

Proxy Methods for Post Hoc Interpretability

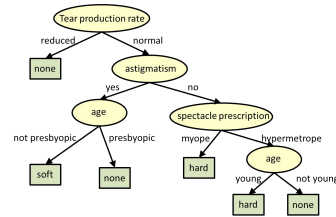
Apr 17, 2020

Dr. Wei Wei, Prof. James Landay

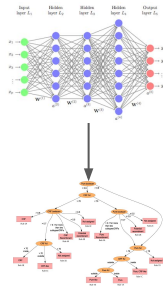
CS 335: Fair, Accountable, and Transparent (FAcCT) Deep Learning
Stanford University

Recap

Previous Lecture: Intrinsic Methods for Interpretability



This Lecture: Post Hoc Methods for Interpretability



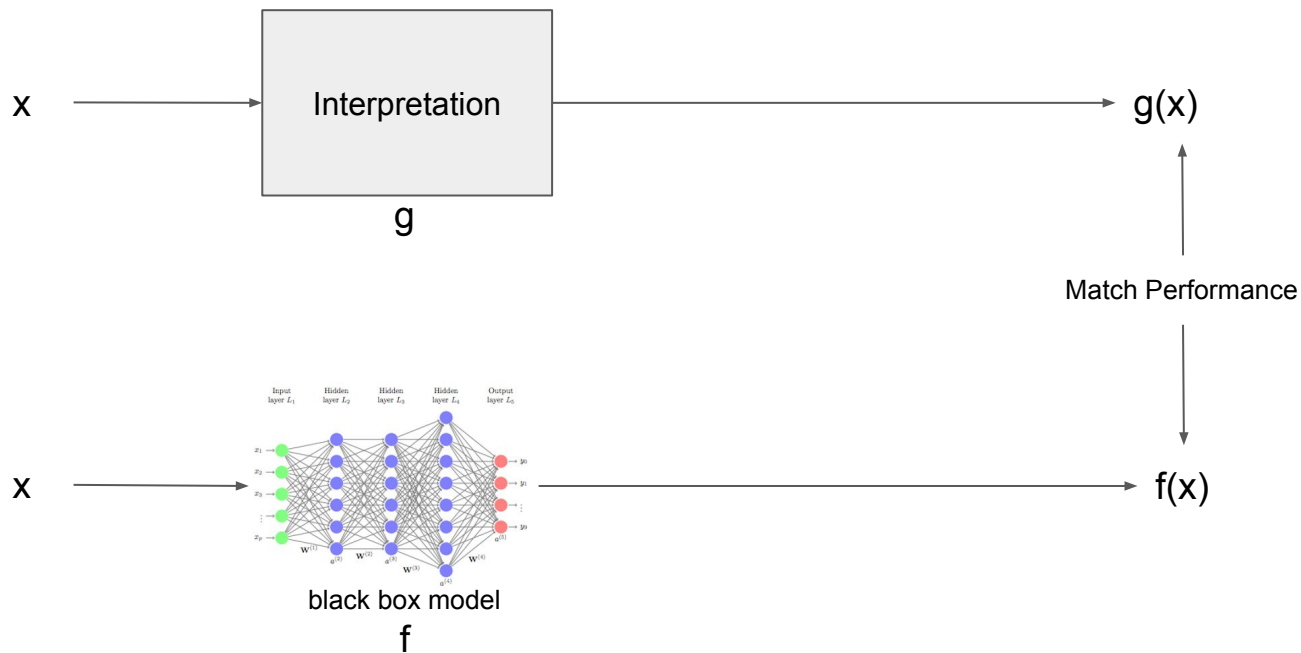
Outline

- Post Hoc Interpretability
 - Proxy Models
- Local Surrogate Methods
 - LIME
- Rule Based Learner
 - Anchors

Post Hoc Interpretability

- Model Agnostic
 - Can be applied across many different black box models
 - Multiple techniques can be applied at the same time
- Availability
 - Do not require training data
 - Do not require model training/fine-tuning
- No Performance Degeneration
 - Will not alter the black box model

Proxy Models for Post Hoc Interpretability



Outline

- Post Hoc Interpretability
 - Proxy Models
- Local Surrogate Methods
 - LIME
- Rule Based Learner
 - Anchors

Local Surrogate Methods

- Local surrogate methods aim at finding explanation g to approximate f around x based on Model Fidelity

neighbors of x

↓

$$F(f, g, N_x) := \mathbb{E}_{x' \sim N_x} [(g(x') - f(x'))^2]$$

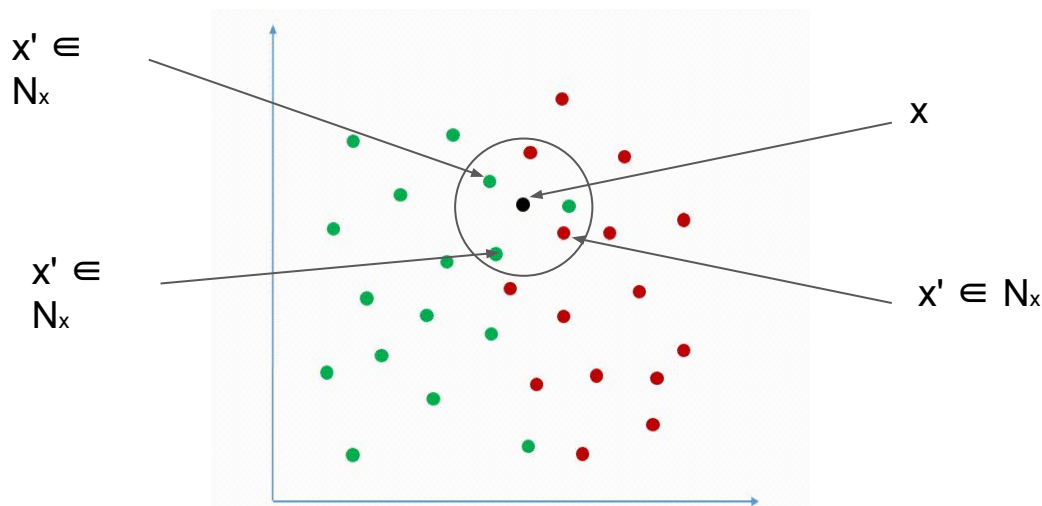
↑ ↑

explanation blackbox

Local Surrogate Methods

- Local surrogate methods aim at finding explanation g to approximate f around x based on Model Fidelity

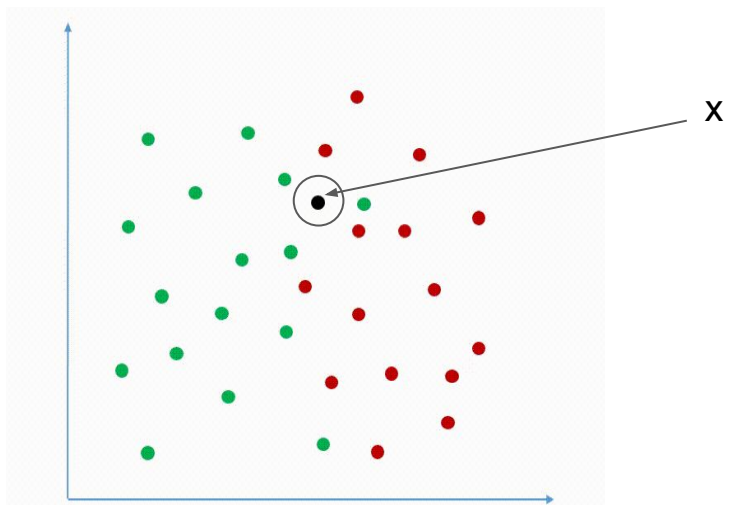
$$F(f, g, N_x) := \mathbb{E}_{x' \sim N_x} [(g(x') - f(x'))^2]$$



Local Surrogate Methods

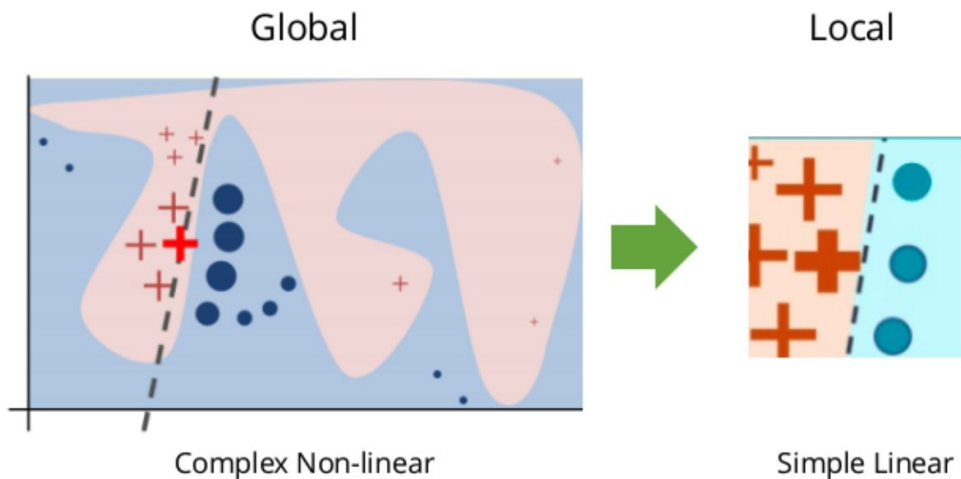
- Local surrogate methods aim at finding explanation g to approximate f around x based on Model Fidelity

$$F(f, g, N_x) := \mathbb{E}_{x' \sim N_x} [(g(x') - f(x'))^2]$$



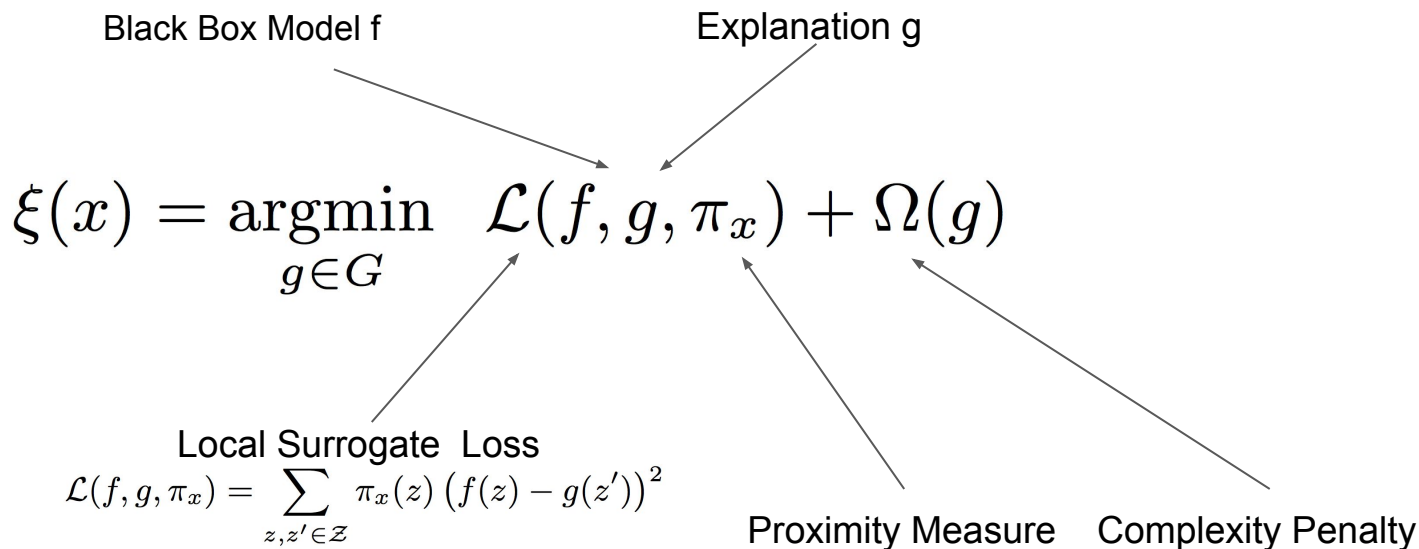
Local Interpretable Model-agnostic Explanations (LIME)

- Deep learning models are usually too complex for global interpretation
 - Instead, we seek for local interpretability using simple interpretable models (e.g. linear models)



LIME

- LIME generates an explainable model that optimizes both model fidelity and explanation



Linear Explainable Model

$$\text{Local Surrogate Loss } \mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in \mathcal{Z}} \pi_x(z) (f(z) - g(z'))^2$$

- Linear Explainable Model
 - We use a linear model for explanation $g(z') = w_g \cdot z'$, $z' \in \{0, 1\}^d$

Linear Explainable Model

$$\text{Local Surrogate Loss } \mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in \mathcal{Z}} \pi_x(z) (f(z) - g(z'))^2$$

- Linear Explainable Model

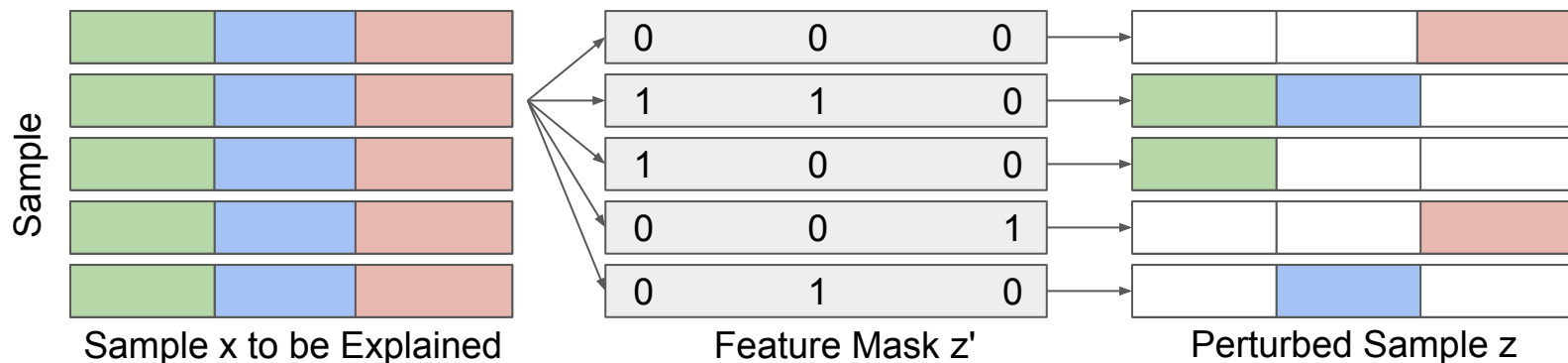
- We use a linear model for explanation $g(z') = w_g \cdot z'$, $z' \in \{0, 1\}^d$
- z' is a feature mask indicating whether a specific input will be included in the explanation
- A perturbed sample z can be recovered from mask z' , $z = h_x(z')$

Linear Explainable Model

$$\text{Local Surrogate Loss } \mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in \mathcal{Z}} \pi_x(z) (f(z) - g(z'))^2$$

- Linear Explainable Model

- We use a linear model for explanation $g(z') = w_g \cdot z'$, $z' \in \{0, 1\}^d$
- z' is a feature mask indicating whether a specific input will be included in the explanation
- A perturbed sample z can be recovered from mask z' , $z = h_x(z')$



Training Objective for LIME

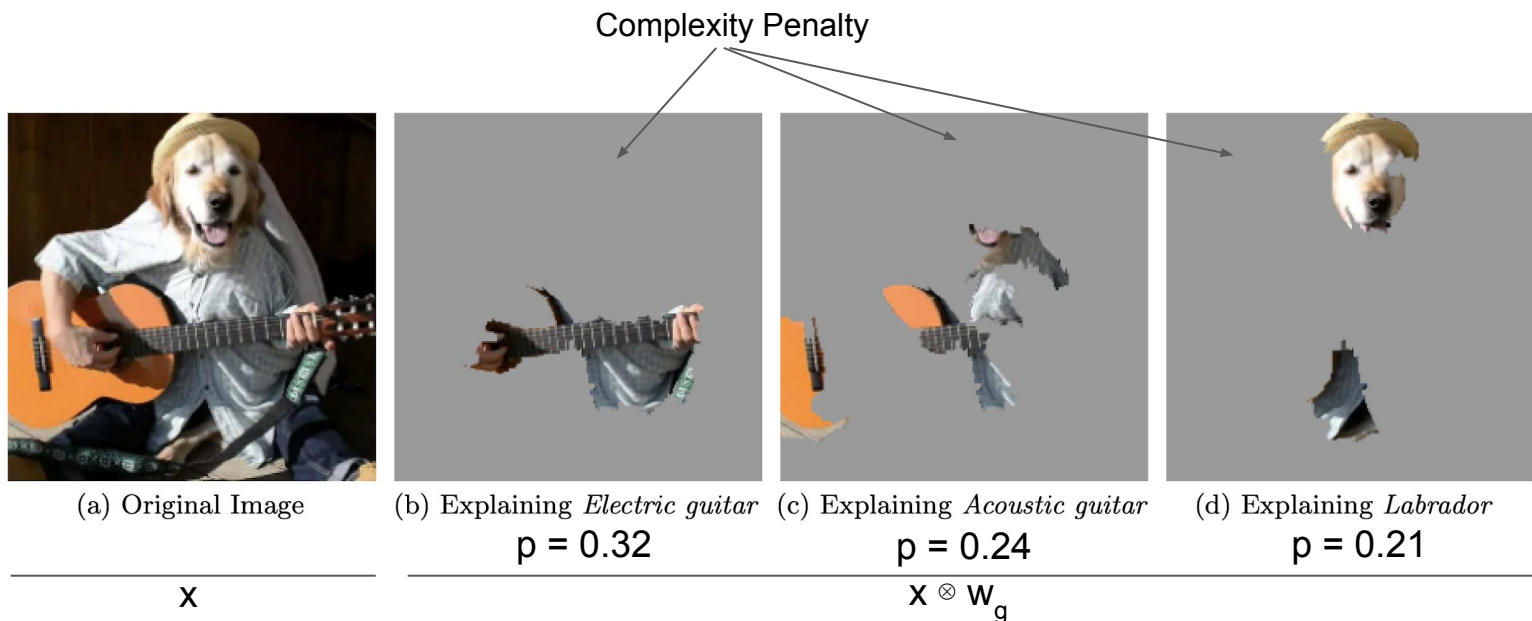
- Loss Function

- Match predictions of the explanation model g with that of the black box model f around x
- We use an exponentially scaled function to measure proximity
 - D = cosine distance for text
 - D = L2 distance for images

$$\pi_x(z) = \exp(-D(x, z)^2 / \sigma^2)$$

$$\epsilon(x) = \arg \min_{g \in G} \underbrace{\sum_{z, z'} \pi_x(z) (f(z) - g(z'))^2}_{\text{Local Surrogate Loss}} + \underbrace{\infty \cdot \mathbb{1}_{\|w_g\|_0 > K}}_{\text{Complexity Penalty}}$$

Explaining Google InceptionNet

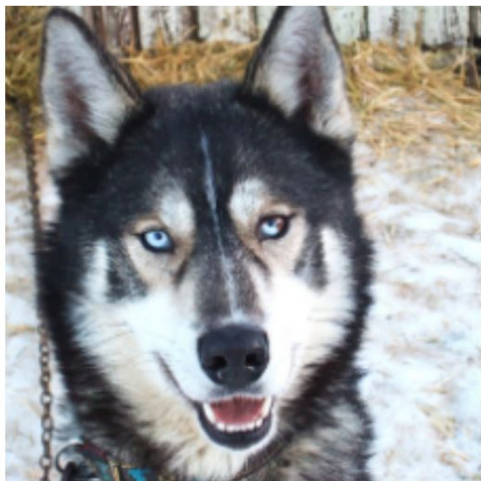


$$\epsilon(x) = \arg \min_{g \in G} \sum_{z, z'} \pi_x(z) (f(z) - g(z'))^2 + \infty \cdot \mathbb{1}_{\|w_g\|_0 > K}$$

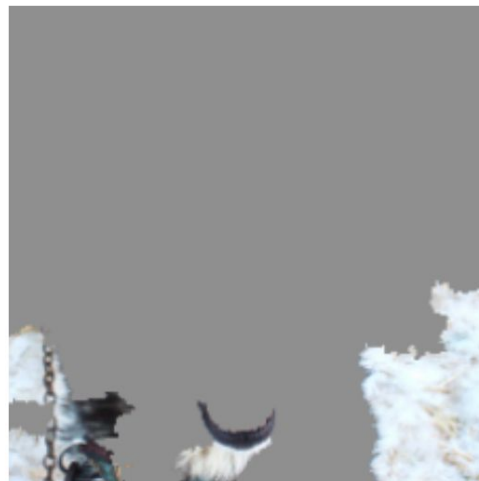
[Ribeiro et al. 2016](#)

Example for Bad ML Predictions

- Explanations on a model that misclassified Husky as Wolf



(a) Husky classified as wolf



(b) Explanation

Explaining Text Classifiers

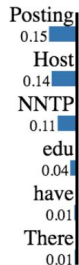
- Explanations for a SVM classifier with 94% accuracy
 - Predictions are made for arbitrary reasons
 - The word “Posting” appears in 22% of examples in the training set
 - 99% of which are samples attribute to class “Atheism”

Prediction probabilities



atheism

christian



Text with highlighted words

From: johnchad@triton.unm.edu (jchadwic)
Subject: Another request for Darwin Fish
Organization: University of New Mexico, Albuquerque
Lines: 11

NNTP-Posting-Host: triton.unm.edu

Hello Gang,

There have been some notes recently asking where to obtain the DARWIN fish.

This is the same question I have and I have not seen an answer on the

net. If anyone has a contact please post on the net or email me.

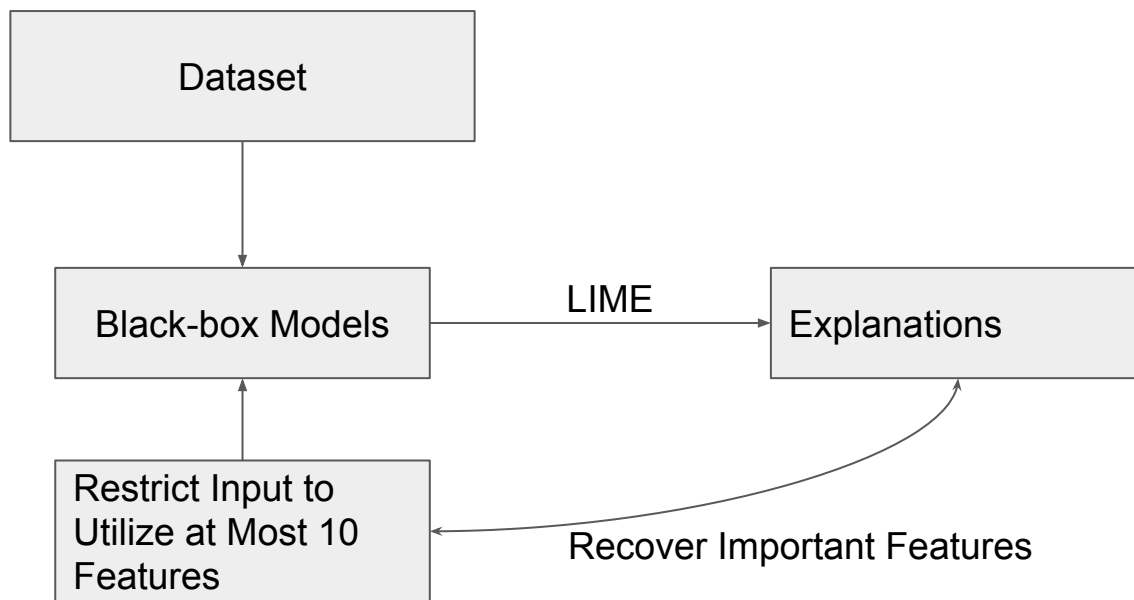
$f(x)$

W_g

$x \otimes W_g$

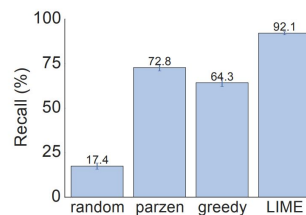
[Ribeiro et al, 2016](#)

Faithfulness of Explanations

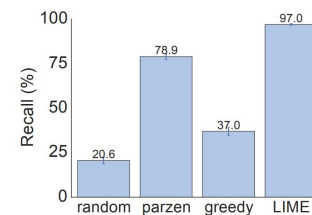


Faithfulness of Explanations

- LIME Achieves Good Faithfulness
- Sentiments classification tasks
 - Books, DVDs
- Classifiers
 - logistic regression with L2 reg. (Sparse LR)
 - decision tree

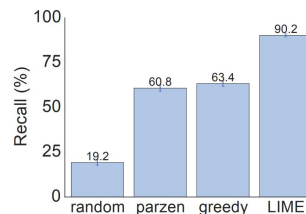


(a) Sparse LR

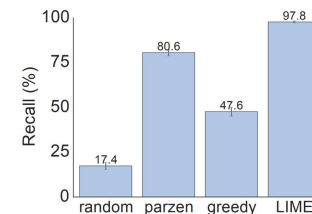


(b) Decision Tree

Books Dataset



(a) Sparse LR



(b) Decision Tree

DVDs Dataset

parzen - [Baehrens et al. 2010](#)
random - randomly pick K features
greedy - remove features contribute most to the classifiers

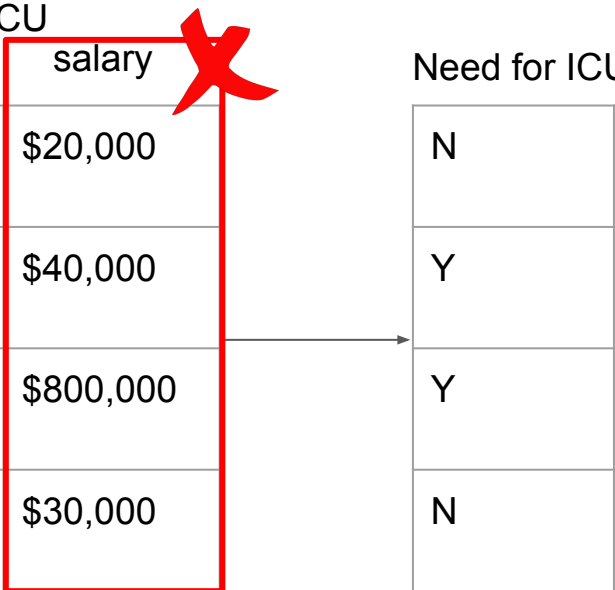
[Ribeiro et al. 2016](#)

Trustworthiness for ML Models

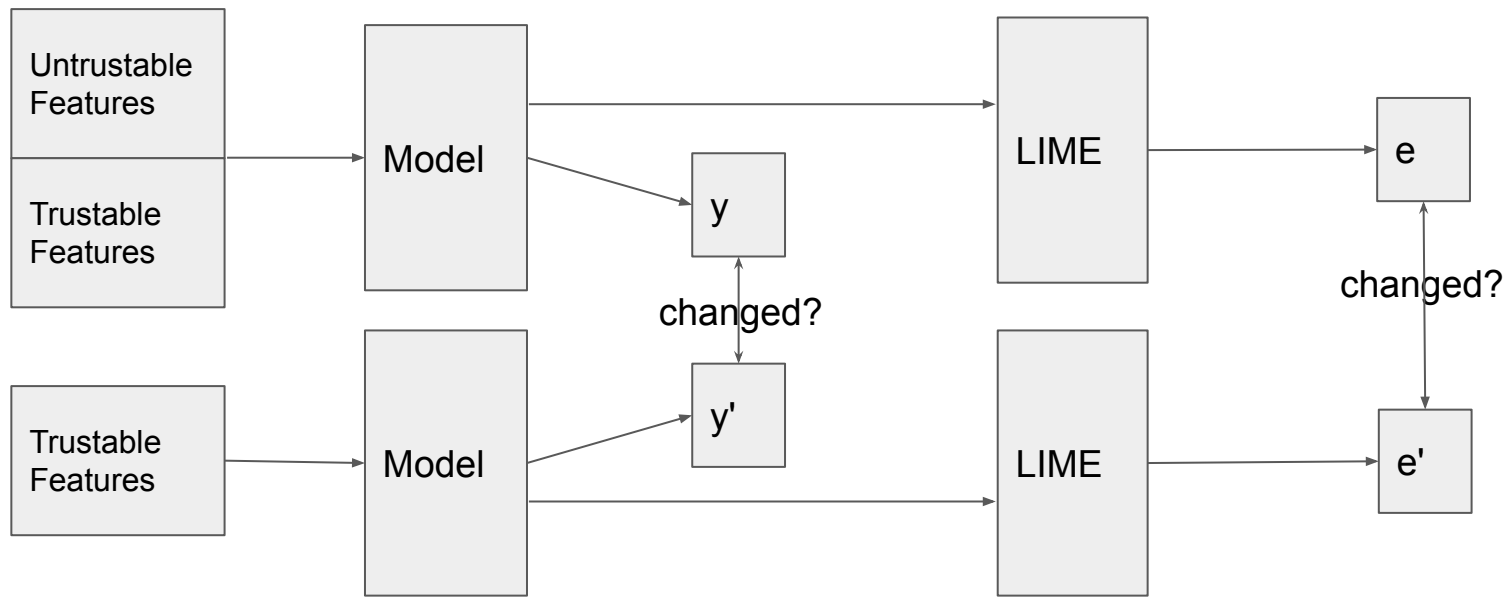
- Human discredits certain features in the learning tasks
- Classifiers that use those features will be considered not trustable.

Predict the need for ICU

heart beat	temperature	salary	Need for ICU?
120 BPM	101 F	\$20,000	N
80 BPM	104.4 F	\$40,000	Y
140 BPM	99 F	\$800,000	Y
110 BPM	100 F	\$30,000	N



Trustworthiness for Explanations



Trustworthiness of Predictions

- Untrustable Features
 - 25% of features are "untrustable features"
- Trustworthiness of Predictions
 - Compares changes of model predictions and the changes of model explanations when unstable features are removed

	Books				DVDs			
	LR	NN	RF	SVM	LR	NN	RF	SVM
Random	14.6	14.8	14.7	14.7	14.2	14.3	14.5	14.4
Parzen	84.0	87.6	94.3	92.3	87.0	81.7	94.2	87.3
Greedy	53.7	47.4	45.0	53.3	52.4	58.1	46.6	55.1
LIME	96.6	94.5	96.2	96.7	96.6	91.8	96.1	95.6

Trustworthiness of LIME with different ML models:

- Logistic Regression with L2 regularization (LR)
- Nearest Neighbors (NN)
- Random Forests (RF)
- Support Vector Machines (SVM)

[Ribeiro et al, 2016](#)

Explaining Multiple Samples

- Explain a set of samples to get a complete picture of the model
 - Each sample $x_i \in X$ will have its interpretation

$$g_{x_i}(z) = w_i \cdot z = \sum_j w_{i,j} \cdot z_j$$

- How do we select samples?
 - Select samples to cover the maximum information about the model

$$I_j = \sqrt{\sum_{x_i \in X} |w_{i,j}|}$$

Explaining Multiple Samples

- How do we select samples?
 - Select samples to cover the maximum information about the model

$$I_j = \sqrt{\sum_{x_i \in X} |w_{i,j}|}$$

- Set function

$$c(V, \mathcal{W}, I) = \sum_{j=1}^{d'} \mathbb{1}_{[\exists i \in V: w_{i,j} > 0]} I_j$$

- We want to get a set of samples V up to B elements that maximize c

$$Pick(\mathcal{W}, I) = \operatorname{argmax}_{V, |V| \leq B} c(V, \mathcal{W}, I)$$

Explaining Multiple Samples

- How do we select samples?
 - We want to get a set of samples V up to B elements such to maximize c

$$c(V, \mathcal{W}, I) = \sum_{j=1}^{d'} \mathbb{1}_{[\exists i \in V: \mathcal{W}_{ij} > 0]} I_j$$

- Optimization
 - Searching for the global optimal set of V is NP-Hard ([Feige, 1998](#))
 - We turn to greedy algorithm as an approximation method

Greedy Algorithm for Sample Selection

- Pick a subset of samples up to B elements from X to maximize c

$$\text{Pick}(\mathcal{W}, I) = \underset{V, |V| \leq B}{\operatorname{argmax}} c(V, \mathcal{W}, I)$$

- Start with an empty set $V_0 = \emptyset$,
- For the i^{th} step
 - Pick the next element $x_i \in X \setminus V_i$, such that x_i maximizes $c(V_i \cup \{x_i\}, \mathcal{W}, I)$
 - repeat until $|V_i| \geq B$

Theoretical Guarantees on Performance

- A function defined on a set is submodular if
 - for every $V_A \subseteq V_B$

$$c(V_A \cup \{x\}) - c(V_A) \geq c(V_B \cup \{x\}) - c(V_B)$$

- Properties of Submodular functions
 - The performance of a greedy algorithm is at least $1-1/e$ (~63%) to the optimum
 - $c(V, \mathcal{W}, I) = \sum_{j=1}^{d'} \mathbb{1}_{\{\exists i \in V: \mathcal{W}_{ij} > 0\}} I_j$ is submodular
 - The performance of a greedy algorithm on c is guaranteed with a lower bound

Human Experiments

- Ask Human to Select the Best Classifier
 - Annotators are shown the explanations
 - Annotators have no knowledge in machine learning

Example #3 of 6 True Class: ● Atheism Instructions Previous Next

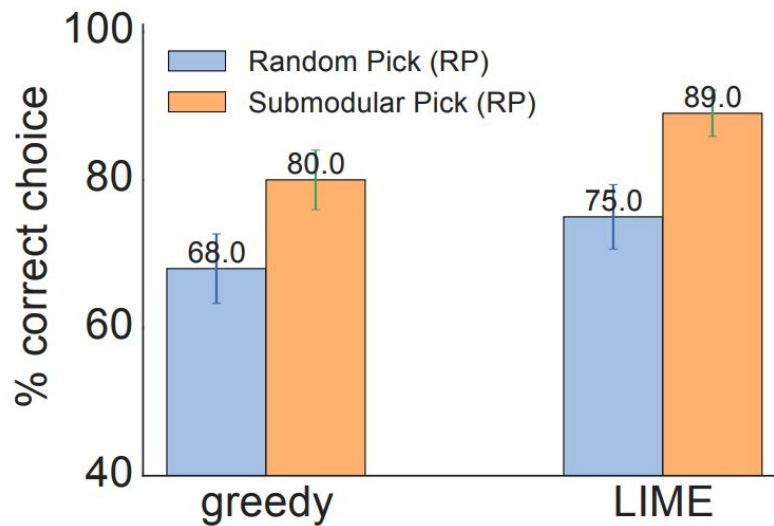
Algorithm 1	Algorithm 2																								
<p>Words that A1 considers important:</p> <table><tr><td>GOD</td><td><div style="width: 80%; height: 10px; background-color: magenta;"></div></td></tr><tr><td>mean</td><td><div style="width: 60%; height: 10px; background-color: magenta;"></div></td></tr><tr><td>anyone</td><td><div style="width: 50%; height: 10px; background-color: green;"></div></td></tr><tr><td>this</td><td><div style="width: 40%; height: 10px; background-color: green;"></div></td></tr><tr><td>Koresh</td><td><div style="width: 20%; height: 10px; background-color: magenta;"></div></td></tr><tr><td>through</td><td><div style="width: 15%; height: 10px; background-color: green;"></div></td></tr></table>	GOD	<div style="width: 80%; height: 10px; background-color: magenta;"></div>	mean	<div style="width: 60%; height: 10px; background-color: magenta;"></div>	anyone	<div style="width: 50%; height: 10px; background-color: green;"></div>	this	<div style="width: 40%; height: 10px; background-color: green;"></div>	Koresh	<div style="width: 20%; height: 10px; background-color: magenta;"></div>	through	<div style="width: 15%; height: 10px; background-color: green;"></div>	<p>Words that A2 considers important:</p> <table><tr><td>Posting</td><td><div style="width: 80%; height: 10px; background-color: magenta;"></div></td></tr><tr><td>Host</td><td><div style="width: 80%; height: 10px; background-color: magenta;"></div></td></tr><tr><td>Re</td><td><div style="width: 40%; height: 10px; background-color: magenta;"></div></td></tr><tr><td>by</td><td><div style="width: 30%; height: 10px; background-color: green;"></div></td></tr><tr><td>in</td><td><div style="width: 30%; height: 10px; background-color: green;"></div></td></tr><tr><td>Nntp</td><td><div style="width: 15%; height: 10px; background-color: magenta;"></div></td></tr></table>	Posting	<div style="width: 80%; height: 10px; background-color: magenta;"></div>	Host	<div style="width: 80%; height: 10px; background-color: magenta;"></div>	Re	<div style="width: 40%; height: 10px; background-color: magenta;"></div>	by	<div style="width: 30%; height: 10px; background-color: green;"></div>	in	<div style="width: 30%; height: 10px; background-color: green;"></div>	Nntp	<div style="width: 15%; height: 10px; background-color: magenta;"></div>
GOD	<div style="width: 80%; height: 10px; background-color: magenta;"></div>																								
mean	<div style="width: 60%; height: 10px; background-color: magenta;"></div>																								
anyone	<div style="width: 50%; height: 10px; background-color: green;"></div>																								
this	<div style="width: 40%; height: 10px; background-color: green;"></div>																								
Koresh	<div style="width: 20%; height: 10px; background-color: magenta;"></div>																								
through	<div style="width: 15%; height: 10px; background-color: green;"></div>																								
Posting	<div style="width: 80%; height: 10px; background-color: magenta;"></div>																								
Host	<div style="width: 80%; height: 10px; background-color: magenta;"></div>																								
Re	<div style="width: 40%; height: 10px; background-color: magenta;"></div>																								
by	<div style="width: 30%; height: 10px; background-color: green;"></div>																								
in	<div style="width: 30%; height: 10px; background-color: green;"></div>																								
Nntp	<div style="width: 15%; height: 10px; background-color: magenta;"></div>																								
<p>Predicted:</p> <p>● Atheism</p> <p>Prediction correct:</p> <p>✓</p>	<p>Predicted:</p> <p>● Atheism</p> <p>Prediction correct:</p> <p>✓</p>																								
<p>Document</p> <p>From: pauld@verdix.com (Paul Durbin) Subject: Re: DAVID CORESH IS! GOD! Nntp-Posting-Host: sarge.hq.verdix.com Organization: Verdix Corp Lines: 8</p>	<p>Document</p> <p>From: pauld@verdix.com (Paul Durbin) Subject: Re: DAVID CORESH IS! GOD! Nntp-Posting-Host: sarge.hq.verdix.com Organization: Verdix Corp Lines: 8</p>																								

Classification of Atheism/Christian in the 20 newsgroups dataset

[Ribeiro et al, 2016](#)

Human Experiments - Select the Best Classifier

- Original model: SVM trained on the dataset with original features
- Cleaned model: SVM trained on the dataset with "cleaned features"



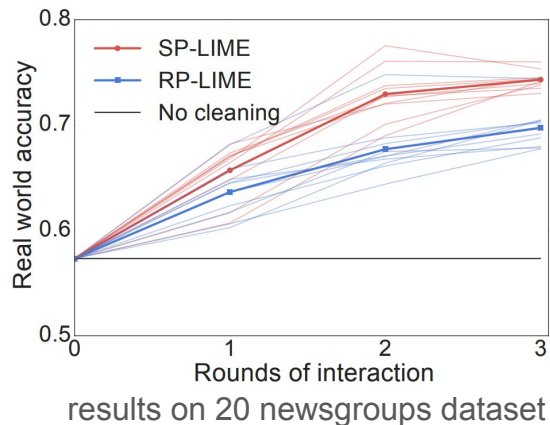
greedy - remove features
contribute most to the classifiers

[Ribeiro et al, 2016](#)

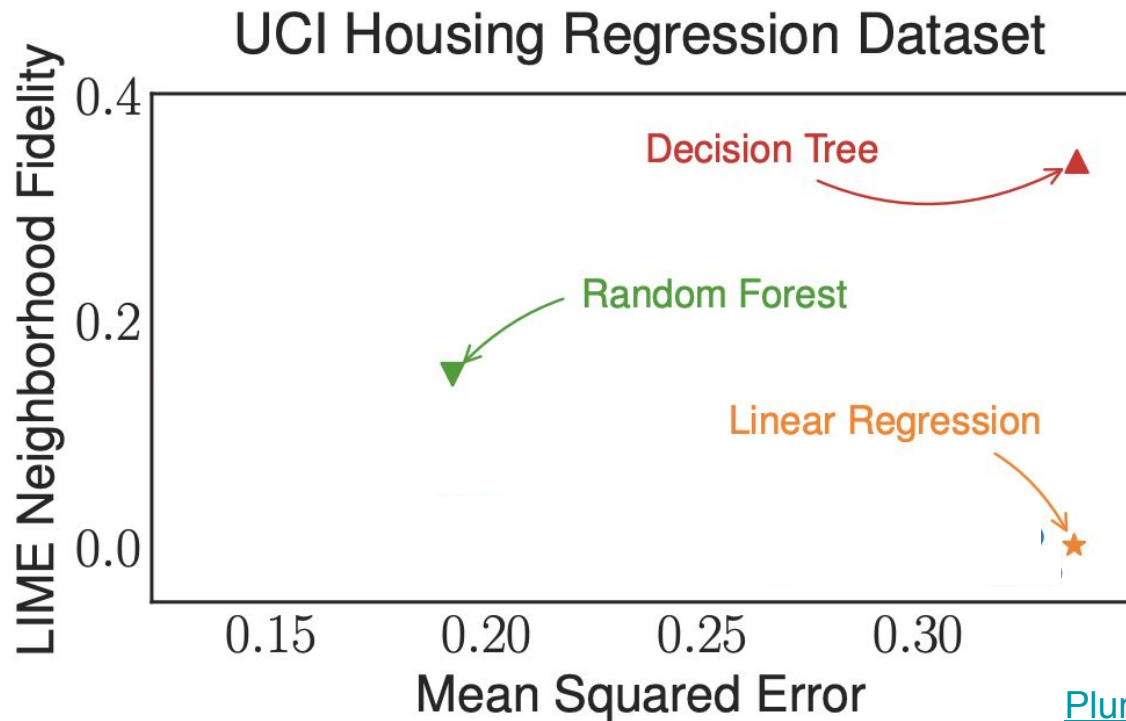
results on 20 newsgroups dataset

Improving Models Through ML Interpretability

- Improving ML Models
 - Human raters are shown model interpretability
 - They are asked to improve the model by masking out unnecessary features
 - Which words from the explanations should be removed from subsequent training
 - SP - select samples by random
 - RP - select samples by greedy algorithm



Faithfulness of Model Explanations



Outline

- Post Hoc Interpretability
 - Proxy Models
- Local Surrogate Methods
 - LIME
- Rule Based Explainers
 - Anchors

Rule Based Explainers

- Explain the Predictions of Deep Learning Models Using Rules
 - How do we find the set of rules for a particular predictor?

	If	Predict
adult	No capital gain or loss, never married	$\leq 50K$
	Country is US, married, work hours > 45	$> 50K$
rcdv	No priors, no prison violations and crime not against property	Not rearrested
	Male, black, 1 to 5 priors, not married, and crime not against property	Re-arrested
lending	FICO score ≤ 649	Bad Loan
	$649 \leq \text{FICO score} \leq 699$ and $\$5,400 \leq \text{loan amount} \leq \$10,000$	Good Loan

Anchors

- Generate A Set of Feature Predicates Known as Anchors A (i.e., rules)
 - Using anchors to explain the performance of deep model f
 - mimic the decisions of deep models on x , $f(x)$
 - explain a wide range of similar decisions in the dataset

	If	Predict
adult	No capital gain or loss, never married	$\leq 50K$
	Country is US, married, work hours > 45	$> 50K$
rcdv	No priors, no prison violations and crime not against property	Not rearrested
	Male, black, 1 to 5 priors, not married, and crime not against property	Re-arrested
lending	FICO score ≤ 649	Bad Loan
	$649 \leq \text{FICO score} \leq 699$ and $\$5,400 \leq \text{loan amount} \leq \$10,000$	Good Loan

Anchors found in adult income dataset

Anchors

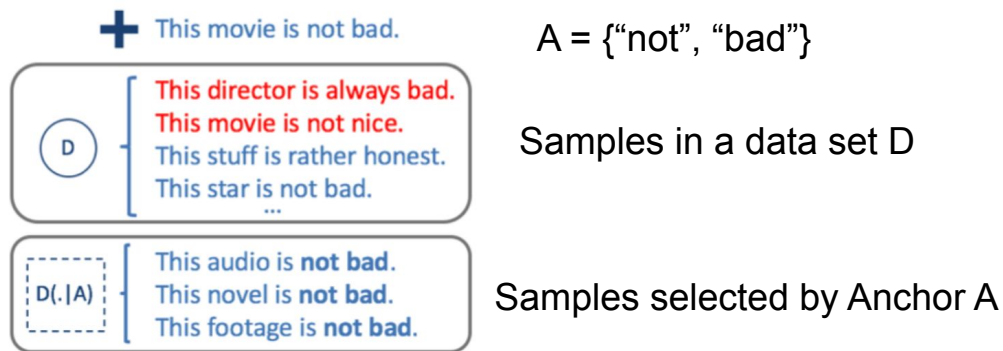
- An Anchor is a set of feature predicates applied to the feature space

$$A = \{\text{"not"}, \text{"bad"}\}$$

- Any text sample x containing both "not" and "bad" will be selected by the anchor

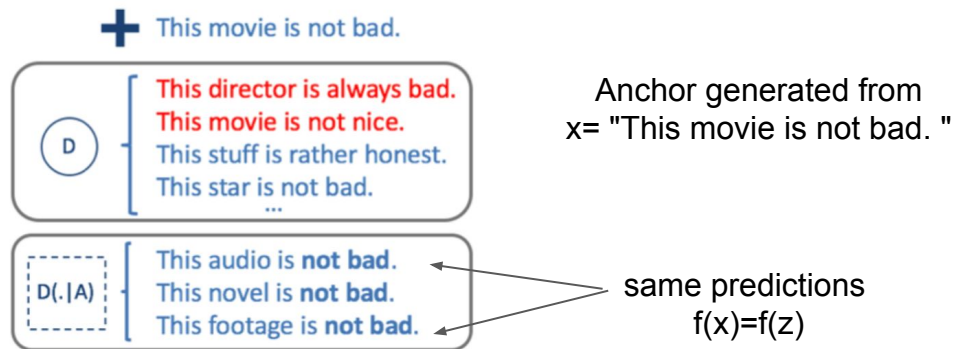
$$A(x) = 1$$

- An anchor can be applied to a dataset D to generate a subset $D(.|A)$



Formal Definitions of Anchors

- Preconditions of Anchors
 - Applies to the sample x being interpreted
 - Precisions
 - Samples covered by the same anchor A need to have the similar predictions
 - i.e., $f(x)=f(z)$ for $z \sim D(.|A)$
 - Coverage
 - A significant portion of the data needs to be covered by Anchor A .



Formal Definitions of Anchors

- Preconditions of Anchors

- Applies to the sample x being interpreted

$$A(x) = 1$$

- Precision

- Samples covered by the same anchor A need to have similar predictions

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

- Coverage

- A significant amount of data needs to be covered by one anchor A .

$$\mathbb{E}_{\mathcal{D}(z)}A(z) \geq c$$

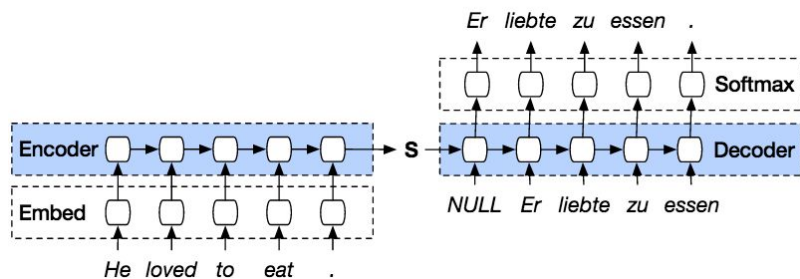
Anchors for Part of Speech Tagging

Instance	If	Predict
I want to play(V) ball.	previous word is PARTICLE	play is VERB.
I went to a play(N) yesterday.	previous word is DETERMINER	play is NOUN.
I play(V) ball on Mondays.	previous word is PRONOUN	play is VERB.

Anchors for Machine Translation

- Group Predictions of Words with Similar Meanings
 - "esta" (feminine of word "this")
 - "este" (masculine of word "this")
 - "isso" (if its referent is not in the sentence)

English	Portuguese
This is the question we must address	Esta é a questão que temos que enfrentar
This is the problem we must address	Este é o problema que temos que enfrentar
This is what we must address	É isso que temos de enfrentar



Anchors for Image Classification (InceptionV3)



original image



Anchors for "beagle"

Anchors for Visual Question Answering (VQA)



What animal is featured in this picture ? **dog**

What floor is featured in this picture? dog

What toenail is paired in this flowchart ? dog

What animal is shown on this depiction ? dog

Anchor for predicting "dog"

Where is the dog? on the floor

What color is the wall? white

When was this picture taken? during the day

Why is he lifting his paw? to play

Other Anchors

Generating Anchors

- Preconditions

- Precision

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

- Coverage

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

- Challenges in Generating the Optimal A

- Calculating precision and coverage is computationally intensive
 - will need to iterate through the predictions of f over the entire dataset
- Usually difficult to apply white box optimization techniques (e.g., gradient descent)

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

Generating Anchors

- Optimization Target

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

- Searching for the Optimal A
 - for each step t,
 - 1) Construct a set of candidate solutions with the best coverage
 - Candidate solutions need to satisfy $\text{cov}(A) \geq c$
 - 2) Pick top-k candidates with the best precision
 - Candidates need to have $\text{prec}(A) \geq \tau$ with confidence at least $1 - \delta$
 - 3) Update the optimal Anchor A^*

Generating Anchors - Optimizing Coverage

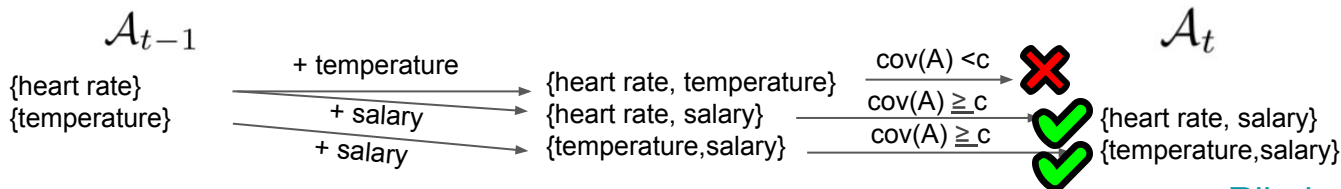
- Searching for the Optimal A
 - 1) Optimizing coverage with $\text{cov}(A) \geq c$
 - 2) Optimizing precision with $\text{prec}(A) \geq \tau$ and confidence at least $1-\delta$
 - 3) Update the optimal solution A^*

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

heart beat	temperature	salary
120 BPM	101 F	\$20,000
80 BPM	104.4 F	\$40,000
140 BPM	99 F	\$800,000

- Optimizing Coverage

- Start with $\mathcal{A}_0 = \emptyset$
- Expand \mathcal{A}_{t-1} by one element to get \mathcal{A}_t



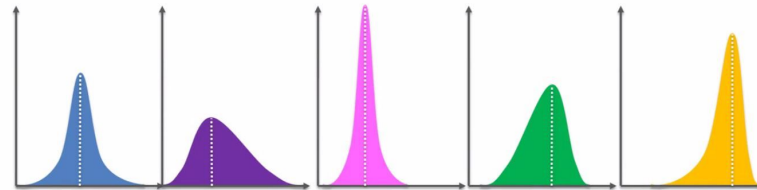
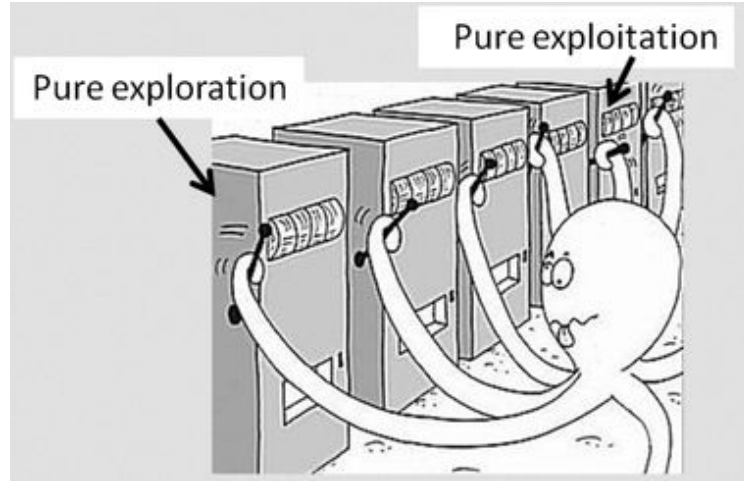
Generating Anchors - Optimizing Precisions

- Searching for the Optimal A
 - 1) Optimizing coverage with $\text{cov}(A) \geq c$
 - 2) Optimizing precision with $\text{prec}(A) \geq \tau$ and confidence at least $1-\delta$
 - 3) Update the optimal solution A^*

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

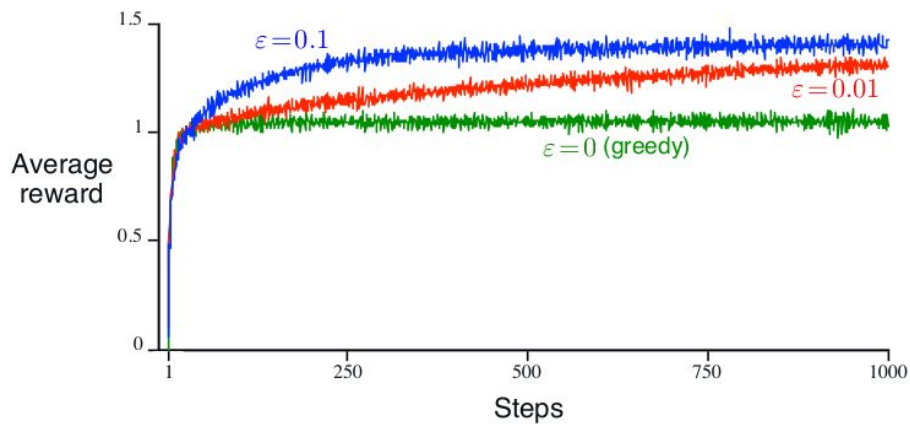
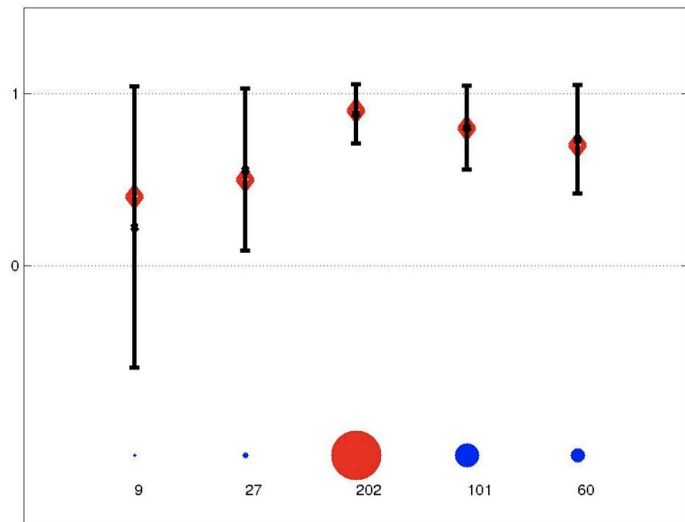
- Optimizing Precisions
 - Formulate it as a Multi-armed bandit optimization problem

Multi-Armed Bandit Problem



Reward Distribution of Each Arm

Exploration and Exploitation Trade-offs



Generating Anchors - Optimizing Precisions

- Searching for the Optimal A
 - 1) Optimizing coverage with $\text{cov}(A) \geq c$
 - 2) Optimizing precision with $\text{prec}(A) \geq \tau$ and confidence at least $1-\delta$
 - 3) Update the optimal solution A^*

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

- Optimizing Precisions
 - Formulate it as a Multi-armed bandit optimization problem
 - Find out candidates with $\text{Prec}(A) \geq \tau$
 - Using minimal costs (number of pulls of the arms)
 - Each candidate solution A is an arm
 - $\text{Prec}(A)$ of a single sample is the latent reward

Generating Anchors - Optimizing Precisions

- Searching for the Optimal A
 - 1) Optimizing coverage with $\text{cov}(A) \geq c$
 - 2) Optimizing precision with $\text{prec}(A) \geq \tau$ and confidence at least $1-\delta$
 - 3) Update the optimal solution A^*

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

- Optimizing Precisions
 - Formulate it as a Multi-armed bandit optimization problem
 - Find out candidates with $\text{Prec}(A) \geq \tau$
 - Using minimal costs (number of pulls of the arms)
 - Each candidate solution A is an arm
 - $\text{Prec}(A)$ of a single sample is the latent reward
 - Return the top K arms (i.e., A) with the highest reward ($\text{Prec}(A)$) that satisfies conditions
 - $\text{Prec}(A) \geq \tau, P(\text{Prec}(A) \geq \tau) \geq 1-\delta$

Generating Anchors - Update Optimal A^*

- Searching for the Optimal A
 - 1) Optimizing coverage with $\text{cov}(A) \geq c$
 - 2) Optimizing precision with $\text{prec}(A) \geq \tau$ and confidence at least $1-\delta$
 - 3) Update the optimal solution A^*
- Update A^*
 - For the top-k A returned in step 2)
 - Find the best A^* based on the Coverage criteria
if $\text{cov}(A) > \text{cov}(A^*)$ then $A^* \leftarrow A$
 - Loop into the next step

Precision and Coverage

- Precision

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

- Coverage

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

- Limes

- lime-n - Naive LIME algorithm
- lime-t - Make predictions only when its predictive probability is above a threshold

		Precision		Coverage	
		anchor	lime-n	anchor	lime-t
adult	logistic	<u>95.6</u>	<u>81.0</u>	<u>10.7</u>	<u>21.6</u>
	gbt	<u>96.2</u>	<u>81.0</u>	<u>9.7</u>	<u>20.2</u>
	nn	<u>95.6</u>	<u>79.6</u>	<u>7.6</u>	<u>17.3</u>
rcdv	logistic	<u>95.8</u>	<u>76.6</u>	<u>6.8</u>	<u>17.3</u>
	gbt	<u>94.8</u>	<u>71.7</u>	<u>4.8</u>	<u>2.6</u>
	nn	<u>93.4</u>	<u>65.7</u>	<u>1.1</u>	<u>1.5</u>
lending	logistic	<u>99.7</u>	<u>80.2</u>	<u>28.6</u>	<u>12.2</u>
	gbt	<u>99.3</u>	<u>79.9</u>	<u>28.4</u>	<u>9.1</u>
	nn	<u>96.7</u>	<u>77.0</u>	<u>16.6</u>	<u>5.4</u>

logistic: logistic regression, gbt: gradient boosted trees
nn: two layers of 50 units MLP

User Study

- Ask Users to Guess the Outcomes of A ML Model After Explanations
 - Human annotators are 26 students who took a machine learning course
 - Calculate precision and coverage of the users' performance
 - Human mark "I don't know" when they are not certain, which makes coverage the perceived one.

Method	Precision				Coverage (perceived)				Time/pred (seconds)			
	adult	rcdv	vqa1	vqa2	adult	rcdv	vqa1	vqa2	adult	rcdv	vqa1	vqa2
No expls	<u>54.8</u>	<u>83.1</u>	<u>61.5</u>	<u>68.4</u>	<u>79.6</u>	<u>63.5</u>	<u>39.8</u>	<u>30.8</u>	<u>29.8 ±14</u>	<u>35.7±26</u>	<u>18.7±20</u>	<u>13.9±20</u>
LIME(1)	<u>68.3</u>	98.1	<u>57.5</u>	<u>76.3</u>	<u>89.2</u>	<u>55.4</u>	<u>71.5</u>	<u>54.2</u>	<u>28.5±10</u>	<u>24.6±6</u>	<u>8.6±3</u>	<u>11.1±8</u>
Anchor(1)	<u>100.0</u>	97.8	<u>93.0</u>	<u>98.9</u>	<u>43.1</u>	<u>24.6</u>	<u>31.9</u>	<u>27.3</u>	<u>13.0±4</u>	<u>14.4±5</u>	<u>5.4±2</u>	<u>3.7±1</u>
LIME(2)	89.9	<u>72.9</u>	-	-	<u>78.5</u>	<u>63.1</u>	-	-	<u>37.8±20</u>	24.4±7	-	-
Anchor(2)	87.4	<u>95.8</u>	-	-	<u>62.3</u>	<u>45.4</u>	-	-	<u>10.5±3</u>	19.2±10	-	-

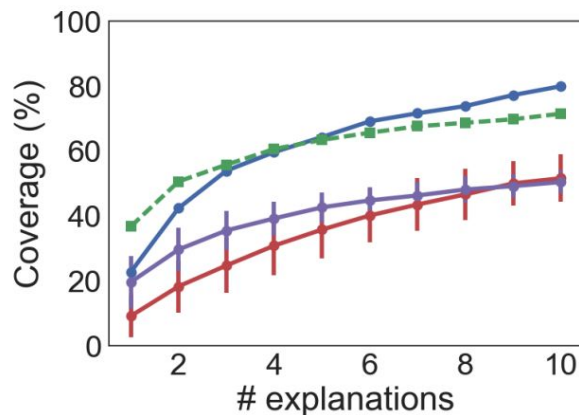
LIME(n): results after n LIME explanations

Anchor(n): results after n Anchor explanations

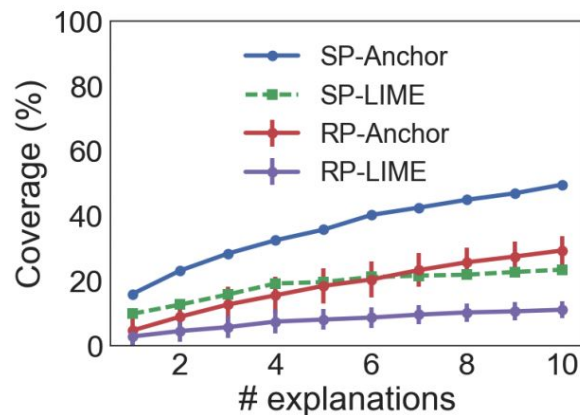
[Ribeiro et al, 2018](#)

User Study Results

- Coverage change with number of explanations seen by the same user.
 - gradient boosted trees(gb)
 - SP - Submodular Pick
 - RP - Random Pick



(a) *adult* dataset

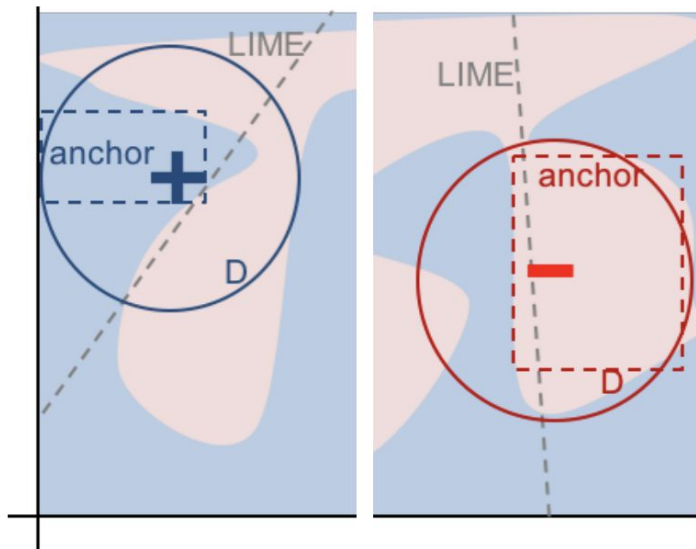


(b) *rcdv* dataset

Comparisons to LIME

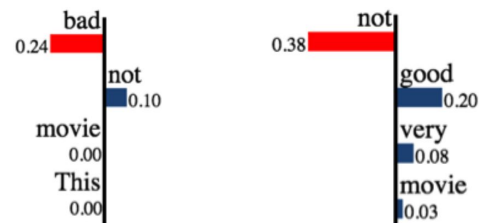
	LIME	Anchors
Explanations	$g(z') = w_g \cdot z'$	Anchors A
Optimization Target	$\mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in \mathcal{Z}} \pi_x(z) (f(z) - g(z'))^2$	$\max_{A \text{ s.t. } P(\text{prec}(A) \geq \tau) \geq 1 - \delta} \text{cov}(A)$

Comparisons to LIME



+ This movie is not bad. - This movie is not very good.

(a) Instances



(b) LIME explanations

{"not", "bad"} → Positive {"not", "good"} → Negative

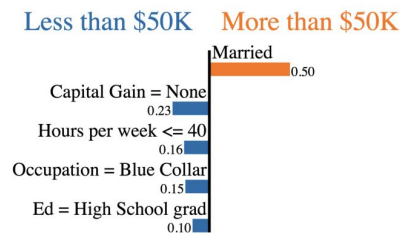
(c) Anchor explanations

Overly Specific Anchors

28 < Age ≤ 37
Workclass = Private
Education = High School grad
Marital Status = Married
Occupation = Blue-Collar
Relationship = Husband
Race = White
Sex = Male
Capital Gain = None
Capital Loss = Low
Hours per week ≤ 40.00
Country = United-States

$P(\text{Salary} > \$50\text{K}) = 0.57$

(a) Instance and prediction



(b) LIME explanation

IF Country = United-States **AND** Capital Loss = Low
AND Race = White **AND** Relationship = Husband
AND Married **AND** 28 < Age ≤ 37
AND Sex = Male **AND** High School grad
AND Occupation = Blue-Collar
THEN PREDICT Salary > \$50K

(c) An *anchor* explanation

Project Review

- Project Proposal Due Apr 22
 - Up to 1.5 pages
 - The problem you are solving
 - Datasets
 - Metrics
 - Baselines
- Use the Slack Channel to Find Partners
 - <https://cs335-2020sp.slack.com/archives/C0120BNJJHW>
- Google Cloud Credits

Required Reading

Molnar: Ch 5.7, Ch 5.8

Reading Assignments (Pick One)

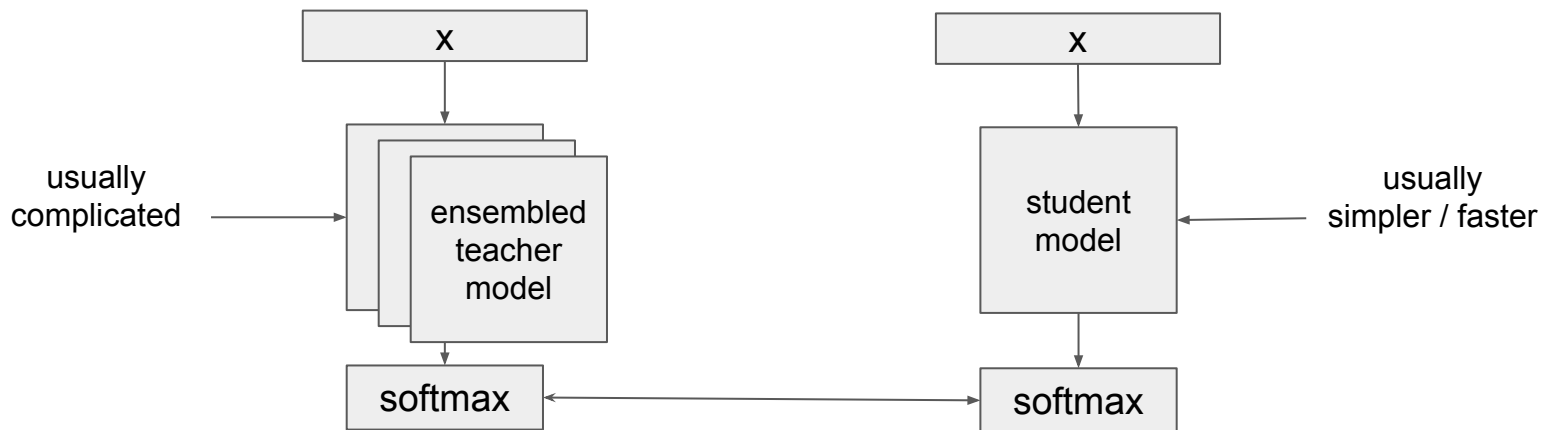
- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "Why should i trust you?" Explaining the predictions of any classifier, SIGKDD 2016
- Lakkaraju, Himabindu, Stephen H. Bach, and Jure Leskovec. Interpretable decision sets: A joint framework for description and prediction, SIGKDD 2016
- Che, Z., Purushotham, S., Khemani, R., & Liu, Y. Distilling knowledge from deep networks with applications to healthcare domain, Arxiv 2015
- Plumb, Gregory, Denali Molitor, and Ameet S. Talwalkar. Model agnostic supervised local explanations, NeurIPS 2018
- Robnik-Šikonja, M., & Kononenko, I. Explaining classifications for individual instances. IEEE Transactions on Knowledge and Data Engineering, 2008

Next Lecture

Feature Interactions for Interpretability

Knowledge Distillation ([Hinton et al, 2015](#))

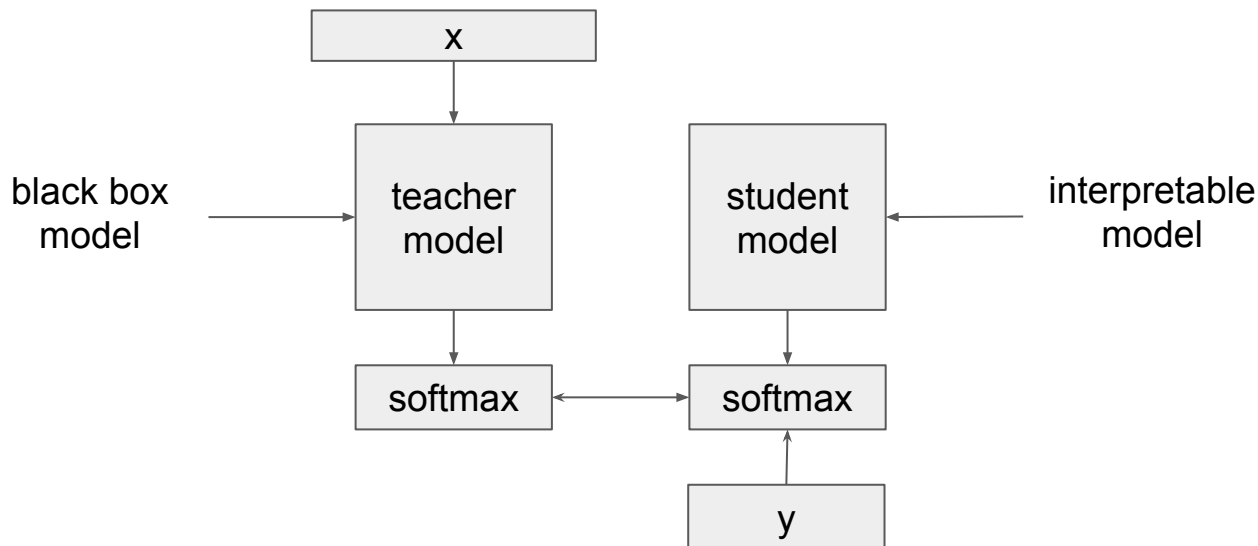
- Distillation of Neural Networks
 - use a simple network to approximate the more complicated ones
 - applications: improve performance (matching ensemble models), improve latency



DarkSight ([Xu et al, 2018](#))

- Teacher - Student Architecture

- match the softmax output between the teacher model and the student model
- $P_T(k | x) \sim P_S(k | y, \Theta)$



DarkSight ([Xu et al, 2018](#))

- Optimization
 - match the distribution of the softmax layer
 - D is a divergence measure

$$L(Y, \theta) = \frac{1}{N} \sum_{i=1}^N D(P_T(\cdot|x_i), P_S(\cdot|y_i; \theta))$$

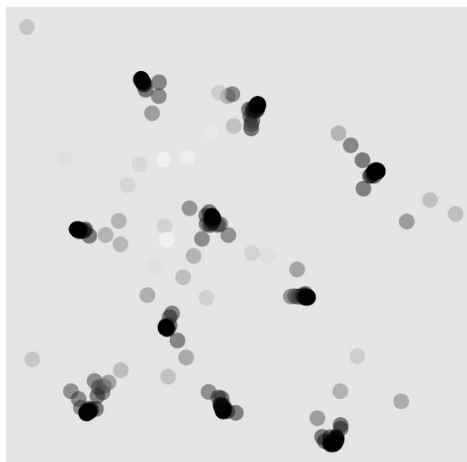
- One such choice can be the symmetric KL

$$KL_{sym}(P, Q) = \frac{1}{2}(KL(P, Q) + KL(Q, P))$$

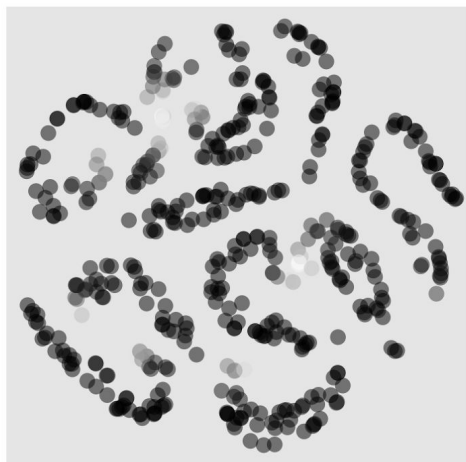
DarkSight ([Xu et al, 2018](#))

- Interpretable Model
 - Naive Bayes Classifier

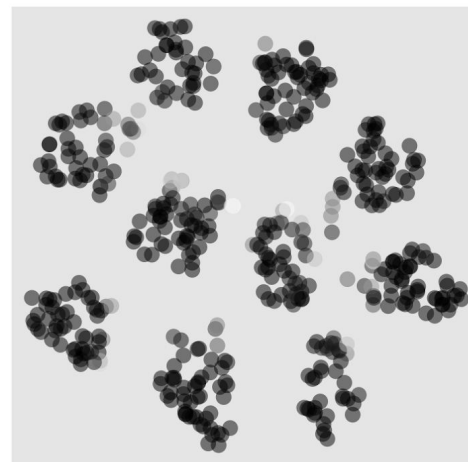
$$P_S(c_i = k|y_i; \theta) = \frac{P(y_i|c_i = k; \theta_c)P(c_i = k; \theta_p)}{P(y_i|\theta)}$$



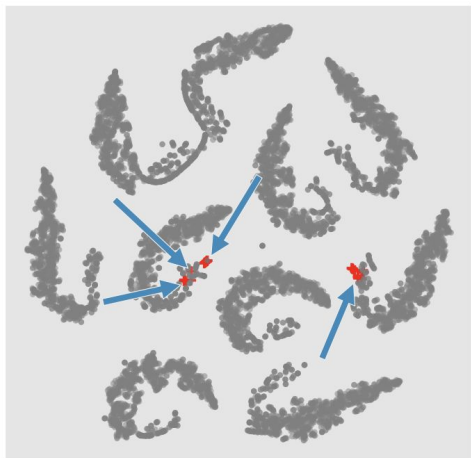
(a) DarkSight



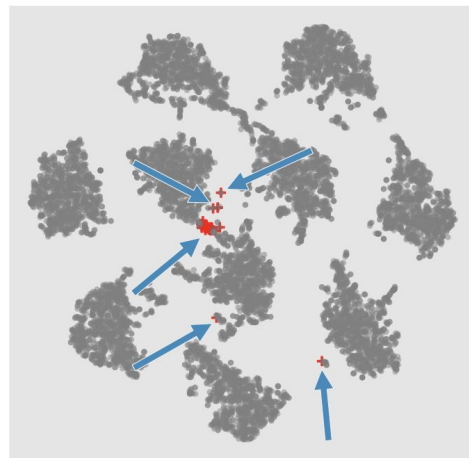
(b) t-SNE prob



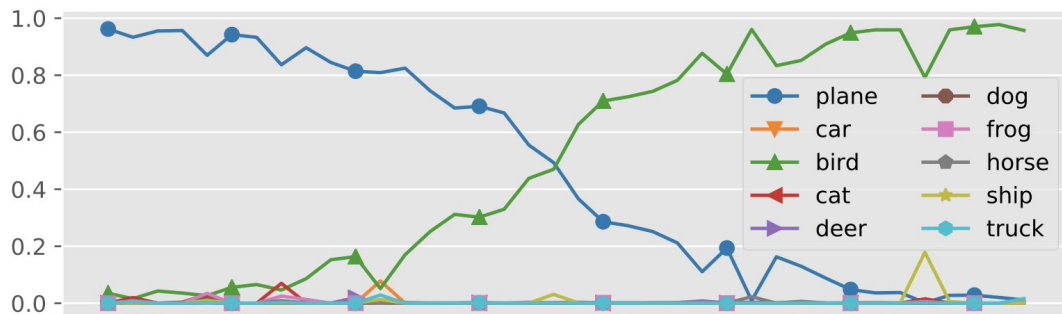
(c) t-SNE logit



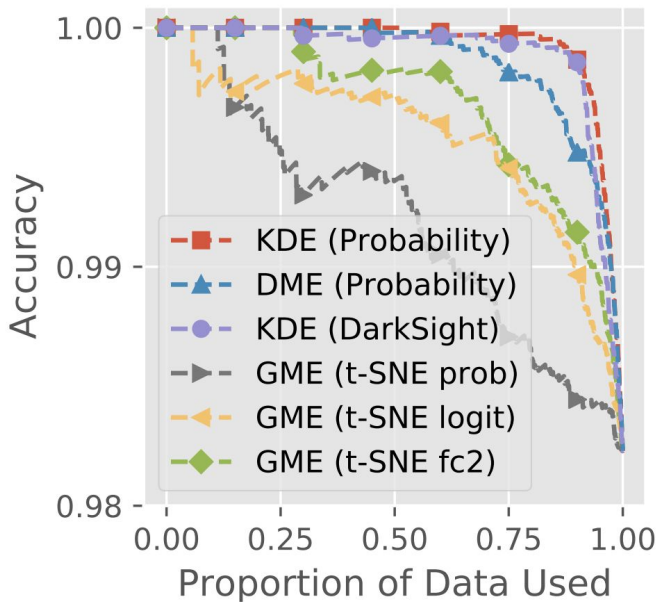
(a) t-SNE prob



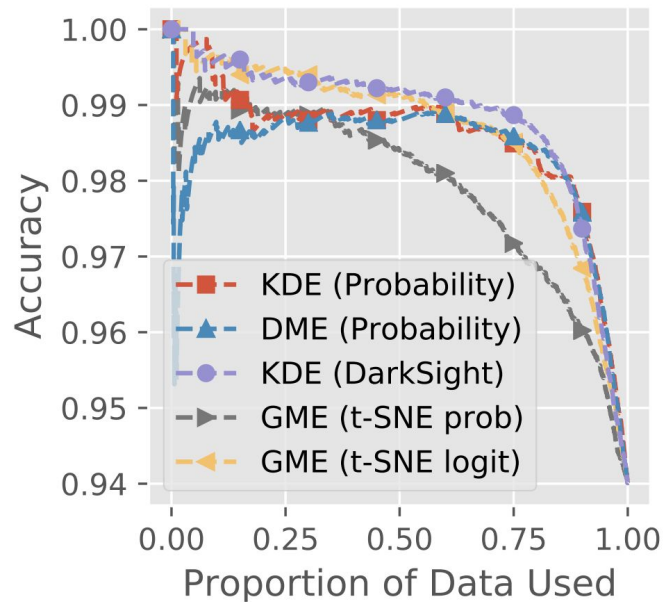
(b) t-SNE logit



(c) Predictive probabilities of points in the black box



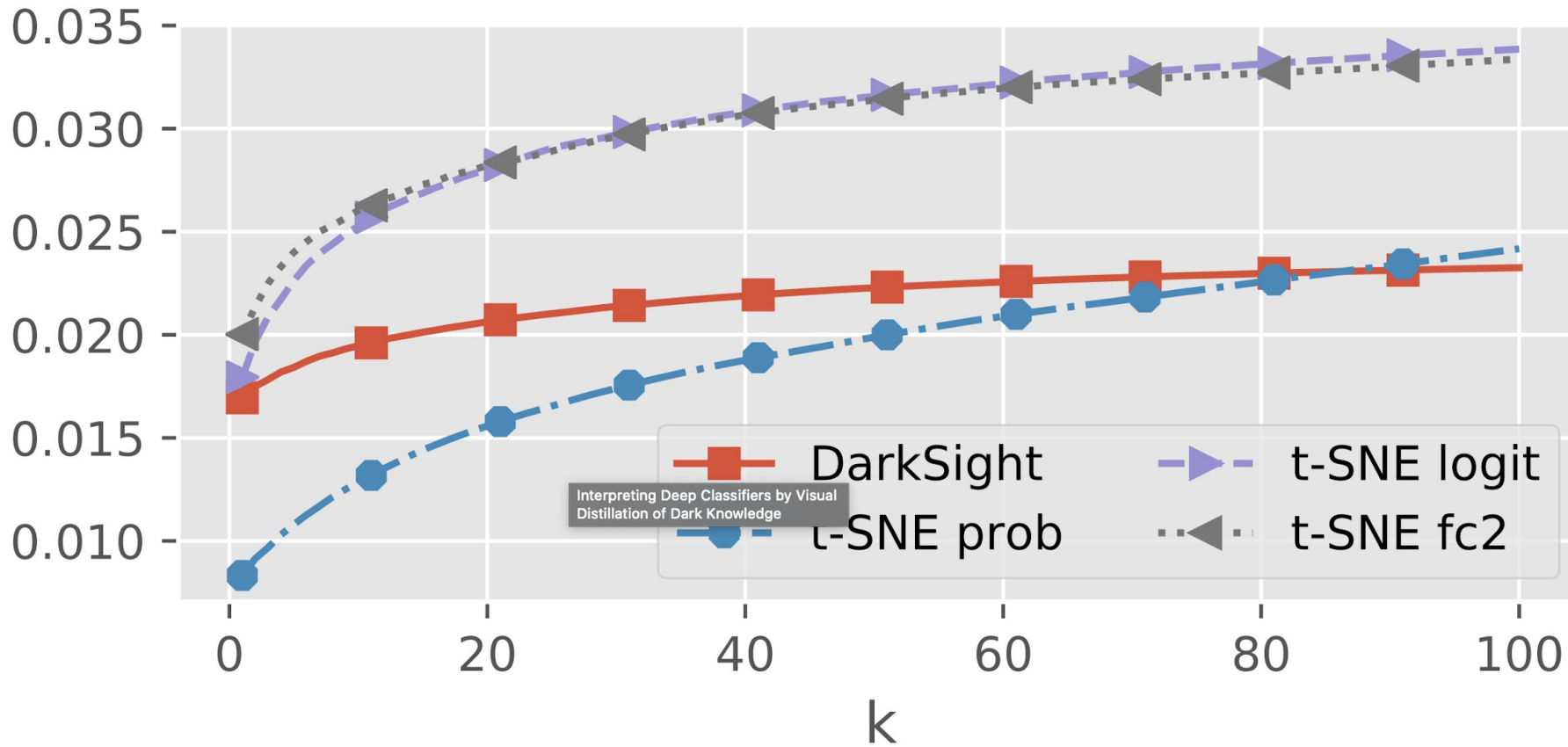
(a) MNIST



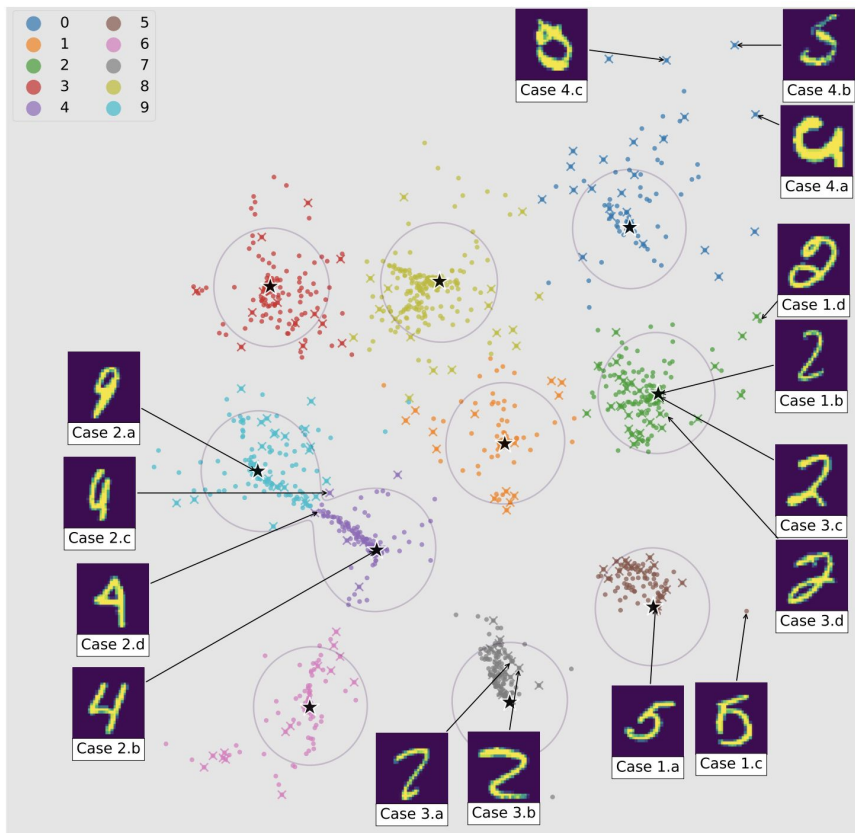
(b) Cifar10

$$M_k(Y) = \frac{1}{N} \sum_{i=1}^N \frac{1}{k} \sum_{j \in \text{NN}_k(y_i)} JSD(p_i, p_j)$$

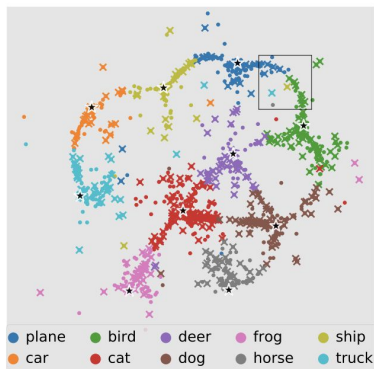
Avg. JSD for point & kNNs



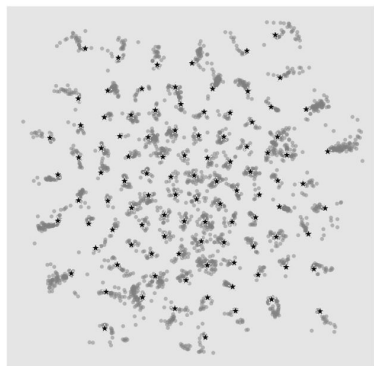
Interpreting Deep Classifiers by Visual Distillation of Dark Knowledge



(a) LeNet on MNIST



(b) VGG16 on Cifar10



(c) Wide-ResNet on Cifar100