

# Anchors



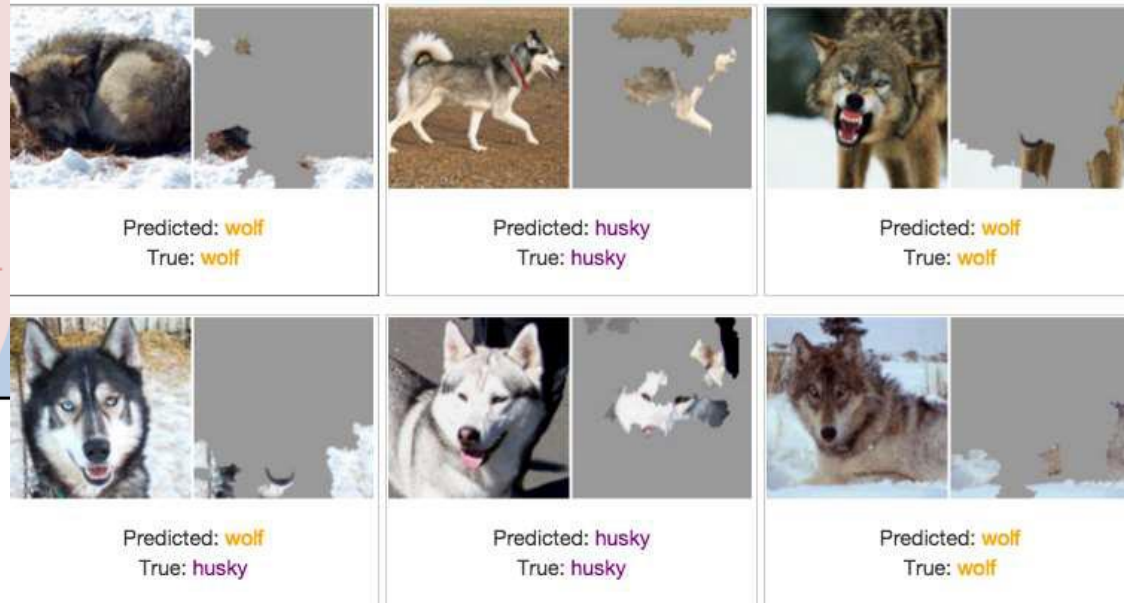
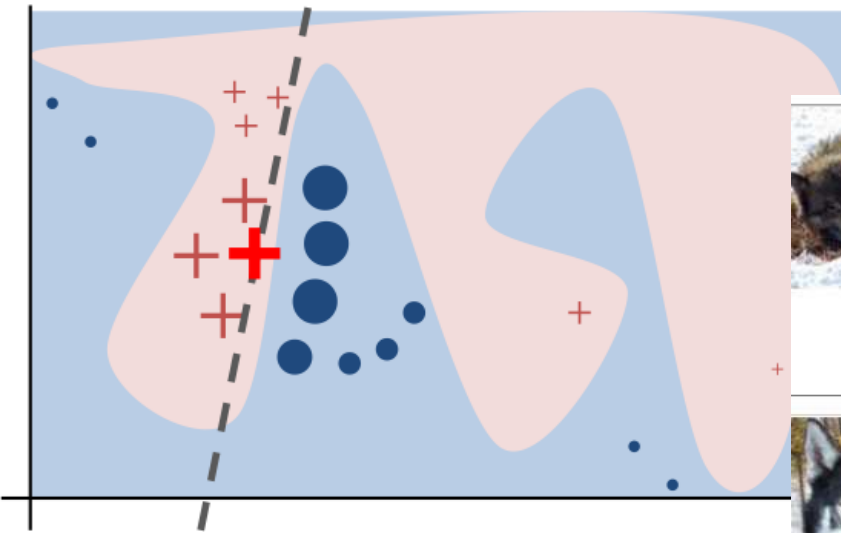
**Alicja Gosiewska**

**Wydział Matematyki i Nauk Informacyjnych  
Politechniki Warszawskiej**

# LIME



## Local Interpretable Model-Agnostic Explanations

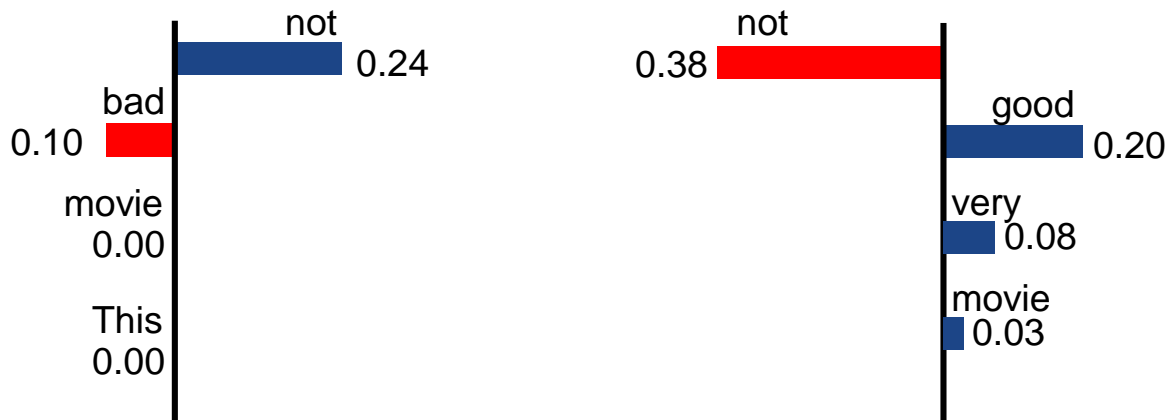


+ This movie is not bad.

— This movie is not very good.



## LIME explanations



## Anchor explanations

{"not", "bad"} →

Positive

{"not", "good"} →

Negative

$f : X \rightarrow Y$  - black box model

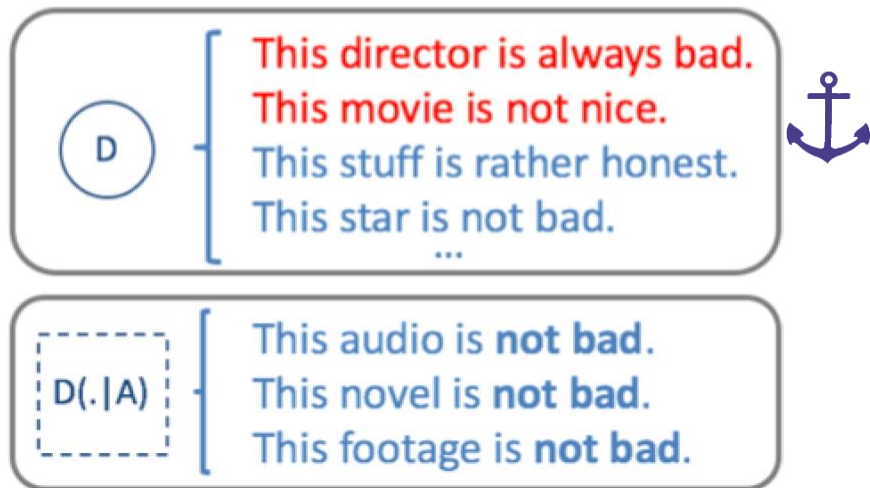
$x \in X$  - an instance to be explained

$\mathcal{D}_x$  - perturbation distribution

$A$  - a rule (set of predicates)

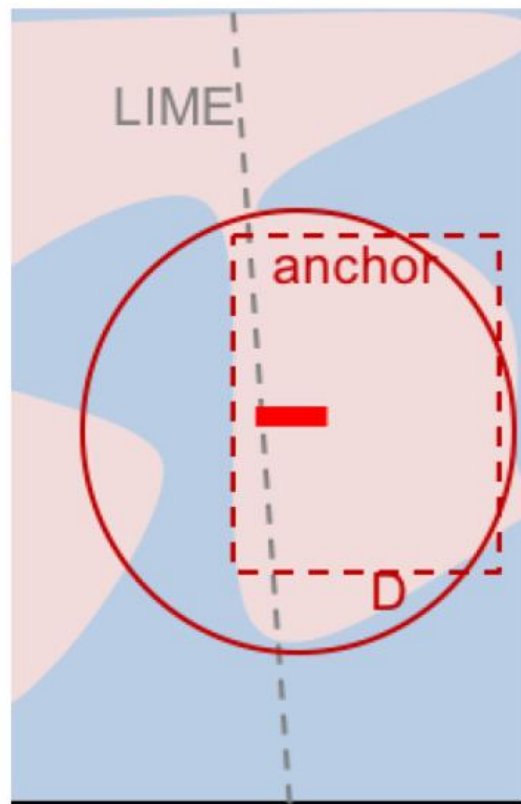
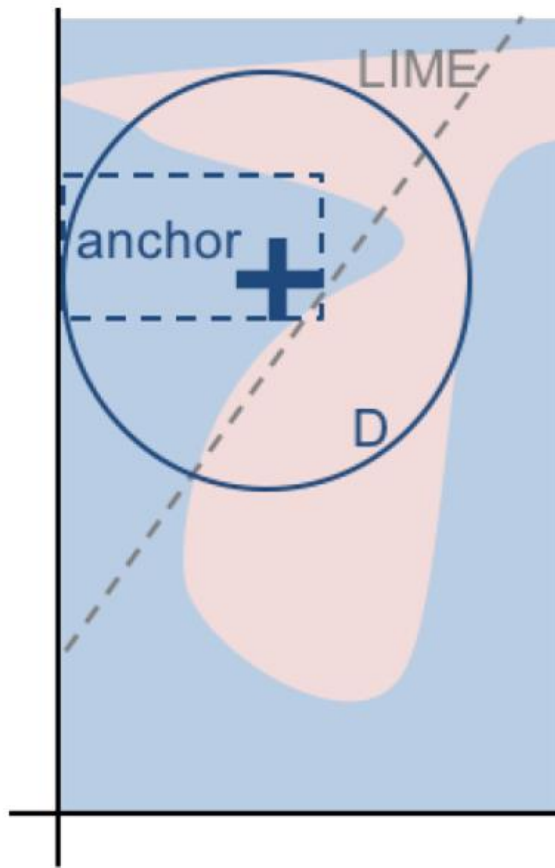
$\mathcal{D}_x(\cdot|A)$  - conditional distribution when the rule  $A$  applies

**+** This movie is not bad.



$A$  is an anchor if:

$$\mathbb{E}_{\mathcal{D}}(z|A) [\mathbb{1}_{f(x)=f(z)}] \geq \tau$$



	<b>If</b>	<b>Predict</b>
<b>adult</b>	No capital gain or loss, never married	$\leq 50K$
	Country is US, married, work hours $> 45$	$> 50K$



An anchor  $A$  is a set of feature predicates that achieves  $\text{prec}(A) \geq \tau$ , where

$$\text{prec}(A) = \mathbb{E}_{\mathcal{D}(z|A)} [\mathbb{1}_{f(x)=f(z)}]$$

Probabilistic definition:

$$P(\text{prec}(A) \geq \tau) \geq 1 - \delta$$

Coverage of an anchor:

$$\text{cov}(A) = \mathbb{E}_{\mathcal{D}(z)} [A(z)]$$



Search for an anchor is the following combinatorial optimization problem:

$$\max_{A \text{ s.t. } P(\text{prec}(A) \geq \tau) \geq 1 - \delta} \text{cov}(A)$$

---

**Algorithm 1** Identifying the *Best* Candidate for Greedy

---

**function** GenerateCands( $\mathcal{A}, c$ )
$$\mathcal{A}_r = \emptyset$$
**for all**  $A \in \mathcal{A}; a_i \in x, a_i \notin A$  **do**

**if**  $\text{cov}(A \wedge a_i) > c$  **then**                      {Only high-coverage}

$$\mathcal{A}_r \leftarrow \mathcal{A}_r \cup (A \wedge a_i) \quad \{\text{Add as potential anchor}\}$$

**return**  $\mathcal{A}_r$                       {Candidate anchors for next round}

**function** BestCand( $\mathcal{A}, \mathcal{D}, \epsilon, \delta$ )

**initialize**  $\text{prec}, \text{prec}_{ub}, \text{prec}_{lb}$  estimates  $\forall A \in \mathcal{A}$

$$A \leftarrow \arg \max_A \text{prec}(A)$$
$$A' \leftarrow \arg \max_{A' \neq A} \text{prec}_{ub}(A', \delta) \quad \{\delta \text{ implicit below}\}$$
**while**  $\text{prec}_{ub}(A') - \text{prec}_{lb}(A) > \epsilon$  **do**

**sample**  $z \sim \mathcal{D}(z|A), z' \sim \mathcal{D}(z'|A')$       {Sample more}

**update**  $\text{prec}, \text{prec}_{ub}, \text{prec}_{lb}$  for  $A$  and  $A'$ 
$$A \leftarrow \arg \max_A \text{prec}(A)$$
$$A' \leftarrow \arg \max_{A' \neq A} \text{prec}_{ub}(A')$$
**return**  $A$ 



Method	Precision		Coverage		Time/pred	
	adult	rcdv	adult	rcdv	adult	rcdv
No expls	<u>54.8</u>	<u>83.1</u>	<u>79.6</u>	<u>63.5</u>	<u>29.8</u> $\pm 14$	<u>35.7</u> $\pm 26$
LIME(1)	<u>68.3</u>	98.1	<u>89.2</u>	<u>55.4</u>	<u>28.5</u> $\pm 10$	<u>24.6</u> $\pm 6$
Anchor(1)	<u>100.0</u>	97.8	<u>43.1</u>	<u>24.6</u>	<u>13.0</u> $\pm 4$	<u>14.4</u> $\pm 5$



Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin.

[Anchors: High-Precision Model-Agnostic Explanations](#)

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin.

[Nothing Else Matters: Model-Agnostic Explanations  
By Identifying Prediction Invariance](#)

[The anchor\\_exp python package](#)

