

WHY AI NEEDS MATHS AND STATS

- lessons from working in a CS field

Martin Jullum

Godt Hjort, Dec 5th 2023

AI

$$\int \frac{1}{\sqrt{x^2 \pm a^2}} dx = \ln|x + \sqrt{x^2 \pm a^2}| + C$$

$$(a+b)^2 = a^2 + 2ab + b^2$$

$$\vec{A} \cdot (\vec{B} + \vec{C}) = \vec{A} \cdot \vec{B} + \vec{A} \cdot \vec{C}$$

$$y = kx + m$$

$$x \in [3; +\infty)$$

$$\sinh x = -i \sin(ix)$$

$$U = \int_a^b \pi f^2(x) dx$$

$$\lim_{n \rightarrow \infty} \exists N \in \mathbb{N} \forall n > N |x_n - a| < \epsilon$$

$$\sinh(x) = \frac{e^x - e^{-x}}{2}$$

$$\int x^n dx = \frac{x^{n+1}}{n+1} + C$$

$$\log(x)$$

$$\cos B \cos C + \sin B \sin C \cos \alpha$$

$$\log(ab) = \log a + \log b$$

$$S = 4\pi R^2$$

$$V = \frac{4}{3}\pi R^3$$

$$(e^x)' = e^x$$

$$\int_a^b f(x) dx$$

$$\ln(a-b)$$

$$\cos x = \operatorname{Re}\{e^{ix}\}$$

$$x! = 1$$

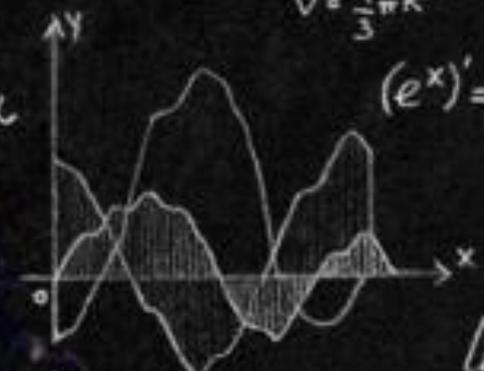
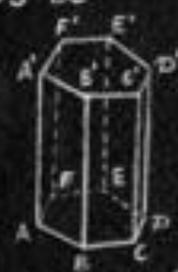
$$\operatorname{tg} \alpha = \frac{\sin \alpha}{\cos \alpha}$$

$$(x^n)' = nx^{n-1}$$

$$(\sqrt{x})' = \frac{1}{2\sqrt{x}}$$

$$(\ln x)' = \frac{1}{x}$$

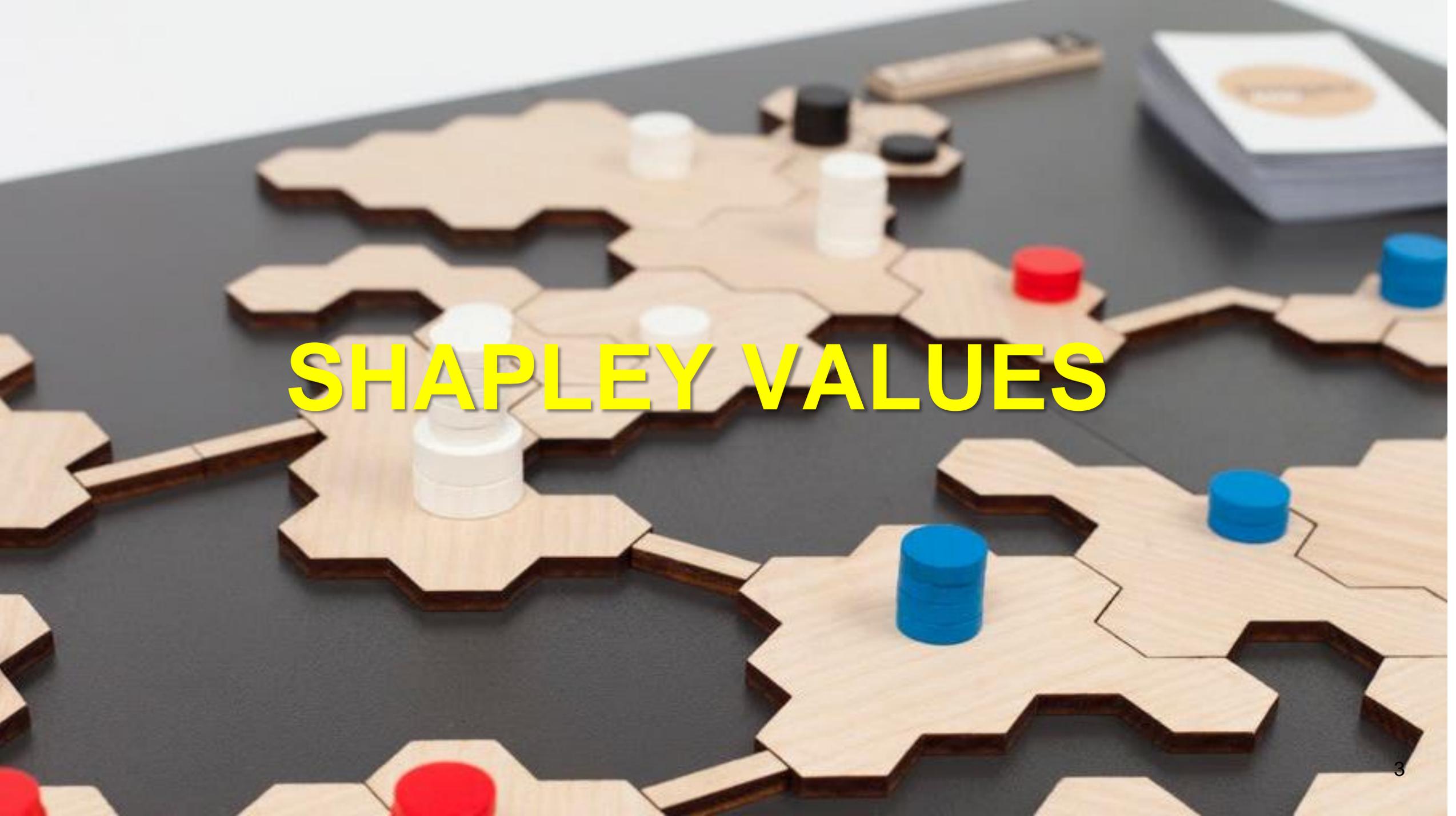
$$\log_a(xy) = \log_a x + \log_a y$$



Explainable AI (XAI)

- Understanding what black box models do
- A field largely driven by computer scientists
- Me and colleagues
 - Work in the subfield restricted to tabular data – i.e. regression: $y = f(x_1, x_2, x_3)$
 - Try to use our statistical mindset to improve/repair the methodology in the field



A wooden board game with a hexagonal grid and various colored pieces (white, black, red, blue) on a dark grey surface. The text "SHAPLEY VALUES" is overlaid in yellow.

SHAPLEY VALUES

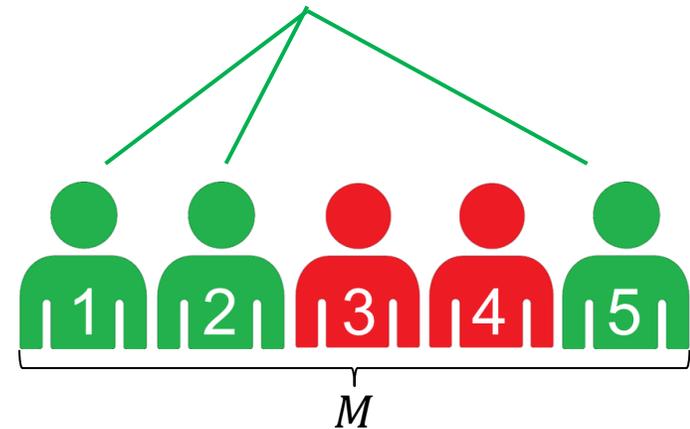
Shapley values

- ▶ Concept from (cooperative) game theory in the 1950s
- ▶ Used to distribute the total payoff to the players
- ▶ Explicit formula for the “fair” payment to every player j :

$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} \frac{|S|! (|M| - |S| - 1)!}{|M|!} (v(S \cup \{j\}) - v(S))$$

$v(S)$ is the payoff with only players in subset S

- ▶ Several mathematical optimality properties



Shapley values for taxi sharing

Costs: \$3/mi

$$v(\{R, B, G\}) = (4 + 6 + 2)mi * \$3 = \$36$$

$$v(\{\}) = \$0$$

$$v(\{R\}) = 4mi * \$3 = \$12$$

$$v(\{B\}) = (5 + 2)mi * \$3 = \$21$$

$$v(\{G\}) = 5mi * \$3 = \$15$$

$$v(\{R, B\}) = (4 + 6)mi * \$3 = \$30$$

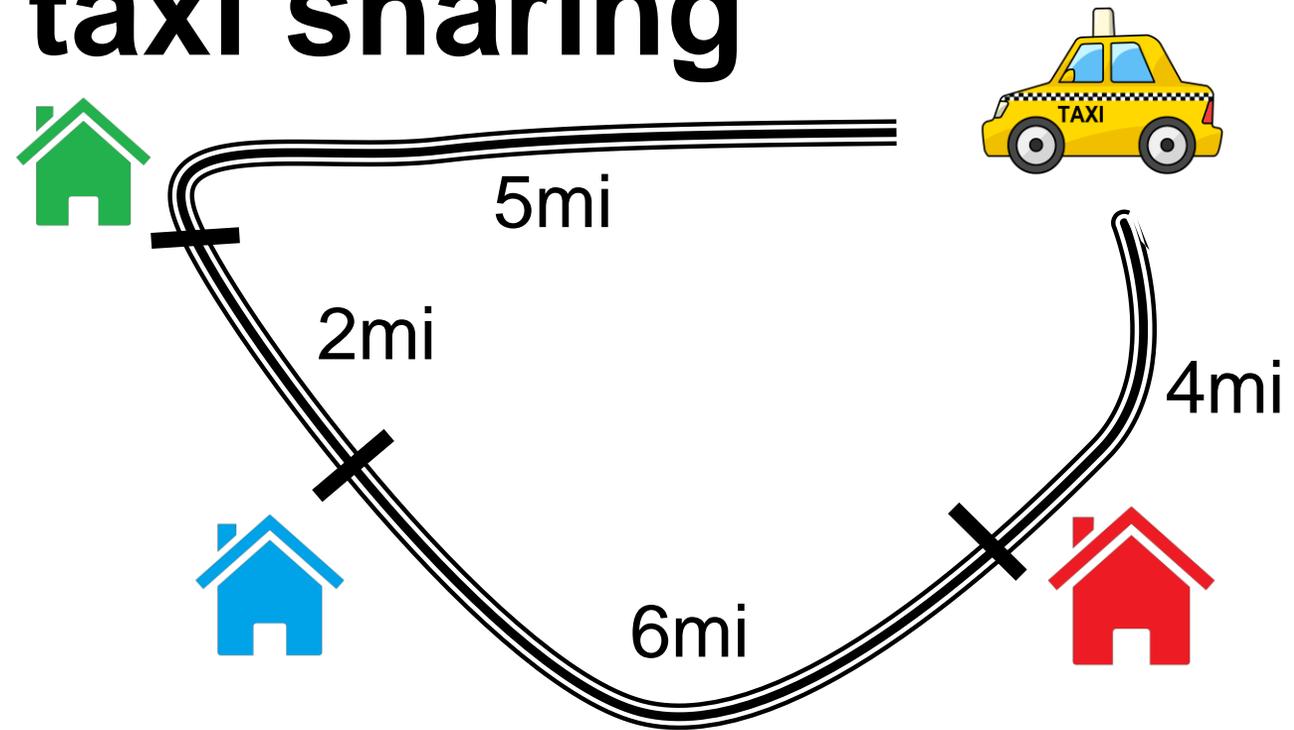
$$v(\{R, G\}) = (4 + 6 + 2)mi * \$3 = \$36$$

$$v(\{B, G\}) = (5 + 2)mi * \$3 = \$21$$

$$\phi_R = \frac{1}{3}(v(\{R, B, G\}) - v(\{B, G\})) + \frac{1}{6}(v(\{R, B\}) - v(\{B\})) + \frac{1}{6}(v(\{R, G\}) - v(\{G\})) + \frac{1}{3}(v(\{R\}) - v(\{\})) = \$14$$

$$\phi_B = \frac{1}{3}(v(\{R, B, G\}) - v(\{R, G\})) + \frac{1}{6}(v(\{R, B\}) - v(\{R\})) + \frac{1}{6}(v(\{B, G\}) - v(\{G\})) + \frac{1}{3}(v(\{B\}) - v(\{\})) = \$11$$

$$\phi_G = \frac{1}{3}(v(\{R, B, G\}) - v(\{R, B\})) + \frac{1}{6}(v(\{R, G\}) - v(\{R\})) + \frac{1}{6}(v(\{B, G\}) - v(\{B\})) + \frac{1}{3}(v(\{G\}) - v(\{\})) = \$11$$



Shapley values for prediction explanation

► Approach popularised by Lundberg & Lee (2017)

- Players = covariates (x_1, \dots, x_M)
- Payoff = prediction ($f(\mathbf{x}^*)$)
- Contribution function: $v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*]$
- Properties

$$\phi_0 + \sum_{j=1}^M \phi_j = f(\mathbf{x}^*)$$

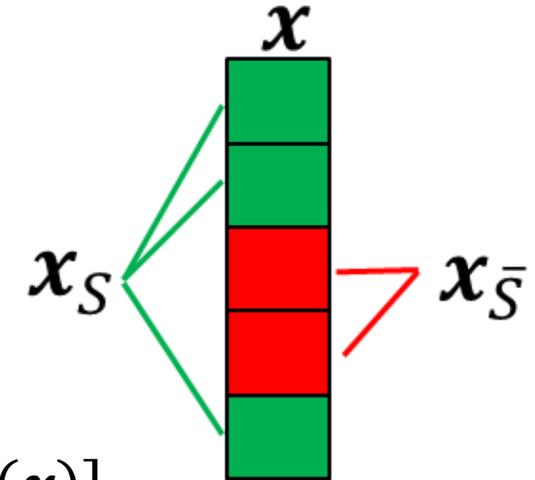
$$\phi_0 = E[f(\mathbf{x})]$$

$$f(\mathbf{x}) \perp\!\!\!\perp x_j$$

implies $\phi_j = 0$

$$x_i, x_j \text{ same contribution}$$

implies $\phi_i = \phi_j$



- Interpretation of ϕ_j : **The prediction change caused by observing the value of x_j – averaged over whether the other covariates were observed or not**

Two main challenges

1. The computational complexity in the Shapley formula is of size 2^M

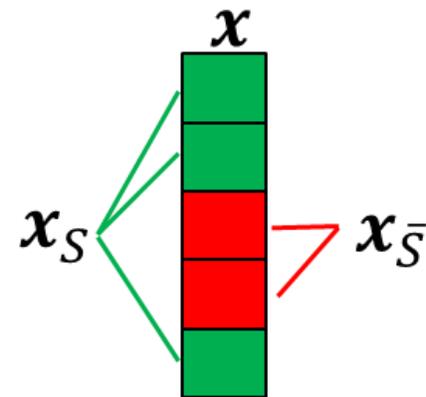
$$\phi_j = \sum_{S \subseteq M \setminus \{j\}} \frac{|S|! (|M| - |S| - 1)!}{|M|!} (v(S \cup \{j\}) - v(S))$$

2. Estimating the contribution function

$$v(S) = E[f(\mathbf{x}) | \mathbf{x}_S = \mathbf{x}_S^*] = \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}} | \mathbf{x}_S = \mathbf{x}_S^*) d\mathbf{x}_{\bar{S}}$$

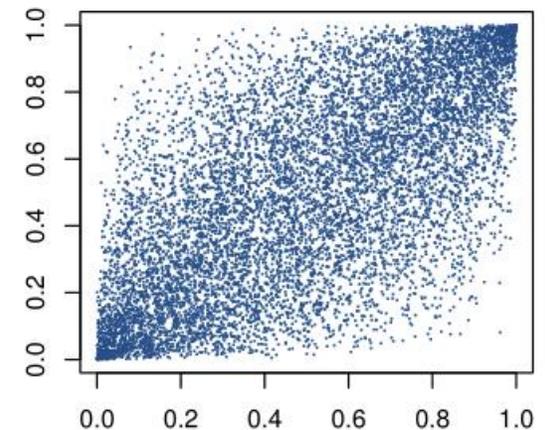
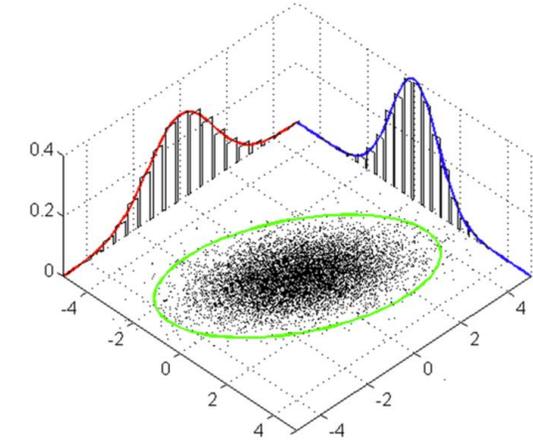
- Lundberg & Lee (2017)
 - Approximates $v(S) \approx \int f(\mathbf{x}_{\bar{S}}, \mathbf{x}_S^*) p(\mathbf{x}_{\bar{S}}) d\mathbf{x}_{\bar{S}}$,
 - Estimates $p(\mathbf{x}_{\bar{S}})$ using the empirical distribution of the training data
 - Monte Carlo integration to solve the integral

This assumes the covariates are independent!



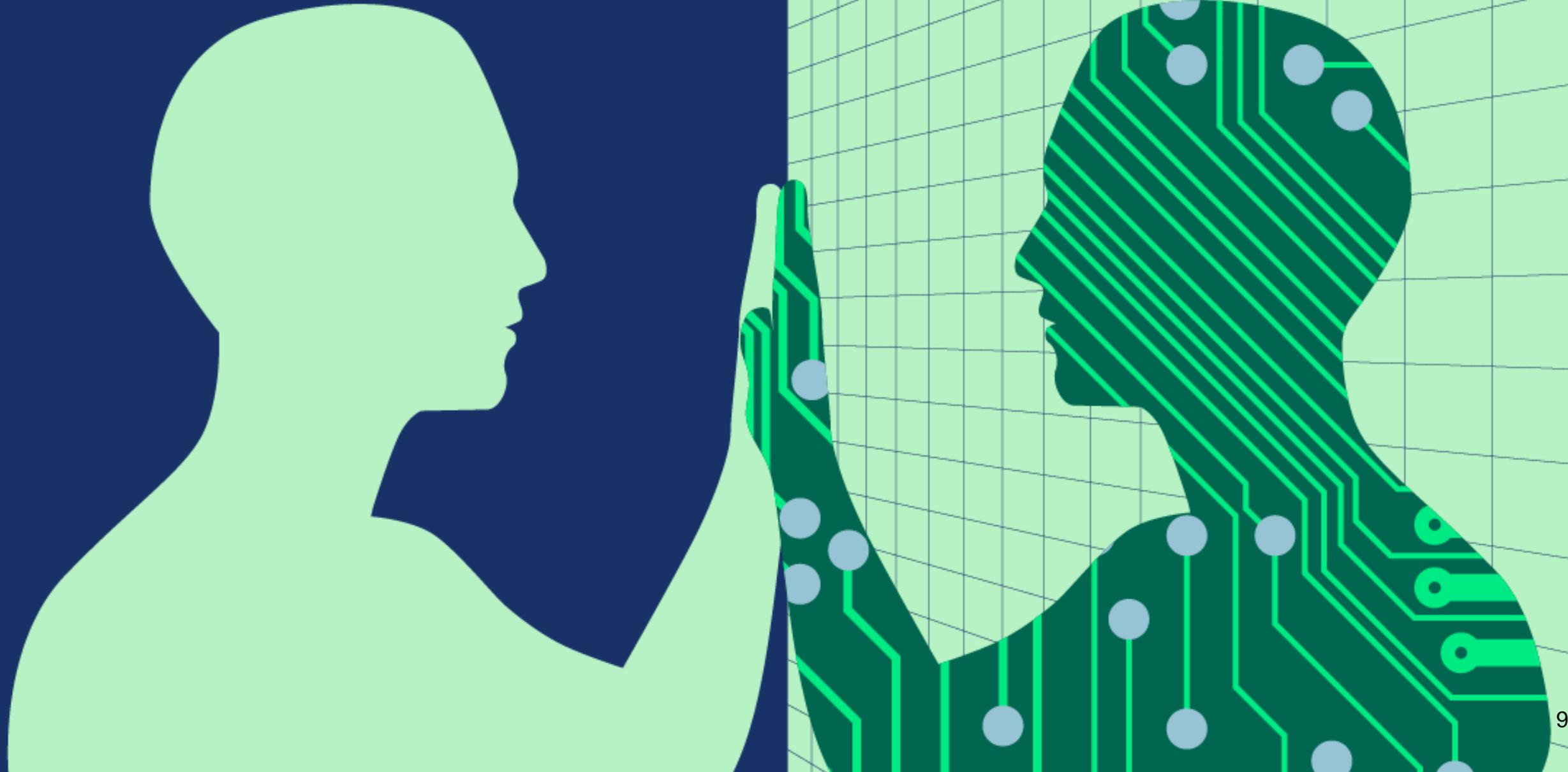
Our contribution

- ▶ Dependence-aware approaches to estimate $v(S) = E[f(\mathbf{x})|\mathbf{x}_S = \mathbf{x}_S^*]$ properly
- ▶ We do this by estimating $p(\mathbf{x}_{\bar{S}}|\mathbf{x}_S = \mathbf{x}_S^*)$ properly
- ▶ Several alternative methods
 - Gaussian distribution
 - Empirical nonparametric method
 - Empirical margins + vine copulas to estimate dependence structure
 - Conditional inference trees (ctree)
 - Variational autoencoders with arbitrary conditioning (VAEAC)
 - Direct regression on $v(S) = E[f(\mathbf{x})|\mathbf{x}_S = \mathbf{x}_S^*] \sim \mathbf{x}_S$
 - Common regression model for any $v(S)$ using masking trick



COUNTERFACTUAL

EXPLANATIONS



Example case

Automatic processing of loaning applications based on default prediction model

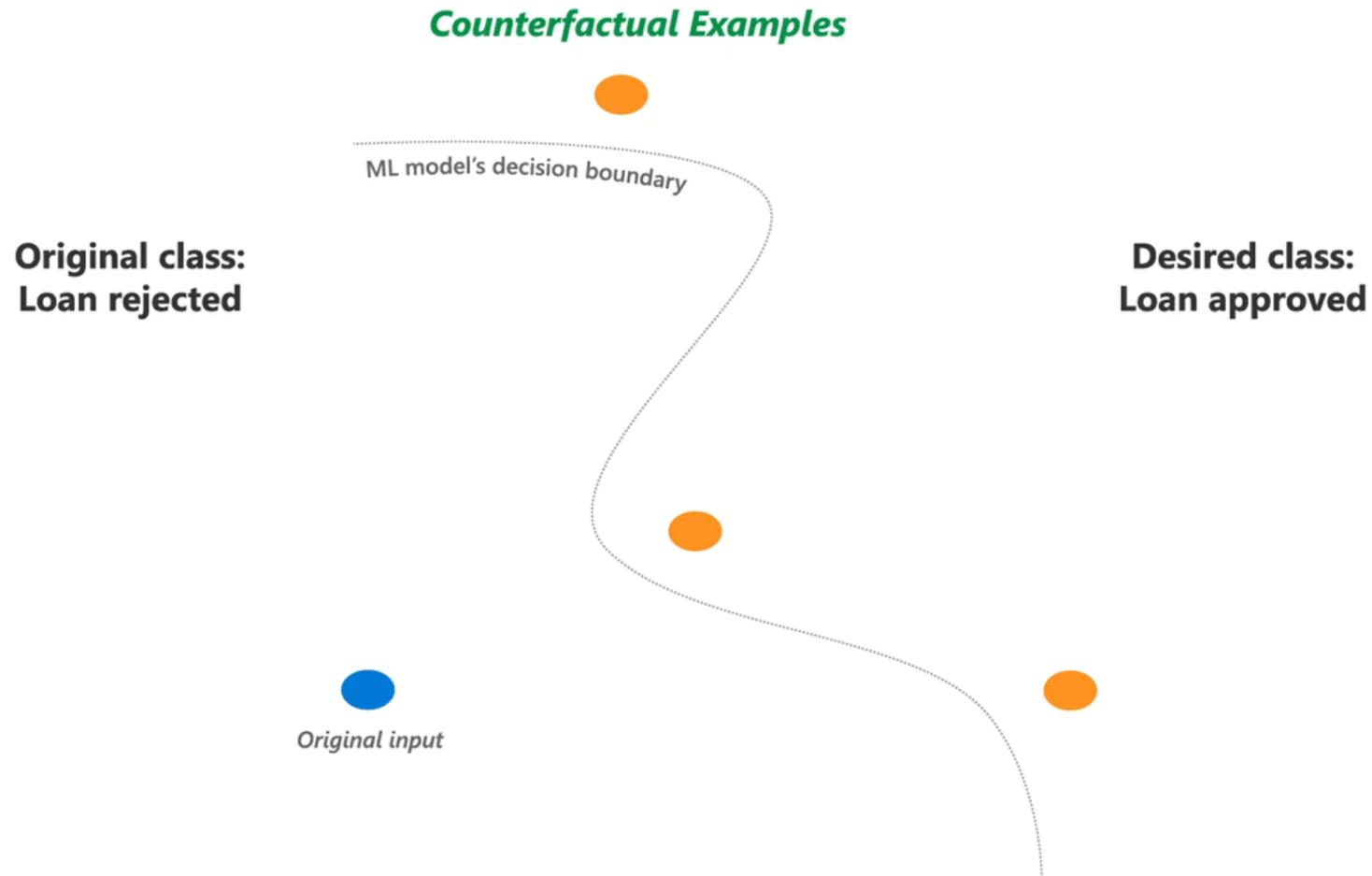
- ▶ Response y : Loan defaulted or not
- ▶ Covariates $\mathbf{x} = (x_1, \dots, x_p)$: Info about the applicant, salary, previous defaults, transactions history, etc
- ▶ Fit regression model f : Model trained to predict probability of default:
$$f(\mathbf{x}) \approx \Pr(y = \text{default} | \mathbf{x})$$
- ▶ *Loan approved if $f(\mathbf{x}) < c = 0.1$*

CASE: Peter has features \mathbf{x}^* , and got his loan application rejected as $f(\mathbf{x}^*) = 0.2 > c$

Question: What can Peter do to receive a loan?

The idea

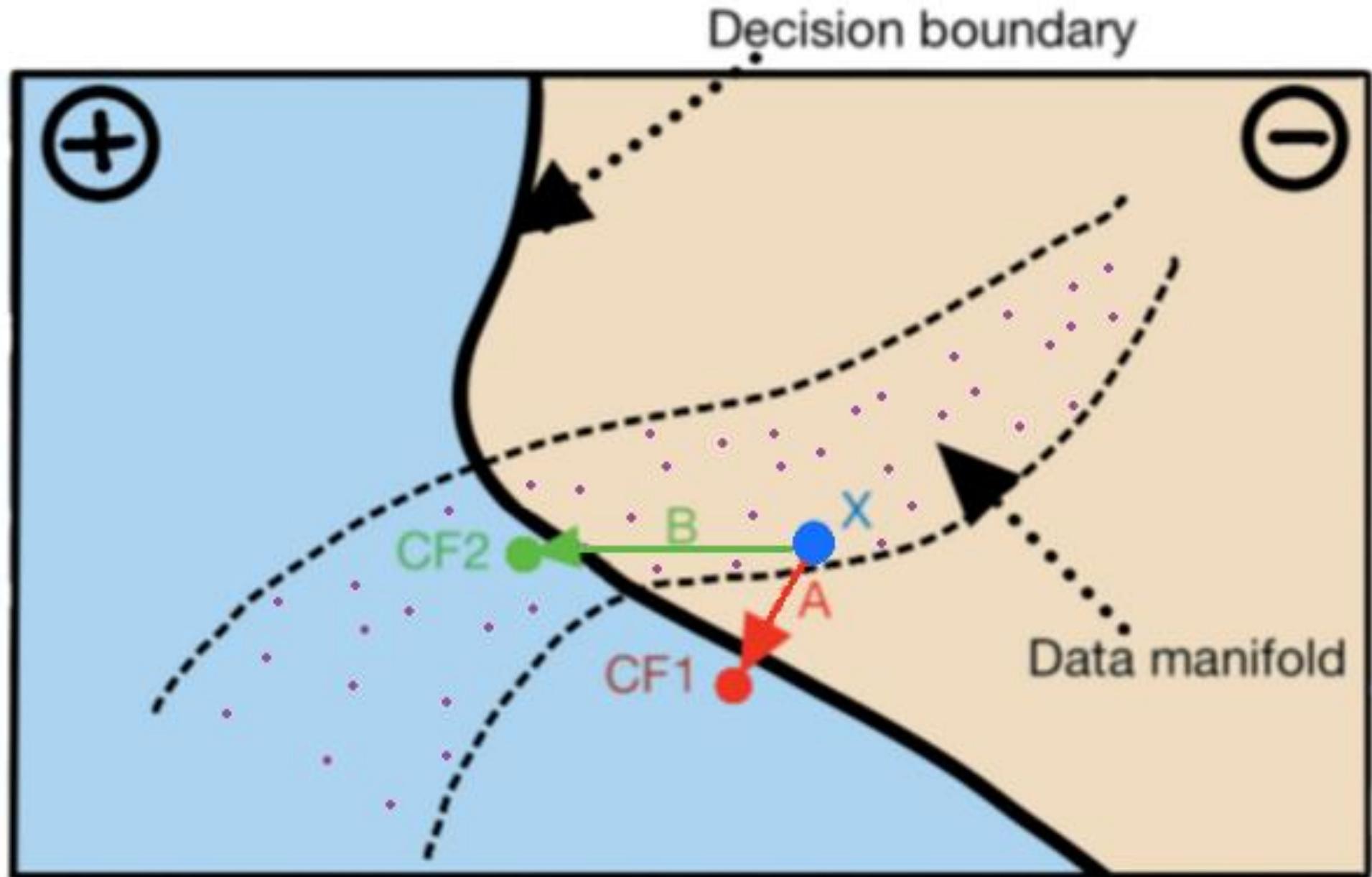
CE solution: Examples of (minimal) changes in covariates which approves the application



Criteria

Desired properties

1. On-manifold
2. Actionable
3. Valid
4. Low cost



Types of CE methods

Optimization based methods

- ▶ Minimize loss (wrt example \mathbf{e}) of the type $L_{\mathbf{x}^*}(\mathbf{e}) = \text{dist}_1(f(\mathbf{e}), c) + \lambda \cdot \text{dist}_2(\mathbf{x}^*, \mathbf{e})$
 - Often require differentiable f
 - Not necessarily on-manifold
 - Categorical covariates more troublesome

Heuristic search-based methods

- ▶ Optimization with heuristic search strategies

Instance-based methods

- ▶ Finds counterfactuals by searching for instances in a reference distribution/dataset

Our simple method: MCCE

MCCE: Monte Carlo sampling of valid and realistic counterfactual explanations

3-step procedure to produce a counterfactual example e

1. **Model**: Model the joint distribution of mutable covariates, given the fixed covariates and *the decision*
2. **Generate**: Generate a large number K of samples from the modelled distribution with the specified fixed covariates x^{*f} and desired decision
3. **Post-process**: Discard the invalid samples, and choose the one “nearest” to x^*

Step 1: Model

- ▶ Utilize the standard probability property:

$$p(\mathbf{X}^m \mid \mathbf{X}^f, Y') = p(X_1^m \mid \mathbf{X}^f, Y') \prod_{i=2}^q p(X_i^m \mid \mathbf{X}^f, Y', X_1^m, \dots, X_{i-1}^m)$$

STEP 1: MODEL

Training data

		Features			
		Immutable	Mutable		
Age	Sex	Salary	Def. last year	f(x)	Decision
30	M	\$ 3500	yes	0.24	0
28	F	\$ 8000	no	0.12	0
42	M	\$ 7500	no	0.04	1
27	F	\$ 9500	yes	0.21	0
39	M	\$ 5000	no	0.09	1
28	F	\$ 4000	no	0.08	1
32	F	\$ 7300	no	0.12	0
⋮	⋮	⋮	⋮	⋮	⋮
23	M	\$ 4300	yes	0.31	0

➔

Tree 1

Salary ~ Age, Sex, Decision

Age < 45

Sex = F

Decision = 1

Tree 2

Def. last y ~ Age, Sex, Decision, Salary

Decision = 1

Age < 30

Salary < \$6000

Salary < \$ 5000

Step 2: Generation

STEP 2: GENERATION

Observations to explain

Explain ID	Features				$f(x)$	Decision
	Immutable		Mutable			
	Age	Sex	Salary	Def. last year		
1	30	F	\$ 6000	yes	0.18	0
2	25	M	\$ 4500	no	0.30	0



	Age	Sex	Decision	Salary	Def. last year
D_1	30	F	1	-	-
	30	F	1	-	-
	30	F	1	-	-
	30	F	1	-	-
D_2	25	M	1	-	-
	25	M	1	-	-
	25	M	1	-	-
	25	M	1	-	-

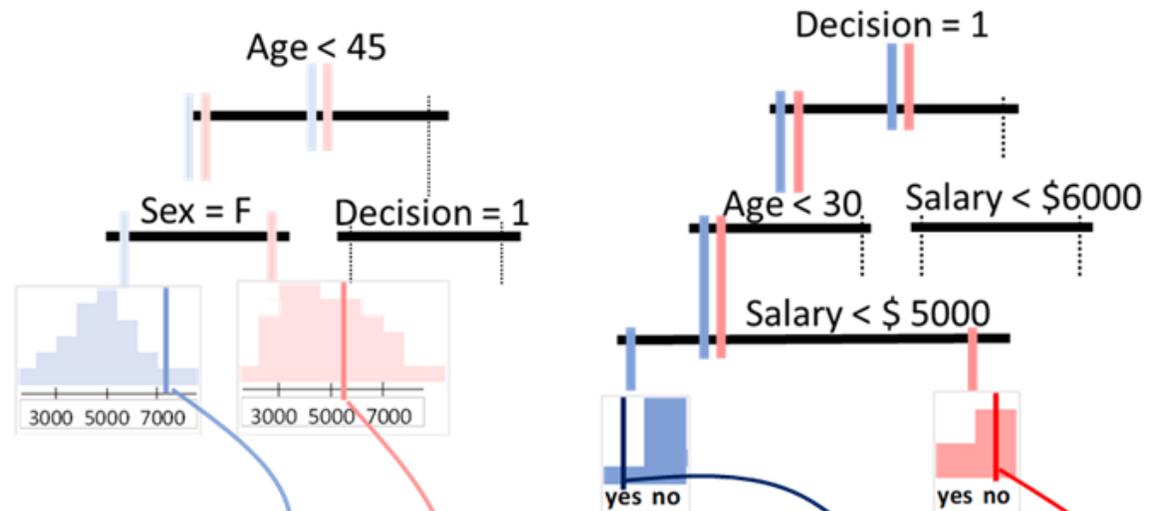
K

Age	Sex	Decision	Salary	Def. last year
30	F	1	\$ 4500	-
30	F	1	\$ 6000	-
30	F	1	\$ 7500	-
30	F	1	\$ 3800	-
25	M	1	\$ 6000	-
25	M	1	\$ 4800	-
25	M	1	\$ 5300	-
25	M	1	\$ 4600	-

Update D_1 & D_2

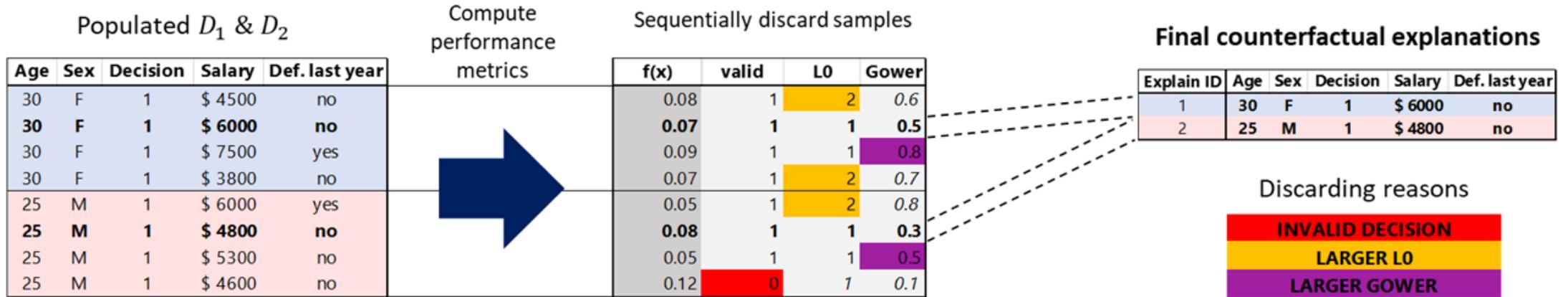
Age	Sex	Decision	Salary	Def. last year
30	F	1	\$ 4500	no
30	F	1	\$ 6000	no
30	F	1	\$ 7500	yes
30	F	1	\$ 3800	no
25	M	1	\$ 6000	yes
25	M	1	\$ 4800	no
25	M	1	\$ 5300	no
25	M	1	\$ 4600	no

Update D_1 & D_2



Step 3: Post-processing

STEP 3: POST-PROCESSING



Our contribution

MCCE

- ▶ Simple, yet effective
- ▶ Flexible
- ▶ Scalable and easy to implement
- ▶ Outperforms competing methods in terms of both accuracy and speed



TAKE HOME

To think differently

To simplify

DARE!

To take a role in an unfamiliar field

To use your statistical mindset

