Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have time
000000000	000	0000	000	0000

Explaining Models or Modelling Explanations Challenging Existing Paradigms in Trustworthy AI

Patrick Altmeyer Arie van Deursen Cynthia C. S. Liem

Delft University of Technology

2025-05-08

Teaching models plausible explanations $_{\rm OOO}$

If we still have time ... 0000

Background

Economist, now PhD CS
How can we make opaque AI more

trustworthy?

Explainable AI, Adversarial ML, Probabilistic ML

✓ Maintainer of Taija (trustworthy AI in Julia)



Figure 1: Scan for slides. Links to www.patalt.org.

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have time
000000000	000	0000	000	0000
Agenda				

What are counterfactual explanations (CE) and algorithmic recourse (AR) and why are they useful?

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have time
000000000	000	0000	000	0000

Agenda

- What are counterfactual explanations (CE) and algorithmic recourse (AR) and why are they useful?
- What dynamics are generated when off-the-shelf solutions to CE and AR are implemented in practice?

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have time
000000000	000	0000	000	0000

Agenda

- What are counterfactual explanations (CE) and algorithmic recourse (AR) and why are they useful?
- What dynamics are generated when off-the-shelf solutions to CE and AR are implemented in practice?
- Can we generate plausible counterfactuals relying only on the opaque model itself?

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we st
000000000	000	0000	000	0000

Agenda

- What are counterfactual explanations (CE) and algorithmic recourse (AR) and why are they useful?
- What dynamics are generated when off-the-shelf solutions to CE and AR are implemented in practice?
- Can we generate plausible counterfactuals relying only on the opaque model itself?
- How can we leverage counterfactuals during training to build more trustworthy models?

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have tin
●000000000	000	0000	000	0000

Background

Background ○●○○○○○○○

ynamics of CE and AR

Plausibility at all cost? 0000 Teaching models plausible explanations 000

If we still have time ... 0000

CE in Five Slides



Figure 2: Cats and dogs in two dimensions.

Patrick Altmeyer, Arie van Deursen, Cynthia C. S. Liem Explaining Models or Modelling Explanations Delft University of Technology

Background 000000000

ynamics of CE and AR

Plausibility at all cost? 0000 Teaching models plausible explanations 000

If we still have time ... 0000

CE in Five Slides

Model Training

Objective:

 $\min_{\boldsymbol{\theta}} \{ \mathsf{yloss}(M_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \}$

Patrick Altmeyer, Arie van Deursen, Cynthia C. S. Liem Explaining Models or Modelling Explanations Delft University of Technology

ynamics of CE and AR

Plausibility at all cost? 0000 Teaching models plausible explanations

If we still have time ... 0000

CE in Five Slides

Model Training

Objective:

 $\min_{\boldsymbol{\theta}}\{\mathsf{yloss}(M_{\boldsymbol{\theta}}(\mathbf{x}),\mathbf{y})\}$

Solution:

$$\begin{split} \theta_{t+1} &= \theta_t - \nabla_{\theta} \{ \mathrm{yloss}(M_{\theta}(\mathbf{x}), \mathbf{y}) \} \\ \theta^* &= \theta_T \end{split}$$

Background 000000000

ynamics of CE and AR

Plausibility at all cost? 0000 Teaching models plausible explanations 000

If we still have time ... 0000

CE in Five Slides



Figure 3: Fitted model. Contour shows predicted probability y= .

Patrick Altmeyer, Arie van Deursen, Cynthia C. S. Liem

Explaining Models or Modelling Explanations

Background 0000000000

ynamics of CE and AR

Plausibility at all cost? 0000 Teaching models plausible explanations 000

If we still have time ... 0000

CE in Five Slides

Counterfactual Search

Objective:

 $\min_{\mathbf{x}}\{\mathrm{yloss}(M_{\theta^*}(\mathbf{x}),\mathbf{y}^+)\}$

Background 000000●00

ynamics of CE and AR

Plausibility at all cost? 0000 Teaching models plausible explanations

If we still have time ... 0000

CE in Five Slides

Counterfactual Search

Objective:

$$\min_{\mathbf{x}}\{\mathsf{yloss}(M_{\theta^*}(\mathbf{x}),\mathbf{y}^+)\}$$

Solution:

$$\begin{split} \mathbf{x}_{t+1} &= \mathbf{x}_t - \nabla_{\theta} \{ \mathrm{yloss}(M_{\theta^*}(\mathbf{x}), \mathbf{y}^+) \} \\ \mathbf{x}^* &= \mathbf{x}_T \end{split}$$



CE in Five Slides

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

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

^a Altmeyer, Deursen, and Liem (2023) @ JuliaCon 2022



Figure 4: Counterfactual explanation for what it takes to be a dog.

Teaching models plausible explanations 000

If we still have time ... 0000

Algorithmic Recourse

Provided CE is valid, plausible and actionable, it can be used to provide recourse to individuals negatively affected by models.

"If your income had been X, then ..."



Figure 5: Counterfactuals for random samples from the Give Me Some Credit dataset (Kaggle 2011). Features 'age' and 'income' are shown.

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we st
000000000	●00	0000	000	0000

Dynamics of CE and AR

Dynamics of CE and AR $\bigcirc \bigcirc \bigcirc$

Plausibility at all cost?

Teaching models plausible explanations 000

If we still have time ... 0000

Hidden Cost of Implausibility



Figure 6: Endogenous Macrodynamics in Algorithmic Recourse.

Figure 7: Illustration of external cost of individual recourse.

^a Altmeyer, Angela, et al. (2023) @ SaTML 2023.

Insight: individual recourse neglects bigger picture.

Plausibility at all cost? 0000 Teaching models plausible explanations 000

Mitigation Strategies

- Incorporate hidden cost in reframed objective.
- Reducing hidden cost is equivalent to ensuring plausibility.

$$\begin{split} \mathbf{s}' &= \arg\min_{\mathbf{s}' \in \mathcal{S}} \{ \mathsf{yloss}(M(f(\mathbf{s}')), y^*) \\ &+ \lambda_1 \mathsf{cost}(f(\mathbf{s}')) + \lambda_2 \mathsf{extcost}(f(\mathbf{s}')) \} \end{split}$$

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have time
000000000	000	●000	000	0000

Plausibility at all cost?

Teaching models plausible explanations 000

Pick your Poison

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

Which one would you pick?



Figure 8: 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).

Teaching models plausible explanations 000

Faithful First, Plausible Second

Counterfactuals as plausible as the model permits¹.



Figure 9: KDE for training data.



Figure 10: KDE for model posterior.

¹ Altmeyer, Farmanbar, et al. (2023) @ AAAI 2024. [blog]

Patrick Altmeyer, Arie van Deursen, Cynthia C. S. Liem

Explaining Models or Modelling Explanations

Plausibility at all cost?

Teaching models plausible explanations 000

If we still have time ... 0000

Faithful Counterfactuals



Figure 11: Turning a 9 into a 7. *ECCCo* applied to MLP (a), Ensemble (b), JEM (c), JEM Ensemble (d).



Figure 12: Results for different generators (from 3 to 5).

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have time
000000000	000	0000	●00	0000

Teaching models plausible explanations

Teaching models plausible explanations $\bigcirc \bigcirc \bigcirc$

If we still have time .. 0000

Counterfactual Training: Method



Let the model compare its own explanations to plausible ones².

- Contrast faithful counterfactuals with data.
- 2 Use nascent CE as adversarial examples.



Figure 13: Example of an adversarial attack. Source: Goodfellow, Shlens, and Szegedy (2015)



Patrick Altmeyer, Arie van Deursen, Cynthia C. S. Liem

Background 000000000 ynamics of CE and AR

Plausibility at all cost?

Teaching models plausible explanations

If we still have time ... 0000

Counterfactual Training: Results



Figure 14: (a) conventional training, all mutable; (b) CT, all mutable; (c) conventional, *age* immutable; (d) CT, *age* immutable.

- Models trained with CT learn more plausible and (provably) actionable explanations.
- Predictive performance does not suffer, robust performance improves.

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have time
000000000	000	0000	000	●000

If we still have time ...

Plausibility at all cost? 0000 Teaching models plausible explanations 000 If we still have time ... ○●○○

Spurious Sparks of AGI

We challenge the idea that the finding of meaningful patterns in latent spaces of large models is indicative of AGI^a.

^a In Altmeyer et al. (2024) @ ICML 2024



Figure 15: Inflation of prices or birds? It doesn't matter!

Background Dynamics of CE and AR Plausibility at all cost? Teaching models plausible ex 000000000 000 000 000 000 If we still have time ... 00●0

Taija

- Model Explainability (CounterfactualExplanations.jl)
- Predictive Uncertainty Quantification (ConformalPrediction.jl)
- Effortless Bayesian Deep Learning (LaplaceRedux.jl)
- ... and more!

- Work presented @ JuliaCon 2022, 2023, 2024.
- Google Summer of Code and Julia Season of Contributions 2024.
- Total of three software projects@ TU Delft.



Figure 16: Trustworthy AI in Julia: github.com/JuliaTrustworthyAI

Background	Dynamics of CE and AR	Plausibility at all cost?	Teaching models plausible explanations	If we still have time
00000000	000	0000	000	000●

References

- Altmeyer, Patrick, Giovan Angela, Aleksander Buszydlik, Karol Dobiczek, Arie van Deursen, and Cynthia CS Liem. 2023. "Endogenous Macrodynamics in Algorithmic Recourse." In 2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML), 418–31. IEEE.
- Altmeyer, Patrick, Andrew M. Demetriou, Antony Bartlett, and Cynthia C. S. Liem. 2024. "Position Paper: Against Spurious Sparks-Dovelating Inflated AI Claims." https://arxiv.org/abs/2402.03962.

Altmeyer, Patrick, Arie van Deursen, and Cynthia C. S. Liem. 2023. "Explaining Black-Box Models through Counterfactuals." In Proceedings of the JuliaCon Conferences, 1:130.

Altmeyer, Patrick, Mojtaba Farmanbar, Arie van Deursen, and Cynthia C.S. Liem. 2023. "Faithful Model Explanations Through Energy-Constrained Conformal Counterfactuals."

Patrick Altmeyer, Arie van Deursen, Cynthia C. S. Liem