

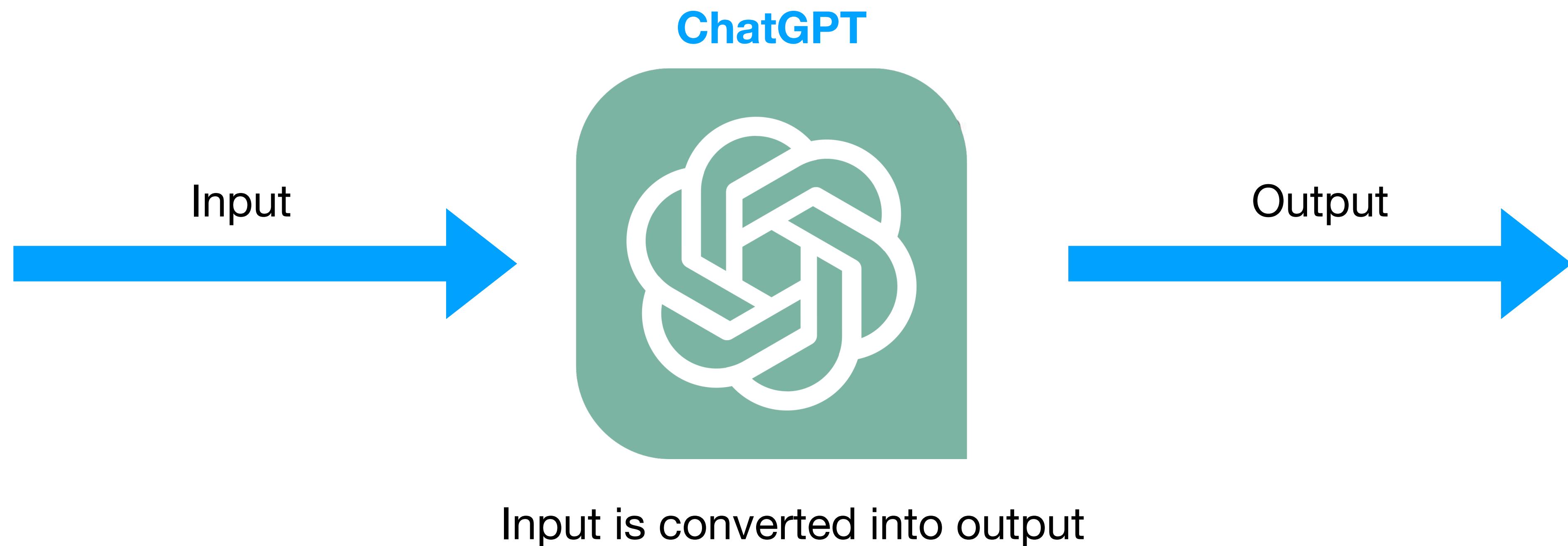
Toward Robust and Reliable Post-hoc Explanations to Machine Learning Models

Ngoc Bui - Graph and Geometric Deep Learning Group - Sep 2023

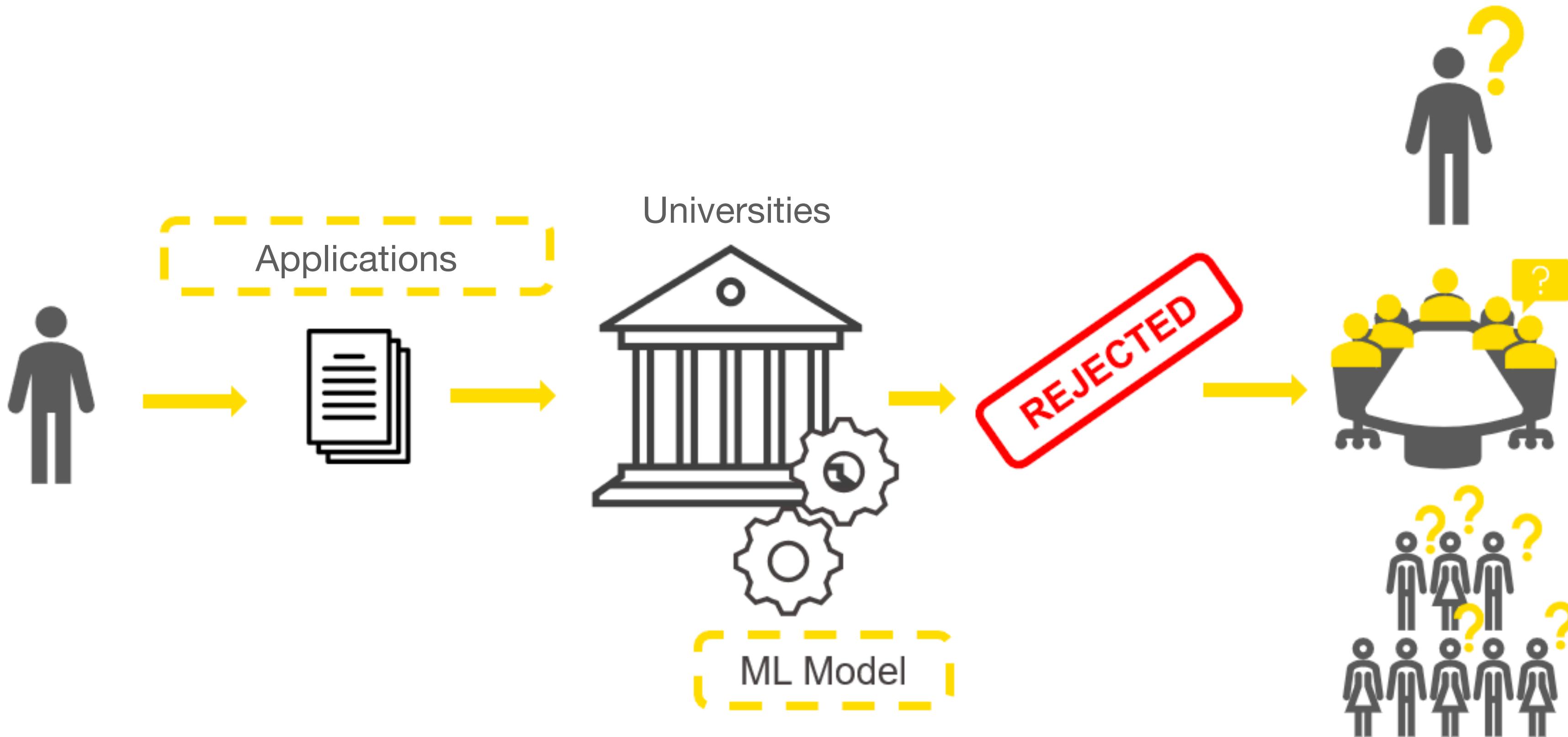
Why's Explainable AI?



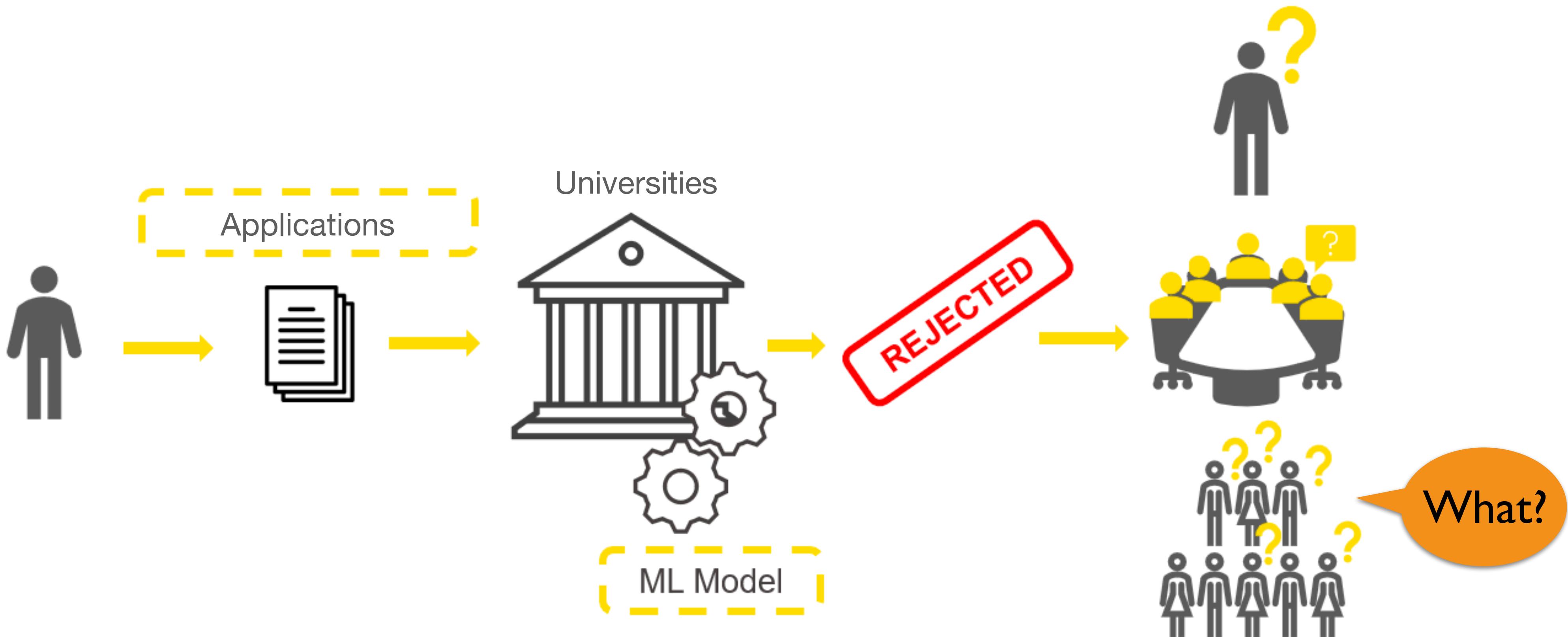
Why's Explainable AI?



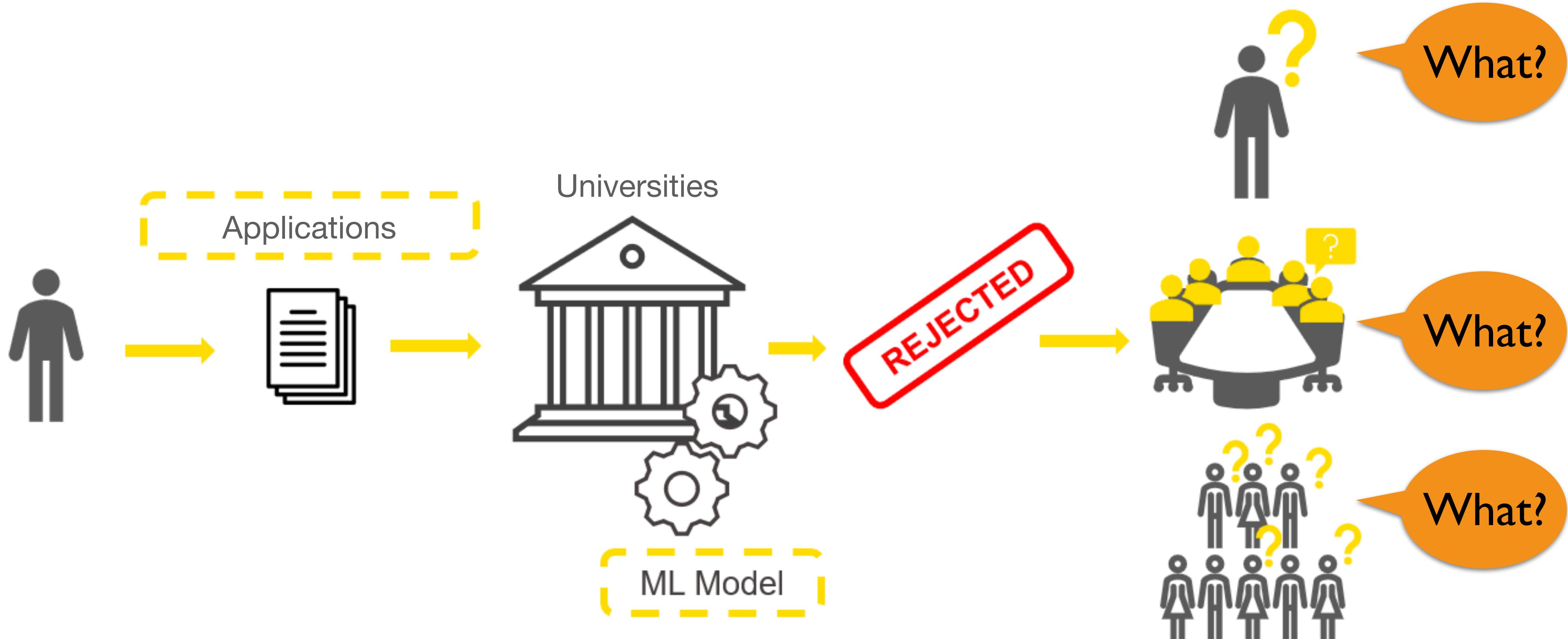
Why's Explainable AI?



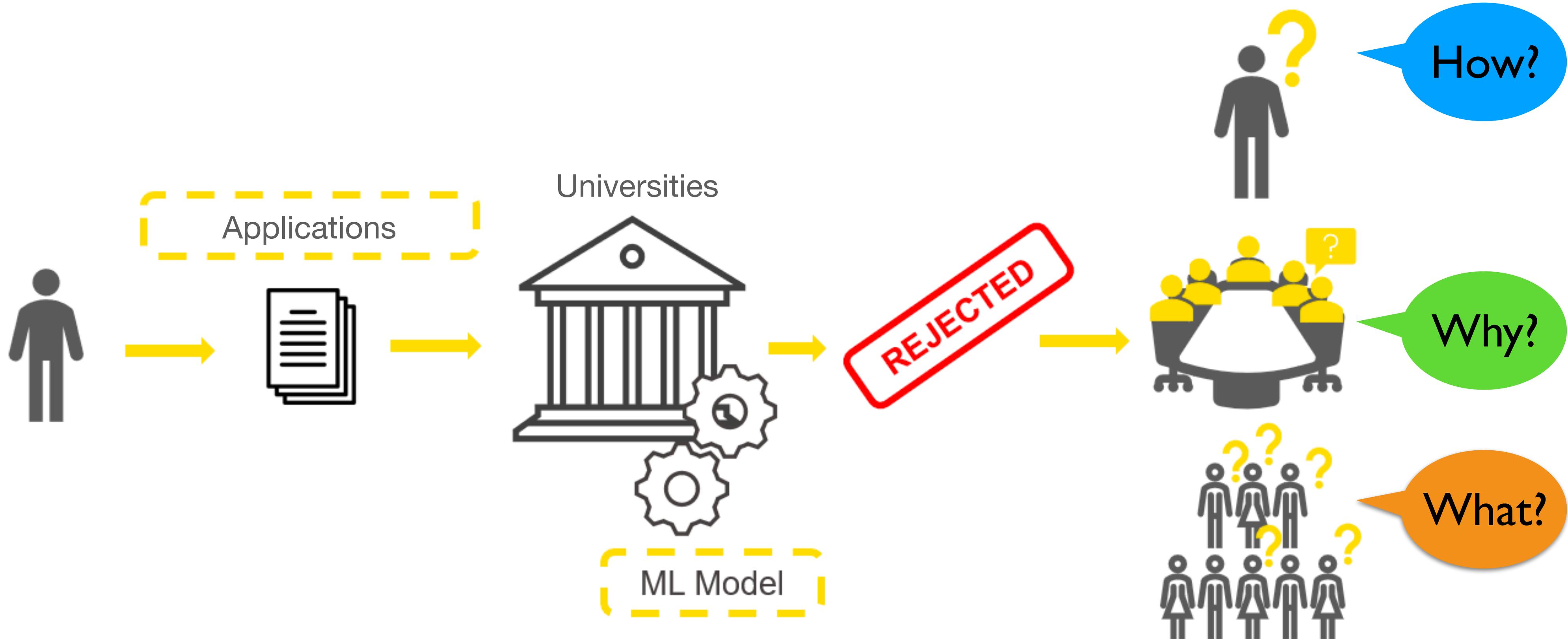
Why's Explainable AI?



Why's Explainable AI?

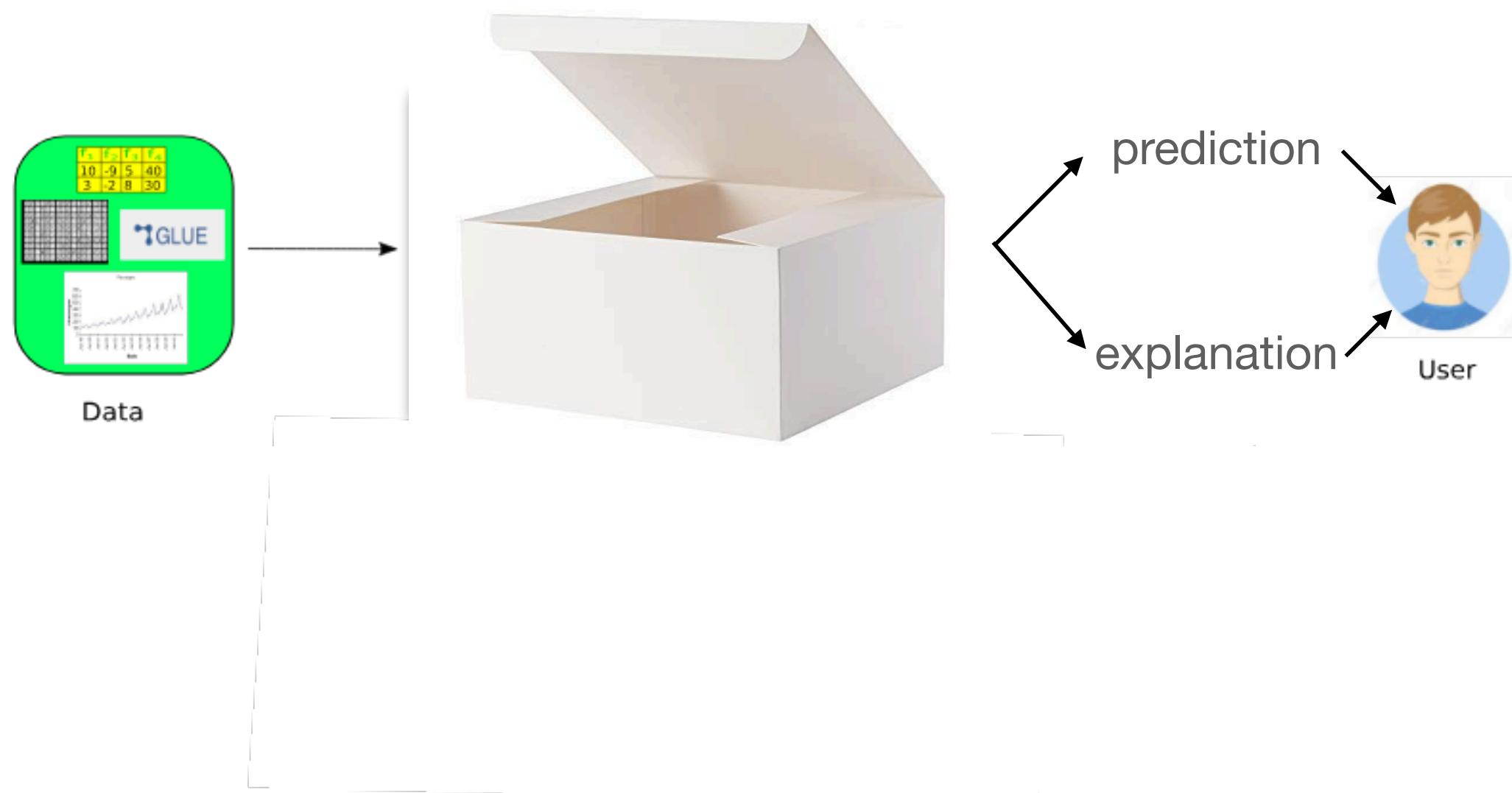


Why's Explainable AI?

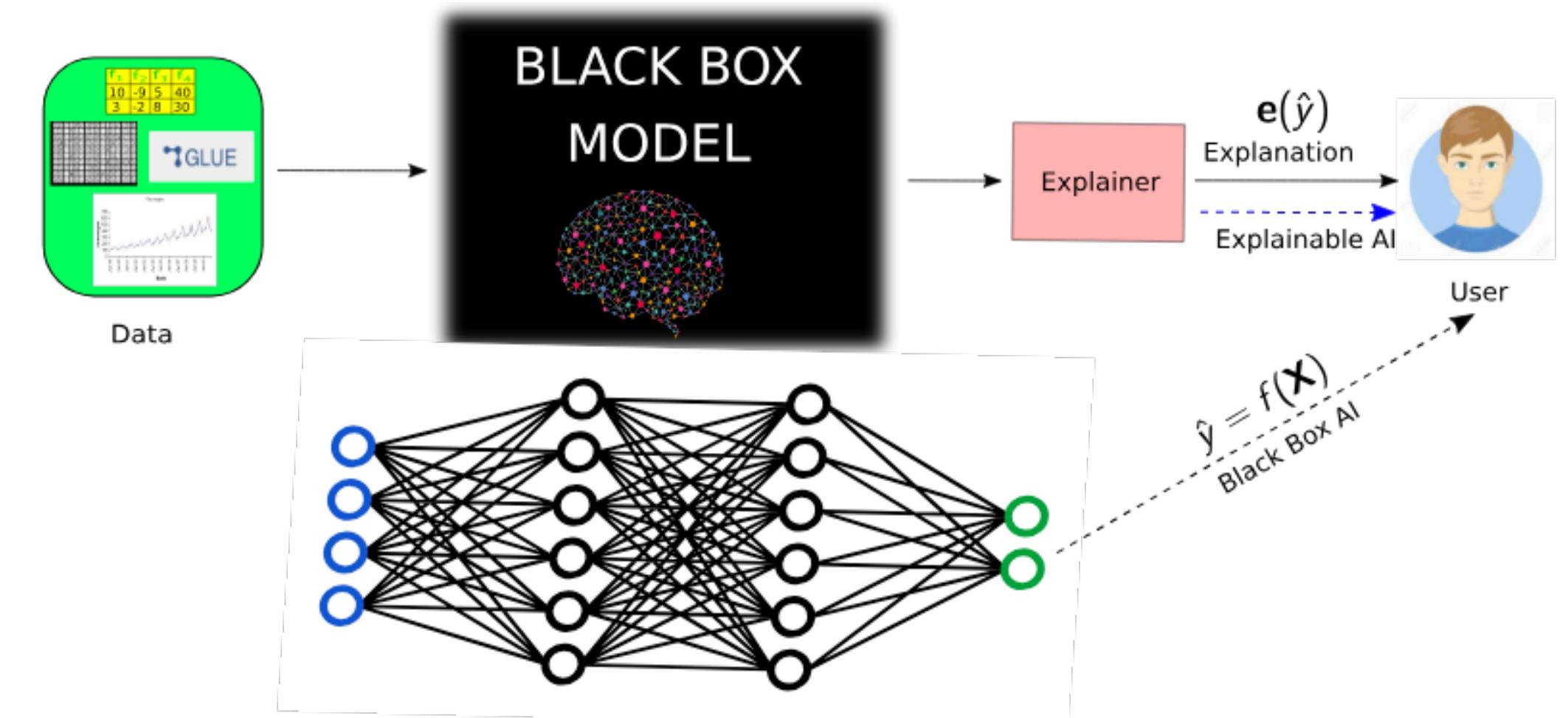


Explanations

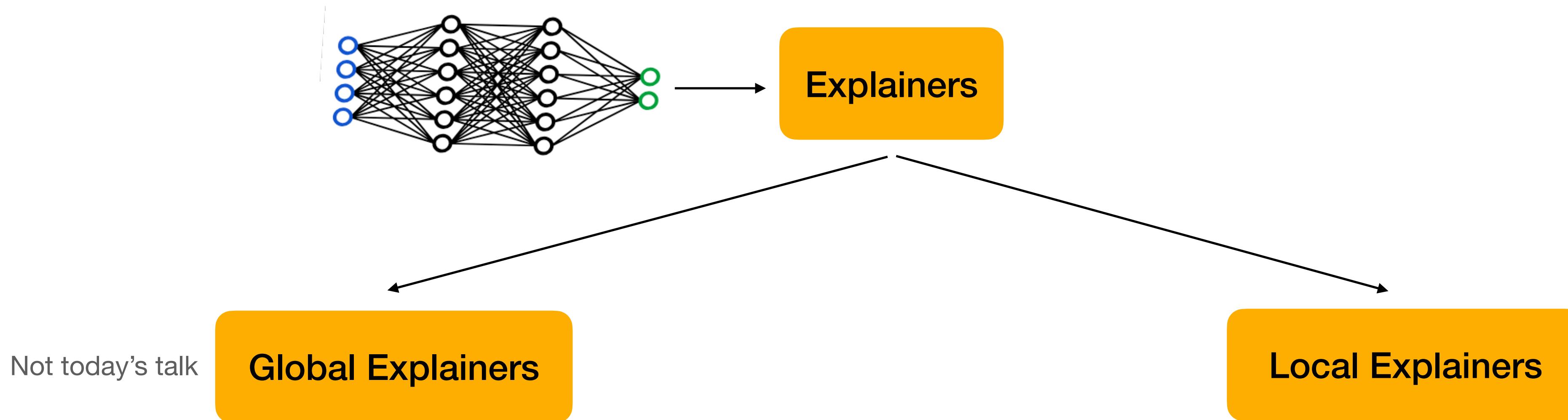
Ante-hoc Explanations



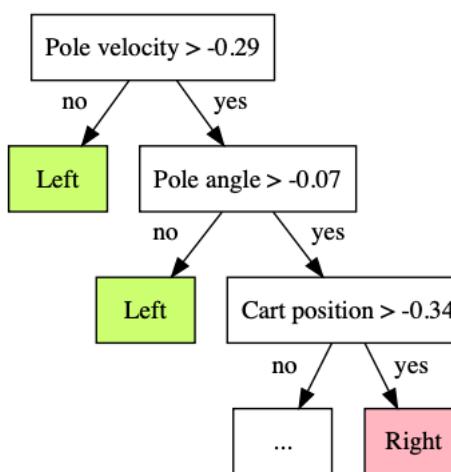
Post-hoc Explanations



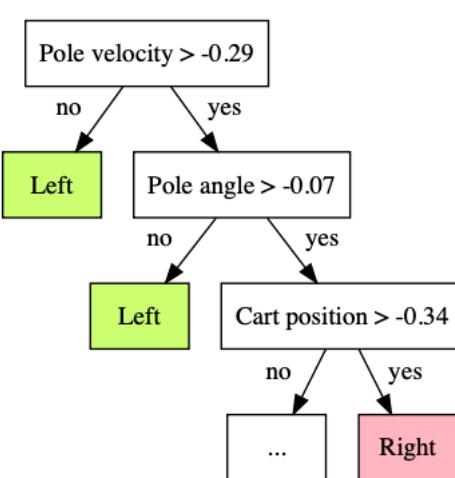
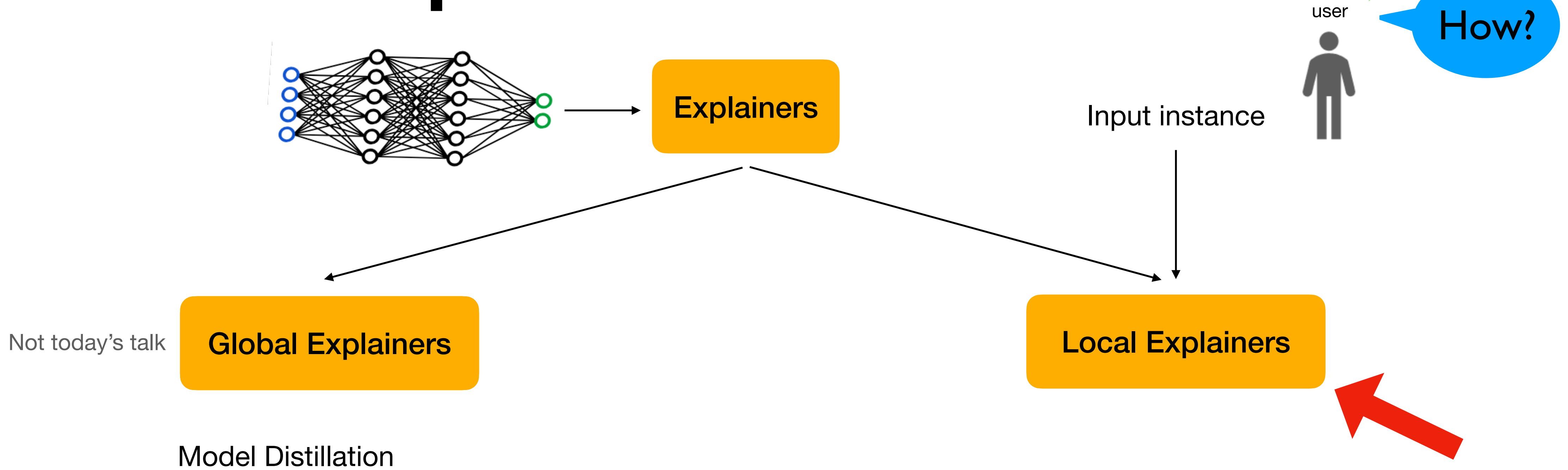
Post-hoc Explanations



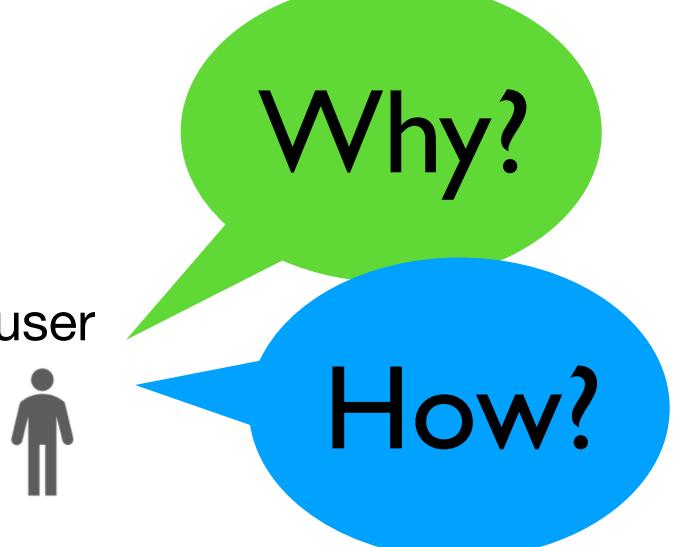
Model Distillation



Post-hoc Explanations



Local Explanations

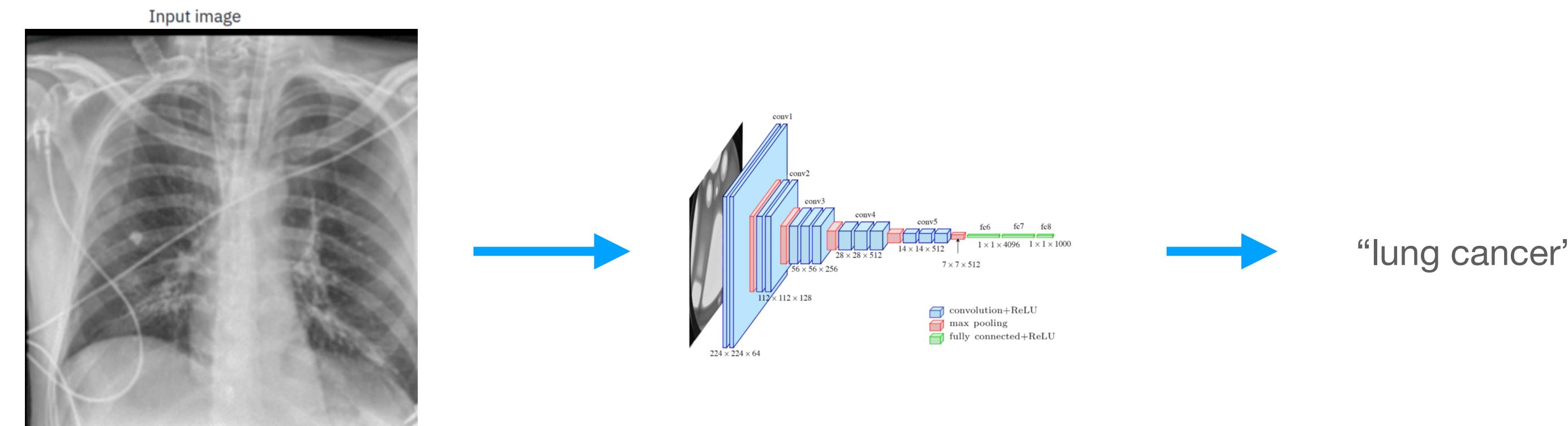


Local Explanations

Why?

user
 |
 |

Which features are important to make such a decision?
Saliency Map (e.g. GradCAM)



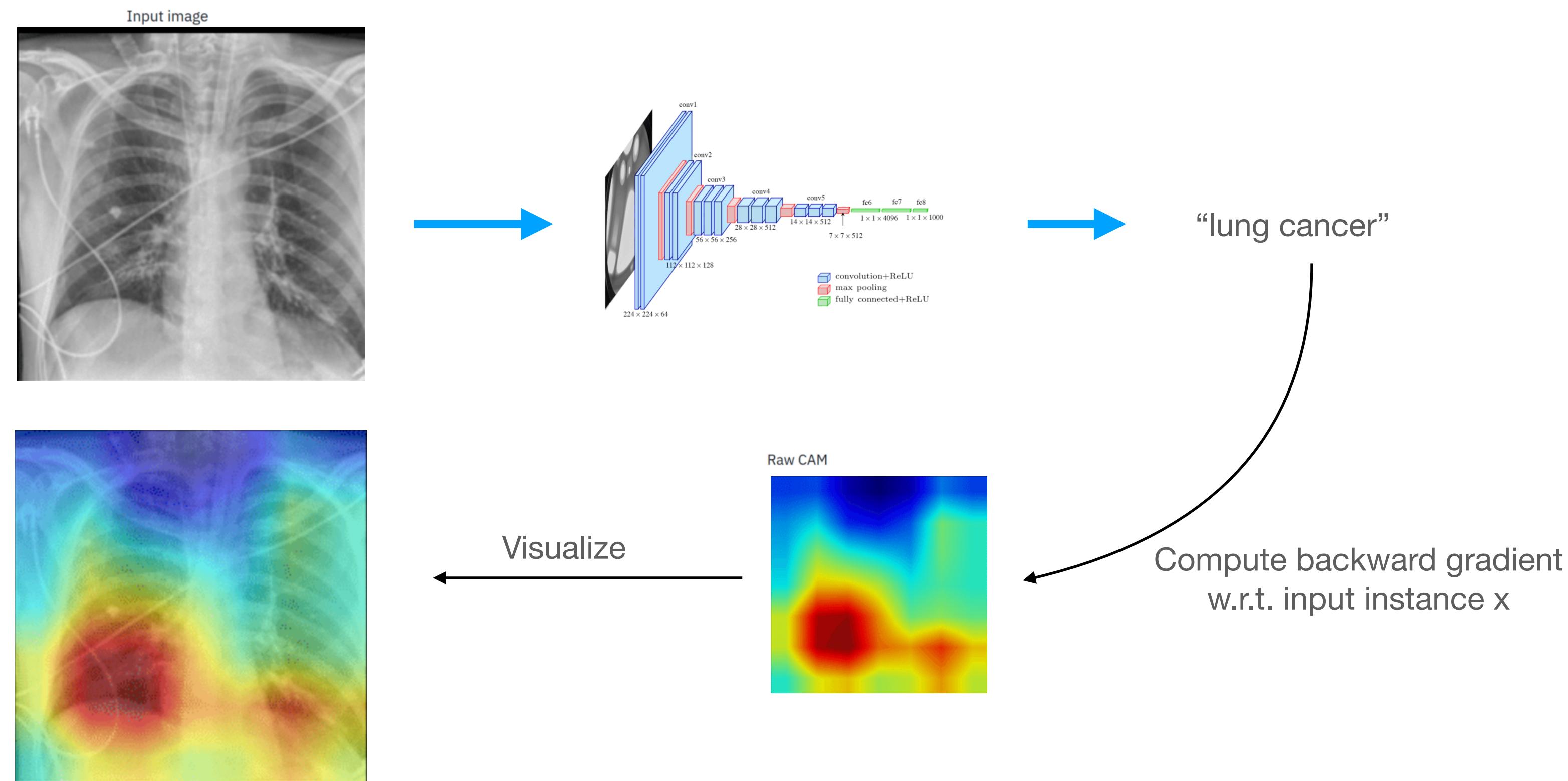
Local Explanations

Why?

user
 |
 |
 |

Which features are important to make such a decision?

Saliency Map (e.g. GradCAM)



Local Explanations

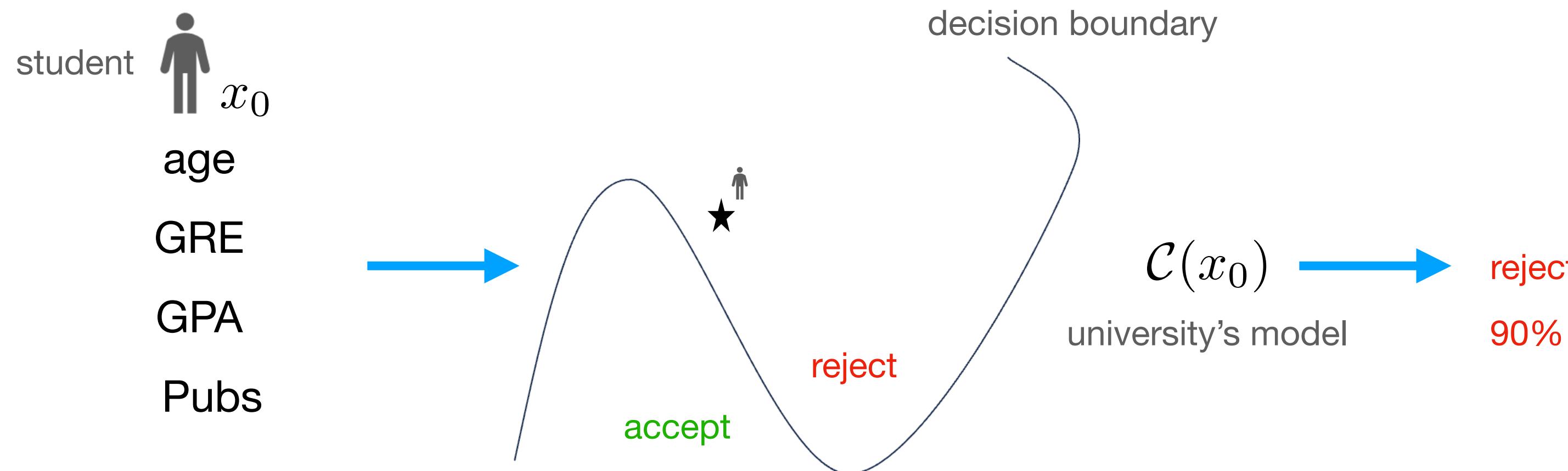
Why?

user
 |
 |
 |
 |
 |

Which features are important to make such a decision?

Saliency Map (e.g. GradCAM)

Feature-based Attribution Methods (e.g. LIME, SHAP)



Local Explanations

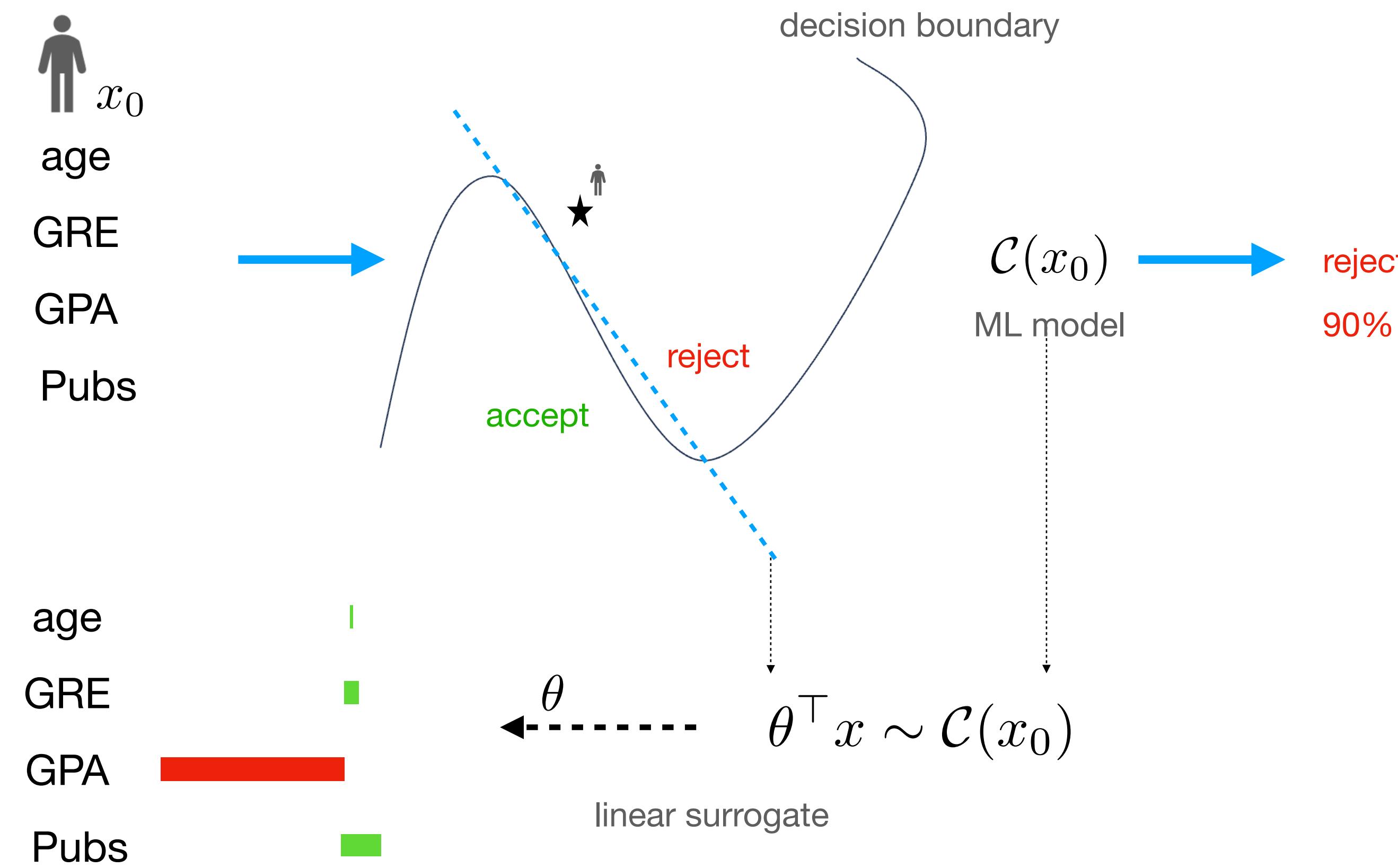
Why?

user
└─┘

Which features are important to make such a decision?

Saliency Map (e.g. GradCAM)

Feature-based Attribution Methods (e.g. LIME, SHAP)

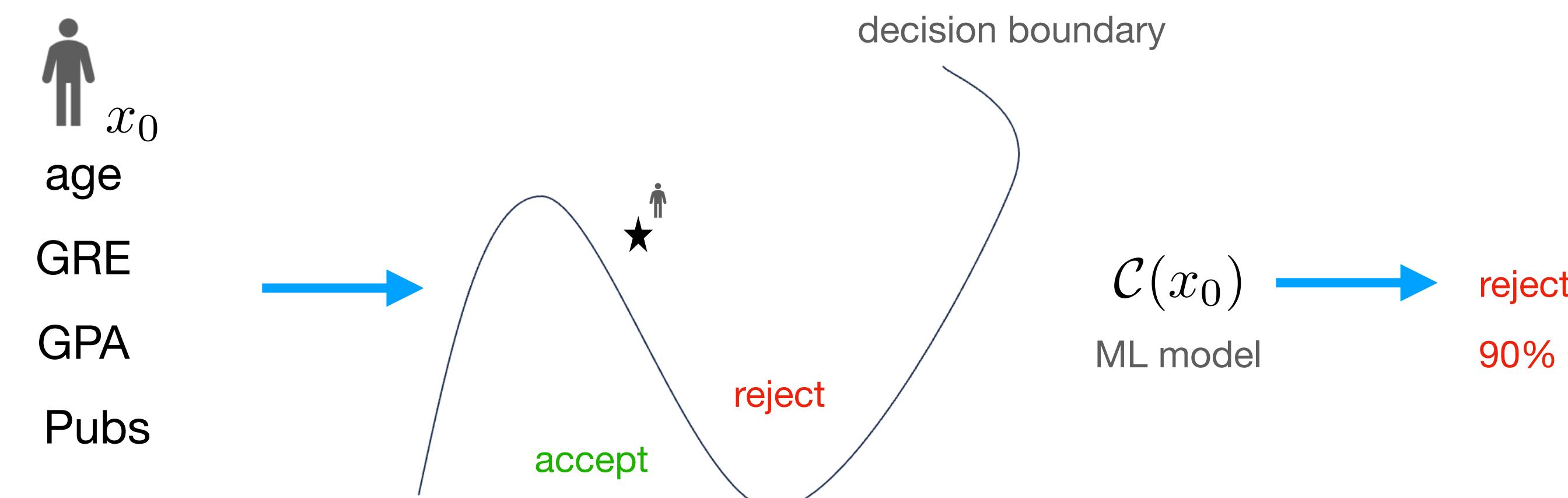


Local Explanations

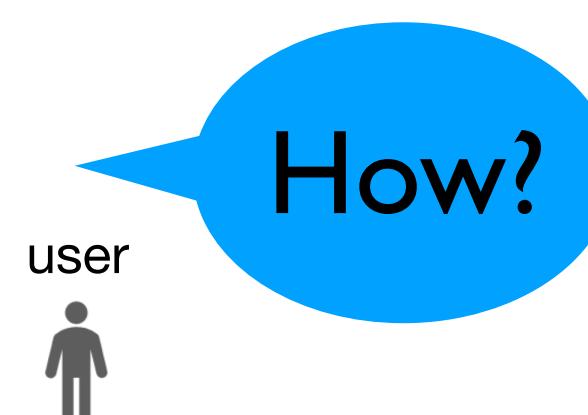
How?

How can I improve to get a favorable result?

user
 └─ person icon



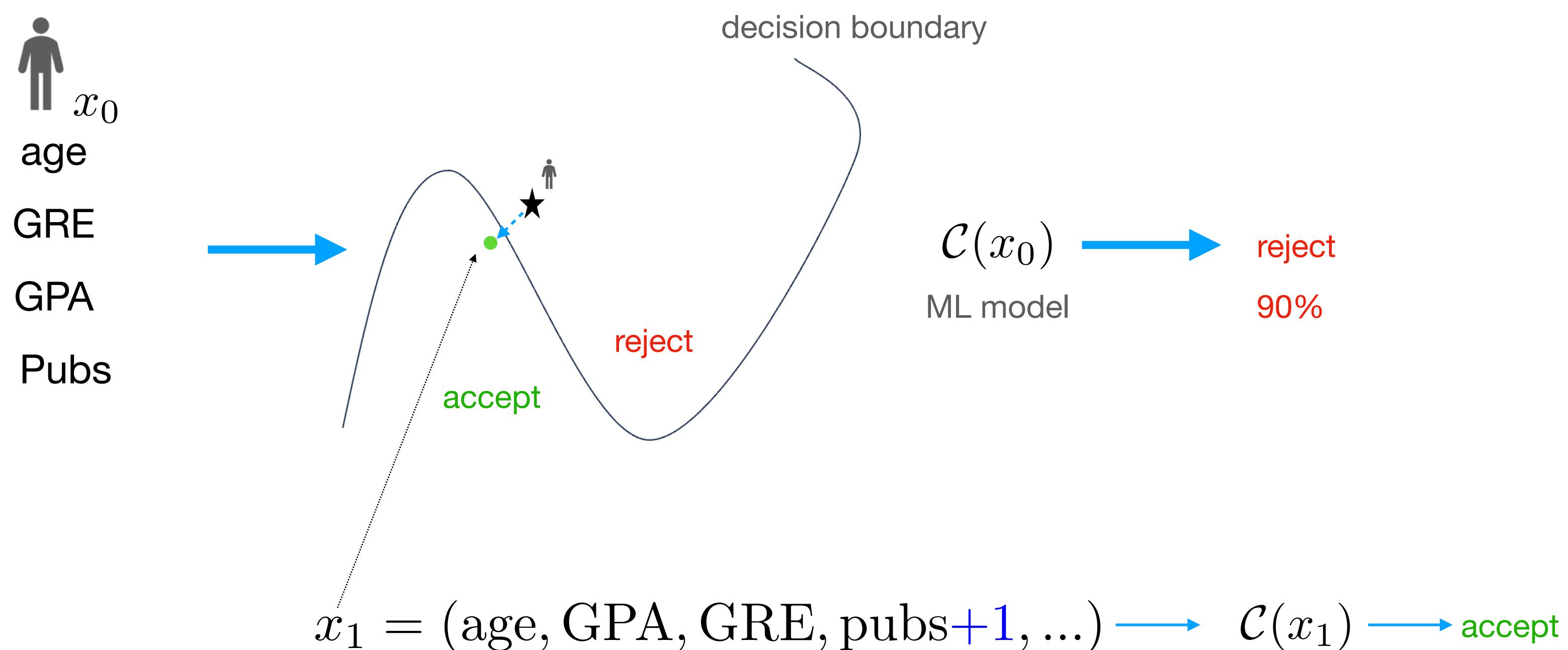
Local Explanations



How?

How can I improve to get a favorable result?

Counterfactual Explanations (a.k.a. Algorithmic Recourse): examples that can flip the model's outcome



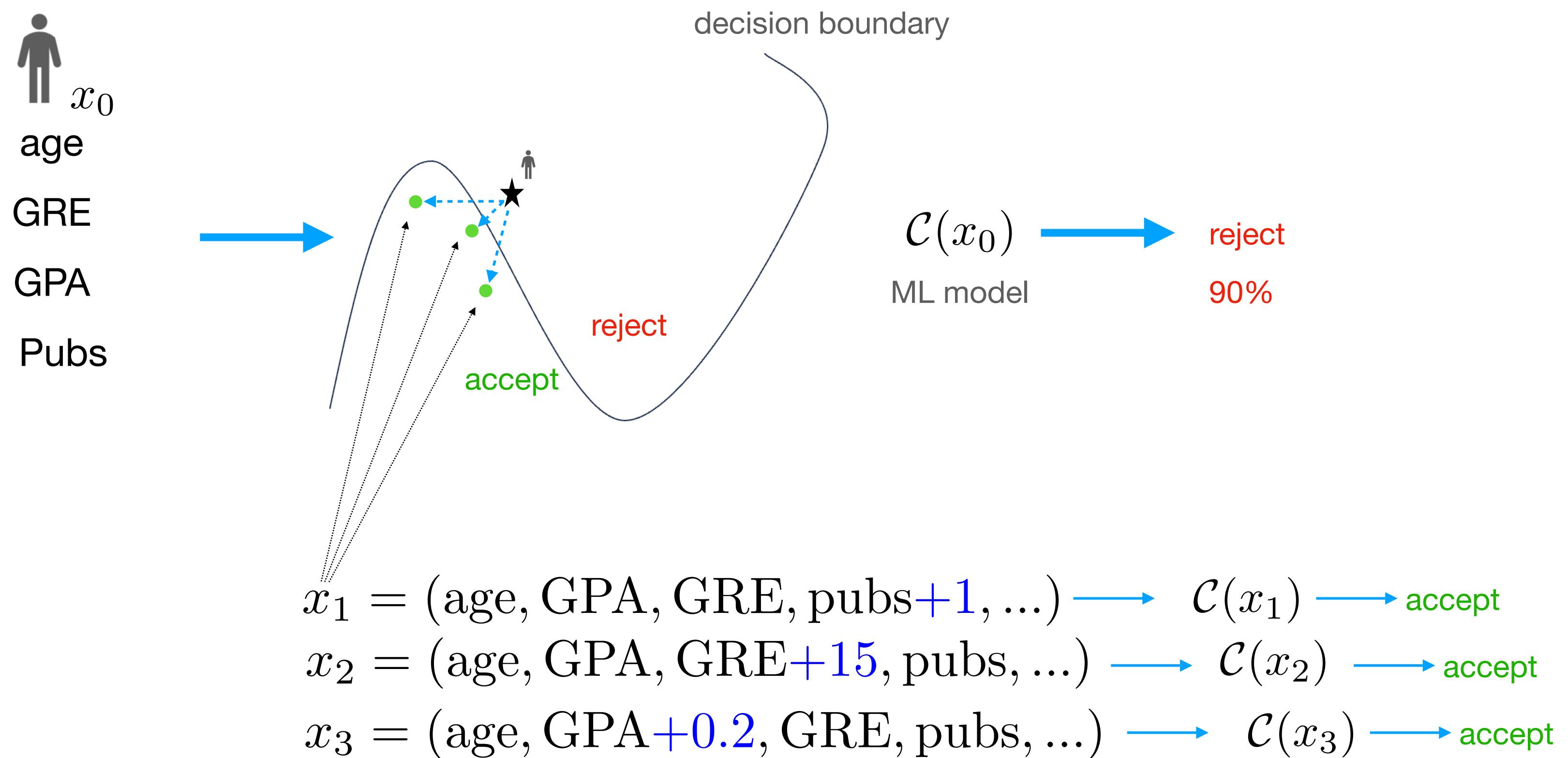
Local Explanations

How?

user
└─┘

How can I improve to get a favorable result?

Counterfactual Plan: a set of counterfactual explanations



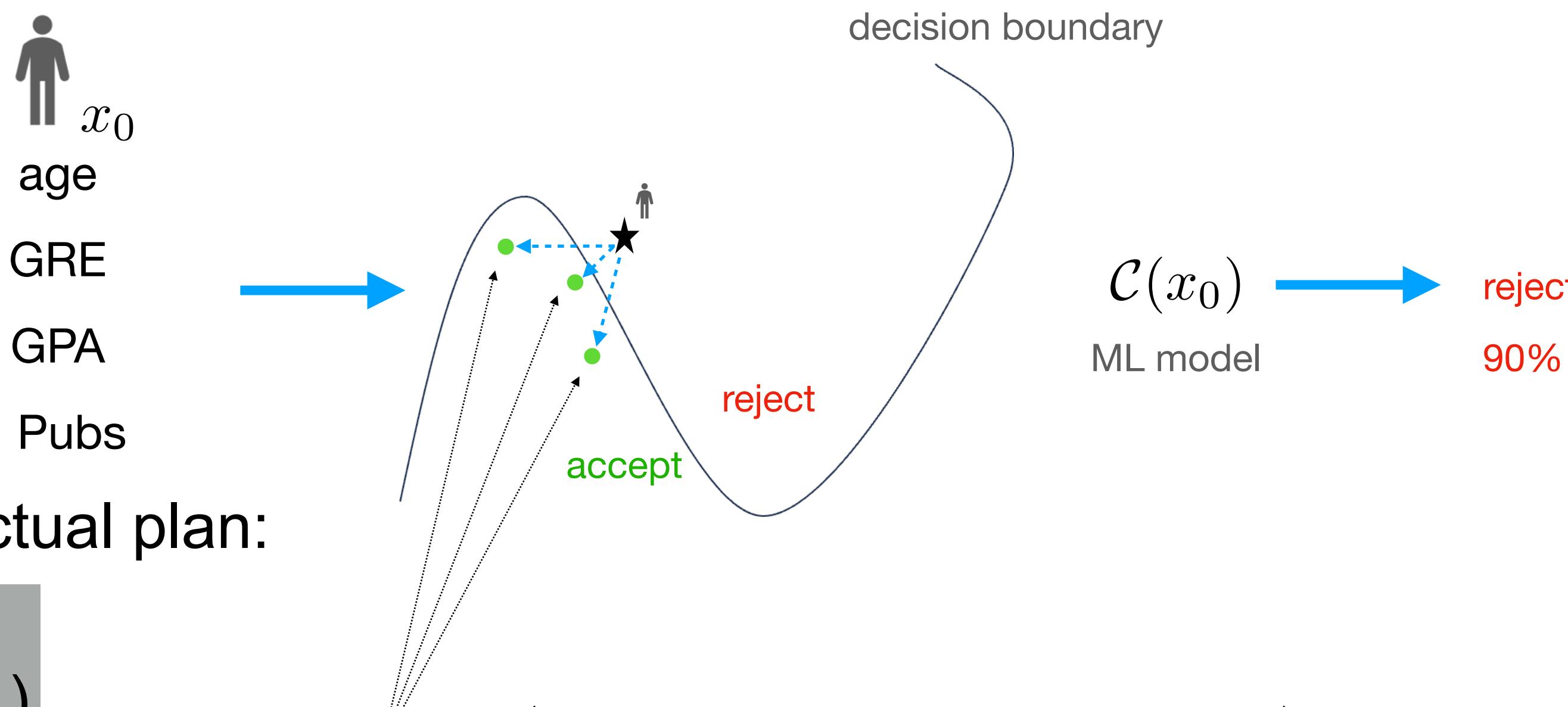
Local Explanations

How?

user
 |
 |
 |
 |
 |

How can I improve to get a favorable result?

Counterfactual Plan: a set of counterfactual explanations



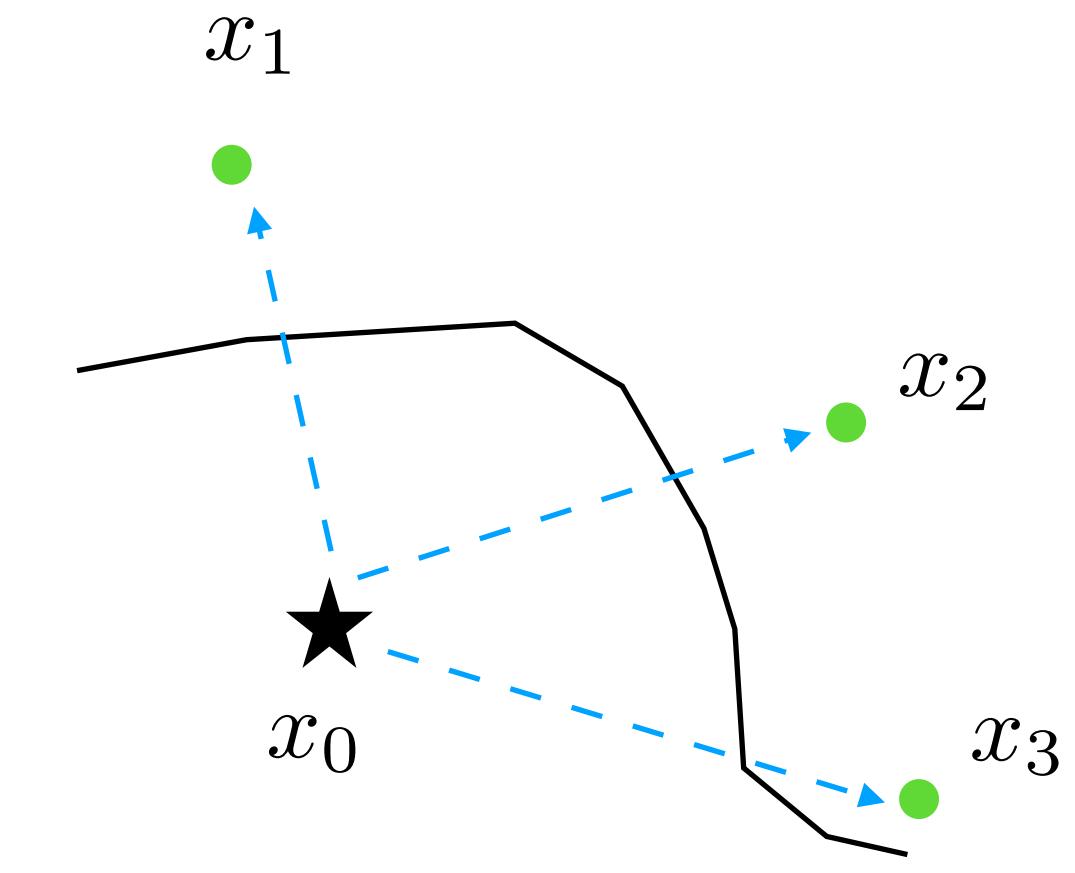
Criteria to construct counterfactual plan:

1. Validity
2. Proximity (closeness to x_0)
3. Diversity
4. Actionability

$$x_1 = (\text{age}, \text{GPA}, \text{GRE}, \text{pubs}+1, \dots) \rightarrow \mathcal{C}(x_1) \rightarrow \text{accept}$$
$$x_2 = (\text{age}, \text{GPA}, \text{GRE}+15, \text{pubs}, \dots) \rightarrow \mathcal{C}(x_2) \rightarrow \text{accept}$$
$$x_3 = (\text{age}, \text{GPA}+0.2, \text{GRE}, \text{pubs}, \dots) \rightarrow \mathcal{C}(x_3) \rightarrow \text{accept}$$

Counterfactual Plan - Construction

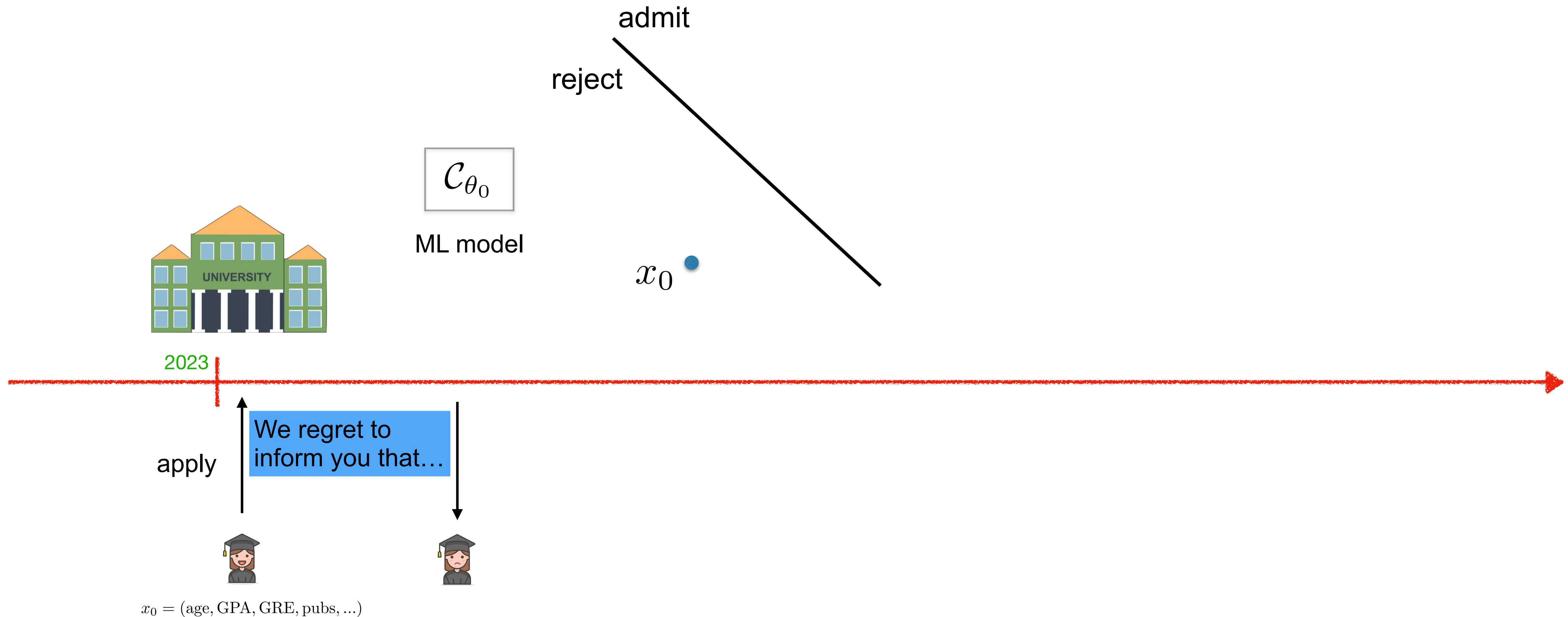
- Mixed-Integer Programming
 - Russell (2019) (for linear models)
 - Similar to Ustun (2018)
 - Promote diversity through diverse sets of constraints
- Multi-objective Evolutionary Algorithm
 - Dandl et al. (2020)
 - Universal but slow
 - Suitable to promote actionability
- Iterative, Gradient-based Optimization
 - Mothilal et al. (2020) (DICE)
 - A weighted sum of three objectives: validity, proximity, and diversity.
 - Solved by projected gradient descent



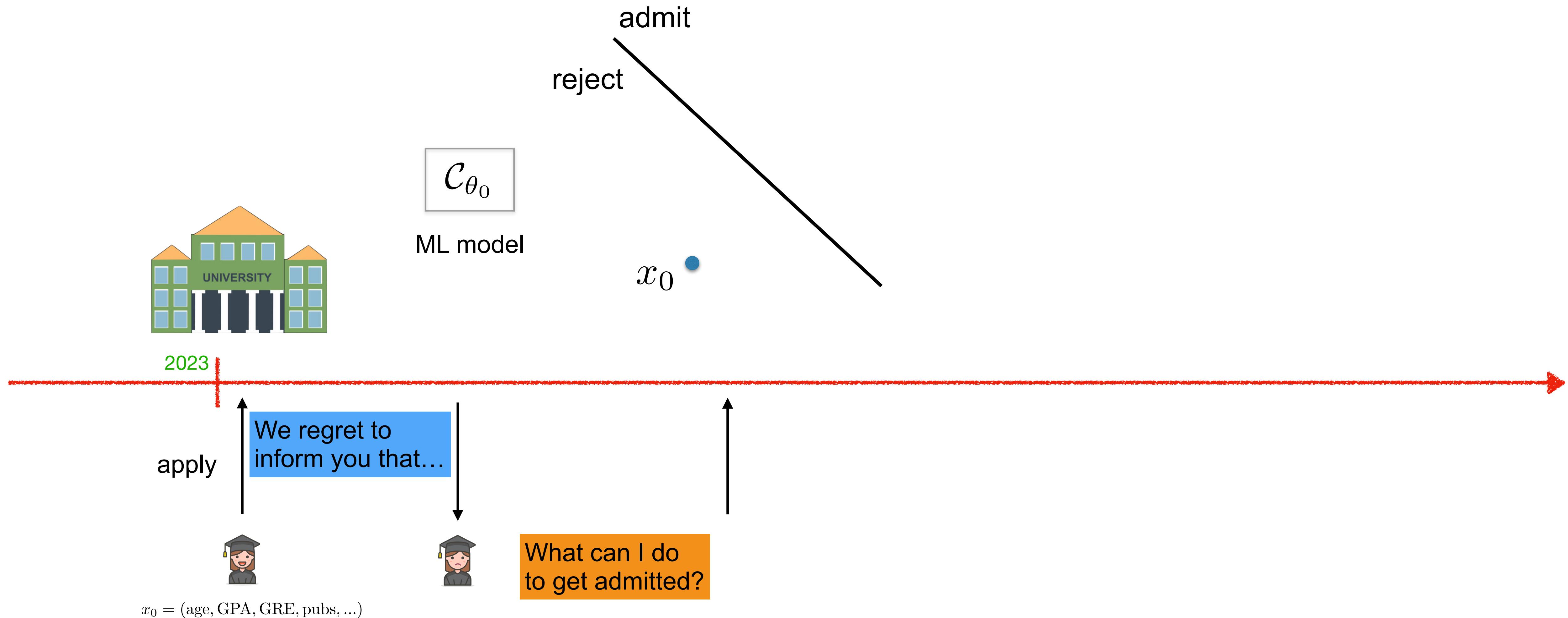
Counterfactual Explanations under Distributional Ambiguity

ICLR'22

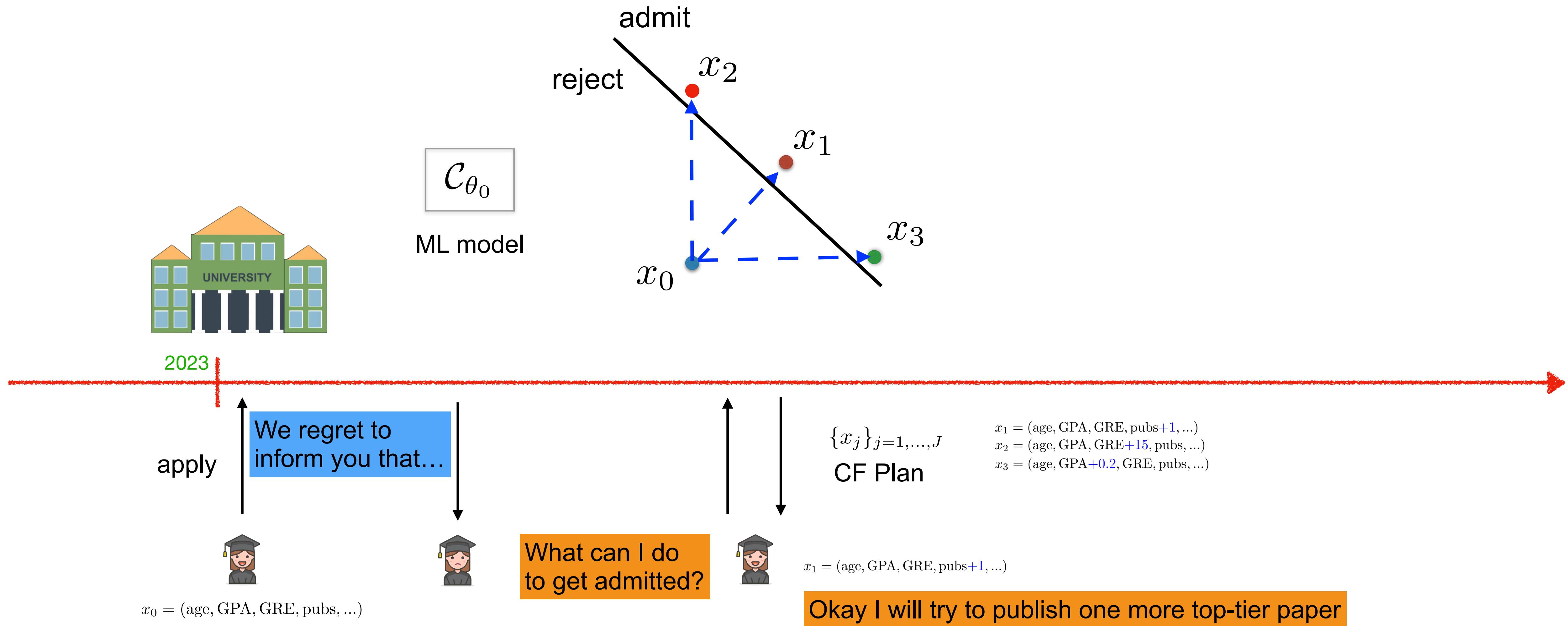
Counterfactual Plan under Model Shifts



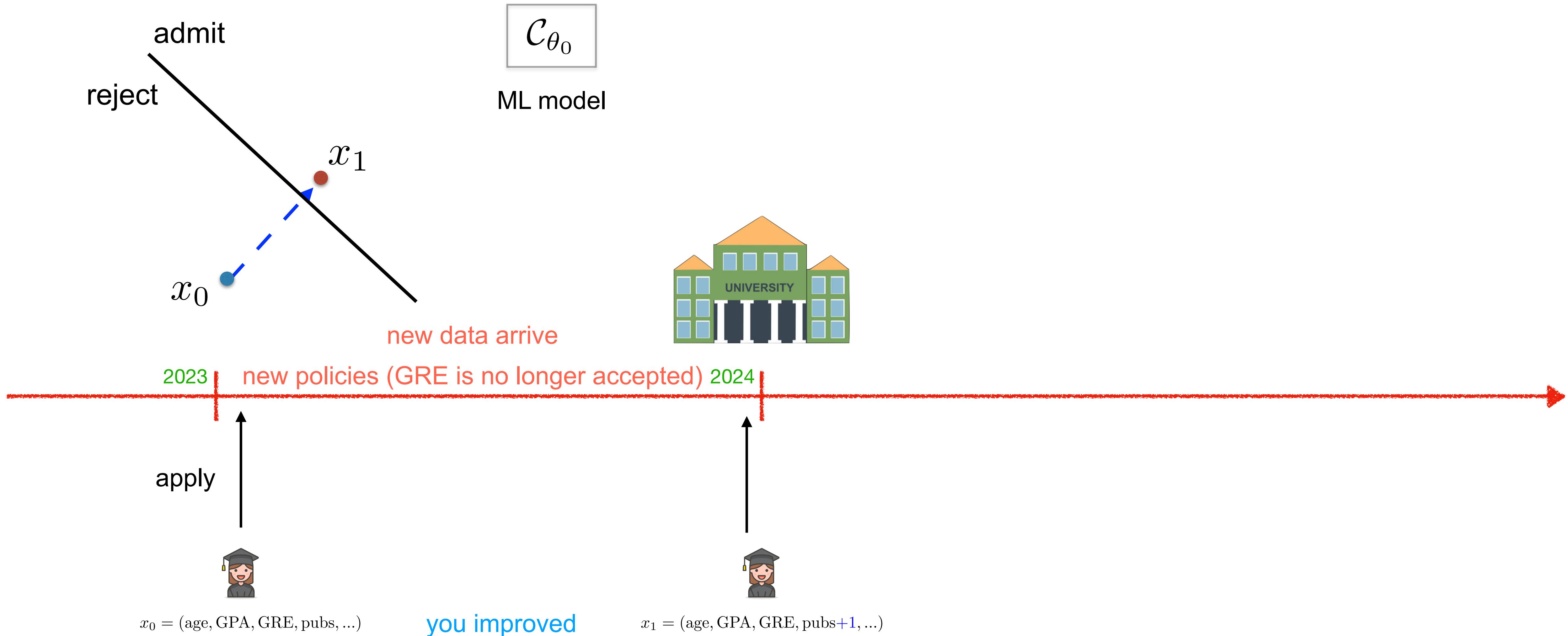
Counterfactual Plan under Model Shifts



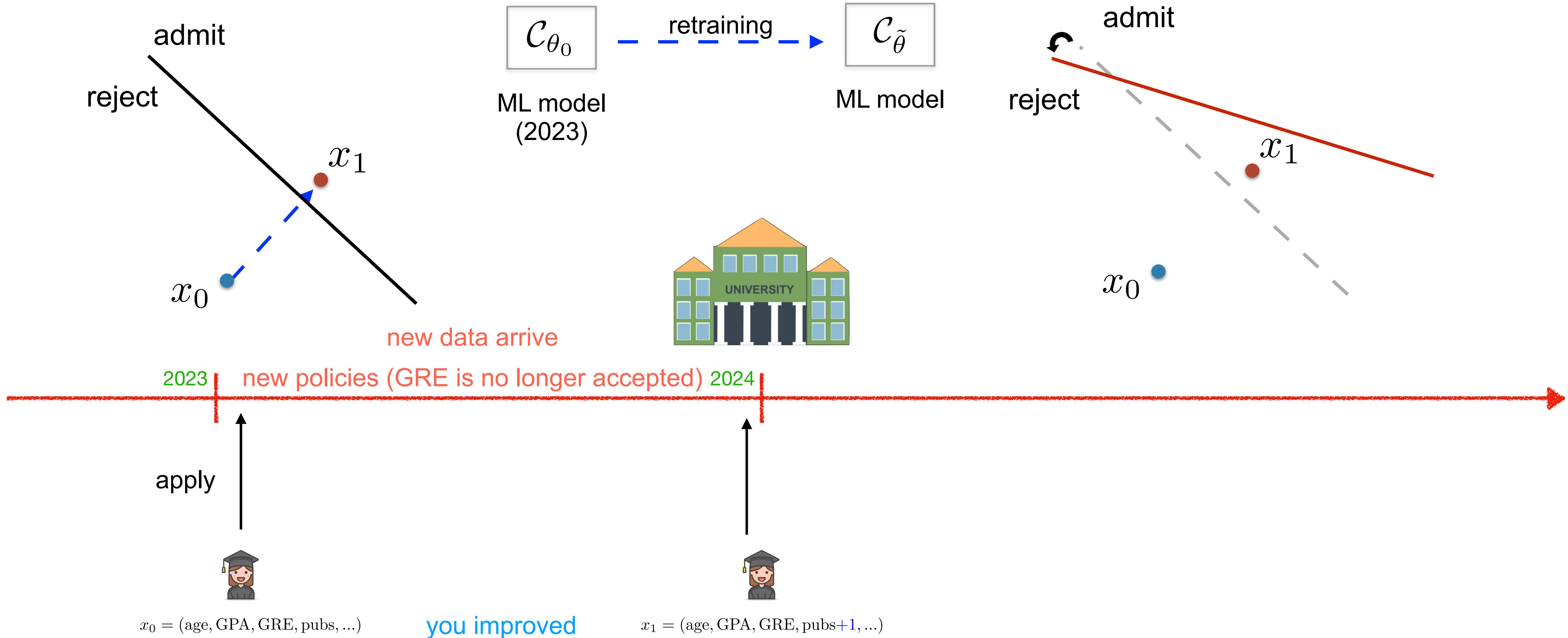
Counterfactual Plan under Model Shifts



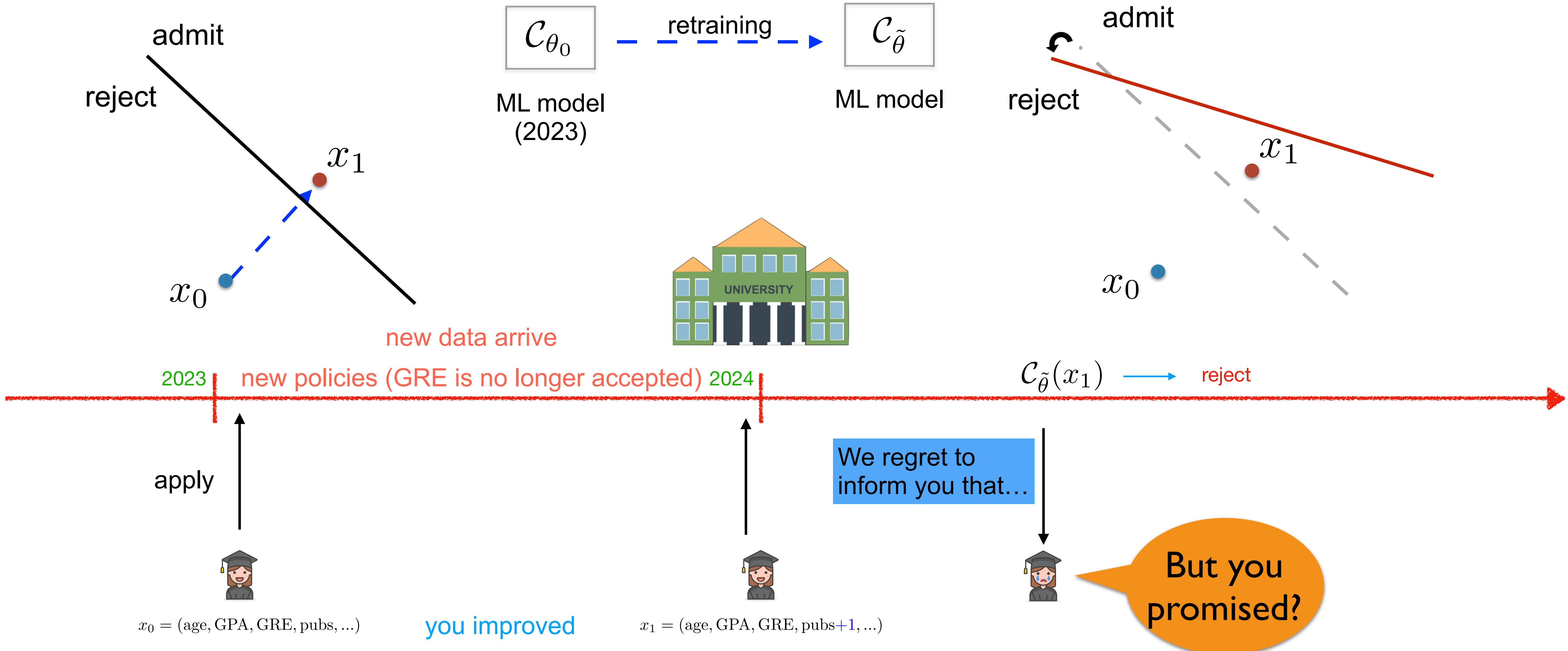
Counterfactual Plan under Model Shifts



Counterfactual Plan under Model Shifts



Counterfactual Plan under Model Shifts



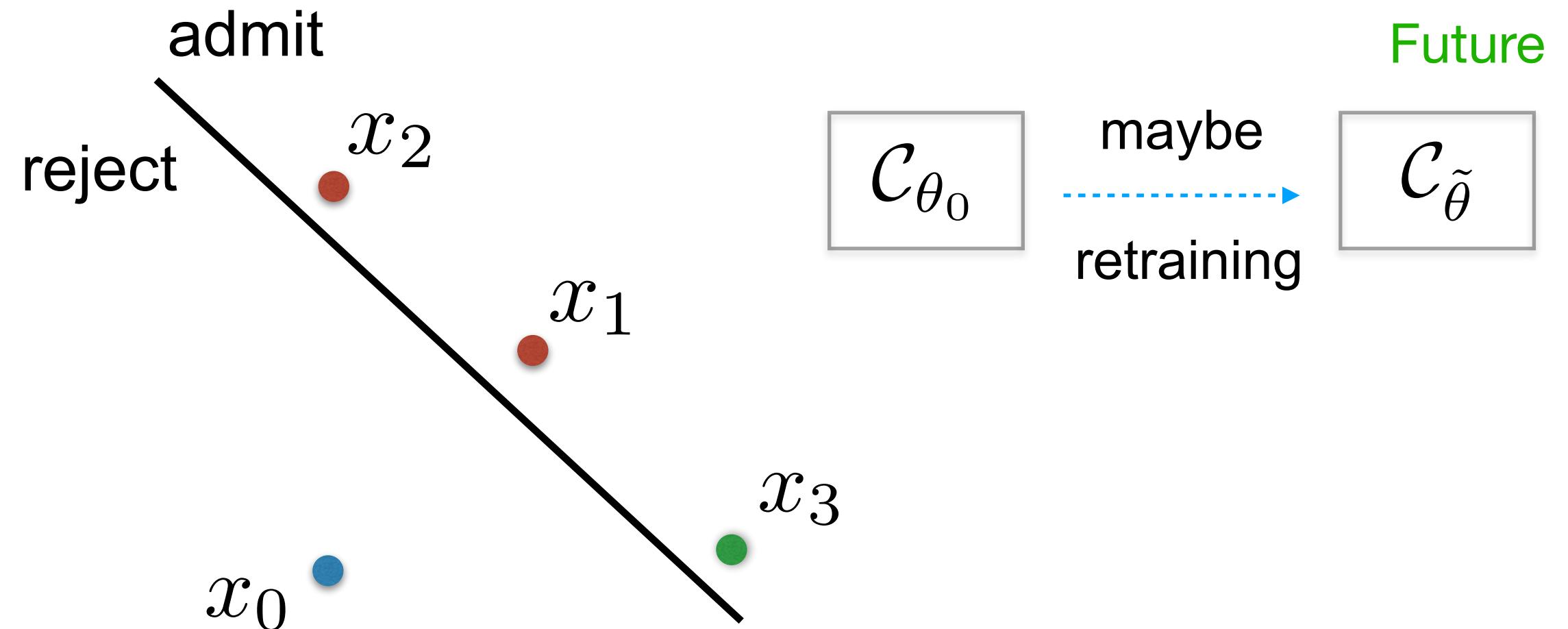
Settings

We are given:

- An ML model (assumed to be linear)

$$\mathcal{C}_{\theta_0}(x) = \begin{cases} 0 \Rightarrow \text{reject (unfavorable)} \\ 1 \Rightarrow \text{accept (favorable)} \end{cases}$$

- A CF plan $\{x_j\}_{j=1,\dots,J}$



Definition: A plan $\{x_j\}_{j=1,\dots,J}$ is valid w.r.t. $\tilde{\theta}$ if

$$\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j = 1, \dots, J$$



$x_0 = (\text{age}, \text{GPA}, \text{GRE}, \text{pubs}, \dots)$

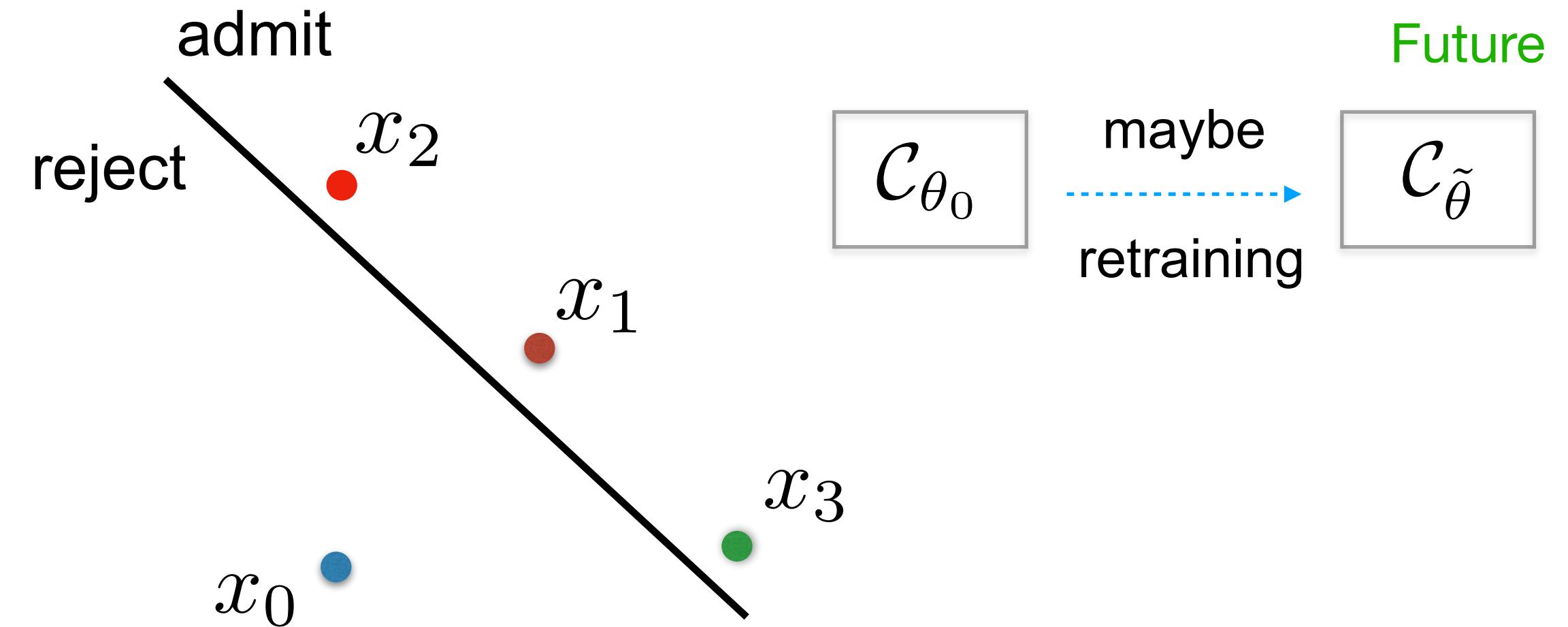
Settings

We are given:

- An ML model (assumed to be linear)

$$\mathcal{C}_{\theta_0}(x) = \begin{cases} 0 \Rightarrow \text{reject (unfavorable)} \\ 1 \Rightarrow \text{accept (favorable)} \end{cases}$$

- A CF plan $\{x_j\}_{j=1,\dots,J}$



Q1. Given $\{x_j\}$, what is the probability that the plan is valid (w.r.t. $\tilde{\theta}$)?

Q2. Can we construct a counterfactual plan that is robust to model shifts?



$x_0 = (\text{age}, \text{GPA}, \text{GRE}, \text{pubs}, \dots)$

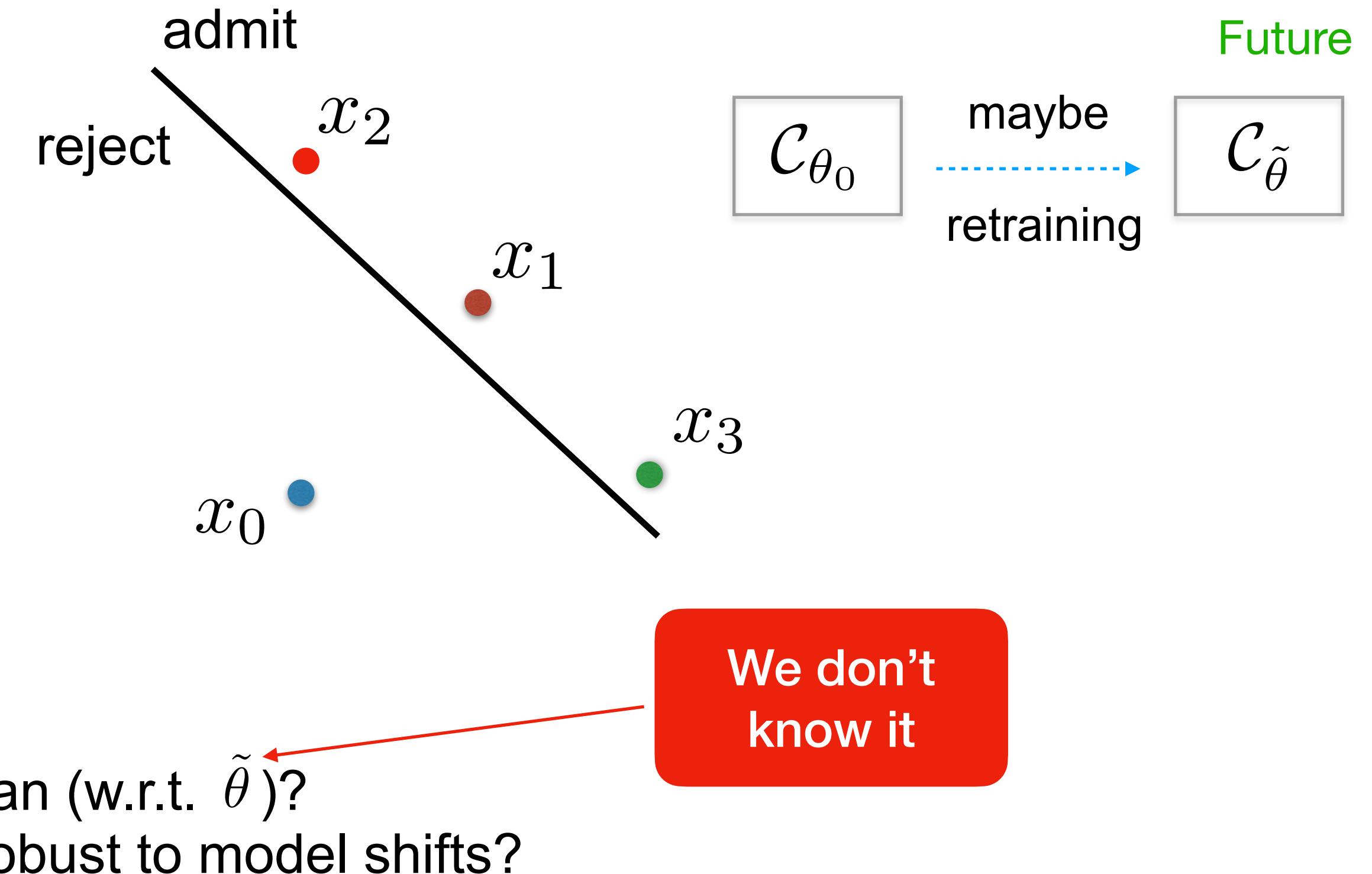
Settings

We are given:

- An ML model (assumed to be linear)

$$\mathcal{C}_{\theta_0}(x) = \begin{cases} 0 \Rightarrow \text{reject (unfavorable)} \\ 1 \Rightarrow \text{accept (favorable)} \end{cases}$$

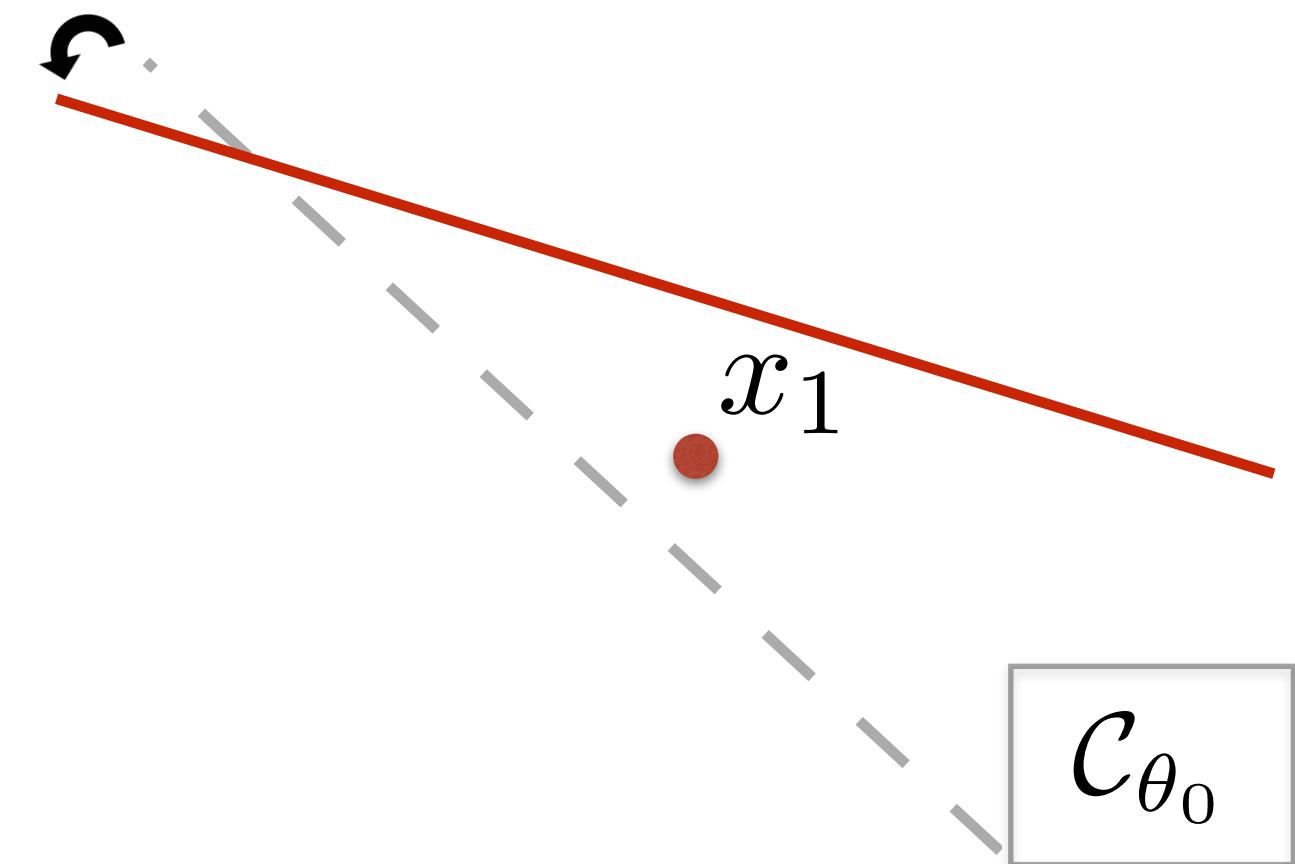
- A CF plan $\{x_j\}_{j=1,\dots,J}$



$x_0 = (\text{age}, \text{GPA}, \text{GRE}, \text{pubs}, \dots)$

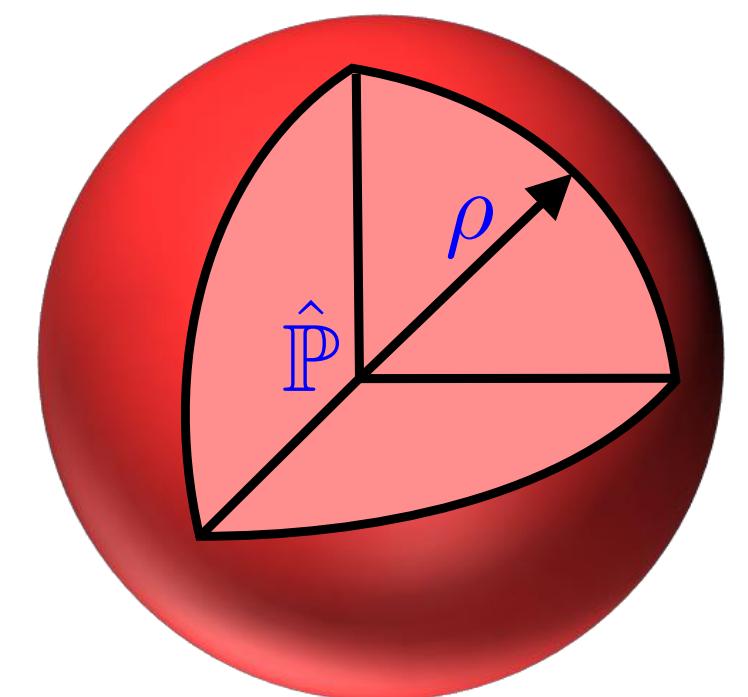
Uncertainty Modeling

Assume that the *future* parameter $\tilde{\theta} \sim \mathbb{P} \in \mathbb{B} = \{\mathbb{P} : \mathbb{G}(\mathbb{P}, \hat{\mathbb{P}}) \leq \rho\}$



Ambiguity set \mathbb{B}

- » neighborhood around $\hat{\mathbb{P}} \sim (\hat{\mu}, \hat{\Sigma})$
- » aims to contain the true distribution \mathbb{P} with a high probability



$\hat{\mathbb{P}} \sim (\hat{\mu}, \hat{\Sigma})$: estimated distribution (e.g. $\hat{\mu} = \theta_0, \hat{\Sigma} = I$)

$\mathbb{G}(\cdot, \cdot)$: distance between 2 distributions

ρ : uncertainty radius (controllable hyperparameter)

Uncertainty Modeling

$$\tilde{\theta} \sim \mathbb{P} \in \mathbb{B} = \{\mathbb{P} : G(\mathbb{P}, \hat{\mathbb{P}}) \leq \rho\}$$

$G(\cdot, \cdot)$: Gelbrich distance between two distributions

Definition: For $\mathbb{Q}_1 \sim (\mu_1, \Sigma_1)$ and $\mathbb{Q}_2 \sim (\mu_2, \Sigma_2)$

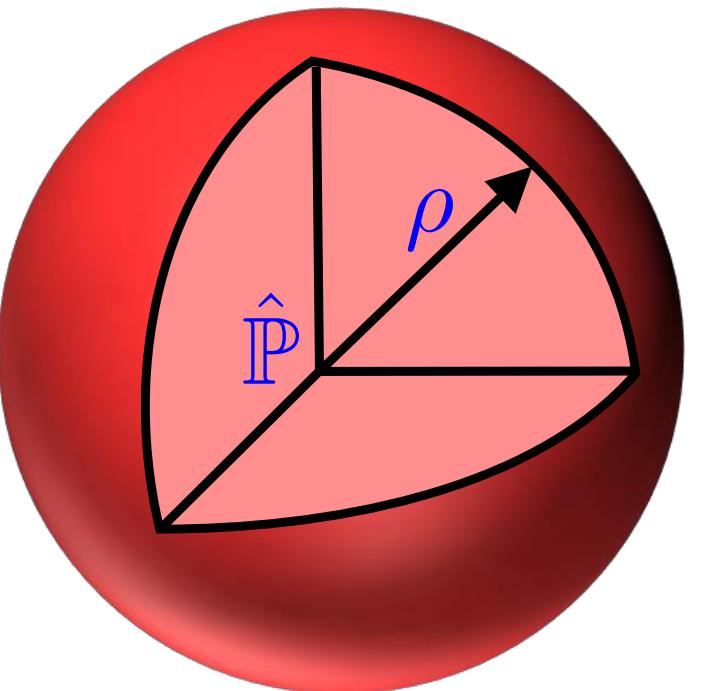
$$G(\mathbb{Q}_1, \mathbb{Q}_2) = \sqrt{\|\mu_1 - \mu_2\|_2^2 + \text{Tr}[\Sigma_1 + \Sigma_2 - 2(\Sigma_2^{\frac{1}{2}} \Sigma_1 \Sigma_2^{\frac{1}{2}})^{\frac{1}{2}}]}$$

Proposition (Connection to Type-2 Wasserstein distance):

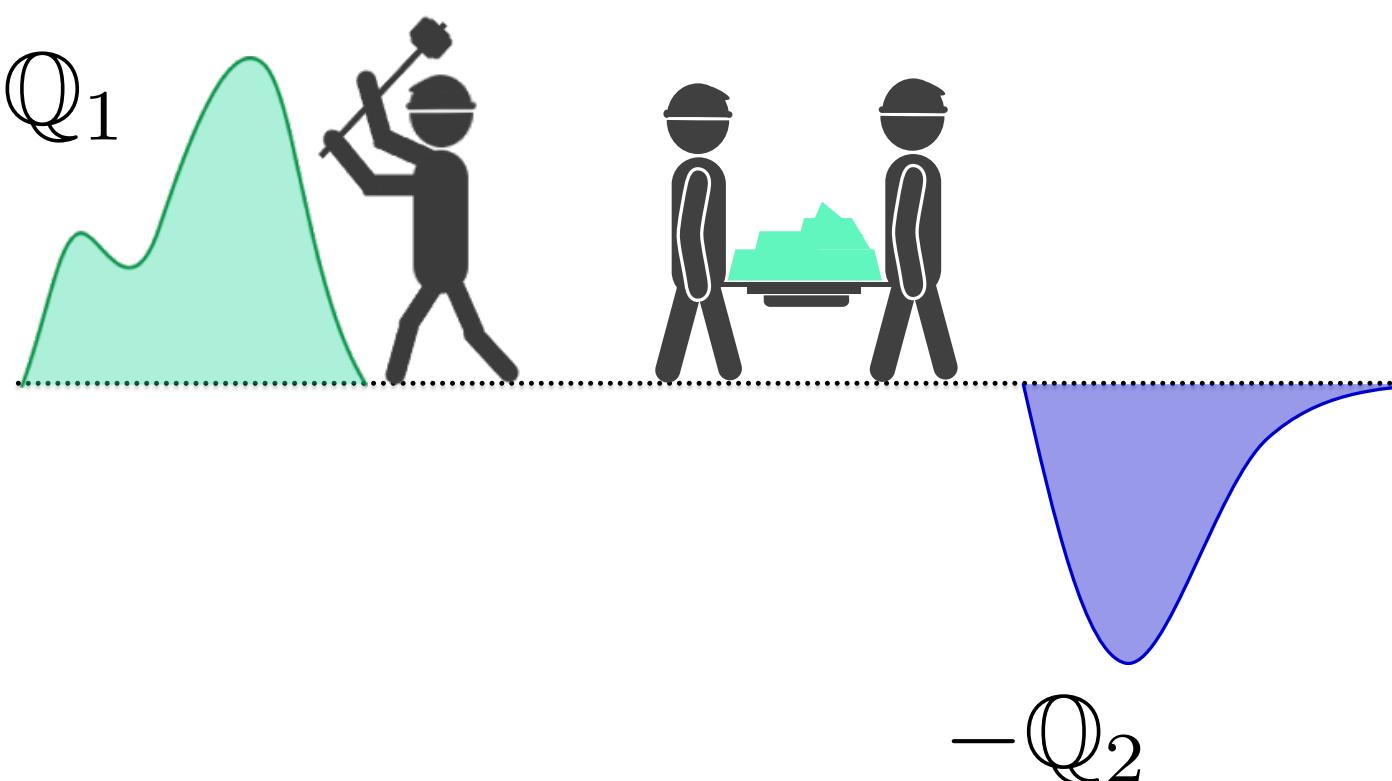
For $\mathbb{Q}_1 \sim (\mu_1, \Sigma_1)$ and $\mathbb{Q}_2 \sim (\mu_2, \Sigma_2)$

$$W(\mathbb{Q}_1, \mathbb{Q}_2) \geq G(\mathbb{Q}_1, \mathbb{Q}_2)$$

Gelbrich distance has a nice connection to Wasserstein distance but is much easier to compute.



$W^2(\mathbb{Q}_1, \mathbb{Q}_2)$ = minimum cost of moving \mathbb{Q}_1 to \mathbb{Q}_2



Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j)$$

True validity

Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j) = \mathbb{P}(\tilde{\theta}^\top x_j \geq 0 \quad \forall j)$$

True validity

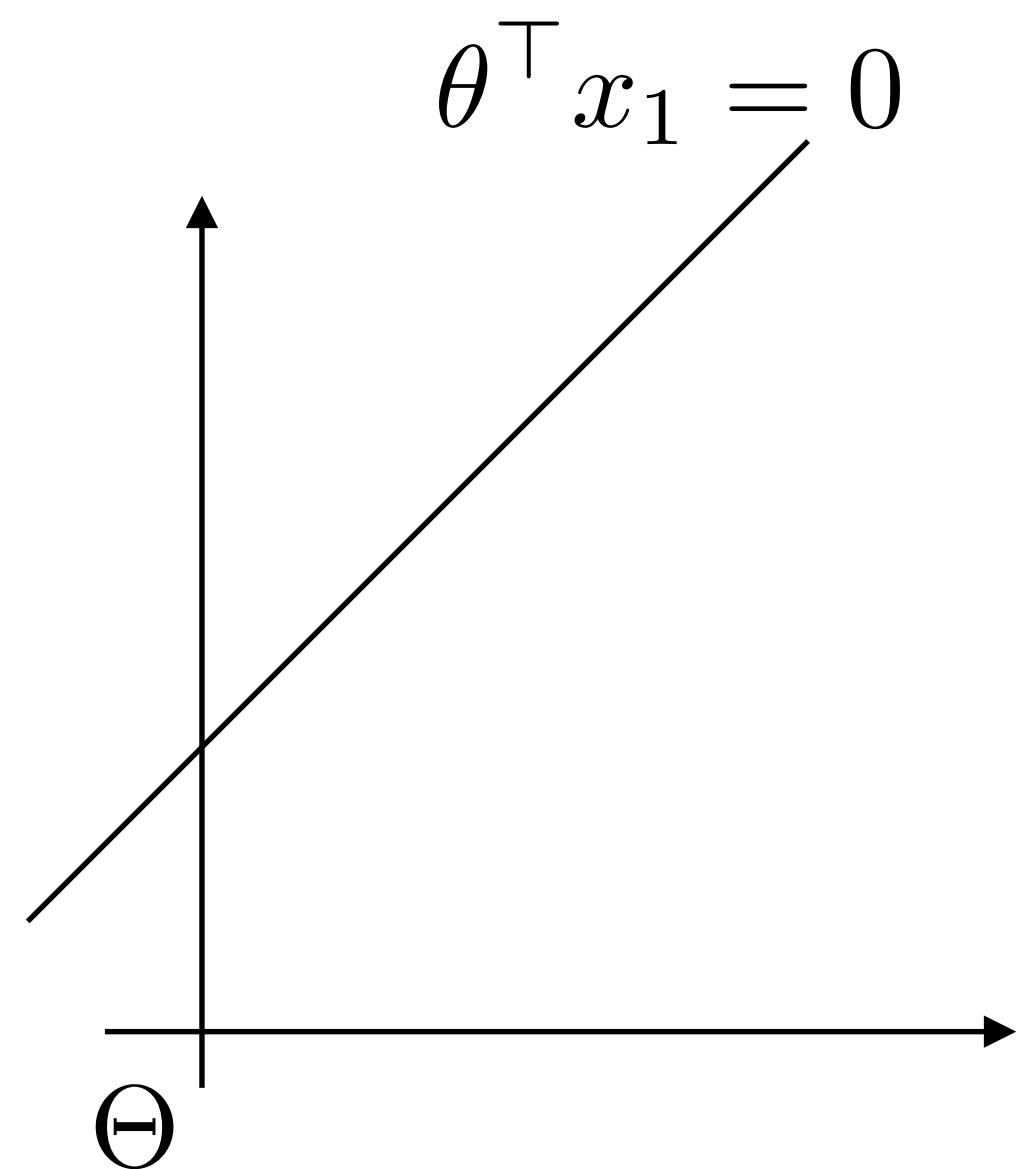
Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j) = \mathbb{P}(\tilde{\theta}^\top x_j \geq 0 \quad \forall j)$$

True validity



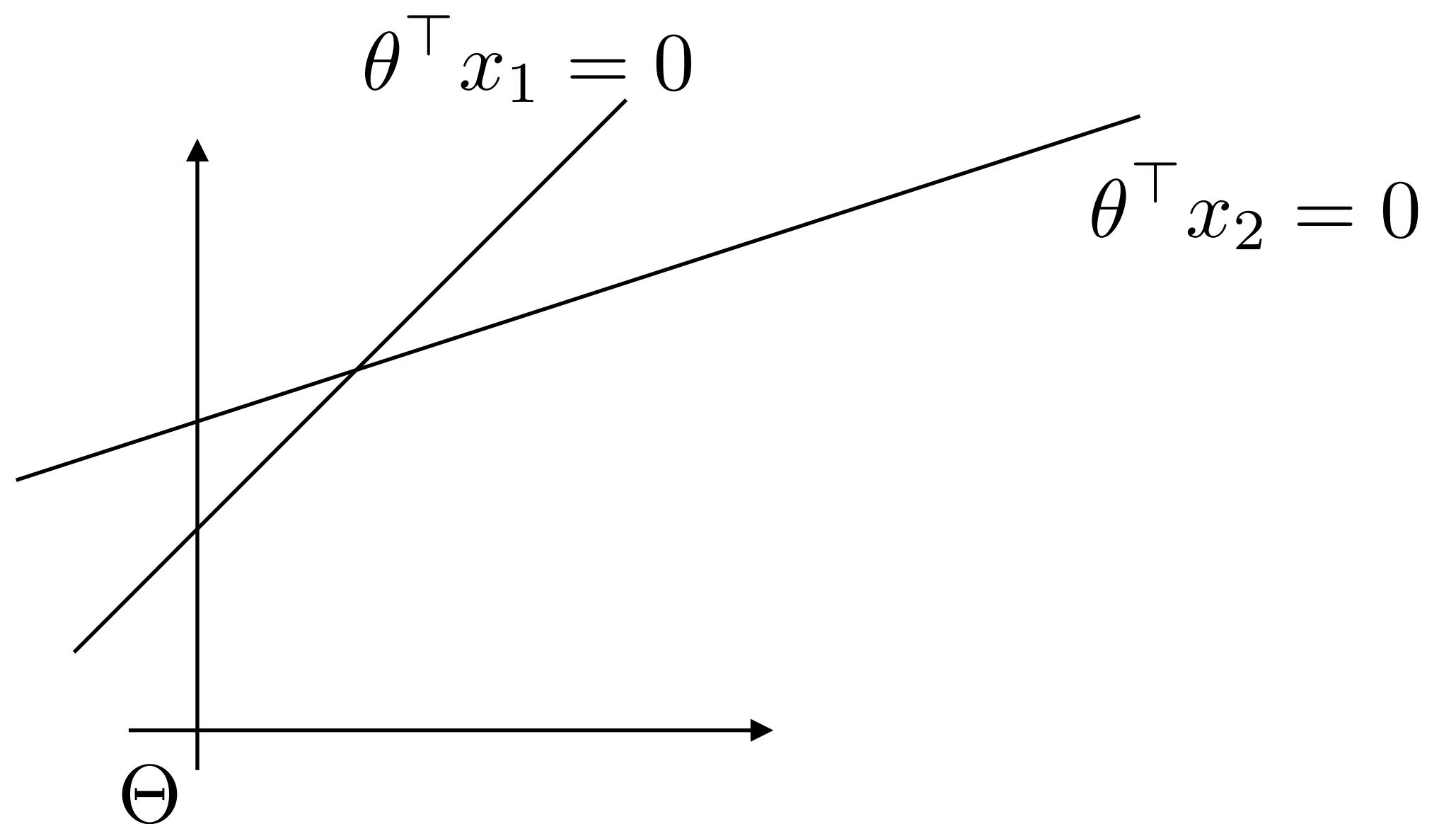
Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j) = \mathbb{P}(\tilde{\theta}^\top x_j \geq 0 \quad \forall j)$$

True validity



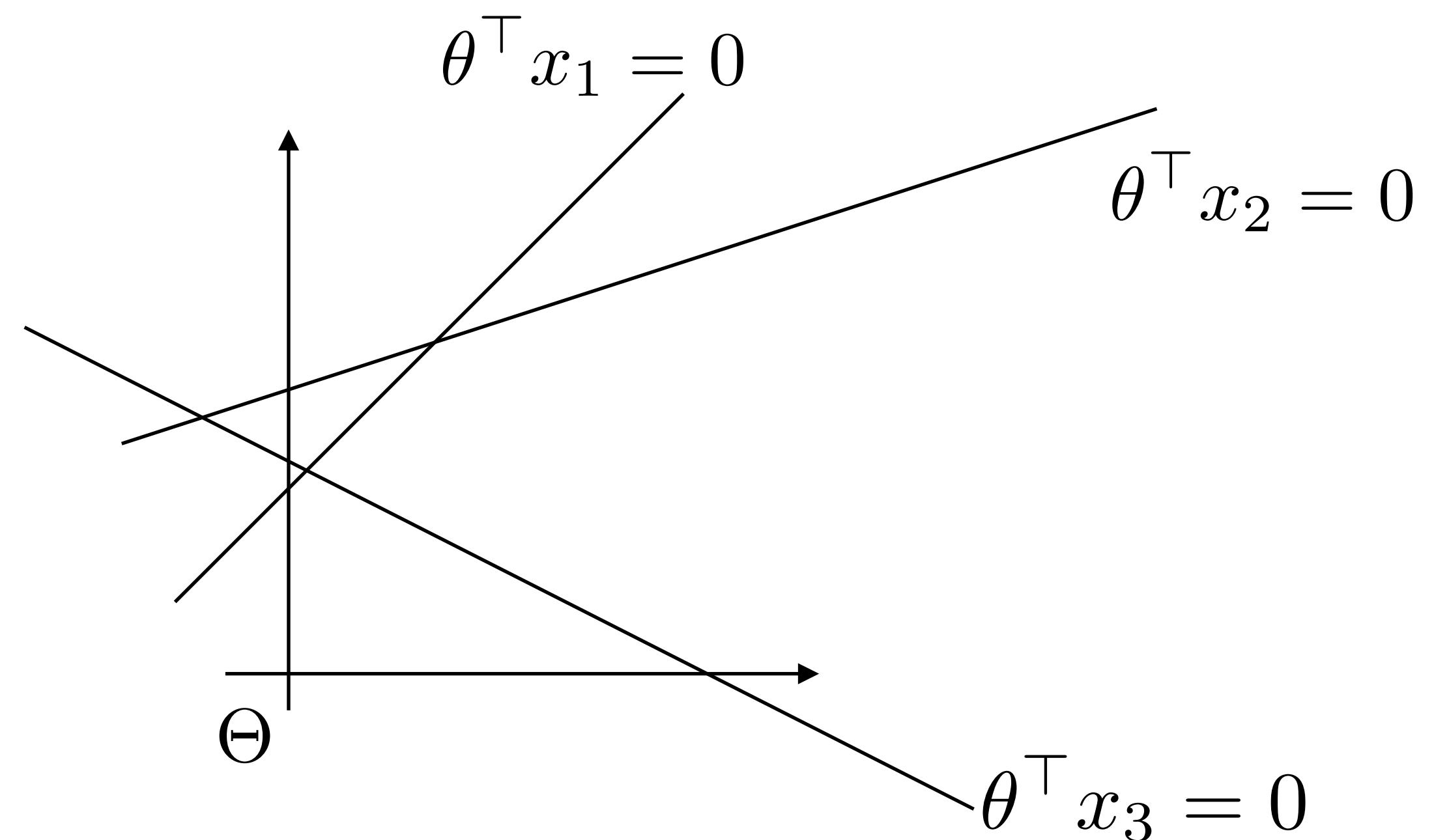
Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j) = \mathbb{P}(\tilde{\theta}^\top x_j \geq 0 \quad \forall j)$$

True validity



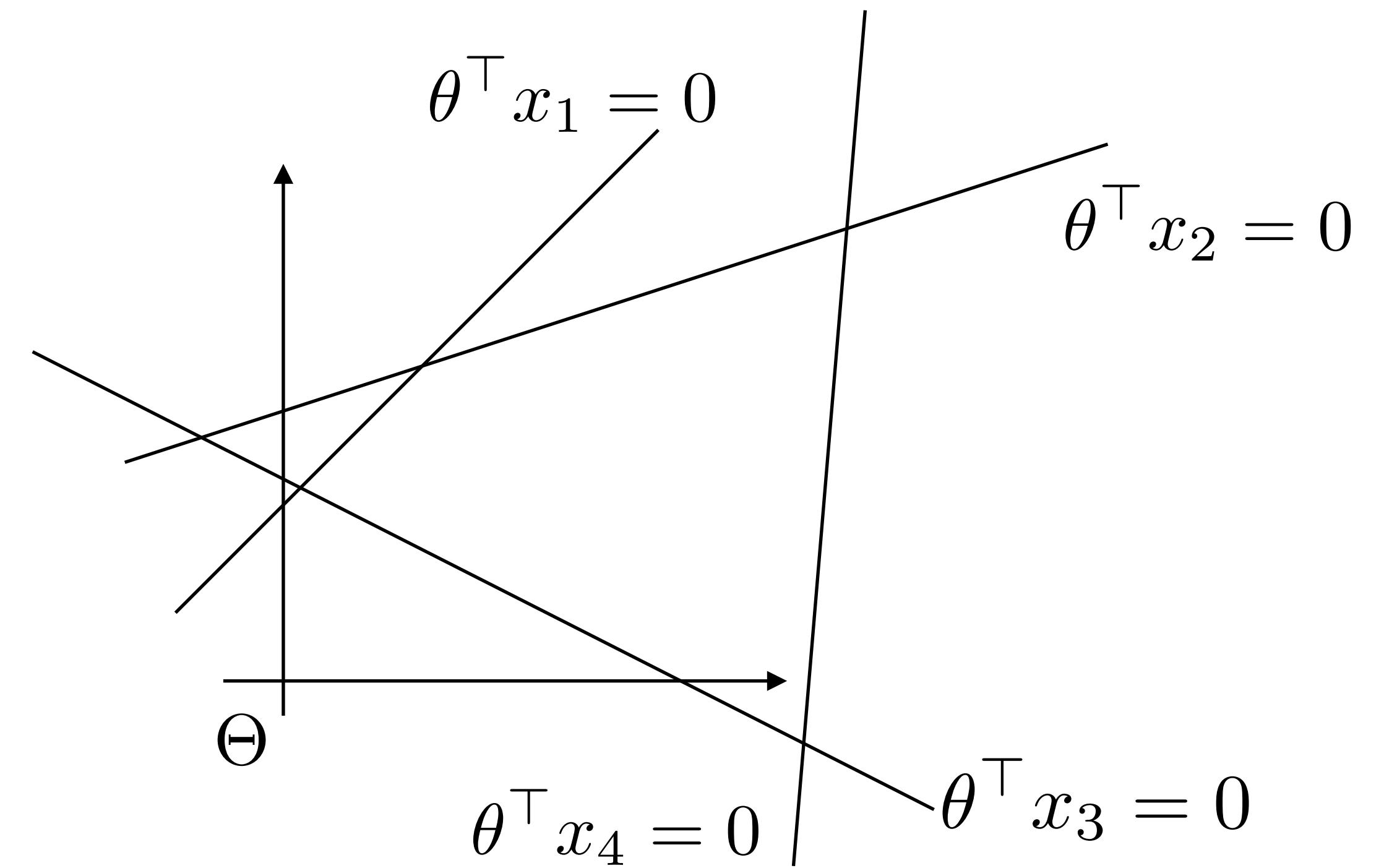
Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j) = \mathbb{P}(\tilde{\theta}^\top x_j \geq 0 \quad \forall j)$$

True validity



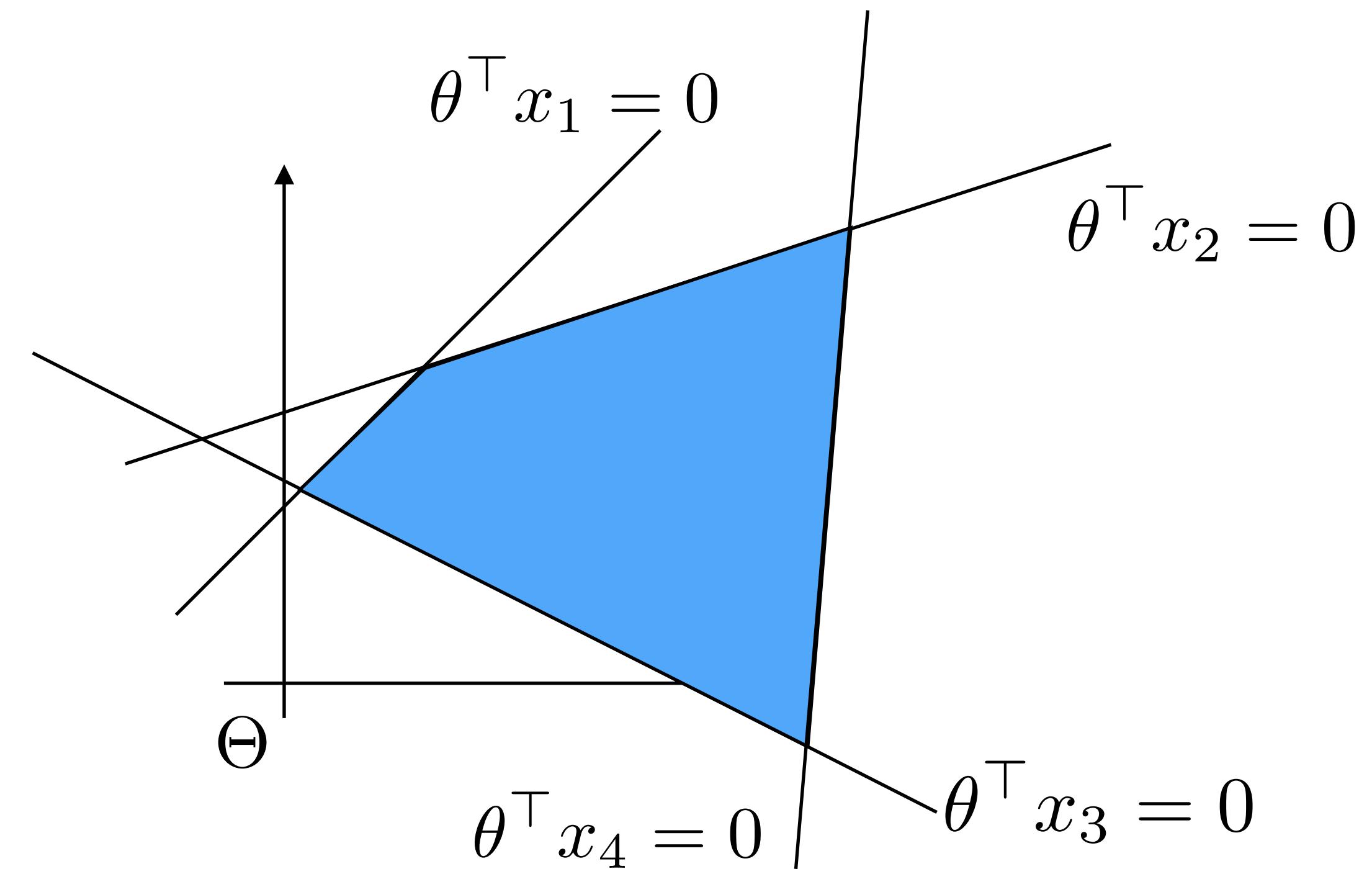
Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j) = \mathbb{P}(\tilde{\theta}^\top x_j \geq 0 \quad \forall j)$$

True validity $= \mathbb{P}(\tilde{\theta} \in \text{triangle})$



Uncertainty Quantification

- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

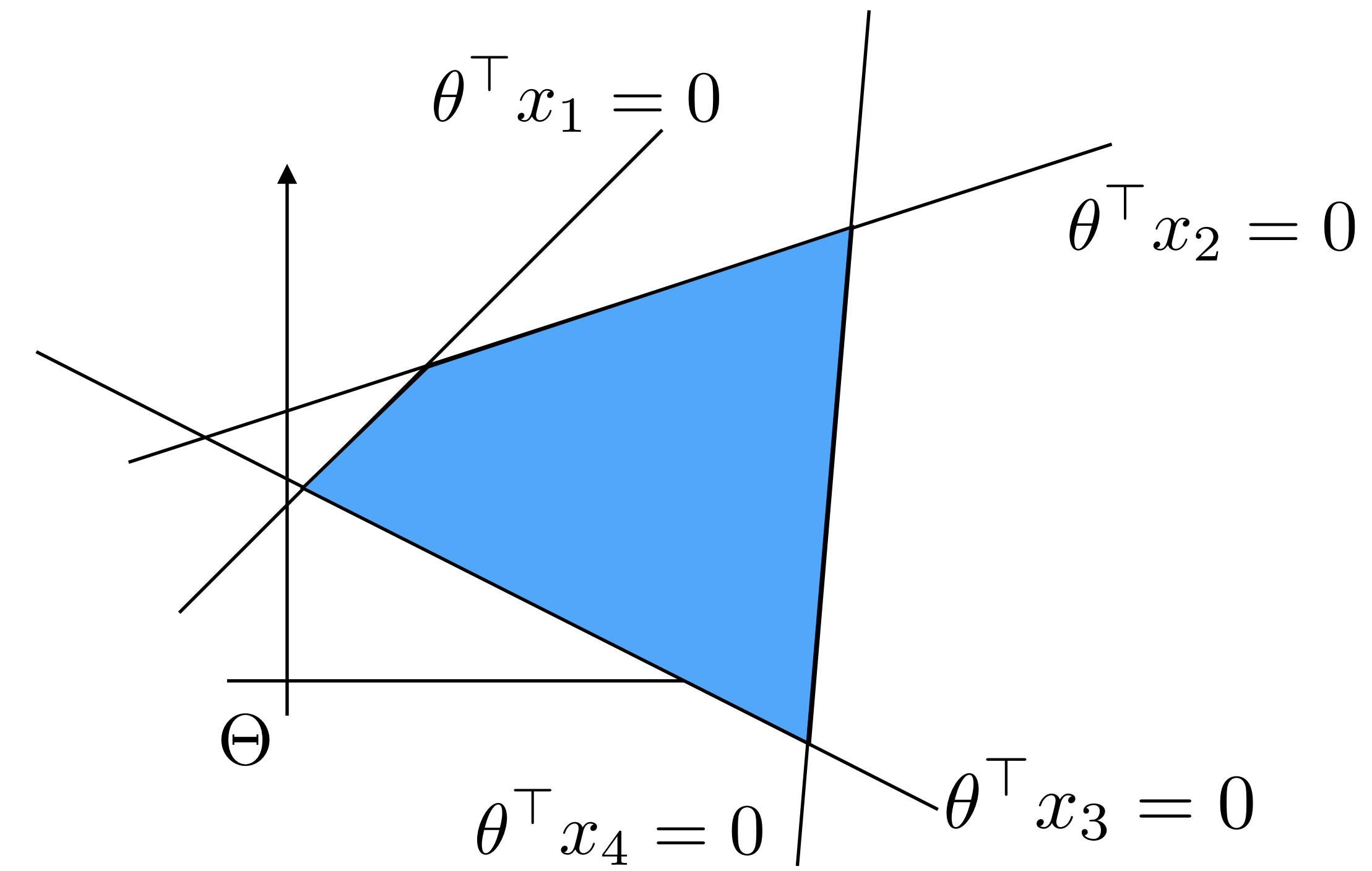
$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j) = \mathbb{P}(\tilde{\theta}^\top x_j \geq 0 \quad \forall j)$$

True validity

$$= \mathbb{P}(\tilde{\theta} \in \triangle)$$

We don't know



Uncertainty Quantification

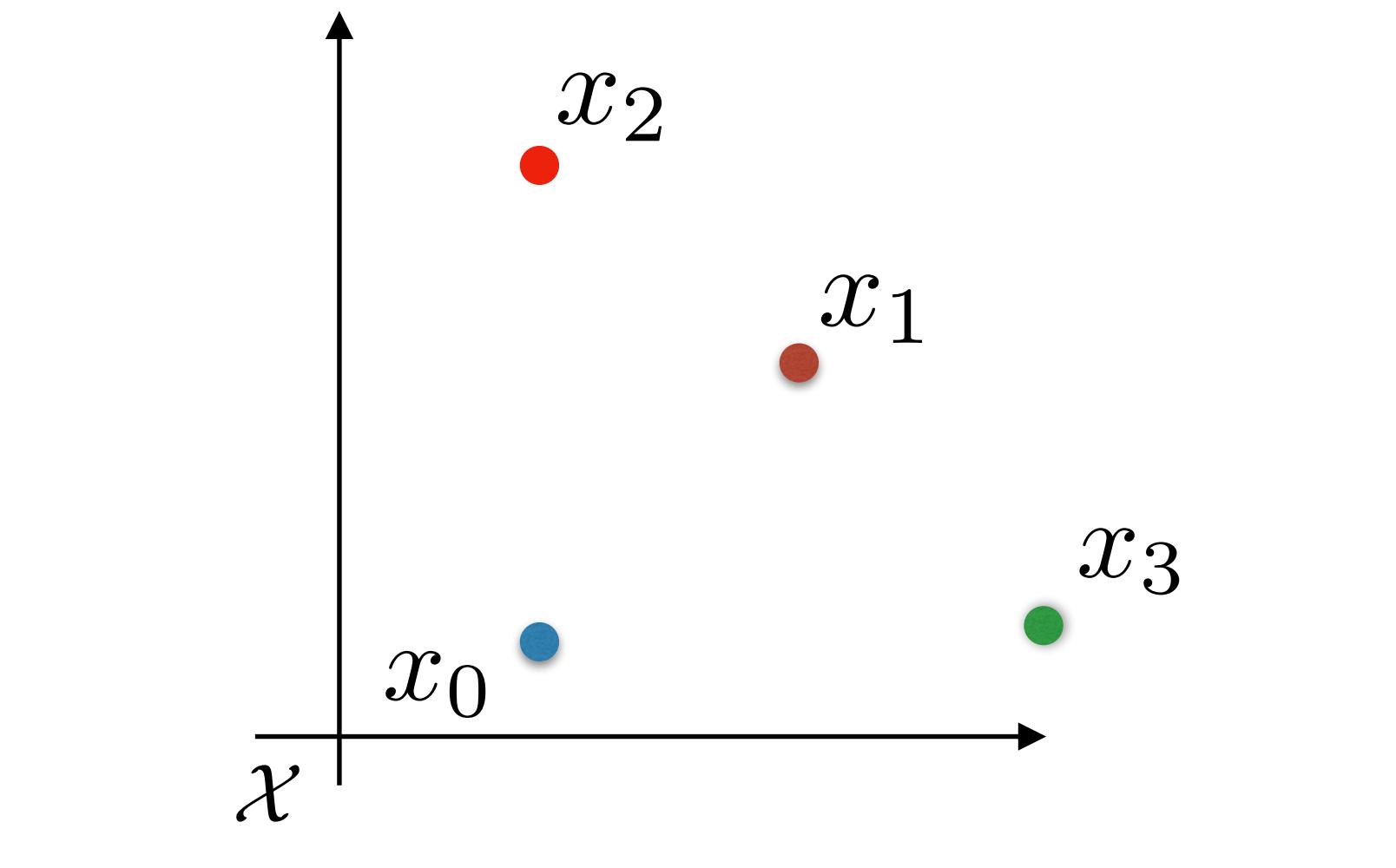
- **Q1.** Given $\{x_j\}$, what is the probability that it is a valid plan (w.r.t. $\tilde{\theta}$)?

$$\mathcal{C}_\theta(x) = \begin{cases} 0 & \text{if } \theta^\top x < 0 \\ 1 & \text{if } \theta^\top x \geq 0 \end{cases}$$

$$\mathbb{P}(\mathcal{C}_{\tilde{\theta}}(x_j) = 1 \quad \forall j) = \mathbb{P}(\tilde{\theta}^\top x_j \geq 0 \quad \forall j)$$

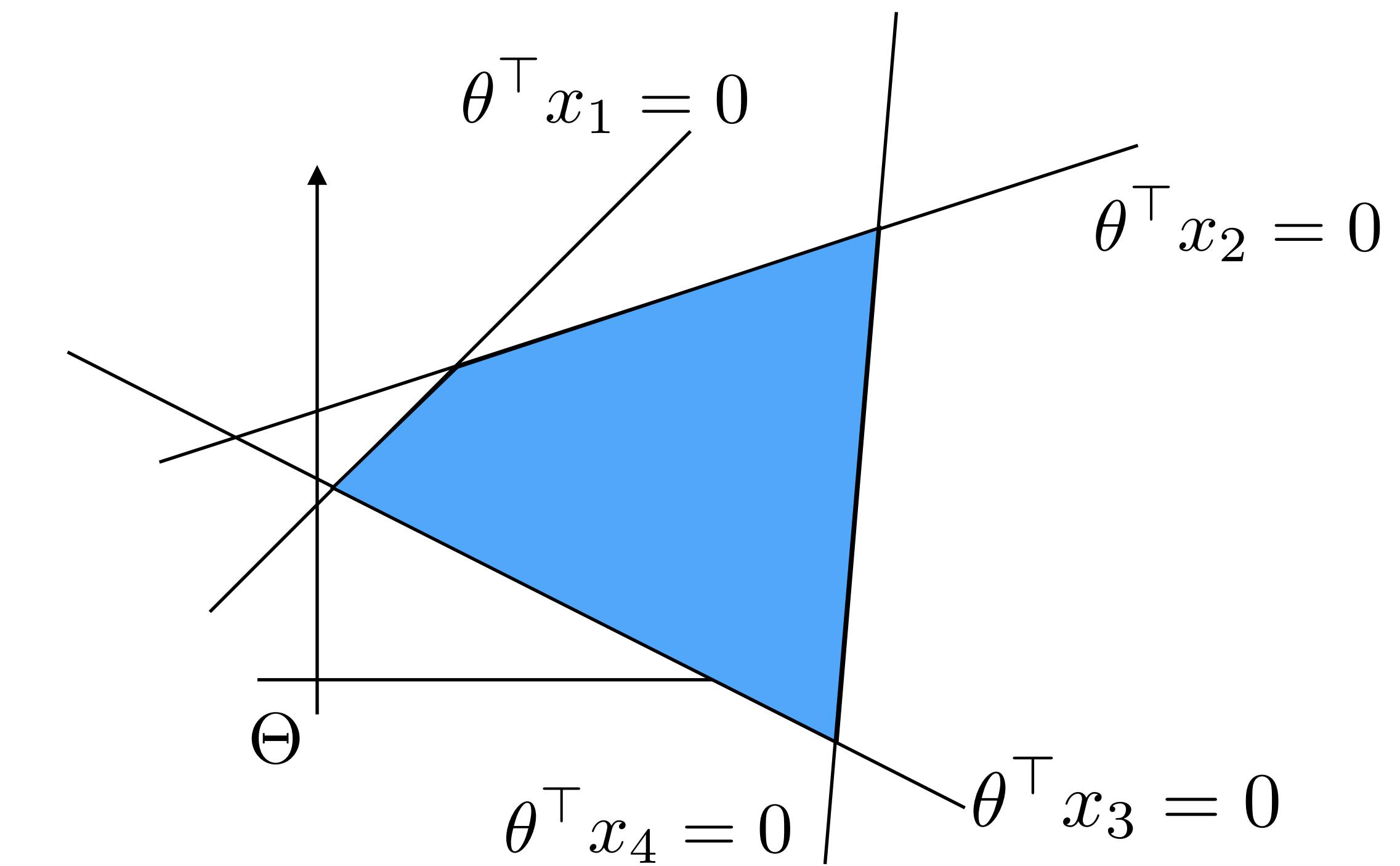
True validity

$$= \mathbb{P}(\tilde{\theta} \in \triangle)$$



Vandenberghe et al. (2007)

~



Uncertainty Quantification

Assume that the *future* parameter $\tilde{\theta} \sim \mathbb{P} \in \mathbb{B} = \{\mathbb{P} : \mathbb{G}(\mathbb{P}, \hat{\mathbb{P}}) \leq \rho\}$

$$\begin{array}{c} \text{Lower bound} \quad \leq \quad \text{True validity} \quad \leq \quad \text{Upper bound} \\ \inf_{\mathbb{Q} \in \mathbb{B}} \mathbb{Q}(\tilde{\theta} \in \text{◀}) \leq \quad \mathbb{P}(\tilde{\theta} \in \text{◀}) \quad \leq \sup_{\mathbb{Q} \in \mathbb{B}} \mathbb{Q}(\tilde{\theta} \in \text{◀}) \\ \uparrow \\ \text{True validity of plan } \{x_j\}_{j=1,\dots,J} \end{array}$$

Uncertainty Quantification

Assume that the *future* parameter $\tilde{\theta} \sim \mathbb{P} \in \mathbb{B} = \{\mathbb{P} : \mathbb{G}(\mathbb{P}, \hat{\mathbb{P}}) \leq \rho\}$

$$L^* \leq \inf_{\mathbb{Q} \in \mathbb{B}} \mathbb{Q}(\tilde{\theta} \in \text{blue cone}) \leq \mathbb{P}(\tilde{\theta} \in \text{blue cone}) \leq \sup_{\mathbb{Q} \in \mathbb{B}} \mathbb{Q}(\tilde{\theta} \in \text{blue cone}) \leq U^*$$

↑
True validity of plan $\{x_j\}_{j=1,\dots,J}$

Lower bound \leq **True validity** \leq **Upper bound**

Semi-definite program **Semi-definite program**

Theorem: L^* and U^* can be computed by solving *linear* semidefinite programs.

- ▶ Computational tractability (solved effectively by MOSEK)
- ▶ Complementary information: either $L^* = 0$ or $U^* = 1$

Uncertainty Quantification

Theorem: We have

$$L^* = \begin{cases} \inf & 1 - \sum_{j \in [J]} \lambda_j \\ \text{s.t.} & \mu \in \mathbb{R}^d, \Sigma \in \mathbb{S}_+^d, C \in \mathbb{R}^{d \times d}, M \in \mathbb{S}_+^d \\ & \lambda_j \in \mathbb{R}, z_j \in \mathbb{R}^d, Z_j \in \mathbb{S}^d \quad \forall j \in [J] \\ & -x_j^\top z_j \geq 0, \begin{bmatrix} Z_j & z_j \\ z_j^\top & \lambda_j \end{bmatrix} \succeq 0 \quad \forall j \in [J] \\ & \sum_{j \in [J]} \begin{bmatrix} Z_j & z_j \\ z_j^\top & \lambda_j \end{bmatrix} \preceq \begin{bmatrix} M & \mu \\ \mu^\top & 1 \end{bmatrix}, \begin{bmatrix} \Sigma & C \\ C^\top & \hat{\Sigma}_1 \end{bmatrix} \succeq 0, \begin{bmatrix} M - \Sigma & \mu \\ \mu^\top & 1 \end{bmatrix} \succeq 0 \\ & \|\hat{\mu}_1\|^2 - 2\hat{\mu}_1^\top \mu + \text{Tr}[M + \hat{\Sigma}_1 - 2C] \leq \rho^2. \end{cases}$$

Proof ideas:

- ▶ Chebyshev bounds on an open polyhedron
- ▶ Rejoin a two-layer optimization problem

Uncertainty Quantification

Theorem: We have

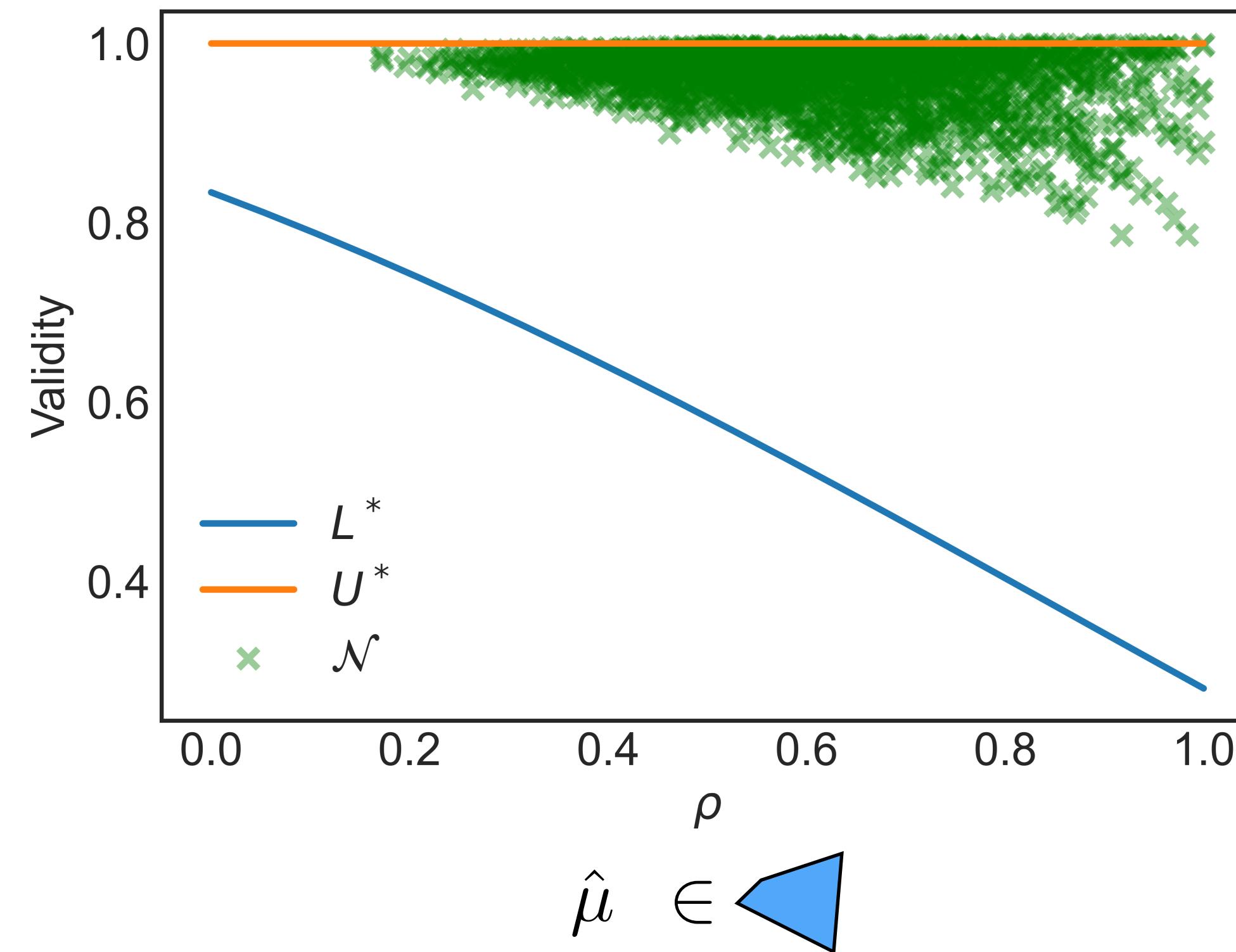
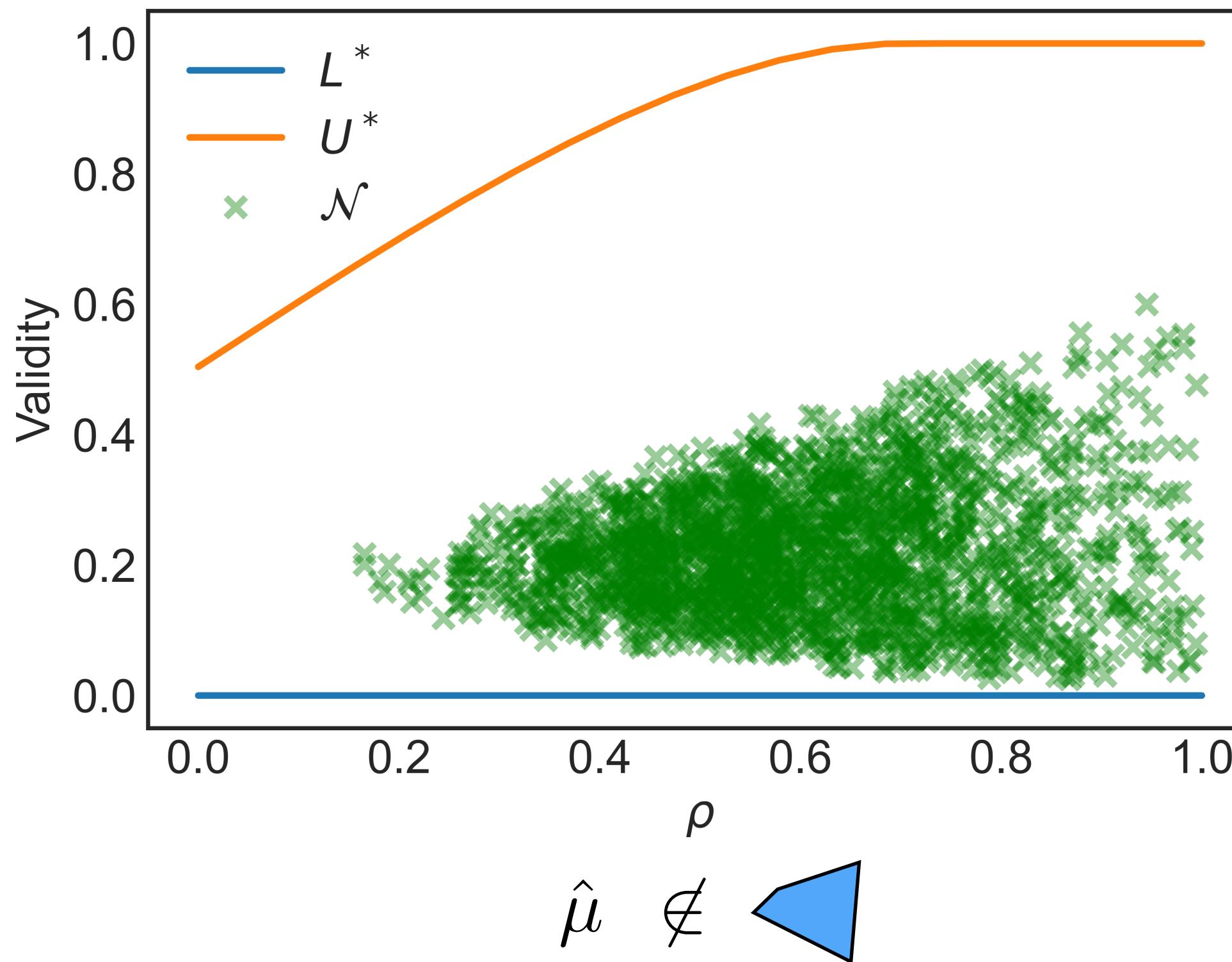
$$U^* = \begin{cases} \inf & z_0 + \gamma(\rho^2 - \|\hat{\mu}_1\|_2^2 - \text{Tr}[\hat{\Sigma}_1]) + q + \text{Tr}[Q] \\ \text{s.t.} & \gamma \in \mathbb{R}_+, z_0 \in \mathbb{R}, z \in \mathbb{R}^d, Z \in \mathbb{S}_+^d, q \in \mathbb{R}_+, Q \in \mathbb{S}_+^d, \lambda \in \mathbb{R}_+^J \\ & \begin{bmatrix} \gamma I - Z & \gamma \hat{\Sigma}_1^{\frac{1}{2}} \\ \gamma \hat{\Sigma}_1^{\frac{1}{2}} & Q \end{bmatrix} \succeq 0, \quad \begin{bmatrix} \gamma I - Z & \gamma \hat{\mu}_1 + z \\ \gamma \hat{\mu}_1^\top + z^\top & q \end{bmatrix} \succeq 0 \\ & \begin{bmatrix} Z & z \\ z^\top & z_0 \end{bmatrix} \succeq 0, \quad \begin{bmatrix} Z & z \\ z^\top & z_0 - 1 \end{bmatrix} \succeq \sum_{j \in [J]} \lambda_j \begin{bmatrix} 0 & \frac{1}{2}x_j \\ \frac{1}{2}x_j^\top & 0 \end{bmatrix}. \end{cases}$$

Proof ideas:

- ▶ Duality techniques from DRO
- ▶ Rejoin with the support function of the uncertainty set

Empirical Results

- ▷ fix a $\{x_j\}_J$
- ▷ fix $\hat{\mu} = \theta_0, \hat{\Sigma} = 0.5I$
- ▷ evaluate L^* and U^*
- ▷ sample $\tilde{\theta} \sim \mathcal{N}(\mu_g, \Sigma_g)$, (μ_g, Σ_g) is chosen randomly



Counterfactual Plan - Construction

$$\min_{x_1, \dots, x_J} \text{Proximity}(\{x_j\}, x_0) - \lambda_1 \text{FutureValidity}(\{x_j\}) - \lambda_2 \text{Diversity}(\{x_j\})$$

Average L-1 distance

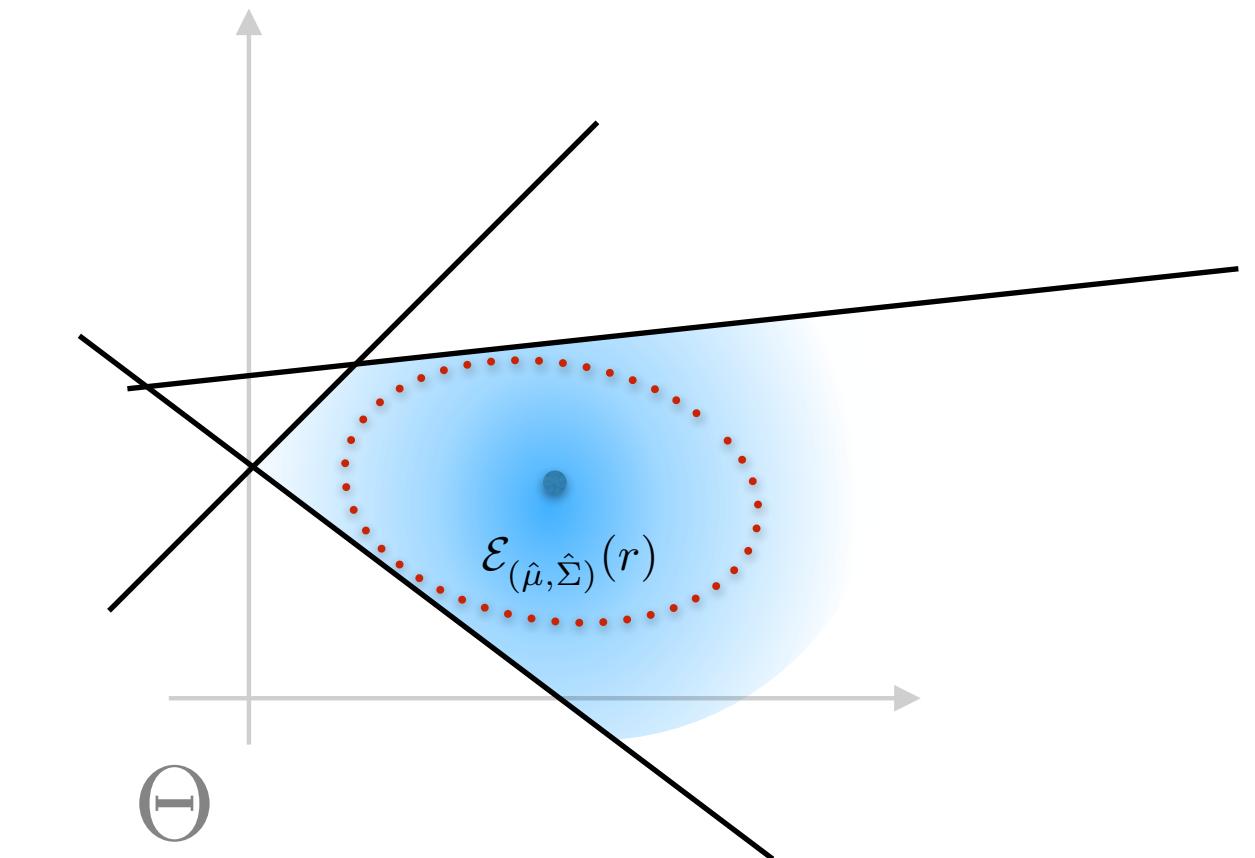
where $\text{Proximity}(\{x_j\}, x_0) \triangleq \frac{1}{J} \sum_{j=1}^J \|x_j - x_0\|_1.$

The maximum volume of the ellipsoid inscribed in a set of feasible parameters

$$\text{FutureValidity}(\{x_j\}) \triangleq \max\{r : r \geq 0, \mathcal{E}_{(\hat{\mu}, \hat{\Sigma})}(r) \subseteq \Theta(\{x_j\})\}.$$

Determinantal point process (DPP)

$$\text{Diversity}(\{x_j\}) \triangleq \det(K), \text{ where } K_{i,j} = (1 + \|x_i - x_j\|_1)^{-1} \quad \forall 1 \leq i, j \leq J.$$



Solved by projected gradient descent

Empirical Results

- each dataset contains two sets of data (present data and future data).
- generate counterfactuals w.r.t. the classifier trained on the present data.
- evaluate empirical validity w.r.t. the classifiers trained on the future data.

Dataset	Method	Proximity	Diversity	L^*	Empirical Validity
German credit	DiCE	0.986 ± 0.324	0.072 ± 0.050	0.649 ± 0.073	0.996 ± 0.008
	MahalanobisCrr	1.002 ± 0.323	0.064 ± 0.047	0.750 ± 0.064	0.999 ± 0.003
	COPA ($\lambda_1 = 0.2; \lambda_2 = 2.0$)	0.916 ± 0.178	0.017 ± 0.058	0.944 ± 0.168	0.997 ± 0.018
	COPA ($\lambda_1 = 0.5; \lambda_2 = 5.0$)	1.154 ± 0.253	0.114 ± 0.101	0.946 ± 0.040	1.000 ± 0.000
	COPA ($\lambda_1 = 1.0; \lambda_2 = 10.0$)	1.351 ± 0.166	0.225 ± 0.045	0.911 ± 0.022	1.000 ± 0.000
SBA	DiCE	2.037 ± 0.470	0.089 ± 0.057	0.946 ± 0.014	0.801 ± 0.061
	MahalanobisCrr	2.014 ± 0.473	0.085 ± 0.055	0.966 ± 0.007	0.945 ± 0.062
	COPA ($\lambda_1 = 0.2; \lambda_2 = 2.0$)	1.831 ± 0.139	0.253 ± 0.026	0.994 ± 0.000	1.000 ± 0.000
	COPA ($\lambda_1 = 0.5; \lambda_2 = 5.0$)	1.966 ± 0.112	0.363 ± 0.012	0.995 ± 0.000	1.000 ± 0.000
	COPA ($\lambda_1 = 1.0; \lambda_2 = 10.0$)	2.010 ± 0.124	0.380 ± 0.006	0.995 ± 0.000	1.000 ± 0.000
Student performance	DiCE	1.486 ± 0.325	0.136 ± 0.044	0.549 ± 0.307	0.408 ± 0.363
	MahalanobisCrr	1.497 ± 0.325	0.126 ± 0.044	0.864 ± 0.117	0.757 ± 0.284
	COPA ($\lambda_1 = 0.2; \lambda_2 = 2.0$)	1.779 ± 0.352	0.052 ± 0.047	0.998 ± 0.000	1.000 ± 0.000
	COPA ($\lambda_1 = 0.5; \lambda_2 = 5.0$)	1.882 ± 0.353	0.089 ± 0.032	0.998 ± 0.000	1.000 ± 0.000
	COPA ($\lambda_1 = 1.0; \lambda_2 = 10.0$)	1.926 ± 0.349	0.109 ± 0.024	0.997 ± 0.000	1.000 ± 0.000

Follow-up Questions

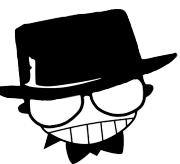
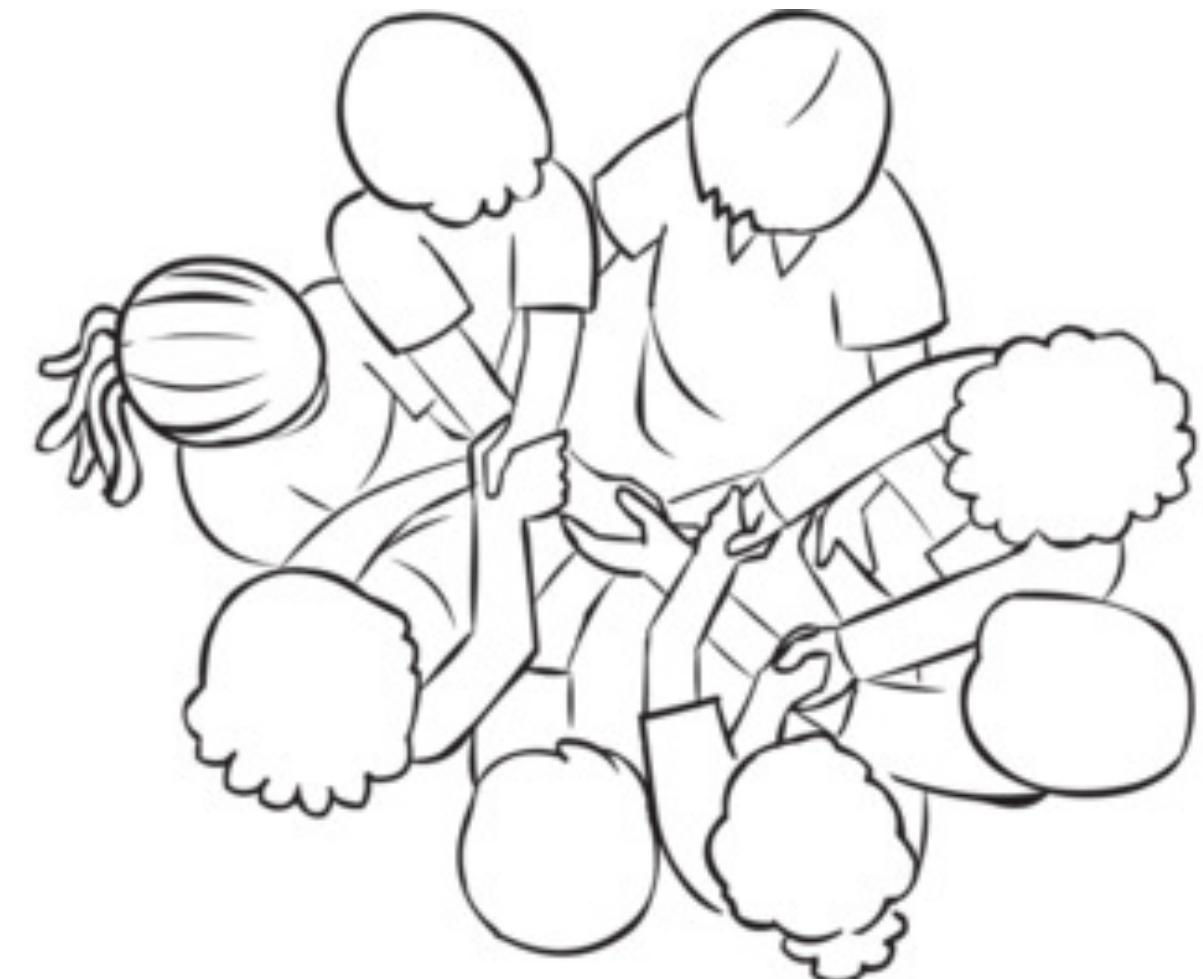
- Can we use the lower bound to construct a *robust* counterfactual plan?
 - No, we have to use a maximum ellipsoid to approximate it.
 - But, if there is only one counterfactual explanation, we can construct it efficiently
-> ICLR'23
- Can we use it for nonlinear models?
 - Yes, but we must use a linear model (e.g., LIME) to construct a surrogate decision boundary before applying our technique.
- How does the linear surrogate affect the constructed explanation?
(e.g., unfaithful to the decision boundary)
 - Analyzing and proposing a more faithful linear surrogate -> Rejected (4 times)
- I don't wanna use a linear surrogate.
 - Robust Bayesian Recourse -> UAI'22

The Shapley value

from game theory to machine learning

Cooperative games

- Cooperative games model scenarios where
 - Agent can benefit by cooperating
 - Binding agreements are possible



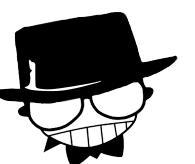
Cooperative games

- How to construct a salary structure in a company or a football club?



Cooperative games

- How to construct a salary structure in a company or a football club?
- How to divide a surplus (profit) to shareholders so that everyone is satisfied?



Cooperative games

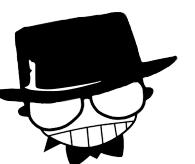
Cooperative game

- Players can benefit by cooperating
- Binding agreements are possible
- Answer the question: **How to divide the surplus when joining the grand coalition?**

Non-cooperative game

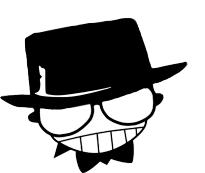
- Players are independent
- No cooperation
- Answer the question: **What is the “good” strategy for each player to maximize their return?**

Nash equilibrium, zero-sum game



Cooperative games

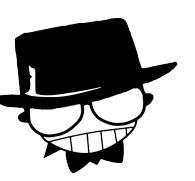
- A cooperative game is described by a tuple (N, F) where $N = \{1, \dots, n\}$ a set of players and $F : 2^N \rightarrow \mathbb{R}$, $F(\emptyset) = 0$ is a characteristic function.
- By von Neumann and Morgenstern:
 - Superadditive game: $F(S \cup T) \geq F(S) + F(T)$ (relaxed by later works)
 - Transferable Utility (TU) game | non-Transferable Utility (NTU) game



Cooperative games

- A cooperative game is described by a tuple (N, F) where $N = \{1, \dots, n\}$ a set of players and $F : 2^N \rightarrow \mathbb{R}$, $F(\emptyset) = 0$ is a characteristic function.
 - \mathbb{G} : set of all possible games
- How to divide a surplus (profit) to players?
 - $\psi : \mathbb{G} \rightarrow \mathbb{R}^n$, $\psi((N, F)) \in \mathbb{R}^n$: a solution concept

When context is clear, I omit the player set N in (N, F)
and use F to denote the game and $\psi(F)$ to denote a solution concept of F



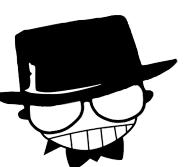
Cooperative games

- A cooperative game is described by a tuple (N, F) where $N = \{1, \dots, n\}$ a set of players and $F : 2^N \rightarrow \mathbb{R}$, $F(\emptyset) = 0$ is a characteristic function.
- \mathbb{G} : set of all possible games
- How to divide a surplus (profit) to players?
 - $\psi : \mathbb{G} \rightarrow \mathbb{R}^n$, $\psi((N, F)) \in \mathbb{R}^n$: a solution concept



$\psi((N, F))$	\$4	\$11	\$9
----------------	-----	------	-----

Payoff vector



Several Solution Concepts

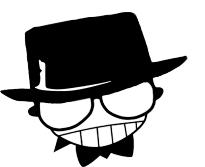
- **Imputation set** (von Neumann & Morgenstern, 1944): Includes all payoff vectors that are *efficient*, i.e., share the whole utility, and *individually rational*, i.e., assign to each player at least her standing alone worth.

$$\psi_{IS}((N, F)) = \left\{ x \in \mathbb{R}^n : \sum_{i \in N} x_i = F(N), x_i \geq v(i), \forall i \in N \right\}$$

- **Core set** (Gillies, 1953): similar to imputation set but considering *coalitional rational* instead of individual rational.

$$\psi_{CS}((N, F)) = \left\{ x \in \mathbb{R}^n : \sum_{i \in S} x_i \geq F(S), \forall S \subseteq N \right\}$$

Could be empty and not unique!



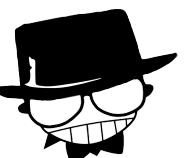
The Shapley Value

- A solution concept for cooperative games that assigns payoffs for each player in the game and that is
 - unique for each game
 - efficient
 - fair


$$\psi((N, F))$$

\$4	\$8	\$7
-----	-----	-----

Payoff vector



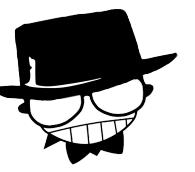
The Shapley Value

Def 1.1 (dividend): The value of the player i is the sum of its dividend share to all possible coalitions to which he could be joined. The dividend of a coalition is shared equally by its members.

$$\phi_i((N, F)) = \sum_{S \subseteq N \setminus i} \frac{\Delta_F(S \cup i)}{s+1}, \quad \forall i \in N,$$

where $\Delta_F(T)$ is Harsanyi dividend of the game F to the coalition T :

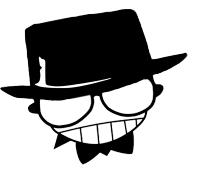
$$\Delta_F(T) = \sum_{R \subseteq T} (-1)^{|T|-|R|} F(R), \quad \forall T \subseteq N$$



The Shapley Value

Def 1.2 (*weighted marginal contribution*): The value of player i is computed by its weighted average marginal contribution over all possible coalitions to which player i could be joined.

$$\phi_i((N, F)) = \frac{1}{n} \sum_{S \subseteq N \setminus i} \frac{1}{\binom{n-1}{s}} [F(S \cup i) - F(S)], \quad \forall i \in N$$

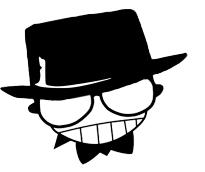


The Shapley Value

Def 1.3 (random order): The grand coalition is formed in a way that each player joins the coalition one-by-one, the value of player i is the average of its marginal contribution over all possible order.

$$\phi_i((N, F)) = \frac{1}{n!} \sum_{\pi \in \Pi(N)} [F(\pi^i \cup i) - F(\pi^i)]$$

where π^i is the set of players that precede player i in the permutation π .



Axiomatic Characterization

We want to design a solution concept ψ that satisfies the following axioms:

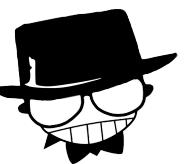
- **Axiom 1** (*Linearity, L*): for any two games F_1 and F_2 ,

$$\psi(F_1 + \alpha F_2) = \psi(F_1) + \alpha\psi(F_2)$$

- **Axiom 2** (*Dummy, D*): if $i \in N$ is a dummy player,
i.e., $F(S \cup i) = F(S) + F(i), \forall S \subseteq N$, then $\psi_i(F) = F(i)$.

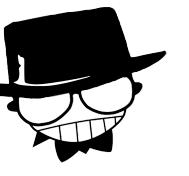
- **Axiom 3** (*Symmetry, S*): For any permutation $\pi \in \Pi(N)$, $\psi_{\pi i}(\pi F) = \psi_i(F)$

- **Axiom 4** (*Efficiency, E*): $\sum_{i \in N} \psi_i(F) = F(N) - F(\emptyset)$

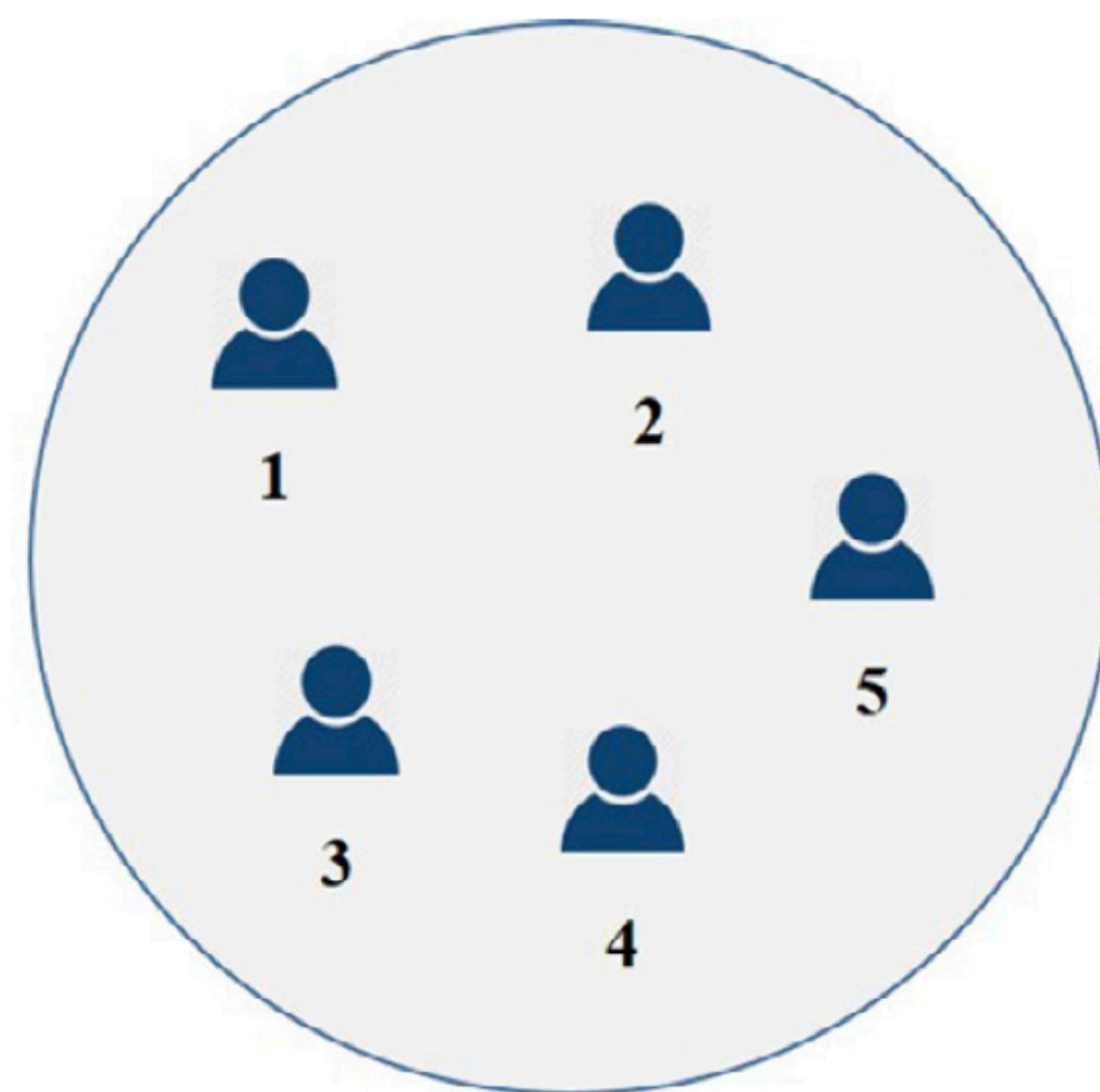


Main result

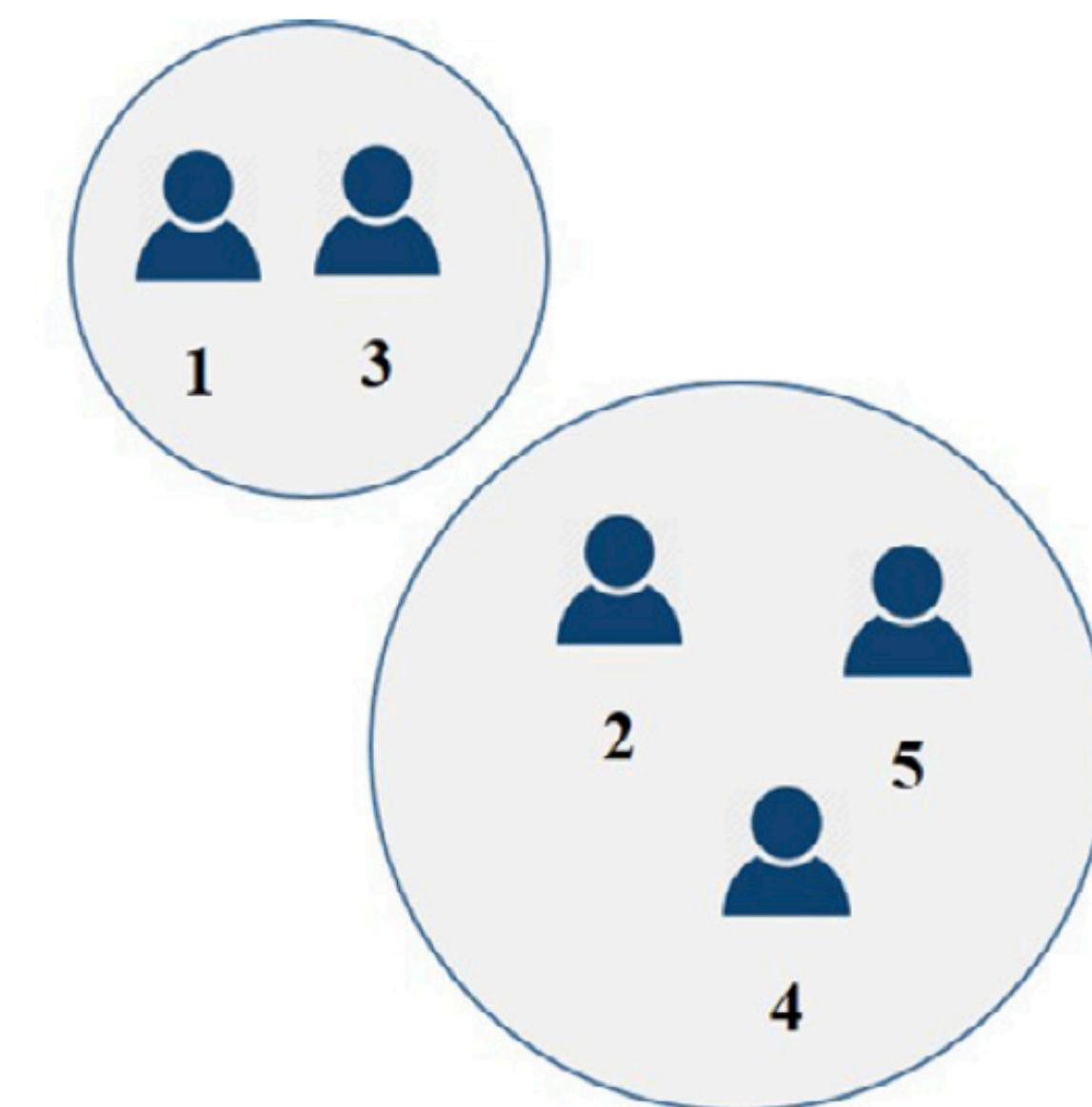
Theorem 1: Shapley value is a **unique** solution concept that satisfies four axioms *linearity, dummy, symmetry, and efficiency*.



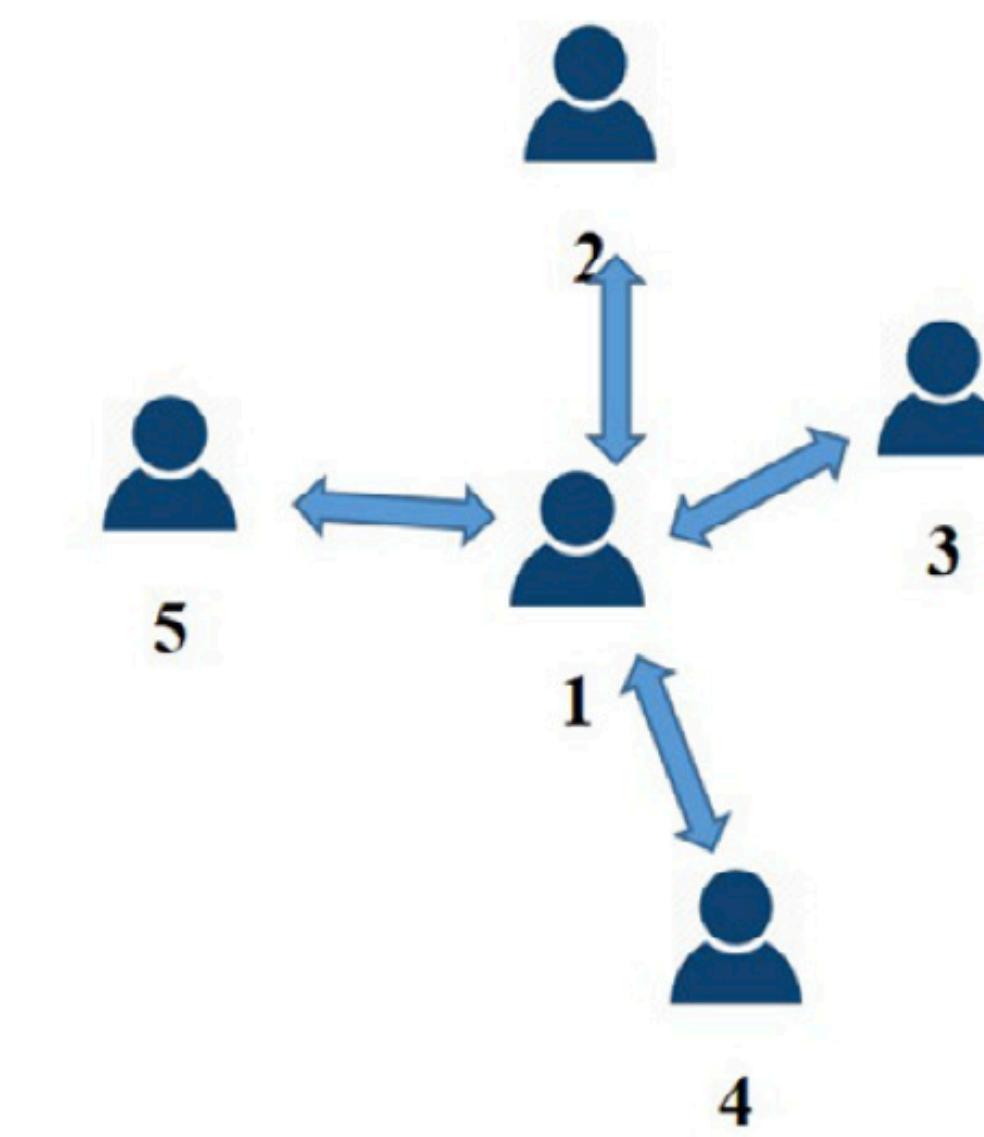
Types of Cooperative Games



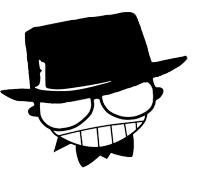
Class I: Canonical Coalitional Games



Class II: Coalition Formation Games



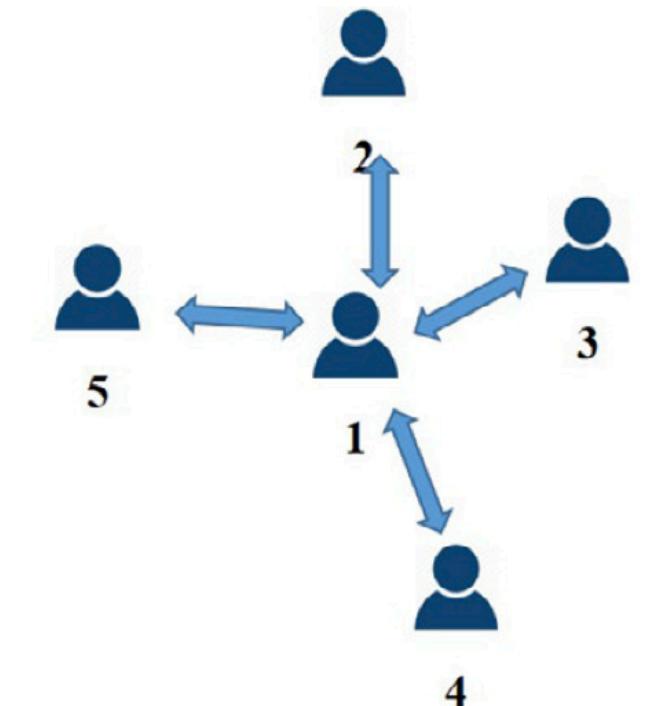
Class III: Coalitional Graph Games



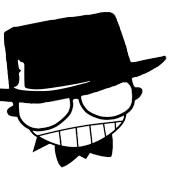
Coalitional Graph Games

- N is a set of nodes, $L = \{ij \mid i, j \in N\}$ is a set of edges. An edge represents a communication connection between two nodes.
- We omit the set N whenever the context is clear.
- $\zeta(N, L) = \{C_1, \dots, C_r\}$ denotes a set of connected components of a graph (N, L) .
- $L|_S$ is a set of edges induced by subset $S \subseteq N$, i.e., $L|_S = \{ij \in L \mid i, j \in S\}$.
- A communication-restricted function $f|_L$ is defined as

$$f|_L(T) = \sum_{R \in \zeta(T, L|_T)} f(R) \quad \forall T \subseteq N.$$



Class III: Coalitional Graph Games



The Myerson Value

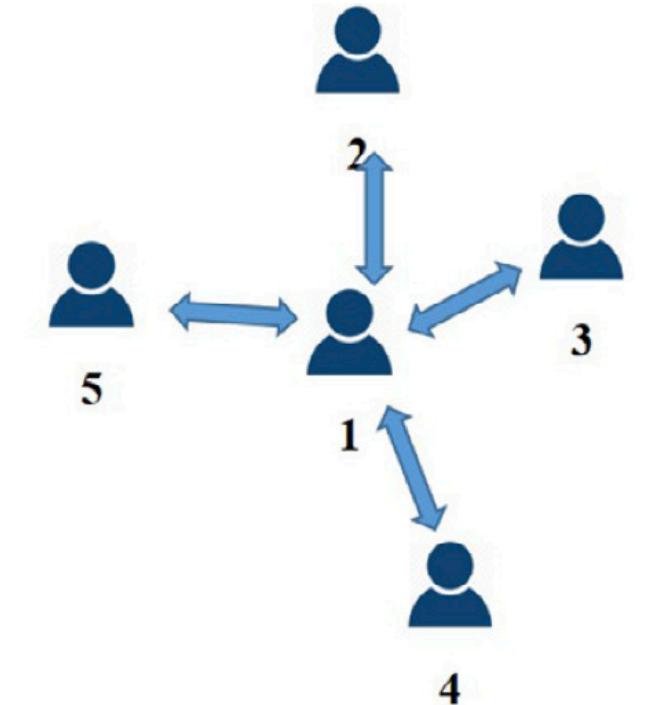
- N is a set of nodes, $L = \{ij \mid i, j \in N\}$ is a set of edges. An edge represents a communication connection between two nodes.
- We omit the set N whenever the context is clear.
- $\zeta(N, L) = \{C_1, \dots, C_r\}$ denotes a set of connected components of a graph (N, L) .
- $L|_S$ is a set of edges induced by subset $S \subseteq N$, i.e., $L|_S = \{ij \in L \mid i, j \in S\}$.
- A communication-restricted function $f|_L$ is defined as

$$f|_L(T) = \sum_{R \in \zeta(T, L|_T)} f(R) \quad \forall T \subseteq N.$$

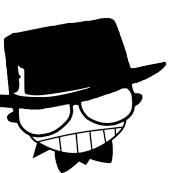
Definition 2.2 (Myerson value ([Myerson, 1977](#))). *Given a graph (N, L) and a characteristic function f , Myerson value on a CO-game (N, f, L) is defined as*

$$\psi(f, L) = \phi(f|_L),$$

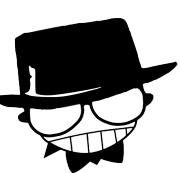
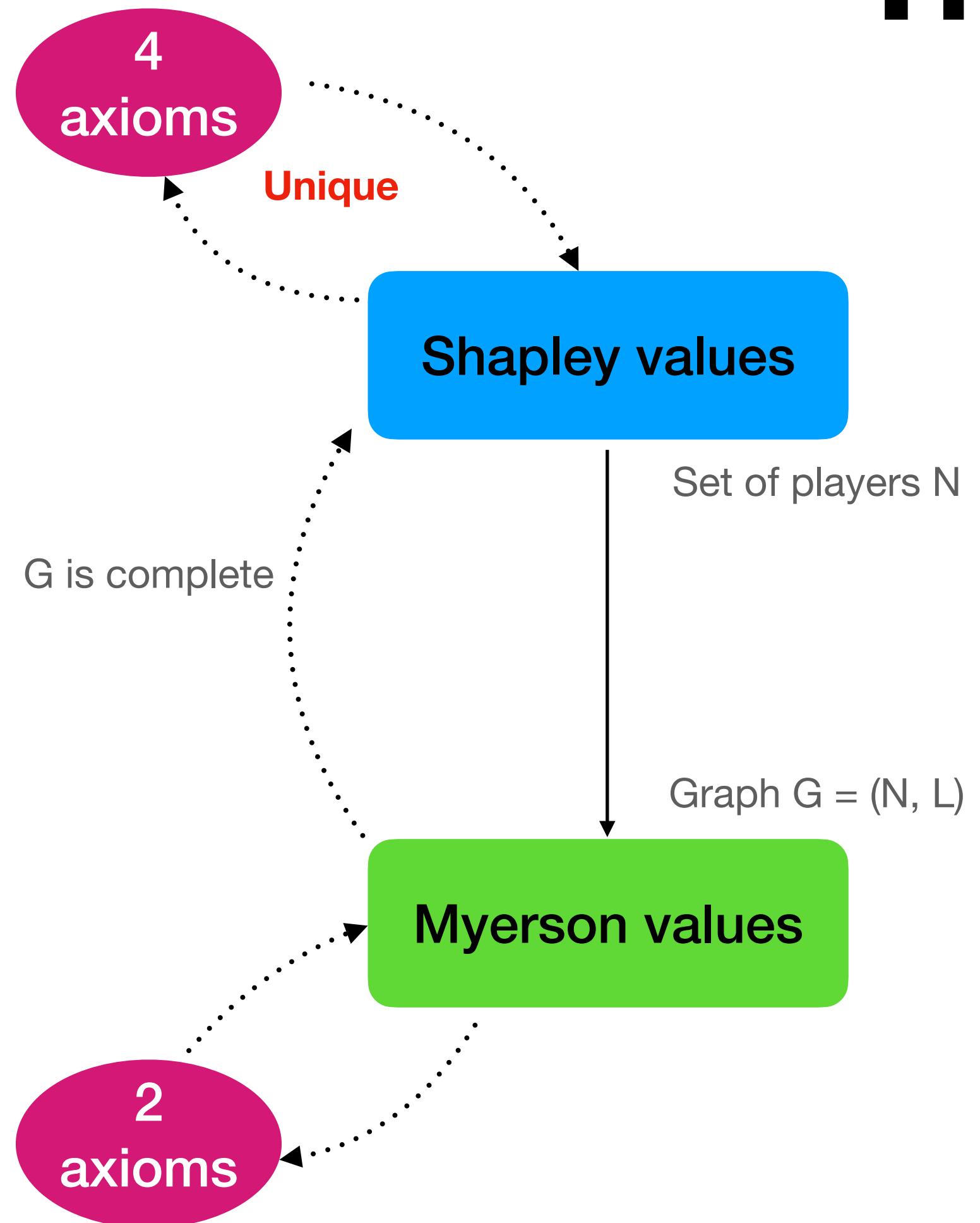
where $\phi(f|_L)$ is Shapley value of a communication-restricted function $f|_L$.



Class III: Coalitional Graph Games



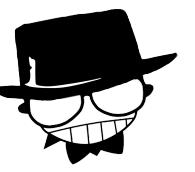
The Myerson Value



The Shapley Taylor Interaction Index

Definition 2.1 (Shapley-Taylor Interaction). *The Shapley-Taylor index is defined as the expectation of $\Phi_{S,\pi}^m(f)$ over an ordering of N players chosen uniformly at random:*

$$\Phi_S^m(f) = \mathbb{E}_\pi(\Phi_{S,\pi}^m(f)) = \frac{m}{n} \sum_{T \subseteq N \setminus S} \delta_S f(T) \frac{1}{\binom{n-1}{t}}.$$



The Myerson Value

