

Finding Regions of Counterfactual Explanations via Robust Optimization

DONATO MARAGNO
JANNIS KURTZ
TABEA RÖBER
ROB GOEDHART
ILKER BIRBIL
DICK DEN HERTOG



UNIVERSITY OF AMSTERDAM
Amsterdam Business School

JANUARY 2023

EXAMPLE: CREDIT SCORING

Labeled Data

Customer ID	Age	Salary	Current Balance	...	Loan Granted
1	28	42.000 EUR	8.200 EUR	...	1
2	56	73.800 EUR	22.300 EUR	...	1
3	42	35.100 EUR	16.900 EUR	...	0
:				:	:

Credit Scoring

Given **individual information** of a new customer, decide if the customer should be **granted a loan or not**.

BINARY CLASSIFICATION PROBLEMS

Classification Problem

Given a data point x in a data space $\mathcal{X} \subseteq \mathbb{R}^n$, classify it as 1 or -1.

Classification Machine Learning Model

- Train a classifier $h : \mathcal{X} \rightarrow [0, 1]$ on labeled data
- Assigns "probability" $h(x)$ to each data points $x \in \mathcal{X}$
- Classify x as

$$h_{\text{class}}(x) = \begin{cases} 1 & \text{if } h(x) \geq \tau \\ -1 & \text{if } h(x) < \tau \end{cases}$$

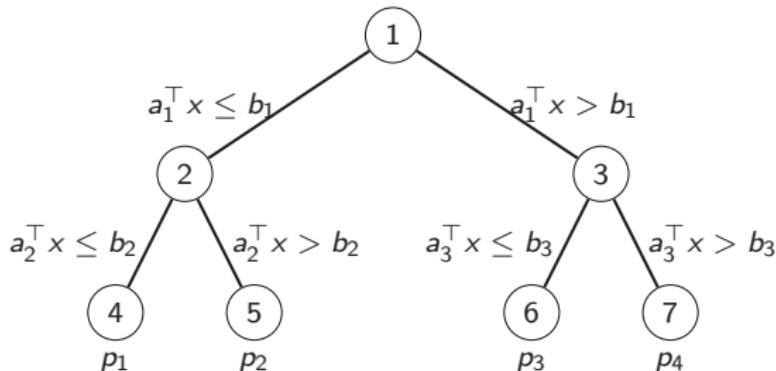
where $\tau \in (0, 1)$ is a given threshold parameter (often $\tau = 0.5$).

Models

Classifier h can be a

- Logistic Regression Function
- Classification Tree
- Neural Network

CLASSIFICATION TREES



Training Classification Trees

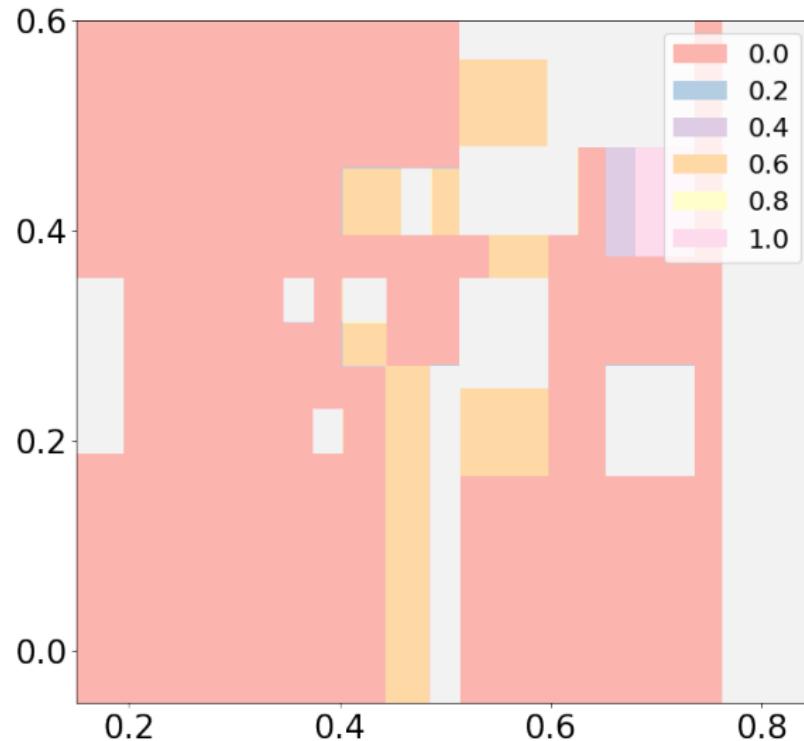
- CART (single feature splits) [Breiman et al. (1984)]
- Optimal Decision Trees (hyperplane splits) [Bertsimas, Dunn (2015)]

Properties

- Each leaf is given by a set of inequalities: $\mathcal{L}_5 = \{x \mid a_1^T x \leq b_1, a_2^T x > b_2\}$
- Classifier h assigns fraction $p_i \in [0, 1]$ to each point x in the leaf.
- p_i is often the fraction of training data in the leaf with label 1

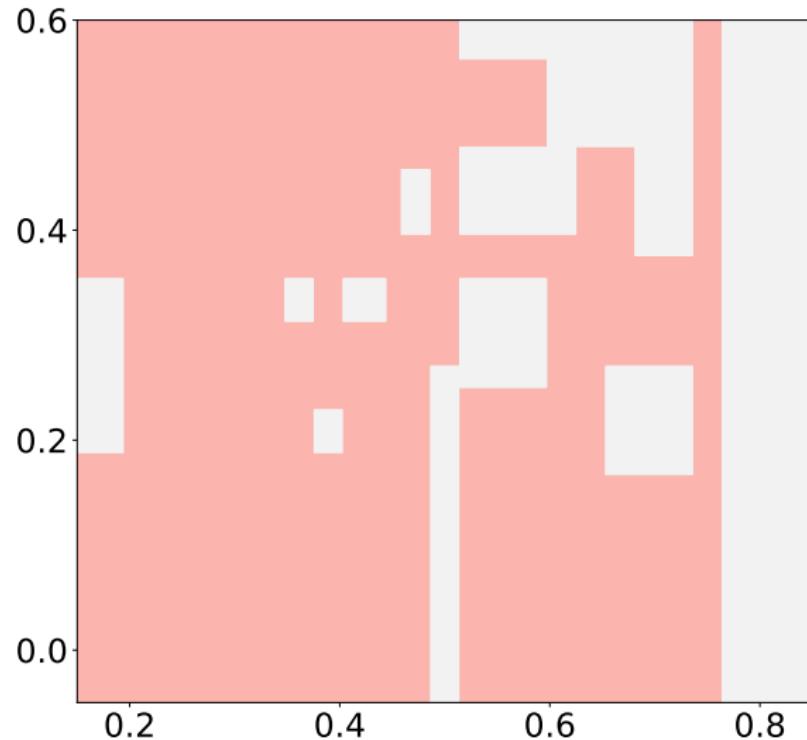
CLASSIFICATION TREES

Prediction values of a classification tree trained by CART on the Iris dataset



CLASSIFICATION TREES

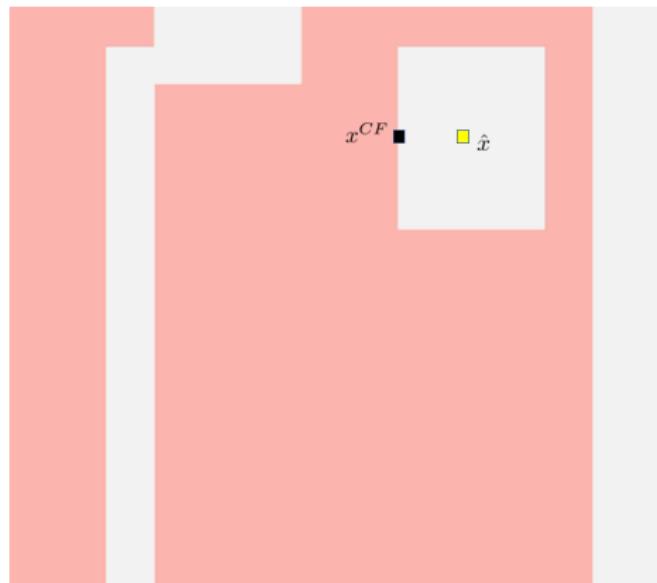
Decision region of a classification tree trained by CART on the Iris dataset



COUNTERFACTUAL EXPLANATIONS

Closest Counterfactual Point

Given a trained classifier h and a **factual point** \hat{x} predicted as -1 , we are looking for the **closest counterfactual point**, i.e. the closest x predicted as 1 .



COUNTERFACTUAL EXPLANATIONS

Optimization Problem

Counterfactual can be calculated by solving

$$\begin{aligned} & \min d(\hat{x}, x) \\ \text{s.t. } & h(x) \geq \tau \\ & x \in \mathcal{X} \end{aligned}$$

Heuristic Method

$$\min_x \max_{\lambda} d(\hat{x}, x) + \lambda(h(x) - \tau)$$

[Wachter et al. (2017)]

EXAMPLE

Credit Scoring

- Assume a person with individual information \hat{x} does **not** get granted a loan.

Age	Salary	Current Balance	...
28	42.000 EUR	8.200 EUR	...

- Calculated counterfactual point:

Age	Salary	Current Balance	...
28	44.500 EUR	10.000 EUR	...

- **Counterfactual Explanation:** If your salary would be 44.500 EUR and your current balance 10.000, then you would have been granted a loan.

RIGHT TO EXPLANATION

The European Union enacted the **right to explanation** in 2016 which was incorporated in the EU General Data Protection Regulation:

[...] In any case, such processing should be subject to suitable safeguards, which should include specific information to the data subject and the right to obtain human intervention, to express his or her point of view, to **obtain an explanation of the decision reached after such assessment** and to challenge the decision. [...]

Motivation

- **Counterfactual Explanation:** If your salary would be 44.500 EUR and your current balance 10.000, then you would have been granted a loan.
- But what if my current balance is 9985 EUR, or 10.099 EUR?
- Find counterfactual points which remain counterfactual points after small changes

LET THE USER DECIDE

CE-OCL Home Demo Questionnaire Documentation About FAQ Contact us

Diabetes

Decision tree

Pregnancies Glucose BloodPressure SkinThickness

1 85 66 29

Insulin BMI DiabetesPedigreeFunction Age

0 26.6 0.351 31

Model ▾

- Logistic regression
- Support vector machine
- Decision tree
- Random forest
- Gradient boosting method
- Neural network

Model prediction 0

Sparcity Data manifold Robustness

Off Off On

Generate counterfactual explanation

Pregnancies	Glucose	BloodPressure	SkinThickness
2.6	99.5002	66	9.2

Insulin	BMI	DiabetesPedigreeFunction	Age
0	26.95007	0.351	31

ROBUST COUNTERFACTUAL EXPLANATIONS

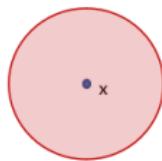
Setup

- Find counterfactual point x^{CE} such that all points in

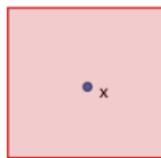
$$x^{CE} + \mathcal{S}$$

are counterfactual points.

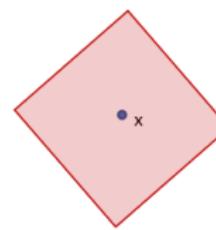
- Uncertainty set $\mathcal{S} = \{s : \|s\| \leq \varepsilon\}$ for a small $\varepsilon > 0$
- E.g. ℓ_1, ℓ_2 or ℓ_∞ -norm can be used



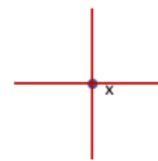
ℓ_2



ℓ_∞

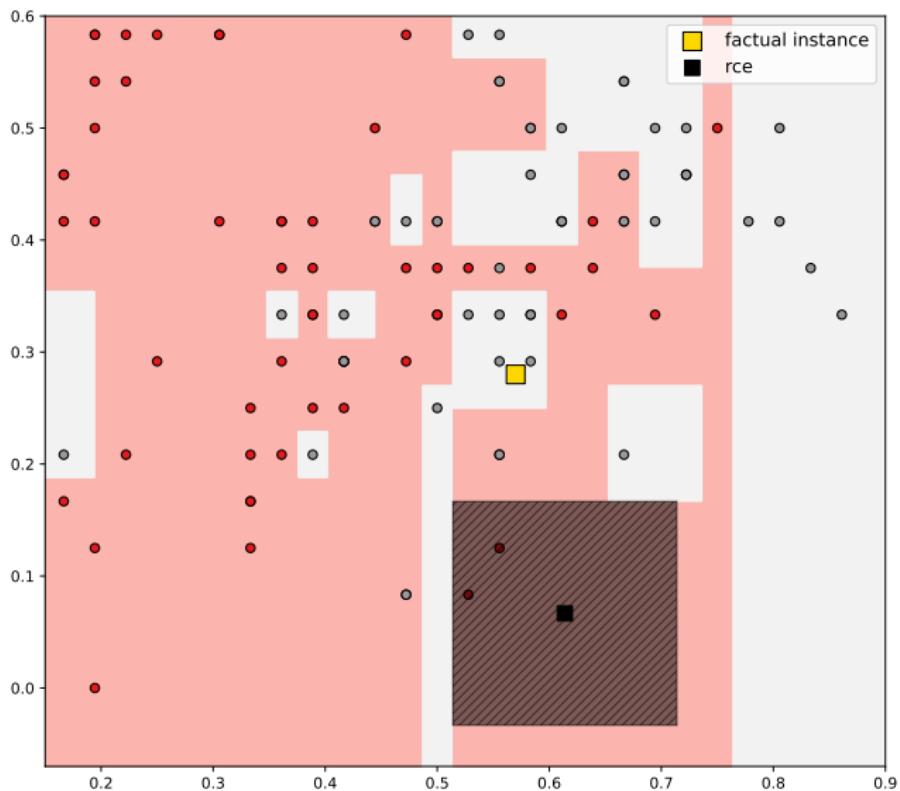


ℓ_1



ℓ_0

EXAMPLE: BOX UNCERTAINTY



ROBUST COUNTERFACTUAL EXPLANATIONS

Literature

- Gradient based heuristic with probabilistic robustness guarantees [Pawelczyk et al. (2022)]
 - ▶ heuristic solution
 - ▶ no full robustness guarantee
 - ▶ not applicable to classification trees
- Gradient based heuristic for robust CE problem [Dominguez-Olmedo et al. (2021)]
 - ▶ heuristic solution
 - ▶ no global optimal solutions, hence no robustness guarantee
 - ▶ not applicable to classification trees

Goal

Derive a method

- which guarantees full robustness
- which is also applicable to decision trees and random forests

ROBUST COUNTERFACTUAL EXPLANATIONS

Optimization Problem

Given a factual instance \hat{x} a **robust counterfactual explanation** can be derived by solving

$$\begin{aligned} & \min d(\hat{x}, x) \\ \text{s.t. } & h(x + s) \geq \tau \quad \forall s \in \mathcal{S} \\ & x \in \mathcal{X} \end{aligned}$$

where \hat{x} is the factual instance.

Difficulties

- Infinitely many constraints
- Duality trick from robust optimization does not work for all h .
- Iterative approach needed: **adversarial approach**

ADVERSARIAL APPROACH

Master Problem

- The master problem (MP) is a relaxation of the problem with a finite number of scenarios $\mathcal{Z} \subset \mathcal{S}$:

$$\begin{aligned} & \min d(\hat{x}, x) \\ \text{s.t. } & h(x + s) \geq \tau \quad \forall s \in \mathcal{Z} \\ & x \in \mathcal{X} \end{aligned}$$

- Provides a **lower bound** for the optimal value.

Adversarial Problem

- The adversarial problem (AP) finds a new scenario in \mathcal{S} which maximally violates the constraints for current MP solution x^* :

$$\begin{aligned} & \max \tau - h(x^* + s) \\ \text{s.t. } & s \in \mathcal{S} \end{aligned}$$

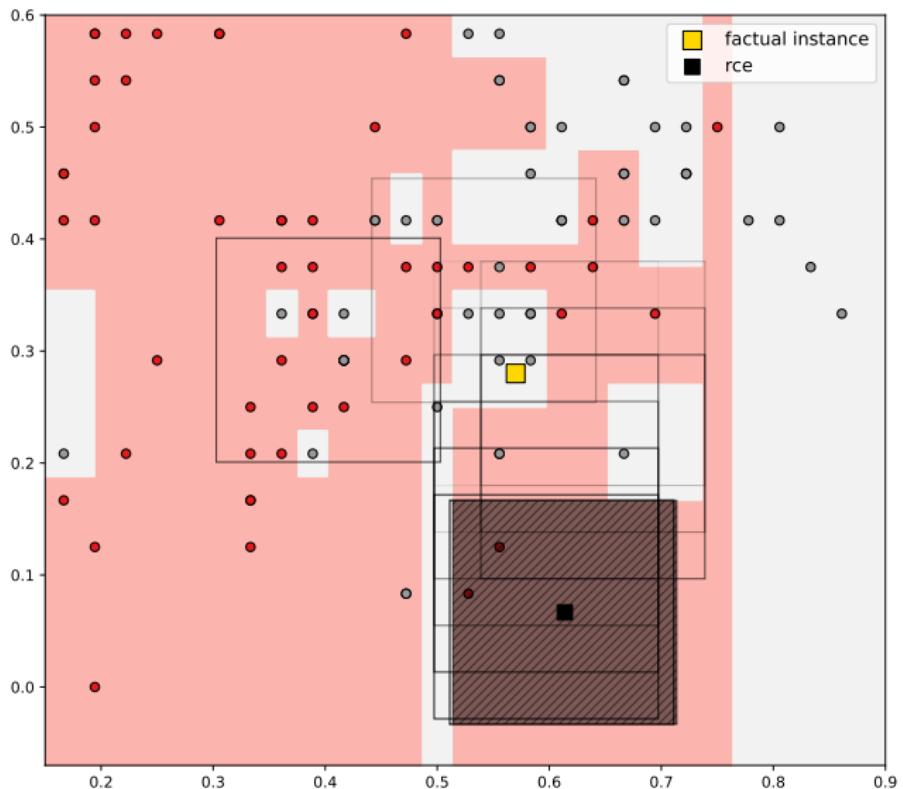
- Current solution x^* is cut-off if optimal value > 0

ADVERSARIAL APPROACH

Algorithm

1. **Input:** \hat{x} , \mathcal{S} , $\varepsilon > 0$
2. $\mathcal{Z} = \{0\}$
3. **Repeat:**
4. $x^* \leftarrow \text{MasterProblem}(\mathcal{Z})$
5. $s^*, \text{opt} \leftarrow \text{AdversarialProblem}(x^*, \mathcal{S})$
6. $\mathcal{Z} \leftarrow \mathcal{Z} \cup \{s^*\}$
7. **Until:** $\text{opt} \leq \varepsilon$
8. **Return:** x^*

ADVERSARIAL APPROACH



CONVERGENCE

Theorem (Mutapcic & Boyd)

If X is bounded and if h is a Lipschitz continuous function, i.e., there exists an $L > 0$ such that

$$|h(x_1) - h(x_2)| \leq L\|x_1 - x_2\|$$

for all $x_1, x_2 \in X$. Then the algorithm terminates after a finite number of steps with a solution x^* such that

$$h(x^* + s) \geq \tau - \varepsilon$$

for all $s \in \mathcal{S}$.

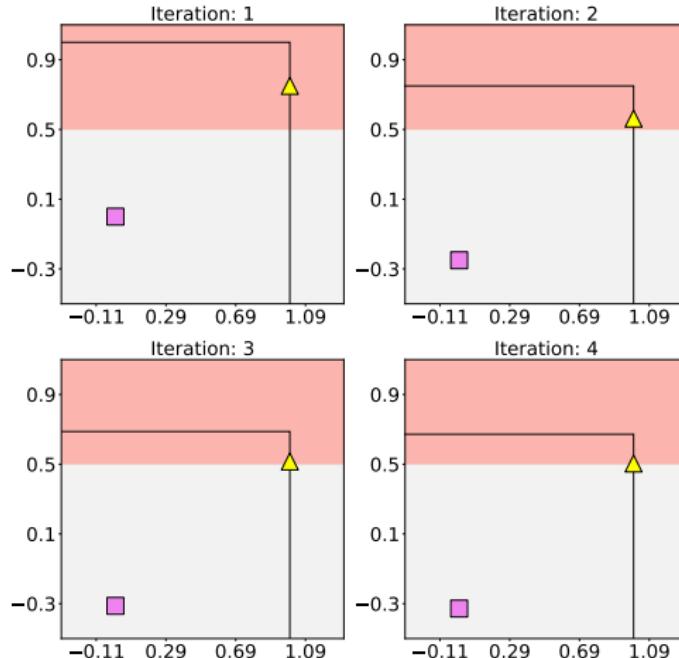
Lipschitz Continuity

Classifier function h of

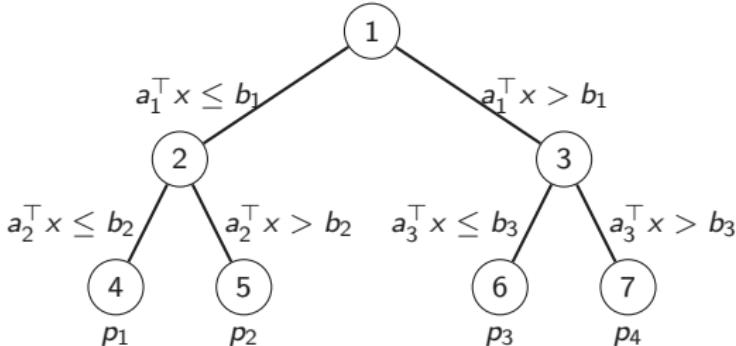
- logistic regression is Lipschitz continuous.
- neural networks with ReLU activation functions is Lipschitz continuous.
- decision trees is not even continuous!

LIPSCHITZ CONTINUITY IS NEEDED

- Classifier $h : \mathbb{R}^2 \rightarrow [0, 1]$ with $h(x) = 0$ if $x_2 \geq \frac{1}{2}$ and $h(x) = 1$ otherwise.
- Factual instance is $\hat{z} = (0, 2)$
- Optimal MP solution in iteration i : $x^i = (0, -\sum_{j=1}^i (\frac{1}{4})^j)$
- Optimal AP solution in iteration i : $s^i = (1, \frac{1}{2} + \sum_{j=1}^i (\frac{1}{4})^j)$



DECISION TREES



Define the continuous function

$$\tilde{h}(x) := \begin{cases} \tau, & x \in \mathcal{L}_i, p_i \geq \tau; \\ \max_{j \in \mathcal{N}_{\leq}^i \cup \mathcal{N}_{<}^i} \{\tau + a_j^T x - b_j\}, & x \in \mathcal{L}_i, p_i < \tau. \end{cases}$$

where $\mathcal{L}_i = \left\{ x : a_j^T x \leq b_j, a_k^T x < b_k, j \in \mathcal{N}_{\leq}^i, k \in \mathcal{N}_{<}^i \right\}$.

Lemma

The function \tilde{h} is Lipschitz continuous and classifies each point which is not on the boundary of a leaf the same as the original decision tree.

MASTER PROBLEM: DECISION TREES

The **master problem** for a trained decision tree can be reformulated as:

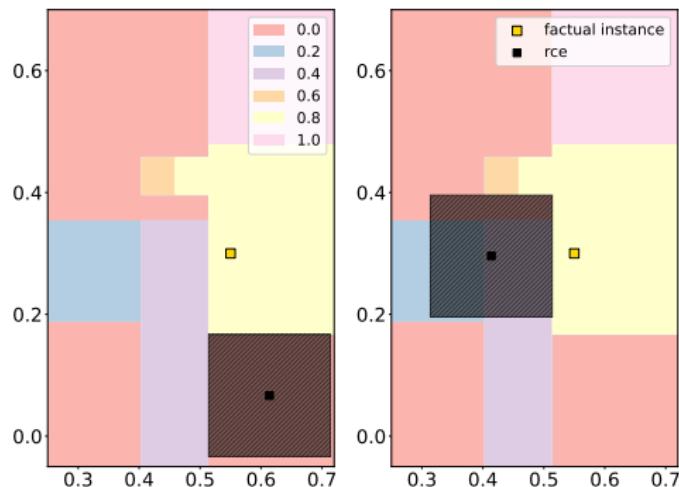
$$\min d(\hat{x}, x)$$

$$\text{s.t. } A^i(x + s) \leq b^i + M(1 - l_i(s)) \\ s \in \mathcal{Z}, \text{ leaf } i$$

$$\sum_{\text{leaves } i} l_i(s) = 1 \quad s \in \mathcal{Z}$$

$$\sum_{\text{leaves } i} l_i(s)p_i \geq \tau \quad s \in \mathcal{Z}$$

$$l_i(s) \in \{0, 1\} \quad s \in \mathcal{Z}, \text{ leaf } i$$

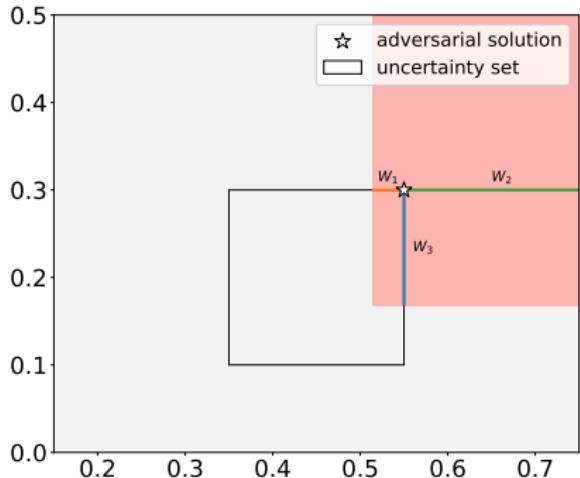


ADVERSARIAL PROBLEM: DECISION TREES

The **adversarial problem** can be solved by solving

$$\begin{aligned} & \max \alpha \\ \text{s.t. } & \alpha \leq w \\ & A^i(x^* + s) + w \leq b^i \\ & s \in S, w \geq 0 \end{aligned}$$

for every leaf i with $p_i < \tau$ and choosing the best optimal value.



ADVERSARIAL APPROACH: DECISION TREES

