
- 128 Guo, Hangzhi, Thanh H. Nguyen, and Amulya Yadav. 2023. “CounterNet: End-to-End Training of Prediction Aware  
129 Counterfactual Explanations.” In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery  
130 and Data Mining*, 577–89. KDD ’23. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3580305.3599290>.

- 132 Joshi, Shalmali, Oluwasanmi Koyejo, Warut Vigitbenjaronk, Been Kim, and Joydeep Ghosh. 2019. “Towards Realistic  
 133 Individual Recourse and Actionable Explanations in Black-Box Decision Making Systems.” <https://arxiv.org/abs/1907.09615>.
- 134 Luu, Hoai Linh, and Naoya Inoue. 2023. “Counterfactual Adversarial Training for Improving Robustness of Pre-  
 135 Trained Language Models.” In *Proceedings of the 37th Pacific Asia Conference on Language, Information and*  
 136 *Computation*, 881–88.
- 137 Pawelczyk, Martin, Chirag Agarwal, Shalmali Joshi, Sohini Upadhyay, and Himabindu Lakkaraju. 2022. “Exploring  
 138 Counterfactual Explanations Through the Lens of Adversarial Examples: A Theoretical and Empirical Analysis.”  
 139 In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, edited by Gustau  
 140 Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, 151:4574–94. Proceedings of Machine Learning Research.  
 141 PMLR. <https://proceedings.mlr.press/v151/pawelczyk22a.html>.
- 142 Ross, Alexis, Himabindu Lakkaraju, and Osbert Bastani. 2024. “Learning Models for Actionable Recourse.” In  
 143 *Proceedings of the 35th International Conference on Neural Information Processing Systems*. NIPS ’21. Red  
 144 Hook, NY, USA: Curran Associates Inc.
- 145 Sauer, Axel, and Andreas Geiger. 2021. “Counterfactual Generative Networks.” <https://arxiv.org/abs/2101.06046>.
- 146 Teney, Damien, Ehsan Abbasnedjad, and Anton van den Hengel. 2020. “Learning What Makes a Difference from  
 147 Counterfactual Examples and Gradient Supervision.” In *Computer Vision–ECCV 2020: 16th European Confer-  
 148 ence, Glasgow, UK, August 23–28, 2020, Proceedings, Part x 16*, 580–99. Springer.
- 149 Wachter, Sandra, Brent Mittelstadt, and Chris Russell. 2017. “Counterfactual Explanations Without Opening the Black  
 150 Box: Automated Decisions and the GDPR.” *Harv. JL & Tech.* 31: 841. <https://doi.org/10.2139/ssrn.3063289>.
- 151 Wu, Tongshuang, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2021. “Polyjuice: Generating Counterfactuals  
 152 for Explaining, Evaluating, and Improving Models.” In *Proceedings of the 59th Annual Meeting of the Associa-  
 153 tion for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing  
 154 (Volume 1: Long Papers)*, edited by Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, 6707–23. Online:  
 155 Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.acl-long.523>.
- 156