What If? Demystifying AI Decisions with Counterfactuals

David Martens

admantwerp.github.io



Advanced AI in our lives



Black Box?

- Deep learning: large artificial neural network with massive number of parameters
 - MobileNetV2: 4 million parameters
 - GPT-4: estimated 100 trillion parameters





Non-linear models

 $y_{SVM} = \operatorname{sign}[\sum_{t=1}^{N} \alpha_t y_t \exp\{-\|\mathbf{x} - \mathbf{x}_t\|_2^2 / \sigma^2\} + b]$

Complex AI models are creeping into our lives Whose decisions are extremely difficult to explain **But Explanations are needed if we want to trust AI**

Counterfactual explanation:

What needs to change in your data, to reach a different decision?

Example: credit scoring using sociodemo and financial data



User x_i: Sam, with income \$ 32,000 and 39 years old.

Sam is denied credit



Counterfactual explanation:

What needs to change in your data, to reach a different decision?

Example: credit scoring using sociodemo and financial data



User x_i: Sam

Sam is denied credit

WHY?

IF Sam would make 8,000\$ more THEN his predicted class would change from *denied* to *granted*



Dieter Brughmans, Pieter Leyman, David Martens (2023) NICE : an algorithm for nearest instance counterfactual explanations. *Data mining and knowledge discovery* p. 1-39

Example: gender prediction using movie viewing data



User x_i: Sam

Sam watched 120 movies Sam is predicted as male **WHY?**

Counterfactual Explanation

IF Sam would not have watched {Taxi driver, The Dark Knight, Die Hard, Terminator 2, Now You See Me, Interstellar}, THEN his predicted class would change from male to female LIME: Linear Interpretable Model-Agnostic Explainer (k=10)



Martens D, Provost F. (2014) *Explaining Data-Driven Document Classifications*. MIS Quarterly 38(1):73-99.

admantwerp.github.io

M.T. Ribeiro, S. Singh, C. Guestrin (2016) *Model-Agnostic Interpretability of Machine Learning* 2016 ICML Workshop on Human Interpretability in Machine Learning (WHI 2016), New York, NY <u>github.com/marcotcr/lime</u>

Ramon Y., Martens D., Provost F., Evgeniou T. (2019) Instance-level explanation algorithms SEDC, LIME, SHAP for behavioral and textual data: a counterfactual-oriented comparison, *Advances in Data Analysis and Classification* 14:4, p. 801-819.

Complex AI models are creeping into our lives Whose decisions are extremely difficult to explain But explanations are needed if we want to trust AI **Counterfactuals explain a decision of a model for an instance**

Counterfactual generating algorithms

Yet more algorithms to explain complex algorithms



 Ill-defined, measures: sparse/short, near, diverse, actionable, feasible, plausible, justified, fast, modelagnostic, data-agnostic, truthful, complete, stable, robust,

						Form	ulatio	n						Solution					200+ 10
Algorithm		Mode	4	Actio	nability	Р	Plausibility		Extra		Data types				Properties			_	"lethod
	тв	KB D	FOT	uncon	d. cond	dom	dens	. prot	o. dive	r. spar		0 8	Tools	Access	opt.	cov.	rtm.	Code	SUS
(2014.03) SEDC [129]		•		0	0	0	0	0	0	0	0	0.	heuristic	query		0	•	•	
(2015.08) OAE [51]		O C	0 0	•		Õ	Õ	Ŏ	Ŏ	Õ	•	Õ Õ	ILP	white-box	ě	Õ	ě	Ō	
(2016.05) HCLS [110, 112]	٠			٠	٠	٠	٠	0	0	0	O	0 0	grad opt/heuristic	gradient/query	٠	0	٠	٠	
(2017.06) Feature Tweaking [186]	٠	00	0 (0	0	0	0	0	0	0	٠	0 0	heuristic	white-box	0	0	٠	٠	
(2017.11) CF Expl. [196]	0	0	0	٠	0	0	0	0	٠	٠	۲	0 0	grad opt	gradient	0	0	0	0	
(2017.12) Growing Spheres [114]	٠	• •		0	0	0	0	0	0	•	0	• 0	heuristic	query	0	0	0	0	
(2018.02) CEM [55]	0	0	0	0	0	•	0	0	0	•	0	• 0	FISTA	class prob.	0	0	0	•	
(2018.02) POLARIS [209]	0	0			0	0	0	0	0	0		• 0	heuristic	gradient	0	0		•	
(2018.05) LORE [80]				8	8		8	8	0	8		00	gen alg + heuristic	query		8			
(2018.00) Local Foll Trees [190] (2018.00) Actionable Recourse [180]	ŏ			ě	ĕ	ĕ	ŏ	ŏ	ŏ	ŏ	ŏ		II P	white-box	ě	ŏ	ŏ		
(2018.11) Weighted CFs [77]	ĕ	ĕ		ō	ŏ	ŏ	ŏ	ŏ	ŏ	ĕ	ŏ	ŏŏ	heuristic	query	ŏ	ŏ	ŏ	ŏ	
(2019.01) Efficient Search [175]	Õ	Õ i	0	•	ŏ	۲	Õ	Õ	•	ě	ě	ΟÖ	MILP	white-box	Õ	Õ	Õ	٠	
(2019.04) CF Visual Expl. [76]	0	0	0	0	0	0	0	0	0	0	0	• 0	greedy search	white-box	٠	0	٠	0	
(2019.05) MACE [99]	٠	• •	•	•	•	٠	0	0	٠	٠	٠	0 0	SAT	white-box	٠	٠	٠	٠	
(2019.05) DiCE [145]	0	0	0	•	0		0	0	•	0	٠	0 0	grad opt	gradient	0	0	0	•	
(2019.05) CERTIFAI [179]	٠	• •		٠	0	•	0	0	•	0	٠	• 0	gen alg	query	0	0	0	0	
(2019.06) MACEM [56]	•			0	0	•	•	0	0	•	•	0 0	FISTA	query	0	0	0	0	
(2019.06) Expl. using SHAP [165]	•	•••	•	0	0	0	0	0	0	0	0	00	heuristic	query	0	0	0	•	
(2019.07) Nearest Observable [201] (2010.07) Quided Prototypes [101]	2						-						brute force	dataset	0	8	-		
(2019.07) Guided Prototypes [191]	ŏ	ŏ 7			ŏ			ŏ	ŏ	ŏ			grad opt/FISTA	gradient/query	ŏ	ŏ		ŏ	
(2019.07) REVISE [95] (2019.08) CLEAR [202]	ĕ	ž	Ĭ	ŏ	ŏ			ŏ	ŏ	ŏ		ŏŏ	heuristic	gradient	ŏ	ŏ	ŏ	ĕ	
(2019.08) MC-BRP [123]	0			ŏ	ŏ	Ō	Ō	ŏ	ŏ	ŏ	õ	ŏŏ	heuristic	query	ŏ	Õ	ŏ	ě	
(2019.09) FACE [162]				٠	•		٠	٠	0	0	۰	• 0	graph + heuristic	query	0	٠	٠	٠	
(2019.09) Equalizing Recourse [83]		•		•	0	0	0	0	0	0	٠	0 0	ILP/heuristic	white-box/query	0	0	0	0	
(2019.10) Action Sequences [163]	0	0	0	•	•	0	0	0	0	0	٠	• 0	program synthesis	class prob.	٠	٠	٠	٠	
(2019.10) C-CHVAE [156]		• •		•	0	•	•	0	0	0	•	00	grad opt + heuristic	query + gradient	0	•	•	•	
(2019.11) FOCUS [124]	•	00	. •	0	0	0	0	0	0	0	0	00	grad opt + heuristic	white-box	•	•	0	•	
(2019.12) Model-based CFs [127]	-	9				0		0		0			grad opt	gradient	0	0			
(2019.12) LIME-C/SHAP-C [104]	÷.			ŏ	ŏ	ŏ	ŏ	ŏ	ŏ	ŏ	Ň		neuristic	query	ŏ	ŏ		ŏ	
(2019.12) EMAF [41] (2019.12) PRINCE [71]	7	1.1		ŏ	ŏ	ŏ	ŏ	ŏ	ŏ	ŏ		ŏŏ	grau opt	dataset/query	ĕ	ĕ		ĕ	
(2019.12) LowProFool [18]	Õ	Õ	0	ŏ	ŏ	ě	ŏ	ŏ	ŏ	ŏ	õ	ŏŏ	grad opt	gradient	ě	Ō	Ō	Ō	
(2020.01) ABELE [79]	٠	0	0	Õ	Õ	Õ	۰	•	Õ	Õ	Õ	• 0	gen alg + heuristic	query + data	Õ	Õ	Õ	٠	
(2020.01) SHAP-based CFs [66]				0	0	0	0	0	0	٠	٠	0 0	heuristic	query	0	0	0	٠	
(2020.02) CEML [11-13]		• •	•	0	0	0	0	٠	0	•	٠	• 0	grad opt/heuristic	gradient/query	٠	0	٠	٠	
(2020.02) MINT [100]	٠			•	•	•	•	0	•	•	•	00	SAT	white-box	•	•	0	•	
(2020.03) ViCE [74]		•		•	0	0	0	0	0	0	0	00	heuristic	query	0	0	0	•	
(2020.03) Plausible CFs [22]	0	0		0	0	0	•	0	0	0	0	• •	grad opt + gen alg	dataset	•	0	0	0	
(2020.04) SEDC-T [193]	÷.											• •	heuristic	query		Ň			
(2020.04) MOC [52] (2020.04) SCOUT [199]	ŏ			ŏ	ŏ	ŏ	ŏ	ŏ	ŏ	ŏ	ŏ	ŏŏ	gen aig grad ont	query	ŏ	ŏ	ĕ	ŏ	
(2020.04) ASP-based CFs [28]				ĕ	ĕ	ŏ	ŏ	ŏ	ŏ	ŏ	ĕ	0 0	answer-set prog	query	ŏ	ŏ	ŏ	ŏ	
(2020.05) CBR-based CFs [103]				ŏ	ŏ	•	۲	ŏ	ŏ	•	ě	õõ	heuristic	query + data	õ	õ	õ	ŏ	
(2020.06) Survival Model CFs [106]	٠			٠	0	0	0	٠	0	0	0	0 0	gen alg	query	٠	0	0	0	
(2020.06) Probabilistic Recourse [101]	٠	• •		•	0	٠	٠	0		0	0	0 0	grad opt/brute force	gradient/query	٠	0	0	٠	
(2020.06) C-CHVAE [155]	٠	0	0	•	0	0	٠	0	0	0	0	00	grad opt	gradient	0	0	0	•	
(2020.07) FRACE [210]	0	0	0	0	0	•	•	0	0	0	0	• 0	grad opt	gradient	0	0	•	0	
(2020.07) DACE [96]	•	0		•		•	•	0	•	0		0 0	MILP	white-box	0	0	•	0	
(2020.07) CRUDS [60]	÷					0		0		0	0	00	grad opt	gradient/data	0	•	0	0	
(2020.07) Gradient Boosted CFs [5] (2020.08) Gradual Construction [07]		0.0		0	0			0		0			neuristic	data class prob	0	0	0		
(2020.08) DECE [44]	õ	ŏ	ŏ		ě	õ		ő	ě	ĕ	ě	õ õ	ared ont	gradient	ŏ	ŏ	ŏ	ĕ	
(2020.08) Time Series CFs [16]	ĕ	ŏ	ŏ	õ	Ő	ŏ	ŏ	ŏ	Ő	Ő	ŏ	ŏŏ	heuristic	query	ŏ	ĕ	ĕ	ŏ	
(2020.08) PermuteAttack [87]			0	ĕ	ŏ	Õ	Õ	Õ	•	õ	ě	0 0	gen alg	query	Õ	Ō	Ō	Õ	
(2020.10) Fair Causal Recourse [195]	٠			•	0	٠	٠	0		0	O	0 0	grad opt/brute force	gradient/query	٠	0	0	٠	
(2020.10) Recourse Summaries [167]	٠	• •		•	0	0	0	0	0	0	٠	0 0	itemset mining alg	query	0	0	0	0	
(2020.10) Strategic Recourse [43]	0	0	0	٠	٠	0	0	0	0	0	٠	0 0	Nelder-Mead	query	٠	٠	0	٠	
(2020.11) PARE [172]	0	0	0	0	0	0	0	0	0	0	0	0.	grad opt + heuristic	query	•	0	0	0	

Karimi et al (2021) A survey of algorithmic recourse: contrastive explanations and consequential recommendations <u>https://arxiv.org/pdf/2010.04050.pdf</u> Complex AI models are creeping into our lives Whose decisions are extremely difficult to explain But Explanations are needed if we want to trust AI Counterfactuals explain a decision of a model for an instance What these should look like is ill-defined and open research

The Counterfactual Explains Wrong Decisions

- To improve the predictive performance of the model
- Example
 - Data: image
 - Task: predict if missile in image



Tom Vermeire, Dieter Brughmans, Sofie Goethals, Raphael Mazzine de Oliveira, David Martens (2022) Explainable Image Classication with Evidence Counterfactual, Pattern Anal Applic 25, 315–335

The Counterfactual Explains Wrong Decisions

- To improve the predictive performance of the model
- Example
 - Data: image
 - Task: predict if missile in image
 - Mainly interested in improving misclassifications
 - Issue: Lighthouse wrongly classified as missile



(a) Original class: *missile*



The Counterfactual Explains Wrong Decisions

- To improve the predictive performance of the model
- Example
 - Data: image
 - Task: predict if missile in image
 - Mainly interested in improving misclassifications
 - Issue: Lighthouse wrongly classified as missile
 - Pattern learnt: line of smoke, resembling exhaust plume behind missiles, indicates missile



(a) Original class: *missile*



(b) Counterfactual explanation



(c) Counterfactual class: *beacon*

Complex AI models are creeping into our lives Whose decisions are extremely difficult to explain But Explanations are needed if we want to trust AI Counterfactuals explain a decision of a model for an instance What these should look like is ill-defined and open research Very useful to explain (incorrect or unfair) decisions

Imagine a world with CF explanations



Volgen

The @AppleCard is such a fucking sexist program. My wife and I filed joint tax returns, live in a community-property state, and have been married for a long time. Yet Apple's black box algorithm thinks I deserve 20x the credit limit she does. No appeals work.

12:34 - 7 nov. 2019



Apple Card 🤣 @AppleCard

Well, if your wife would also have had a 20+ year relationship with our bank, and would have been regarded as Premium customer at some point in time, she would also receive a 20x credit limit



Apple Card 🤣 @AppleCard

Well, if your wife's relationship status would have been "husband" instead of "wife", she would also receive a 20x credit limit

We clearly messed up, we're updating our models now.



Ah, ok, thanks for the additional feedback!



Steve Wozniak 🤣 @stevewoz

Glad you found this and react responsibly. It's how big tech should be in the 21st century.

Q&A

Play around with CF expanations: <u>admantwerp.github.io</u>

Data	Algorithm	Reference
TEXT	SEDC	Martens and Provost (2014)
	LIME-C, SHAP-C	Ramon et al (2020)
BEHAVIORAL	SEDC	Martens and Provost (2014)
	LIME-C, SHAP-C	Ramon et al (2020)
IMAGE	SEDC(-T)	Vermeire and Martens (2020)
STRUCTURED	NICE	Brughmans and Martens (2021)

Book on Data Science Ethics: <u>www.dsethics.com</u>

David Martens (2022) Data Science Ethics: Concepts, Techniques and Cautionary Tales Oxford University Press, 272 pages.

