ECCCos from the Black Box

Faithful Model Explanations through Energy-Constrained Conformal Counterfactuals

Patrick Altmeyer

Mojtaba Farmanbar Cynthia C. S. Liem

Arie van Deursen

Delft University of Technology

2024-05-14

Pick your Poison

All of these counterfactuals are valid explanations for the model's prediction.

Which one would you pick?



Figure 1: Turning a 9 into a 7: Counterfactual explanations for an image classifier produced using *Wachter* (Wachter, Mittelstadt, and Russell 2017), *Schut* (Schut et al. 2021) and *REVISE* (Joshi et al. 2019).

Faithfulness first, plausibility second.

Patrick Altmeyer, Mojtaba Farmanbar, Arie van Deursen, Cynthia C. S. Liem ECCCos from the Black Box

Delft University of Technology

Faithfulness first, plausibility second.

We propose *ECCCo*: a new way to generate faithful model explanations that are as plausible as the underlying model permits.

Faithfulness first, plausibility second.	Reconciling Faithfulness and Plausibility	Results	Questions?
00●0	0000	000	0000

Idea: generate counterfactuals that are consistent with what the model has learned about the data.

- Idea: generate counterfactuals that are consistent with what the model has learned about the data.
- **Method**: constrain the model's energy and predictive uncertainty for the counterfactual.

- Idea: generate counterfactuals that are consistent with what the model has learned about the data.
- **Method**: constrain the model's energy and predictive uncertainty for the counterfactual.
- Result: faithful counterfactuals that are as plausible as the model permits.

- Idea: generate counterfactuals that are consistent with what the model has learned about the data.
- **Method**: constrain the model's energy and predictive uncertainty for the counterfactual.
- Result: faithful counterfactuals that are as plausible as the model permits.
- Benefits: enable us to distinguish trustworthy from unreliable models.

Counterfactual Explanations

$$\min_{\mathbf{Z}' \in \mathcal{Z}^L} \{ \mathsf{yloss}(M_\theta(f(\mathbf{Z}')), \mathbf{y}^+) + \lambda \mathsf{cost}(f(\mathbf{Z}')) \}$$

Counterfactual Explanations (CE)

explain how inputs into a model need to change for it to produce different outputs.



Figure 2: Gradient-based counterfactual search.

Reconciling Faithfulness and Plausibility

Plausibility

Definition (Plausible Counterfactuals)

Let $\mathcal{X}|\mathbf{y}^+ = p(\mathbf{x}|\mathbf{y}^+)$ denote the true conditional distribution of samples in the target class \mathbf{y}^+ . Then for \mathbf{x}' to be considered a plausible counterfactual, we need: $\mathbf{x}' \sim \mathcal{X}|\mathbf{y}^+$.

Why Plausibility?

Plausibility is positively associated with actionability, robustness (Artelt et al. 2021) and causal validity (Mahajan, Tan, and Sharma 2020).

Patrick Altmeyer, Mojtaba Farmanbar, Arie van Deursen, Cynthia C. S. Liem ECCCos from the Black Box



Figure 3: Kernel density estimate (KDE) for the conditional distribution, $p(\mathbf{x}|\mathbf{y}^+)$, based on observed data. Counterfactual path as in Figure 2.

Reconciling Faithfulness and Plausibility

Results

Questions?

Faithfulness

Definition (Faithful Counterfactuals)

Let $\mathcal{X}_{\theta}|\mathbf{y}^{+} = p_{\theta}(\mathbf{x}|\mathbf{y}^{+})$ denote the conditional distribution of \mathbf{x} in the target class \mathbf{y}^{+} , where θ denotes the parameters of model M_{θ} . Then for \mathbf{x}' to be considered a faithful counterfactual, we need: $\mathbf{x}' \sim \mathcal{X}_{\theta}|\mathbf{y}^{+}$.

Trustworthy Models

If the model posterior approximates the true posterior $(p_{\theta}(\mathbf{x}|\mathbf{y}^{+}) \rightarrow p(\mathbf{x}|\mathbf{y}^{+}))$, faithful counterfactuals are also plausible.



Figure 4: KDE for learned conditional distribution, $p_{\theta}(\mathbf{x}|\mathbf{y}^{+})$. Yellow stars indicate conditional samples generated through SGLD for a joint energy model (JEM).

ECCCo

Key Idea

Use the hybrid objective of joint enmodels ergy (JEM) and a model-agnostic penalty for predictive uncertainty: Energy-Constrained (\mathcal{E}_{θ}) Conformal (Ω) Counterfactuals (ECCCo).

ECCCo objective^a:

$$\begin{split} & \min_{\mathbf{Z}' \in \mathcal{Z}^L} \{ L_{\mathsf{clf}}(f(\mathbf{Z}'); M_{\theta}, \mathbf{y}^+) + \lambda_1 \mathsf{cost}(f(\mathbf{Z}')) \\ & + \lambda_2 \mathcal{E}_{\theta}(f(\mathbf{Z}') | \mathbf{y}^+) + \lambda_3 \Omega(C_{\theta}(f(\mathbf{Z}'); \alpha)) \} \end{split}$$



Figure 5: Gradient fields and counterfactual paths for different generators.

^aWe leverage ideas from Grathwohl et al. (2020) and Stutz et al. (2022). See the paper and appendix for a derivation of the objective from first principles.

Results

Visual Evidence



Figure 6: Turning a 9 into a 7. *ECCCo* applied to MLP (a), Ensemble (b), JEM (c), JEM Ensemble (d).

Figure 7: Results for different generators (from 3 to 5).

Faithfulness first, plausibility second.	Reconciling Faithfulness and Plausibility	Results	Questions?
0000	0000	00●	0000

The Numbers

- Large benchmarks on a variety of models and datasets from various domains.
- ECCCo achieves state-of-the-art faithfulness across models and datasets and approaches state-of-the-art plausibility for more trustworthy models.

		C	alifornia Housing			GMSC	
Model	Generator	Unfaithfulness \downarrow	Implausibility \downarrow	Uncertainty \downarrow	Unfaithfulness \downarrow	Implausibility \downarrow	Uncertainty \downarrow
MLP Ensemble	ECCCo ECCCo+ ECCCo (no CP) ECCCo (no EBM) REVISE Schut Wachter	$\begin{array}{c} \textbf{3.69} \pm \textbf{0.08}^{**} \\ \textbf{3.88} \pm \textbf{0.07}^{**} \\ \textbf{3.70} \pm \textbf{0.08}^{**} \\ \textbf{4.03} \pm \textbf{0.07} \\ \textbf{3.96} \pm \textbf{0.07}^{*} \\ \textbf{4.00} \pm \textbf{0.06} \\ \textbf{4.04} \pm \textbf{0.07} \end{array}$	$\begin{array}{c} 1.94 \pm 0.13 \\ 1.20 \pm 0.09 \\ 1.94 \pm 0.13 \\ 1.12 \pm 0.12 \\ \textbf{0.58} \pm \textbf{0.03}^{**} \\ 1.15 \pm 0.12 \\ 1.13 \pm 0.12 \end{array}$	$\begin{array}{c} \textbf{0.09} \pm \textbf{0.01}^{**} \\ 0.15 \pm 0.02 \\ 0.10 \pm 0.01^{**} \\ 0.14 \pm 0.01^{**} \\ 0.17 \pm 0.03 \\ 0.10 \pm 0.01^{**} \\ 0.16 \pm 0.01 \end{array}$	$\begin{array}{c} 3.84 \pm 0.07^{**} \\ \textbf{3.79} \pm 0.05^{**} \\ 3.85 \pm 0.07^{**} \\ 4.08 \pm 0.06 \\ 4.09 \pm 0.07 \\ 4.04 \pm 0.08 \\ 4.10 \pm 0.07 \end{array}$	$\begin{array}{c} 2.13 \pm 0.08 \\ 1.81 \pm 0.05 \\ 2.13 \pm 0.08 \\ 0.97 \pm 0.08 \\ \textbf{0.63 \pm 0.02^{**}} \\ 1.21 \pm 0.08 \\ 0.95 \pm 0.08 \end{array}$	$\begin{array}{c} \textbf{0.23} \pm \textbf{0.01}^{**} \\ 0.30 \pm 0.01^{*} \\ 0.23 \pm 0.01^{**} \\ 0.31 \pm 0.01^{*} \\ 0.33 \pm 0.06 \\ 0.30 \pm 0.01^{*} \\ 0.32 \pm 0.01 \end{array}$
JEM Ensemble	ECCCo ECCCo+ ECCCo (no CP) ECCCo (no EBM) REVISE Schut Wachter	$\begin{array}{c} 1.40 \pm 0.08^{**} \\ \textbf{1.28} \pm \textbf{0.08}^{**} \\ 1.39 \pm 0.08^{**} \\ 1.70 \pm 0.09 \\ 1.39 \pm 0.15^{**} \\ 1.59 \pm 0.10^{*} \\ 1.71 \pm 0.09 \end{array}$	$\begin{array}{c} 0.69 \pm 0.05^{**} \\ 0.60 \pm 0.04^{**} \\ 0.69 \pm 0.05^{**} \\ 0.99 \pm 0.08 \\ \textbf{0.59 \pm 0.04^{**}} \\ 1.10 \pm 0.06 \\ 0.99 \pm 0.08 \end{array}$	$\begin{array}{c} 0.11 \pm 0.00^{**} \\ 0.11 \pm 0.00^{**} \\ 0.11 \pm 0.00^{**} \\ 0.14 \pm 0.00^{*} \\ 0.25 \pm 0.07 \\ \textbf{0.09 \pm 0.00^{**}} \\ 0.14 \pm 0.00 \end{array}$	$\begin{array}{c} 1.20 \pm 0.06 ^{*} \\ 1.01 \pm 0.07 ^{**} \\ 1.21 \pm 0.07 ^{*} \\ 1.31 \pm 0.07 \\ 1.01 \pm 0.07 ^{**} \\ 1.34 \pm 0.07 \\ 1.31 \pm 0.08 \end{array}$	$\begin{array}{c} 0.78 \pm 0.07^{**} \\ 0.70 \pm 0.07^{**} \\ 0.77 \pm 0.07^{**} \\ 0.97 \pm 0.10 \\ \textbf{0.63 \pm 0.04^{**}} \\ 1.21 \pm 0.10 \\ 0.95 \pm 0.10 \end{array}$	$\begin{array}{c} 0.38 \pm 0.01 \\ 0.37 \pm 0.01 \\ 0.39 \pm 0.01 \\ 0.32 \pm 0.01^{**} \\ 0.33 \pm 0.07 \\ \textbf{0.26 \pm 0.01^{**}} \\ 0.33 \pm 0.01 \end{array}$

Table 1: Results for tabular datasets: sample averages +/- one standard deviation across valid counterfactuals. The best outcomes are highlighted in bold. Asterisks indicate that the given value is more than one (*) or two (**) standard deviations away from the baseline (*Wachter*).

Questions?

Patrick Altmeyer, Mojtaba Farmanbar, Arie van Deursen, Cynthia C. S. Liem ECCCos from the Black Box

Delft University of Technology

Faithfulness first, plausibility second.	Reconciling Faithfulness and Plausibility	Results	Questions?
0000	0000	000	0●00



With thanks to my co-authors Mojtaba Farmanbar, Arie van Deursen and Cynthia C. S. Liem.



Faithfulness first, plausibility second.	Reconciling Faithfulness and Plausibility	Results	Questions?
0000	0000	000	○○●○
Code			

The code used to run the analysis for this work is built on top of CounterfactualExplanations.jl.

There is also a corresponding paper, *Explaining Black-Box Models through Counterfactuals*, which has been published in JuliaCon Proceedings.



Figure 8: Trustworthy AI in Julia: github.com/JuliaTrustworthyAI

Faithfulness first, plausibility second. 0000	Reconciling Faithfulness and Plausibility	Results 000	Questions? 000●

References

- Artelt, André, Valerie Vaquet, Riza Velioglu, Fabian Hinder, Johannes Brinkrolf, Malte Schilling, and Barbara Hammer. 2021. "Evaluating Robustness of Counterfactual Explanations." In 2021 IEEE Symposium Series on Computational Intelligence (SSCI), 01–09. IEEE.
- Grathwohl, Will, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. 2020. "Your Classifier Is Secretly an Energy Based Model and You Should Treat It Like One." In International Conference on Learning Representations.
- Joshi, Shalmali, Oluwasanmi Koyejo, Warut Vijitbenjaronk, Been Kim, and Joydeep Ghosh. 2019. "Towards Realistic Individual Recourse and Actionable Explanations in Black-Box Decision Making Systems." https://arxiv.org/abs/1907.09615.

Mahajan, Divyat, Chenhao Tan, and Amit Sharma. 2020. "Preserving Causal Constraints in Counterfactual Explanations for Machine