CounterNet: End-to-End Training of Prediction Aware Counterfactual Explanations

Background

Counterfactual (CF) explanations offer contrast cases for individuals who are *adversely impacted* by algorithmic decisions.



Balancing the **cost-invalidity trade-off** is key to generating high-quality CF explanations.



Post-hoc Explainability. Existing methods are post-hoc (*explain after prediction*) • by assuming underlying models as *black boxes*.

Motivations

1 EU-GDPR enforces the "*Right to Explanation*". • ML developers want to pair explanations with predictions (i.e., the black-box assumption is overly limiting).

- The post-hoc paradigm is *sub-optimal*.
- It fails to properly balance the cost-invalidity trade-off.
- Most of the post-hoc explanation methods run slow.

Research Objective

Can we depart from the dominant post-hoc paradigm in CF explanations by integrating predictive model training and CF explanation generation within an end-to-end pipeline?

Hangzhi Guo¹, Thanh Hong Nguyen², Amulya Yadav¹

¹Pennsylvania State University, ²University of Oregon

CounterNet Architecture



$$\underset{\theta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^{N} \left[\lambda_1 \cdot \underbrace{(y_i - \hat{y}_{x_i})^2}_{\operatorname{Prediction Loss} (\mathcal{L}_1)} + \lambda_2 \cdot \underbrace{(y_i - \hat{y}_{$$

Two issues occur when directly optimizing E.q. 1:

• Issue I: Poor Convergence.

LEMMA I*. $\nabla \mathcal{L}_1 \cdot \nabla \mathcal{L}_2 < 0.$

• Issue II: Adversarial Vulnerability.

LEMMA II*. Minimizing \mathcal{L}_2 w.r.t. θ_f increases the Lipschitz constant of f. *Refer to our paper for details of Lemma I & II.

Block-Wise Coordinate Descent

The Solution to Issue I & II: a block-wise coordinate descent procedure.

For each batch of m data points $\{x^{(i)}, y^{(i)}\}^m$,

 $\boldsymbol{\theta}^{(1)} = \boldsymbol{\theta}^{(0)} - \nabla^{(0)}_{\boldsymbol{\theta}} (\lambda_1 \cdot \mathcal{L}_1) \\ \boldsymbol{\theta}^{(2)}_g = \boldsymbol{\theta}^{(1)}_g - \nabla^{(1)}_{\boldsymbol{\theta}^{(1)}_a} (\lambda_2 \cdot L_2 + \lambda_3 \cdot \mathcal{L}_3)$



Evaluation

lethod	Adult	Credit	HELOC	OULAD
anillaCF	1432.09	1358.26	1340.42	1705.93
iverseCF	4685.39	3898.43	3921.72	5478.17
cotoCF	2348.21	2056.01	1956.71	2823.29
ncertainCF	379.95	60.80	7.91	6.81
-CHVAE	3.28	568.28	2.68	4.79
AE-CF	1.72	1.28	1.48	1.84
ounteRGAN	1.96	1.77	1.59	2.40
CNet	1.39	1.23	1.13	1.81
ounterNet	0.64	0.39	0.44	0.79

/Iodel	Adult	Credit	HELOC	OULAD
Base Model	0.831	0.813	0.717	0.934
CounterNet	0.828	0.819	0.716	0.929

• Post-hoc explainability can be sub-optimal and overly limiting in counterfactual explanations. • CounterNet represents a first step towards developing end-to-end CF explanation systems.

Code: https://github.com/BirkhoffG/counternet



