



Slides and
Resources

Post Hoc, Local and Model-Agnostic Explanations

Yilun Zhou

MIT CSAIL & Amazon



Slides and
Resources

Outline

- Why model explanations?
- How to compute model explanations? -- Definitions
- How to evaluate model explanations? -- Evaluations
- A definition-evaluation duality



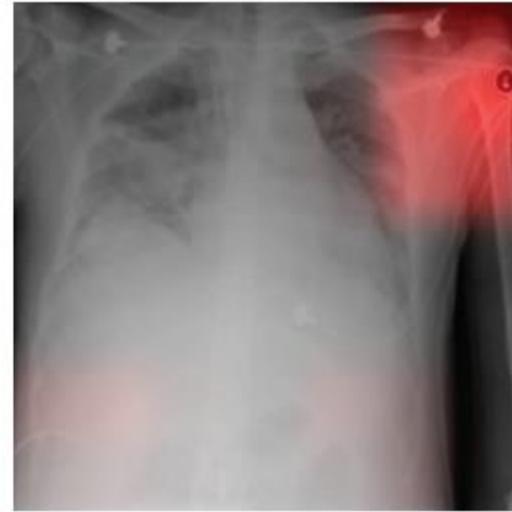
Slides and
Resources

Why Model Explanations?



Slides and
Resources

Why Model Explanations?



Task for DNN	Caption image	Recognise pneumonia
Problem	Describes green hillside as grazing sheep	Fails on scans from new hospitals
Shortcut	Uses background to recognise primary object	Looks at hospital token, not lung

(Geirhos et al. 2020)



Slides and
Resources

Why Model Explanations?

*Why is my model failing?
Because ... (debugging and diagnosis)*

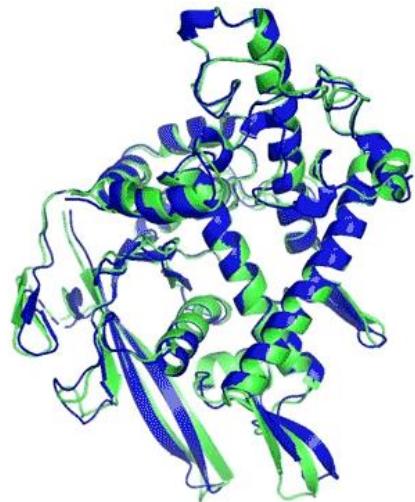
Task for DNN	Caption image	Recognise pneumonia
Problem	Describes green hillside as grazing sheep	Fails on scans from new hospitals
Shortcut	Uses background to recognise primary object	Looks at hospital token, not lung

(Geirhos et al. 2020)

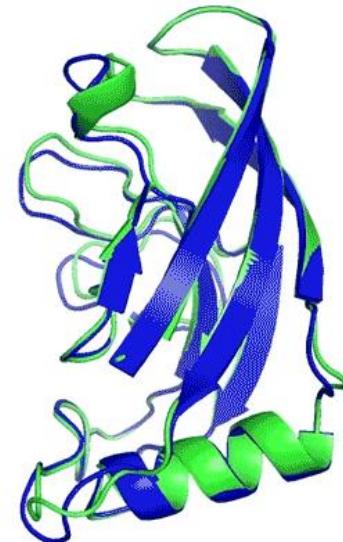


Slides and
Resources

Why Model Explanations?



T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)



T1049 / 6y4f
93.3 GDT
(adhesin tip)

- Experimental result
- Computational prediction



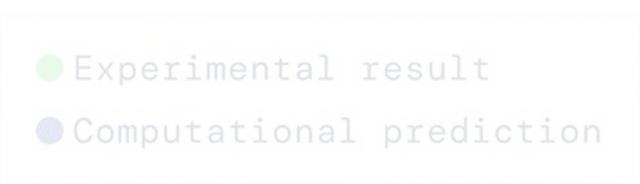
Slides and
Resources

Why Model Explanations?

*How does my model predict this structure?
Because ... (scientific discovery)*

T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

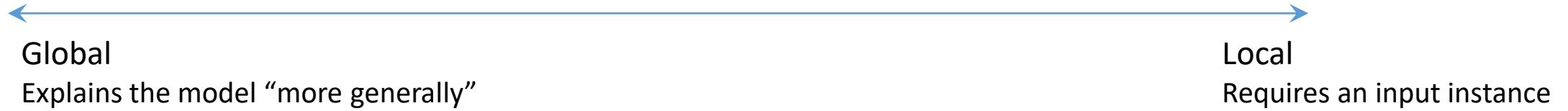
T1049 / 6y4f
93.3 GDT
(adhesin tip)





Slides and
Resources

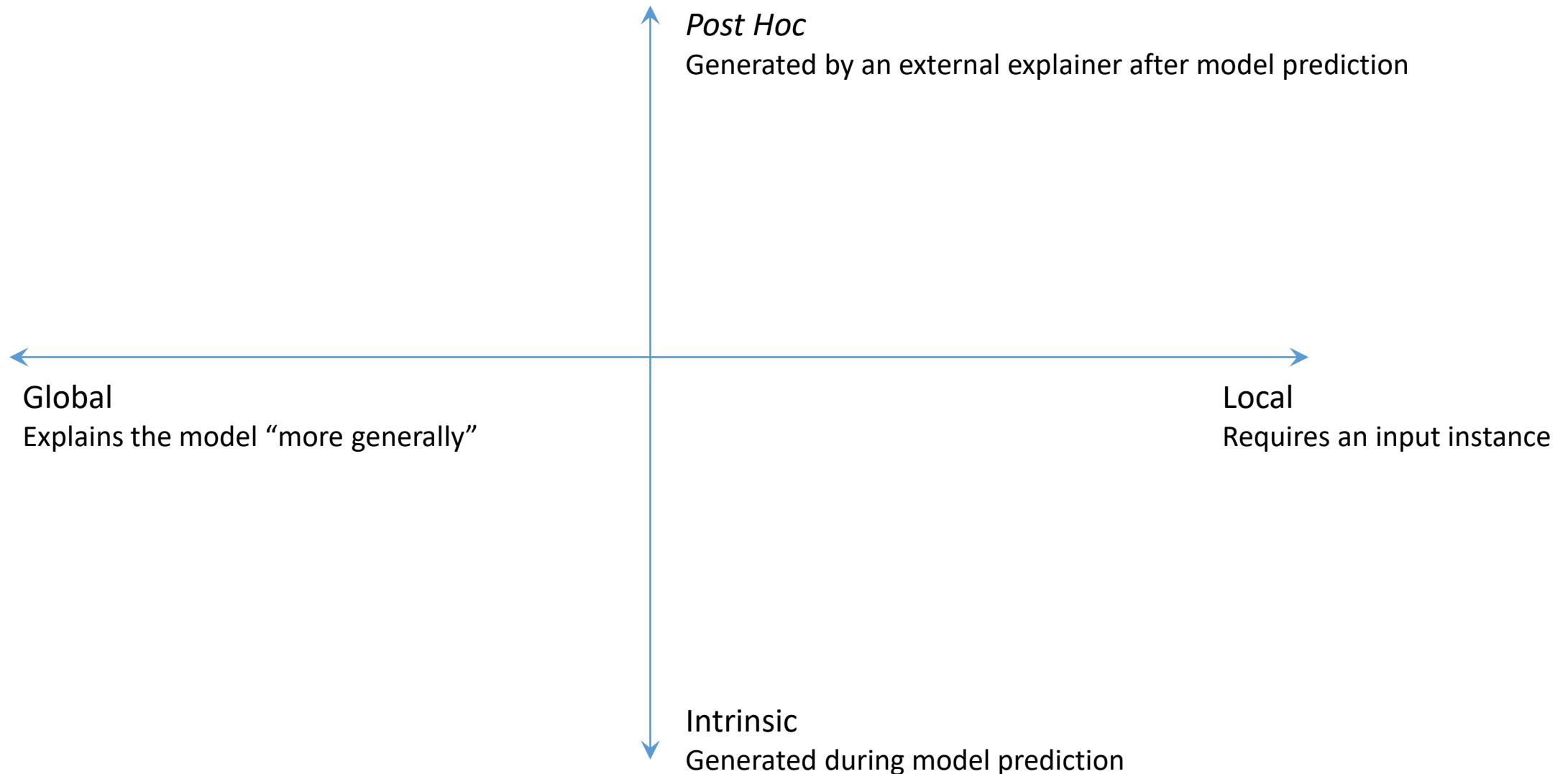
The Two Major Axes of Interpretability





Slides and
Resources

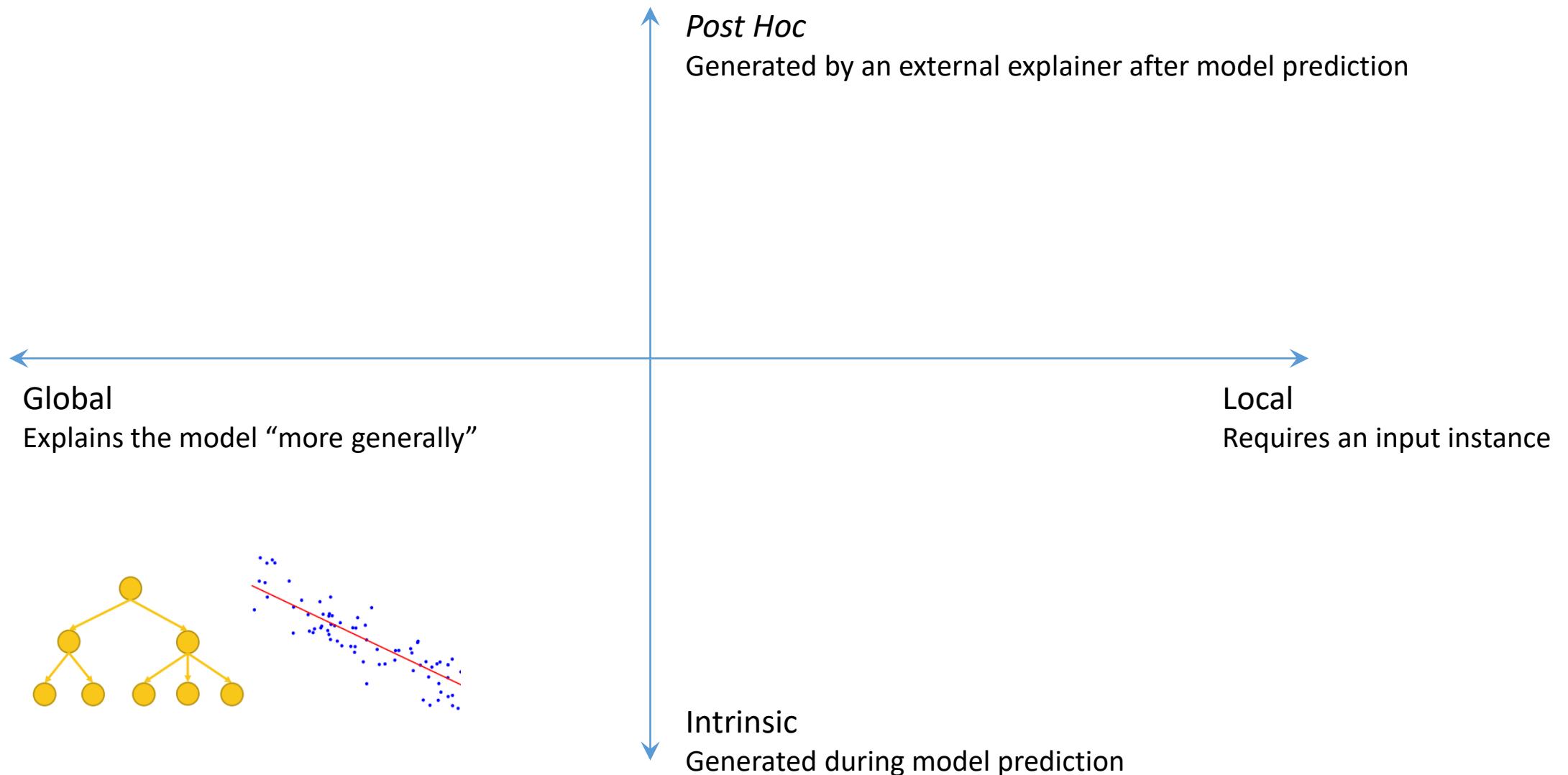
The Two Major Axes of Interpretability





Slides and
Resources

The Two Major Axes of Interpretability

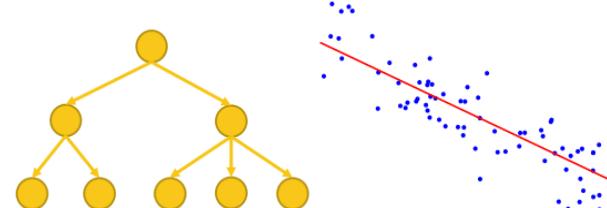




Slides and
Resources

The Two Major Axes of Interpretability

Global
Explains the model “more generally”



Post Hoc

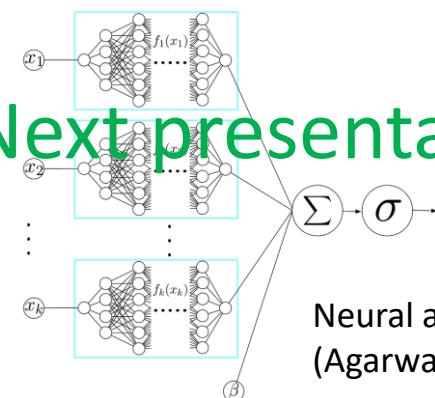
Generated by an external explainer after model prediction

Intrinsic

Generated during model prediction

Local
Requires an input instance

Next presentation



Neural additive model
(Agarwal et al., 2021)



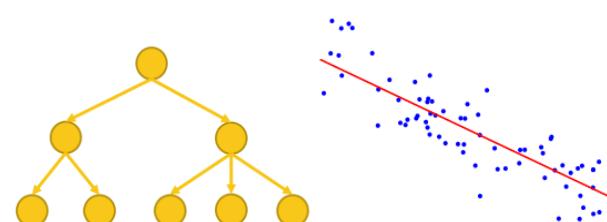
Slides and
Resources

The Two Major Axes of Interpretability



CNN filter visualization
(Olah et al., 2017)

Global
Explains the model “more generally”



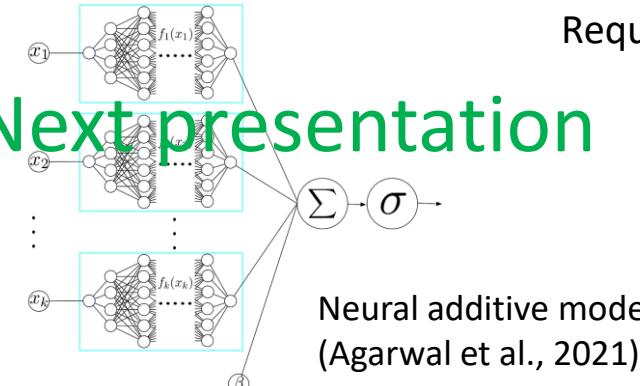
Post Hoc

Generated by an external explainer after model prediction

Intrinsic

Generated during model prediction

Local
Requires an input instance



Neural additive model
(Agarwal et al., 2021)



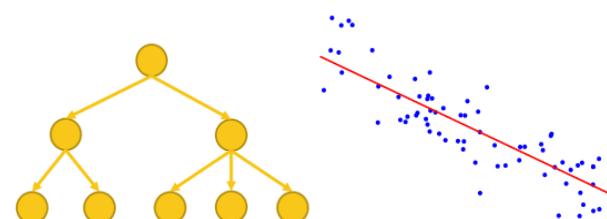
Slides and
Resources

The Two Major Axes of Interpretability



CNN filter visualization
(Olah et al., 2017)

Global
Explains the model “more generally”



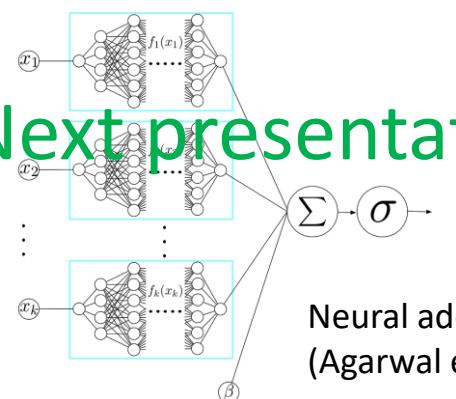
Post Hoc

Generated by an external explainer after model prediction

This presentation

Intrinsic
Generated during model prediction

Local
Requires an input instance



Next presentation

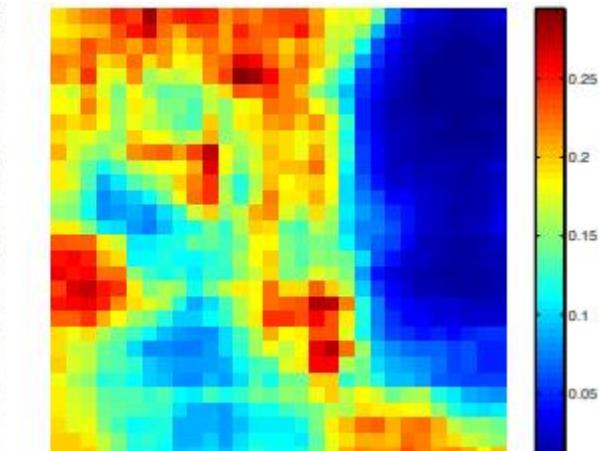
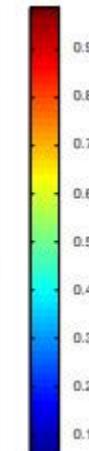
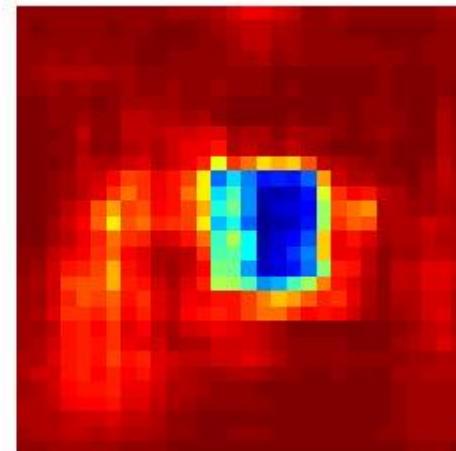
Neural additive model
(Agarwal et al., 2021)



Slides and Resources

Local *Post Hoc* Explanations

- What is a local *post hoc* explanation?
 - Some description of the model's local decision making logic
- If this image patch is grayed out, the model prediction will change to this.



(a)

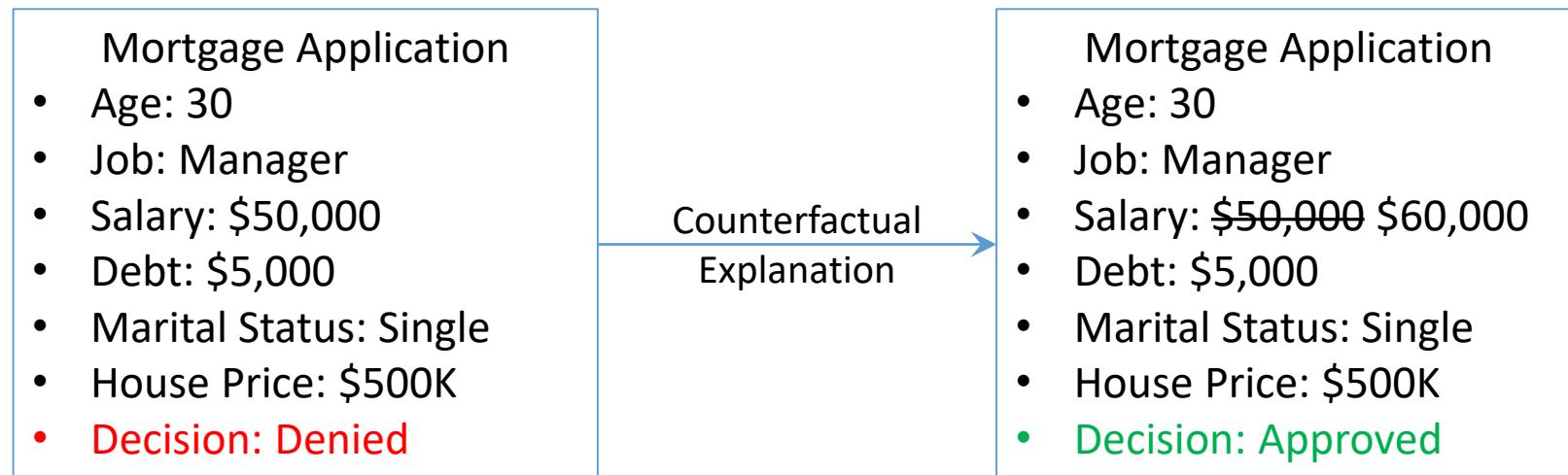
Occlusion Saliency (Zeiler and Fergus, 2014)



Slides and
Resources

Local *Post Hoc* Explanations

- What is a local *post hoc* explanation?
 - Some description of the model's local decision making logic
- If this input feature is changed to this value, the model prediction will be different.

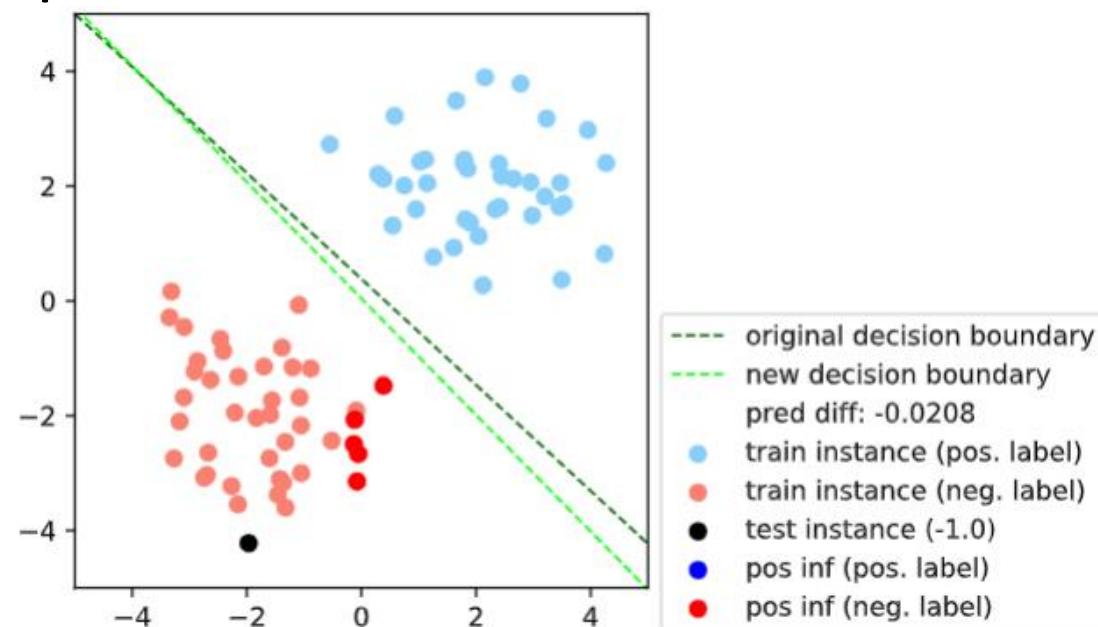




Slides and
Resources

Local *Post Hoc* Explanations

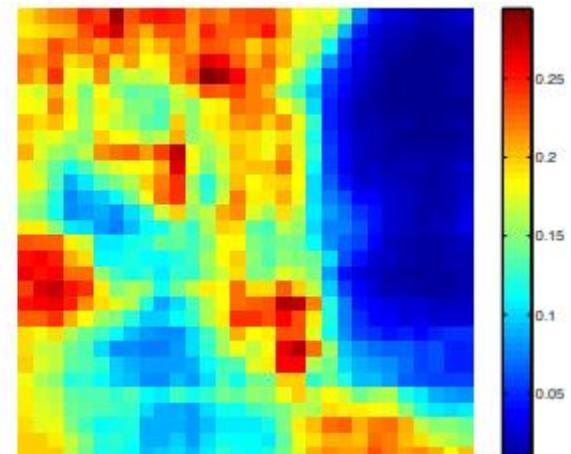
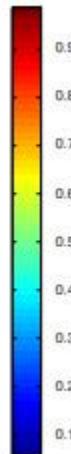
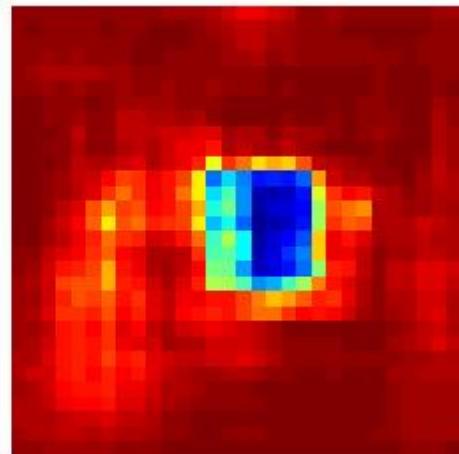
- What is a local *post hoc* explanation?
 - Some description of the model's local decision making logic
- If this training instance is not present in the dataset, the model will make a different prediction.





Slides and
Resources

Feature Attributions (Saliency Maps)



Occlusion Saliency (Zeiler and Fergus, 2014)



Slides and
Resources

Feature Attributions (Saliency Maps)

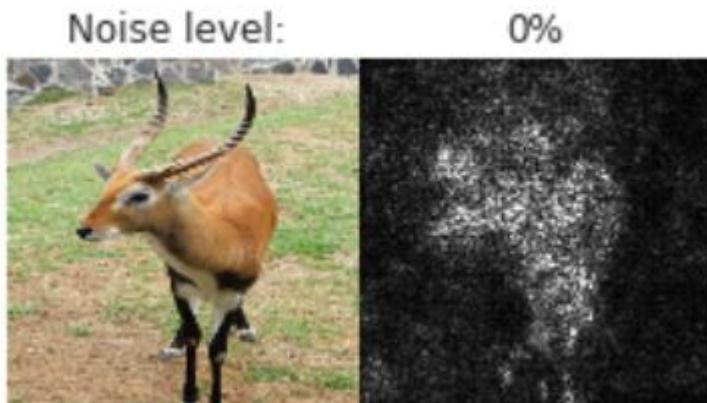
- Vanilla gradient: $s = \nabla_x f(x)$



Slides and
Resources

Feature Attributions (Saliency Maps)

- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_\epsilon[\nabla_x f(x + \epsilon)]$; ϵ_i independent for every x_i

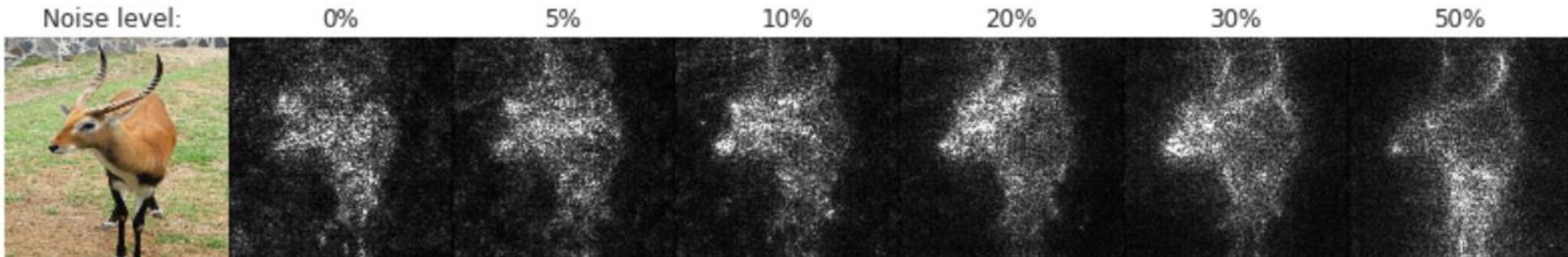




Slides and
Resources

Feature Attributions (Saliency Maps)

- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_\epsilon[\nabla_x f(x + \epsilon)]$; ϵ_i independent for every x_i

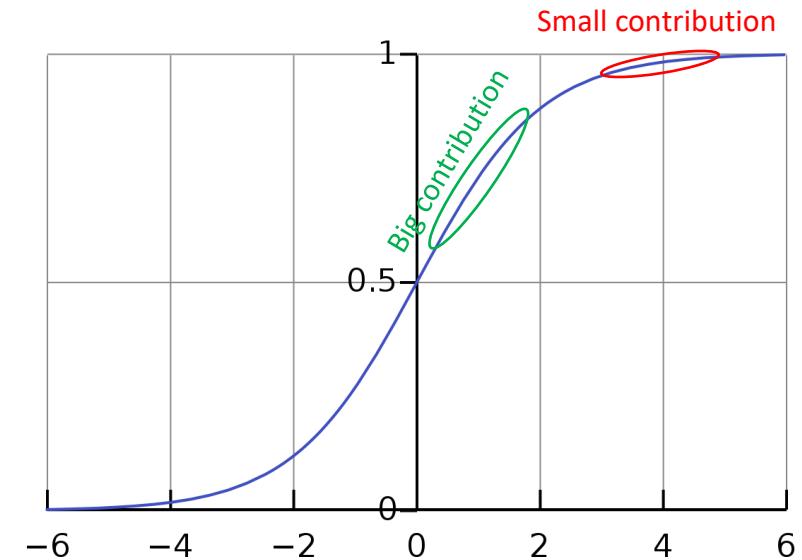




Slides and
Resources

Feature Attributions (Saliency Maps)

- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_\epsilon[\nabla_x f(x + \epsilon)]$; ϵ_i independent for every x_i
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x - x_0)) du$





Slides and
Resources

Feature Attributions (Saliency Maps)

- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_\epsilon[\nabla_x f(x + \epsilon)]$; ϵ_i independent for every x_i
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x - x_0)) du$
- Occlusion: $s_i = f(x) - f(x_{-i})$



Slides and
Resources

Feature Attributions (Saliency Maps)

- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_\epsilon[\nabla_x f(x + \epsilon)]$; ϵ_i independent for every x_i
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x - x_0)) du$
- Occlusion: $s_i = f(x) - f(x_{-i})$
- LIME: $f \sim s^T m + b$

Feature mask m	Masked sentence	$f(\cdot)$
0, 0, 0, 0	""	0.49
0, 0, 0, 1	"beautiful"	0.9
0, 0, 1, 0	"and"	0.52
0, 0, 1, 1	"and beautiful"	0.91
...		
1, 1, 1, 0	"It's good and"	0.89
1, 1, 1, 1	"It's good and beautiful"	0.95



Slides and
Resources

Feature Attributions (Saliency Maps)

- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_\epsilon[\nabla_x f(x + \epsilon)]$; ϵ_i independent for every x_i
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x - x_0)) du$
- Occlusion: $s_i = f(x) - f(x_{-i})$
- LIME: $f \sim s^T m + b$
- SHAP: $s_i = \sum_{S \subseteq F \setminus \{i\}} \frac{(|S|!(|F|-|S|-1)!)}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)]$



Slides and
Resources

Feature Attributions (Saliency Maps)

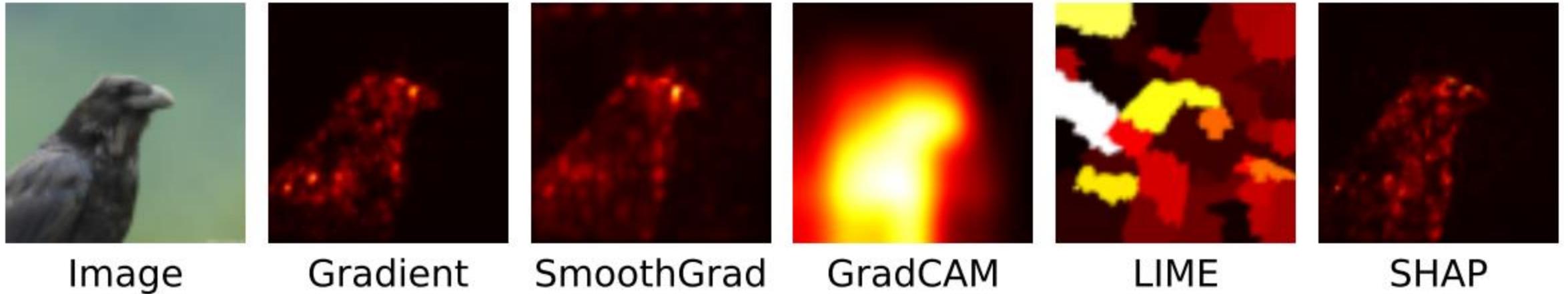
- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_\epsilon[\nabla_x f(x + \epsilon)]$; ϵ_i independent for every x_i
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x - x_0)) du$
- Occlusion: $s_i = f(x) - f(x_{-i})$
- LIME: $f \sim s^T m + b$
- SHAP: $s_i = \sum_{S \subseteq F \setminus \{i\}} \frac{(|S|!(|F|-|S|-1)!)}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)]$

$$e = D(x) \in \mathbb{R}^L$$



Slides and
Resources

Feature Attributions (Saliency Maps)



$$e = D(x) \in \mathbb{R}^L$$



Slides and
Resources

Evaluating Explanations

Problem: don't know
how models work



Slides and
Resources

Evaluating Explanations





Slides and
Resources

Evaluating Explanations

Problem: don't know
how models work

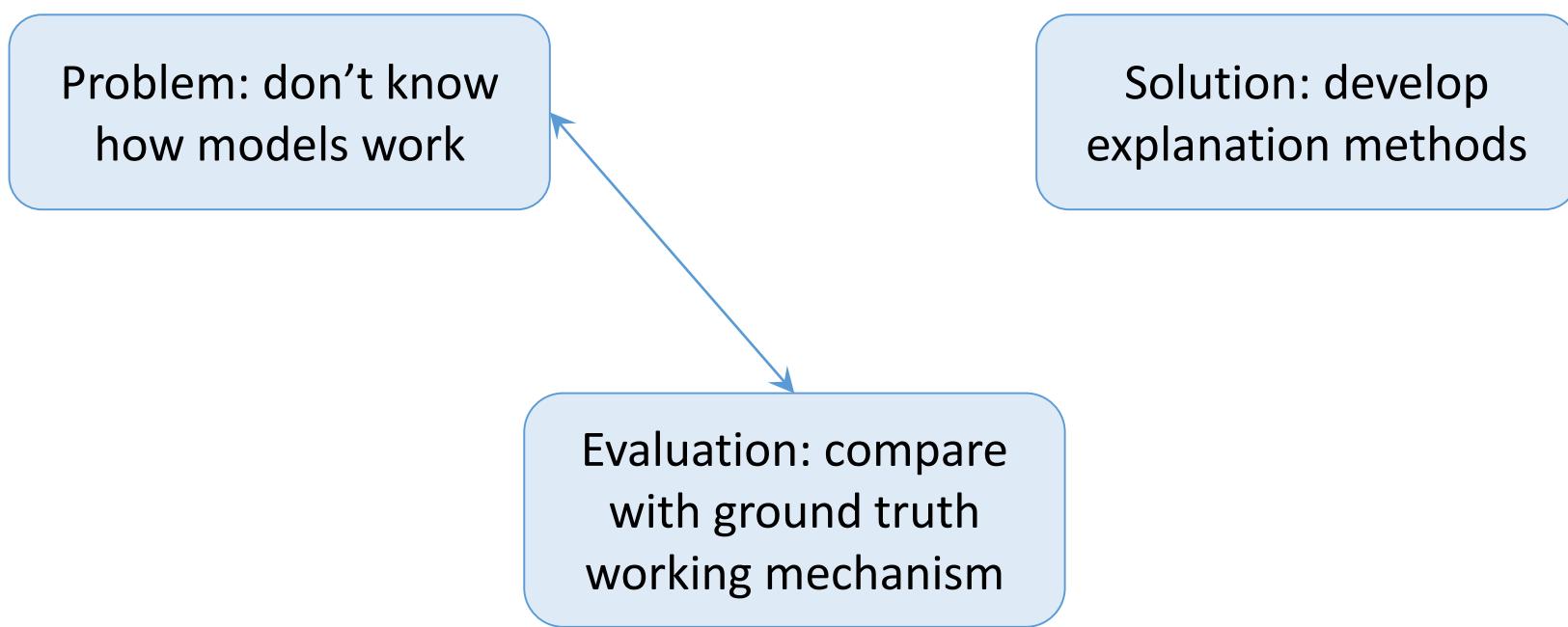
Solution: develop
explanation methods

Evaluation: compare
with ground truth
working mechanism



Slides and
Resources

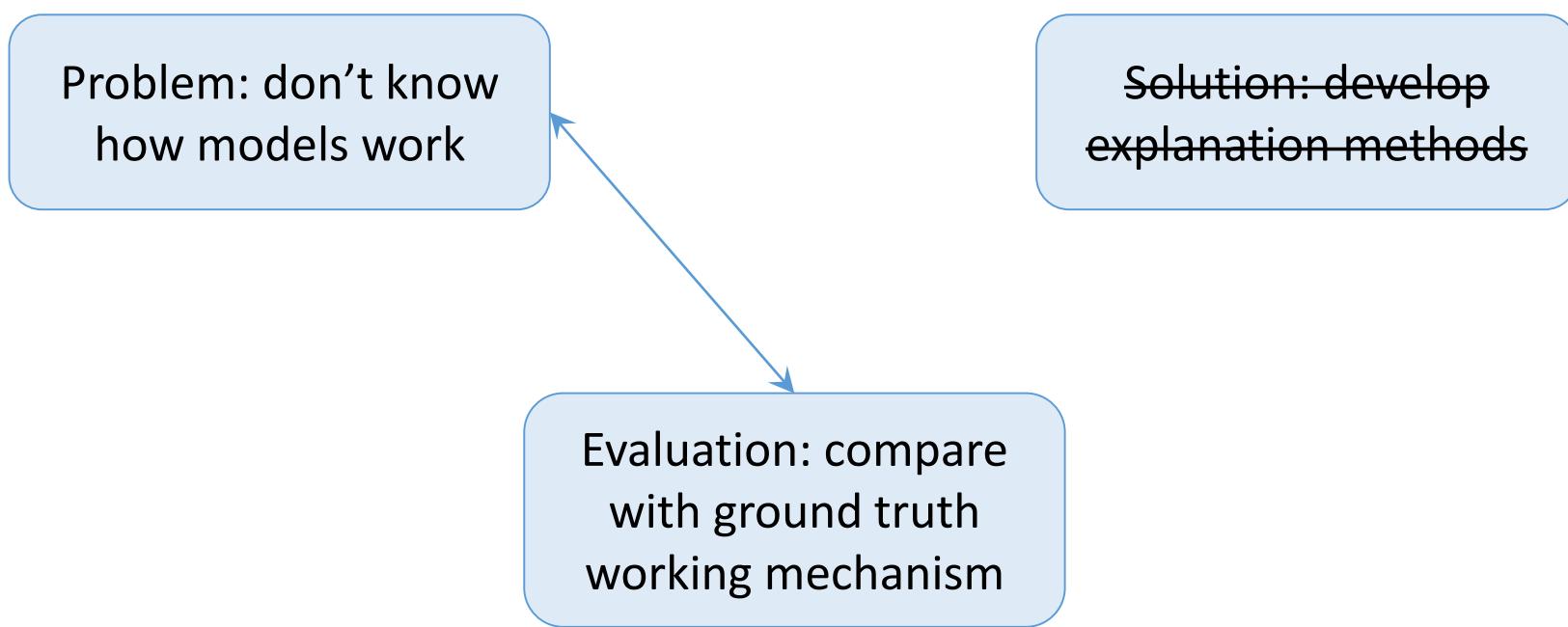
Evaluating Explanations





Slides and
Resources

Evaluating Explanations





Slides and
Resources

Proxy Metrics for Explanation Quality

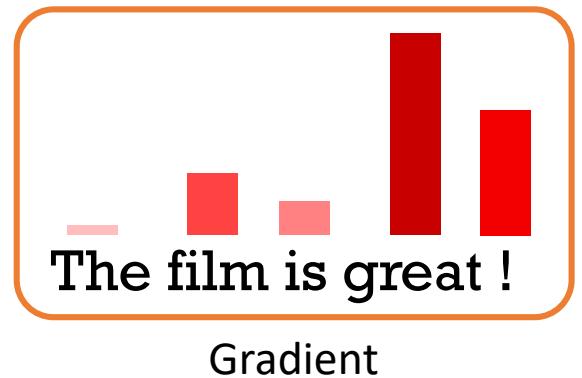
- Feature importance \Leftrightarrow model prediction change with feature removal



Slides and
Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal

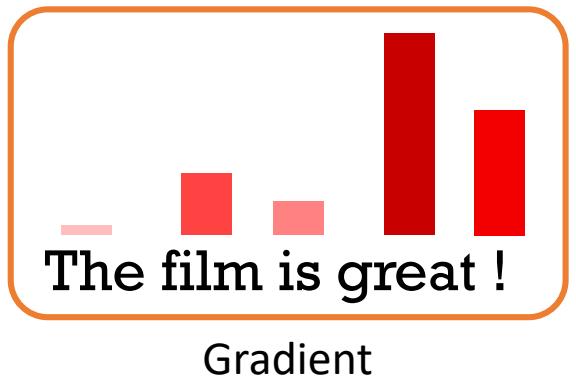




Slides and
Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal



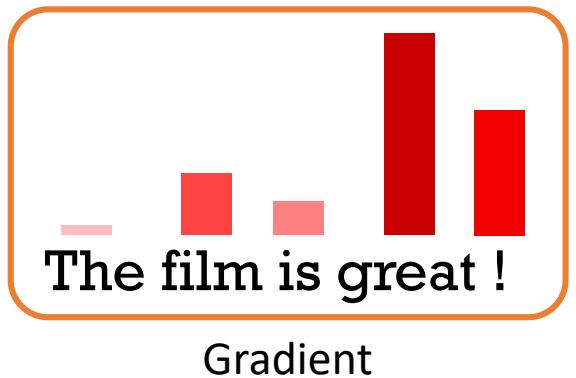
The film is great !



Slides and
Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal



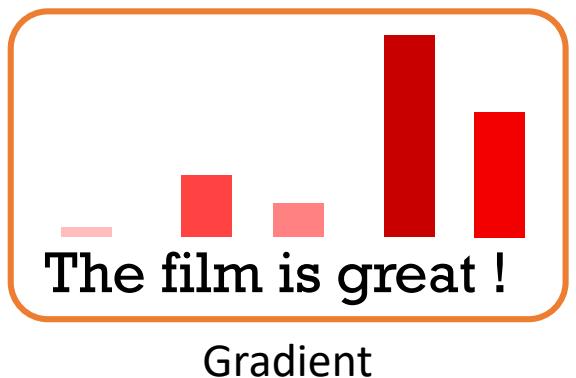
The film is great !
The film is !



Slides and
Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal



The film is great !
The film is !
The film is

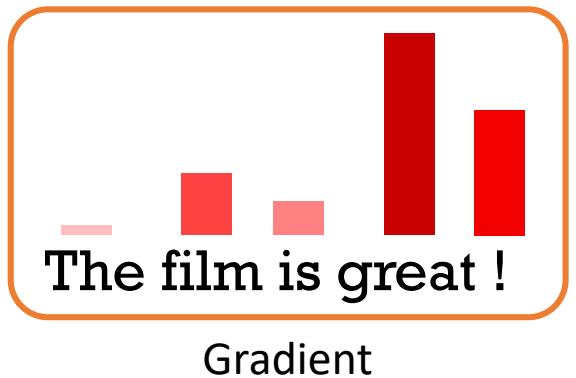


Slides and
Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal

The film is great !
The film is !
The film is
The is
The
(null)

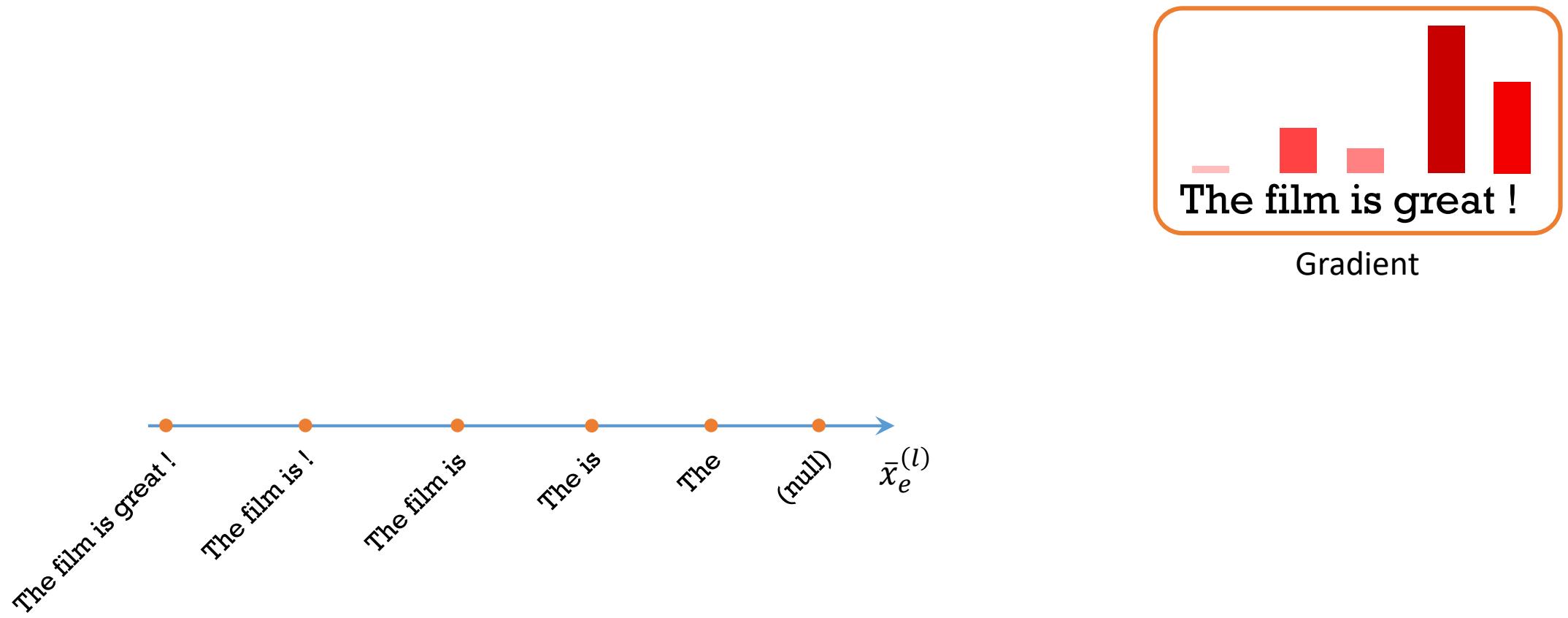




Slides and
Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal

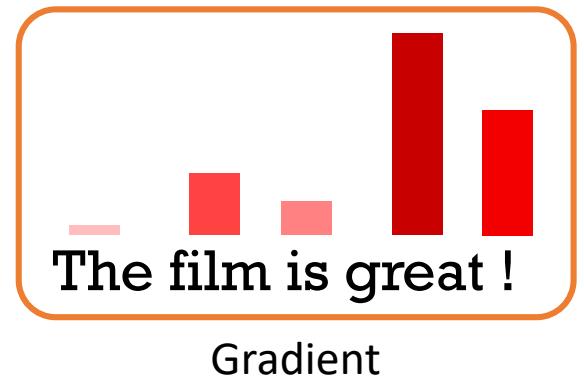
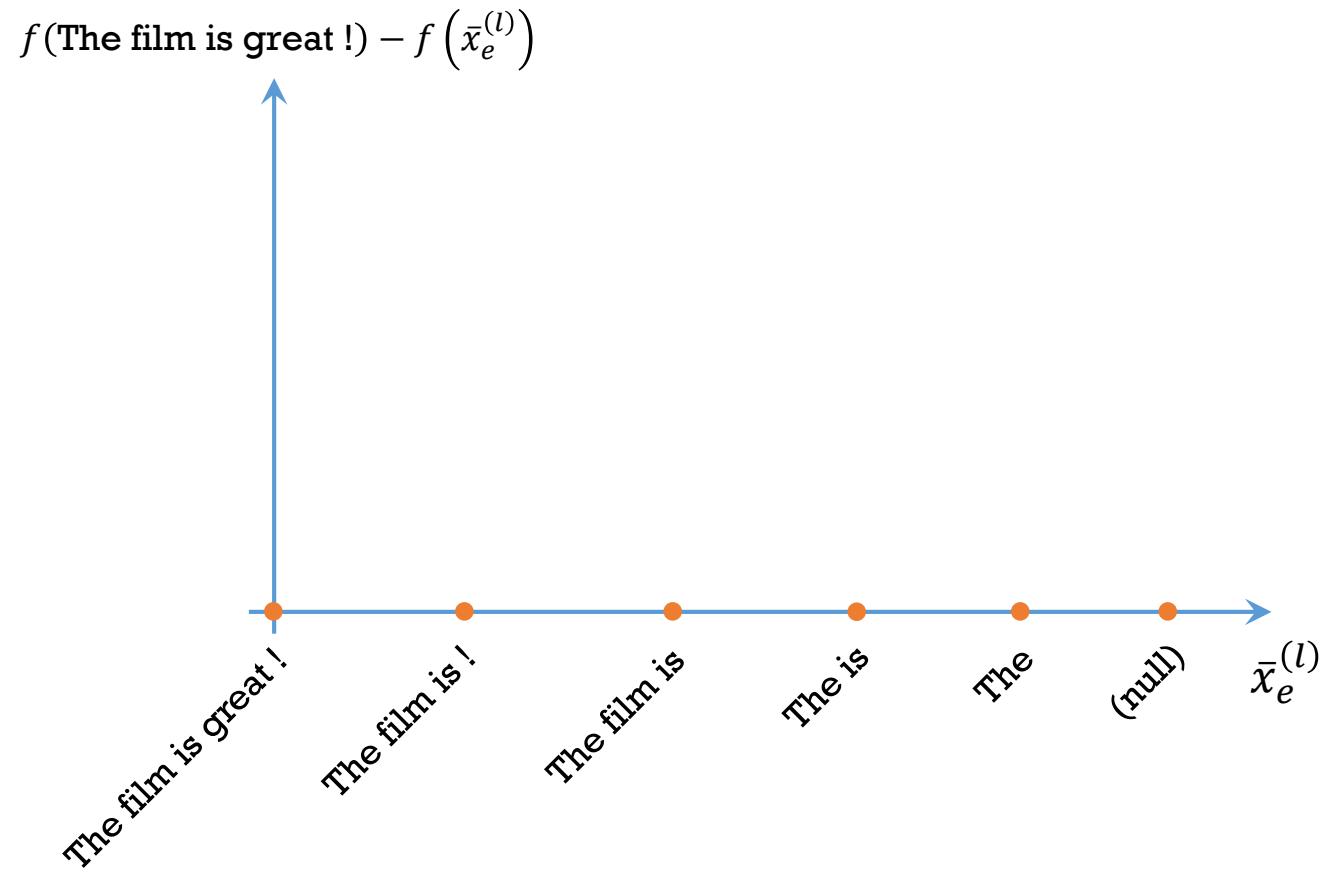




Slides and
Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal

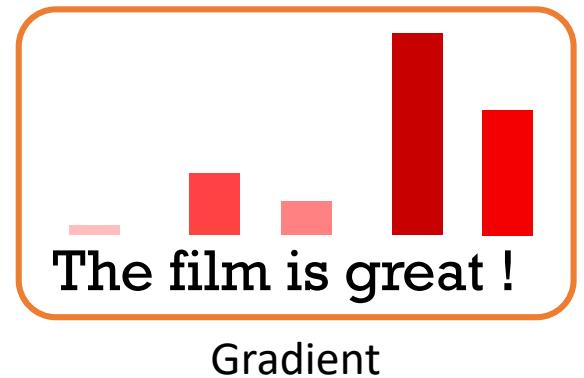
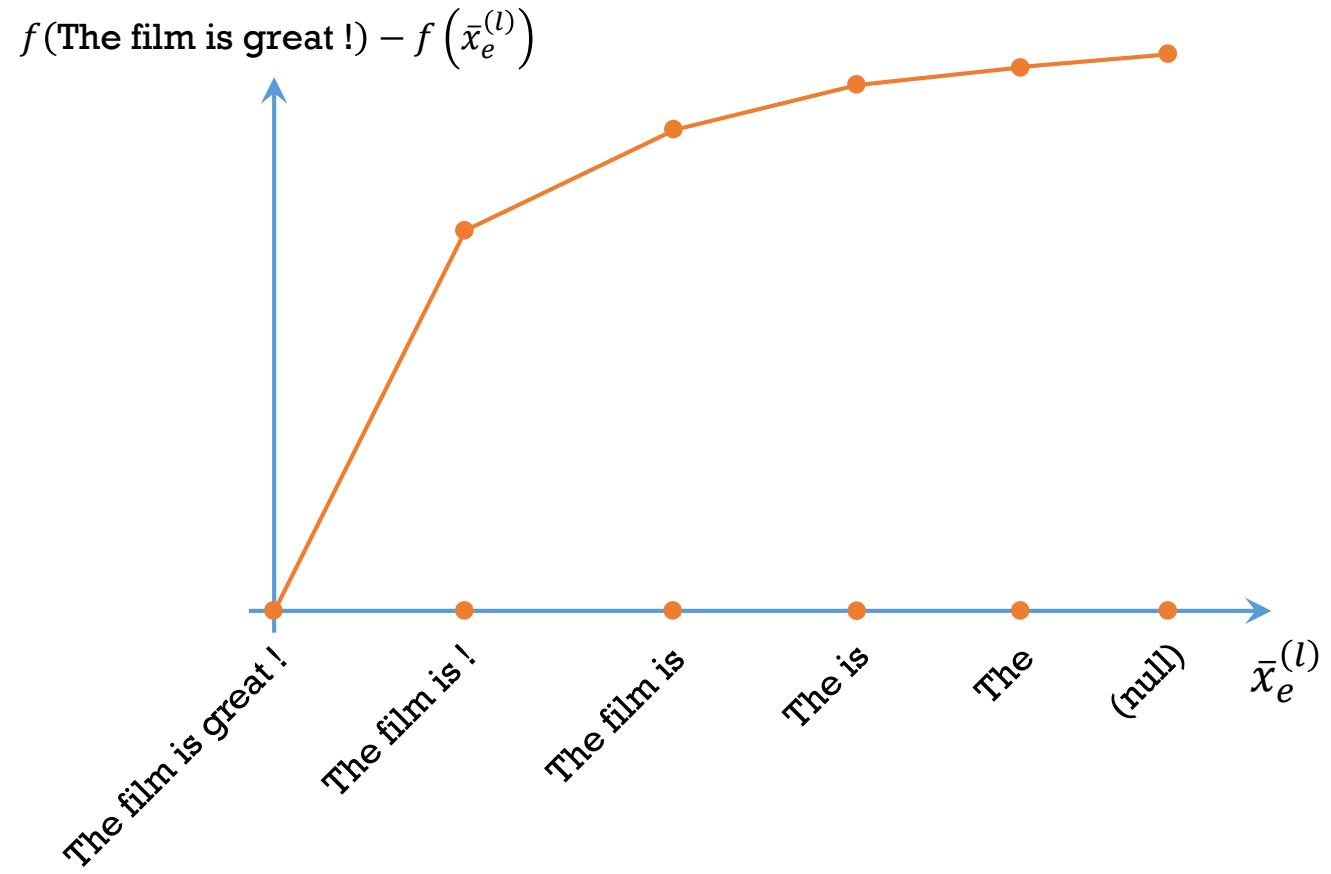




Slides and Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal

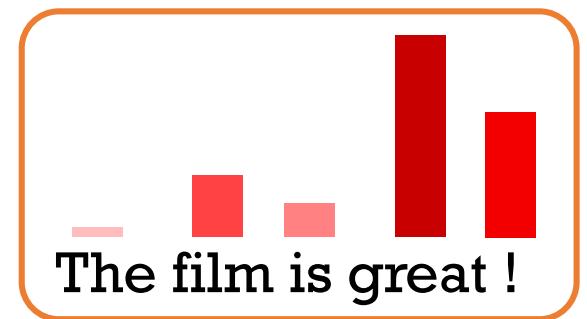
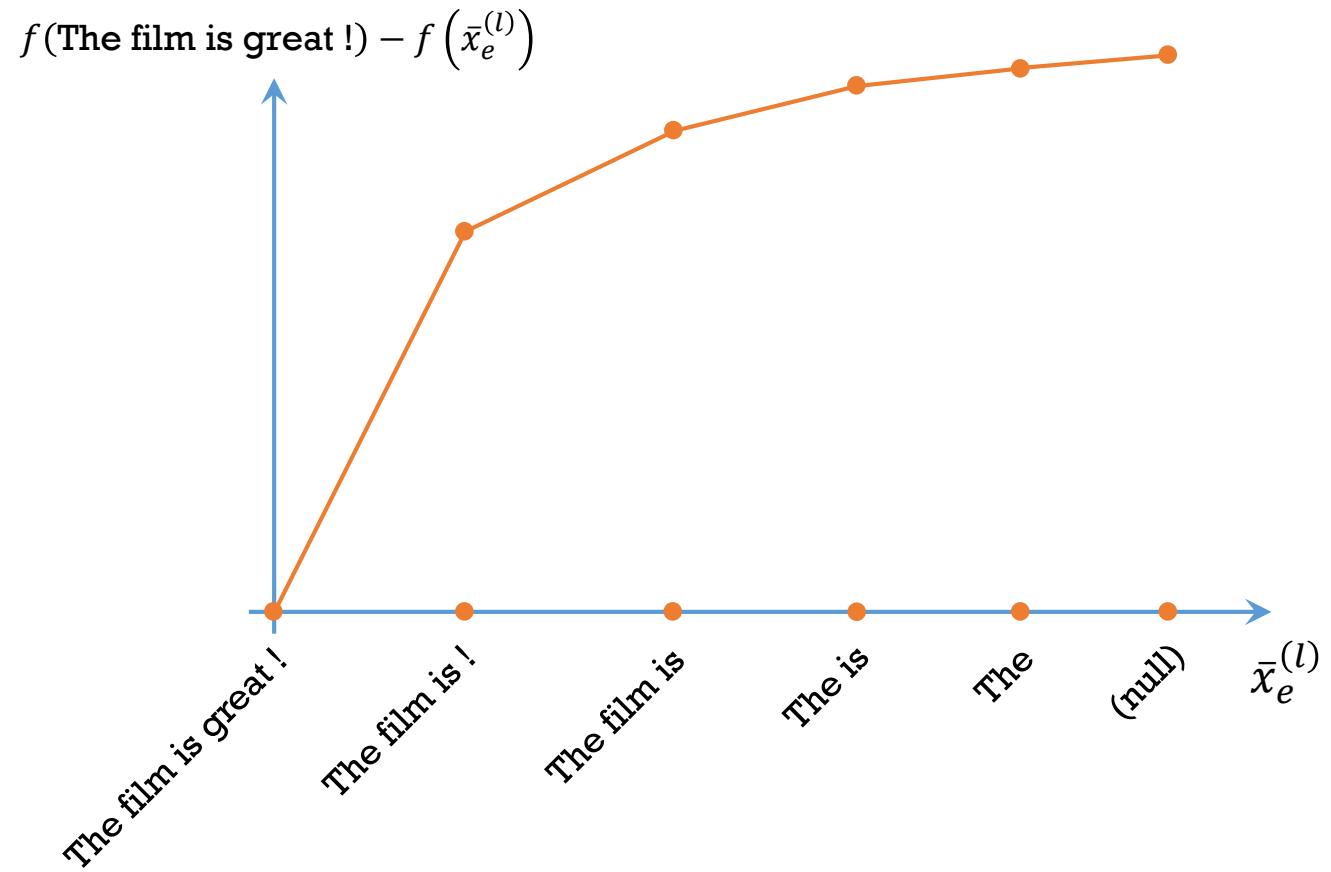




Slides and Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal



Comprehensiveness

$$\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

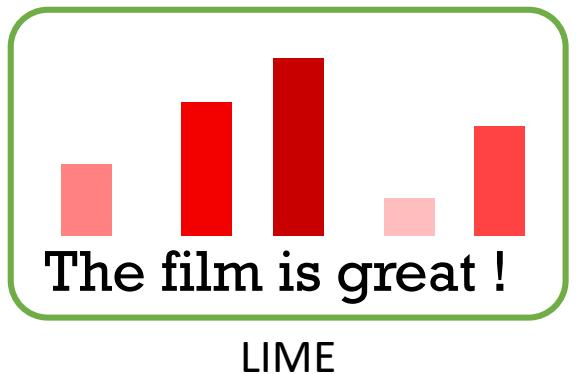
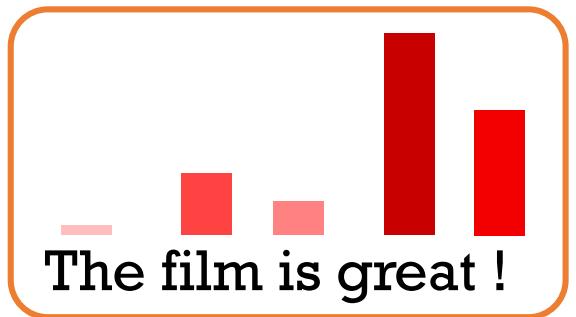
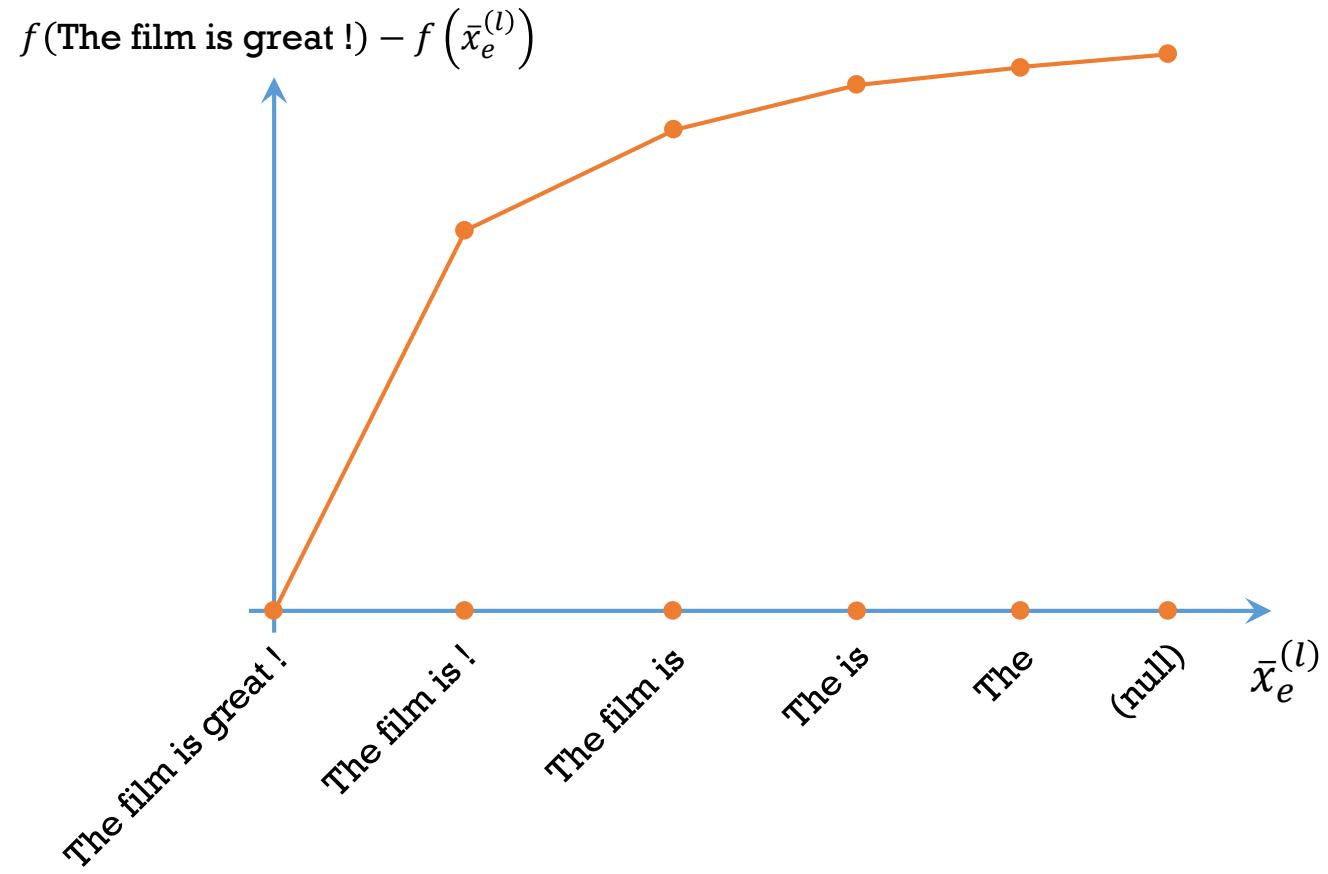
$\bar{x}^{(l)}$: input x with l most important features removed according to e
 $e \in \mathbb{R}^D$ for D -dim input



Slides and Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal

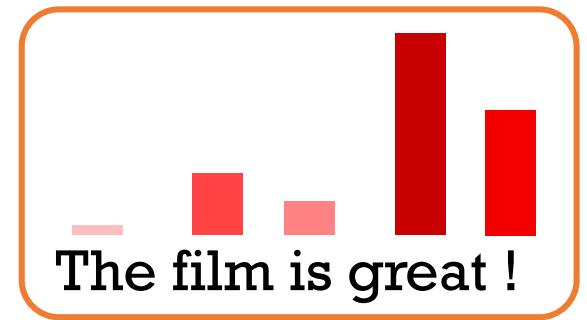
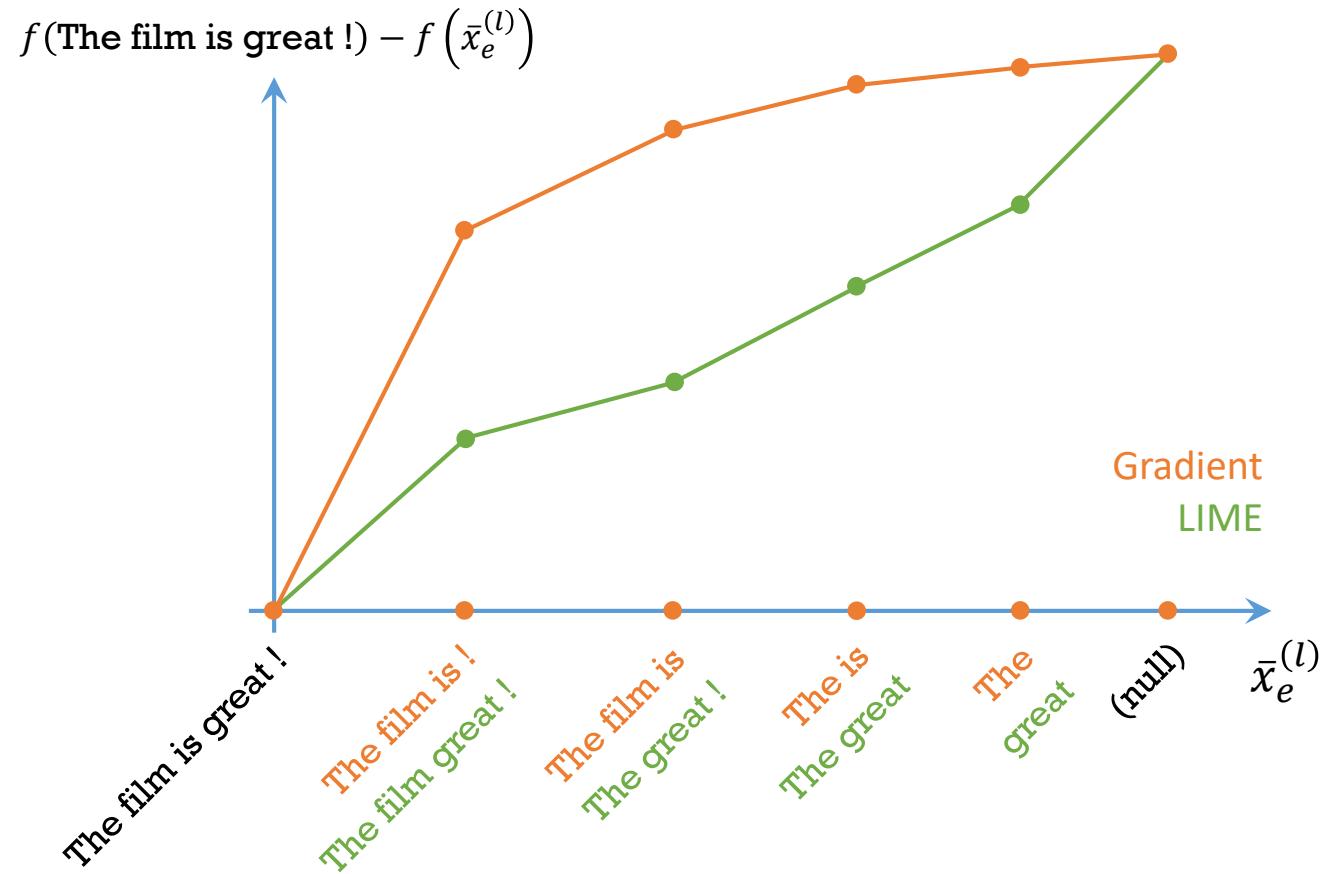




Slides and Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal





Slides and
Resources

Common Evaluation Metrics

- Comprehensiveness and sufficiency
 - Also known as deletion and insertion metrics
 - Comprehensiveness also known as area over perturbation curve (AoPC)



Slides and
Resources

Common Evaluation Metrics

- Comprehensiveness and sufficiency
 - Also known as deletion and insertion metrics
 - Comprehensiveness also known as area over perturbation curve (AoPC)
- Decision flip rate under most important feature removal
- Number of removals required for decision flip
- Prediction change rank correlation
- Etc.



Slides and
Resources

Common Evaluation Metrics

- Comprehensiveness and sufficiency
 - Also known as deletion and insertion metrics
 - Comprehensiveness also known as area over perturbation curve (AoPC)
- Decision flip rate under most important feature removal
- Number of removals required for decision flip
- Prediction change rank correlation
- Etc.

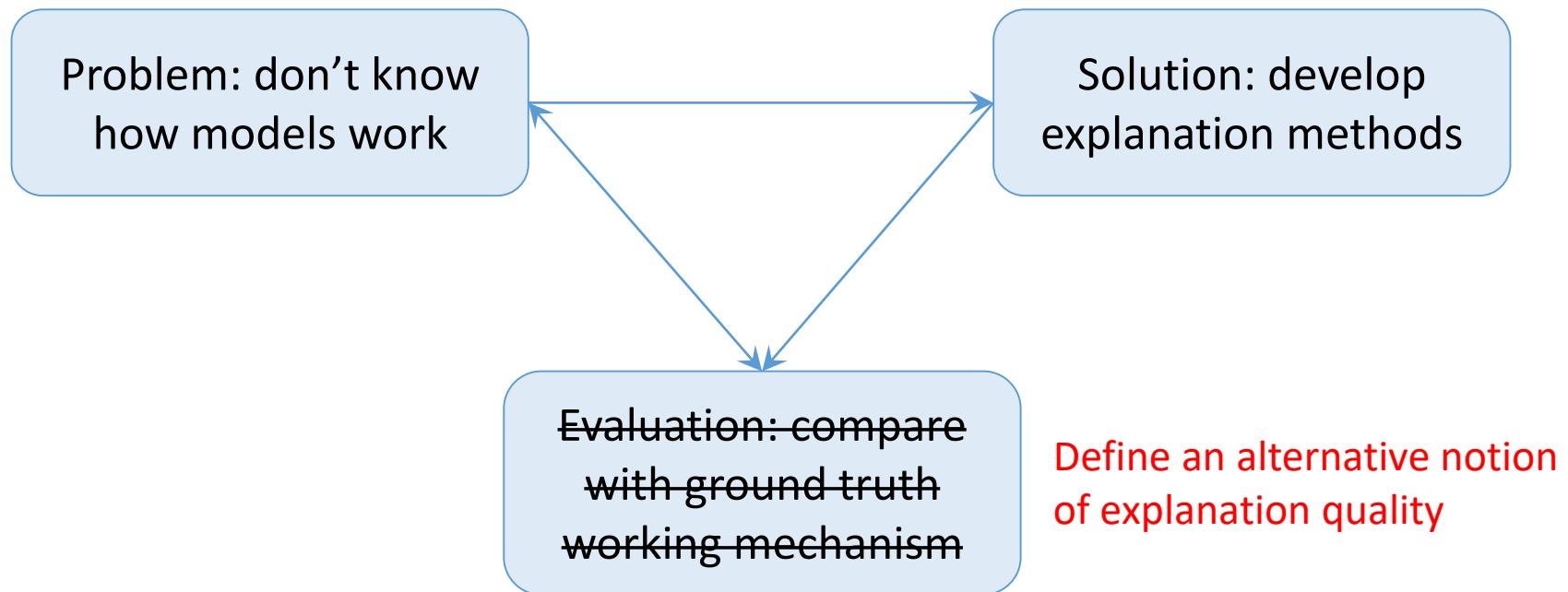
Definition: $e = D(x) \in \mathbb{R}^D$

Evaluation: $q = E(x, e) \in \mathbb{R}$



Slides and
Resources

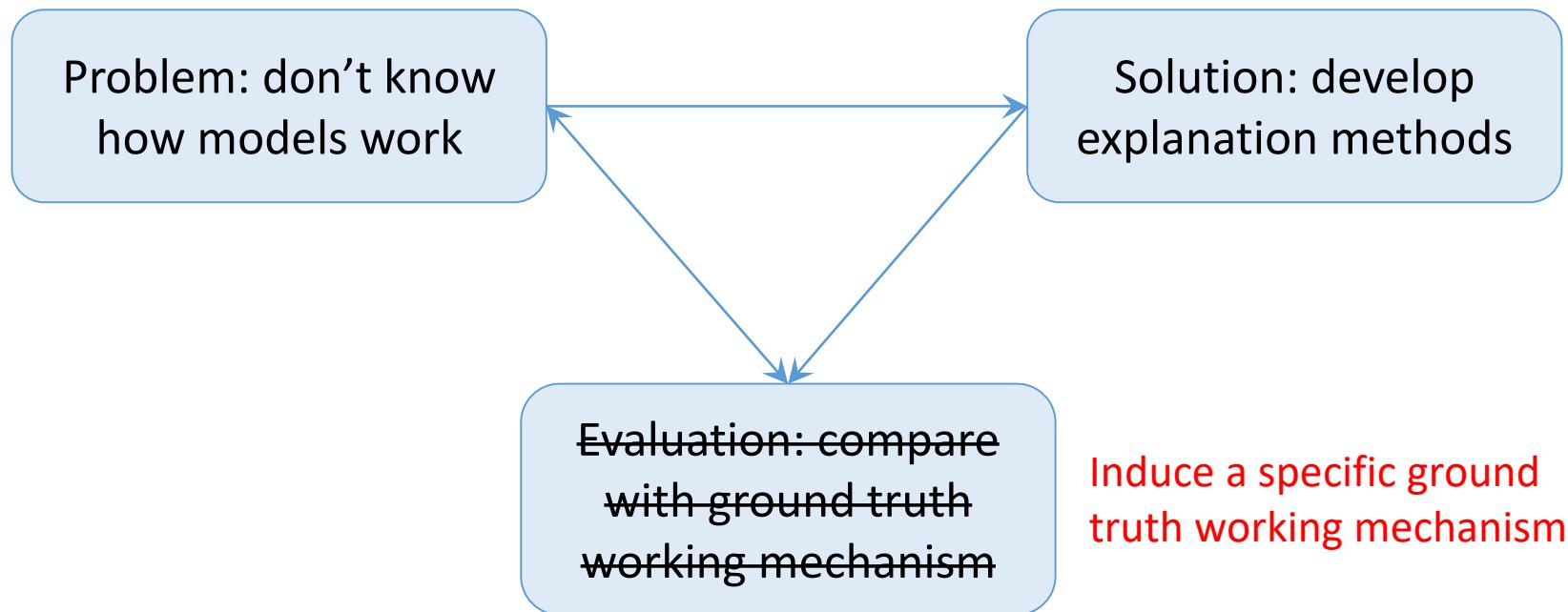
Evaluating Explanations





Slides and
Resources

Evaluating Explanations





Slides and
Resources

Are Explanations (Necessarily) Correct?

*If we **know** that a specific feature is crucial to the model prediction, can feature attribution explanations identify its importance?*



Slides and
Resources

Dataset Modification

X

Y



→ Crow



→ Crow



→ Sparrow

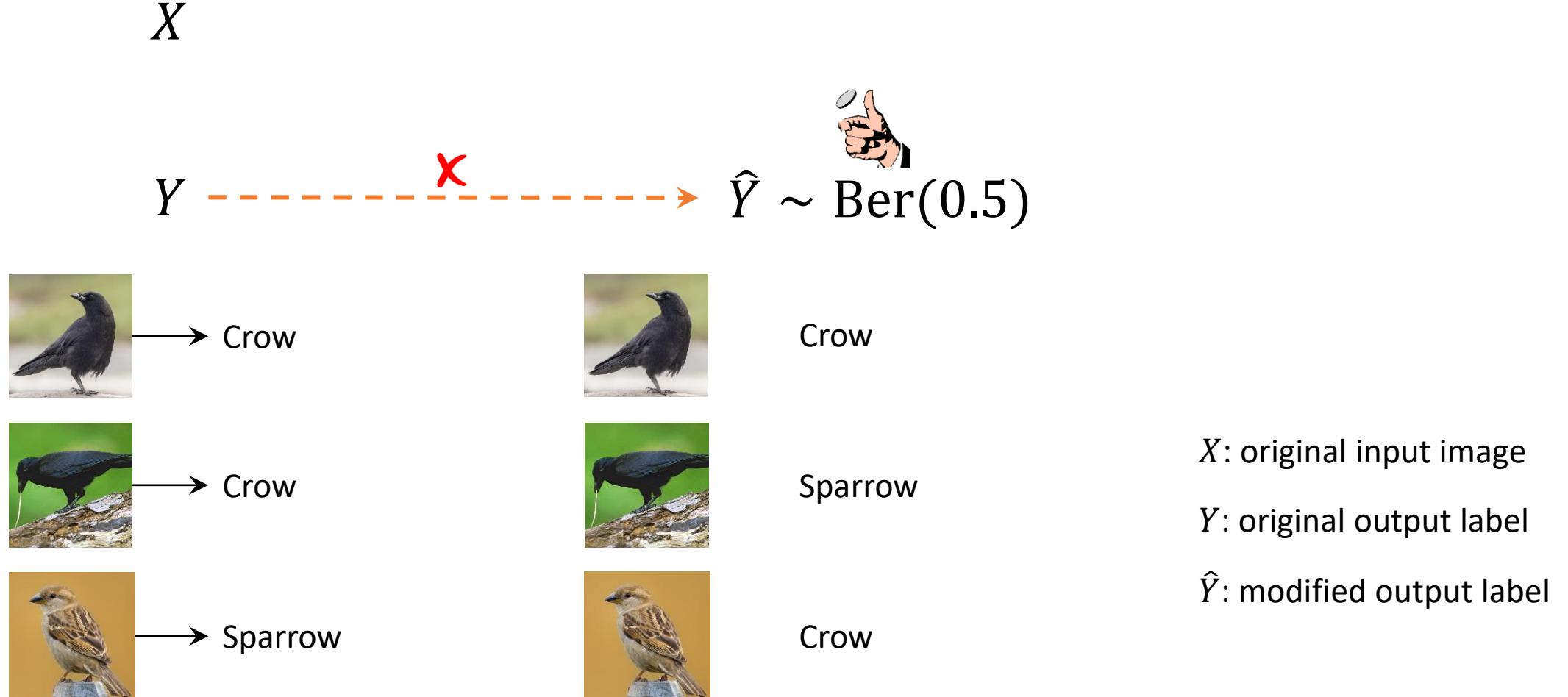
X : original input image

Y : original output label



Slides and
Resources

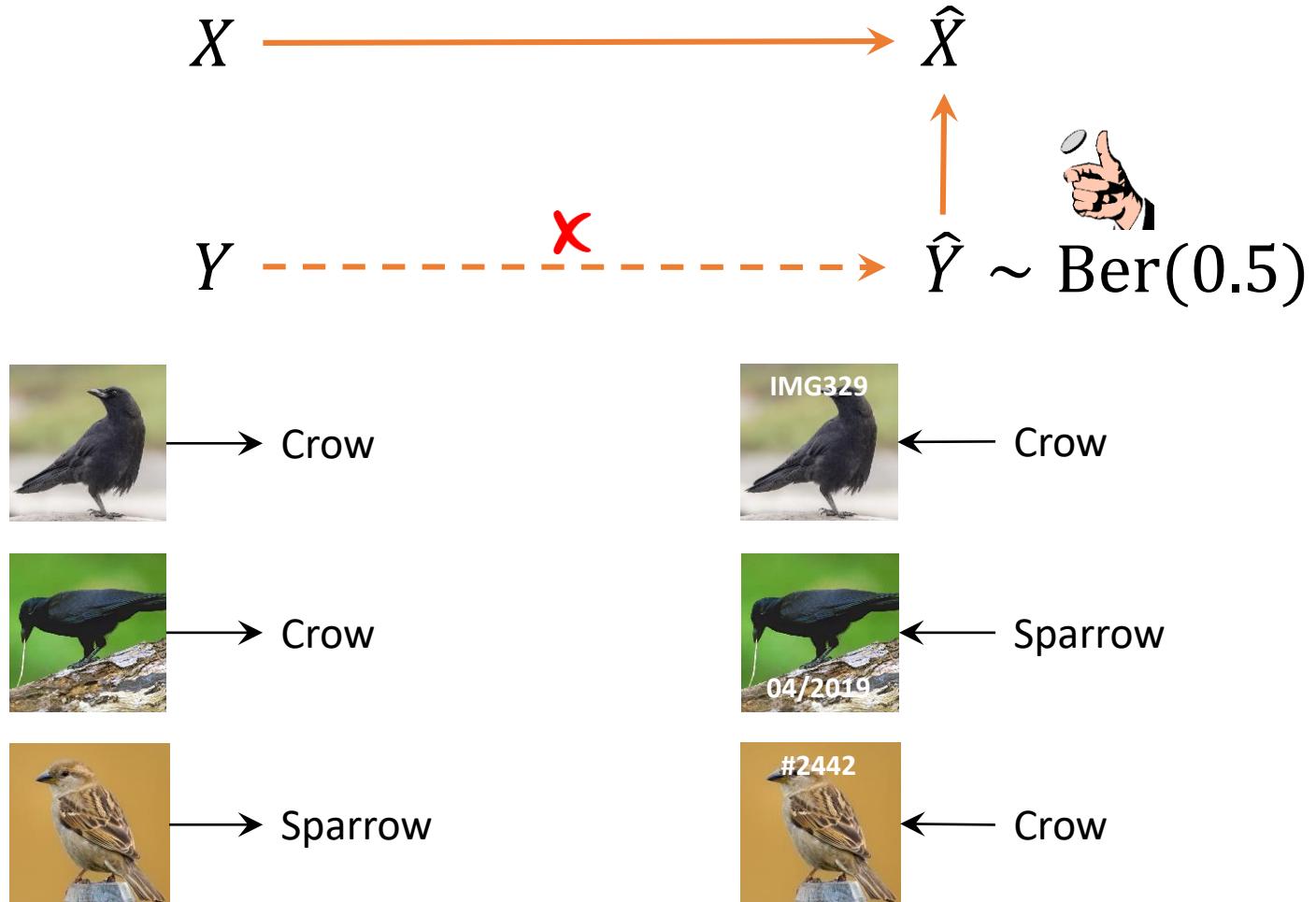
Dataset Modification





Slides and
Resources

Dataset Modification



X : original input image

Y : original output label

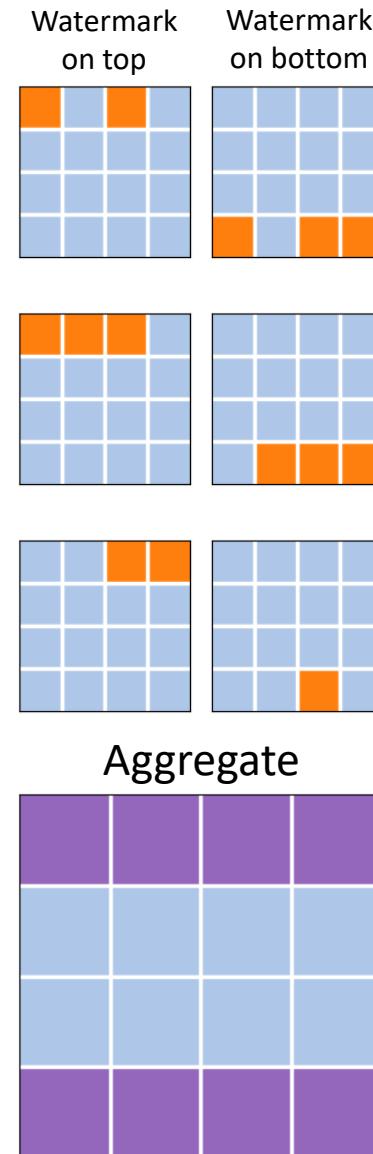
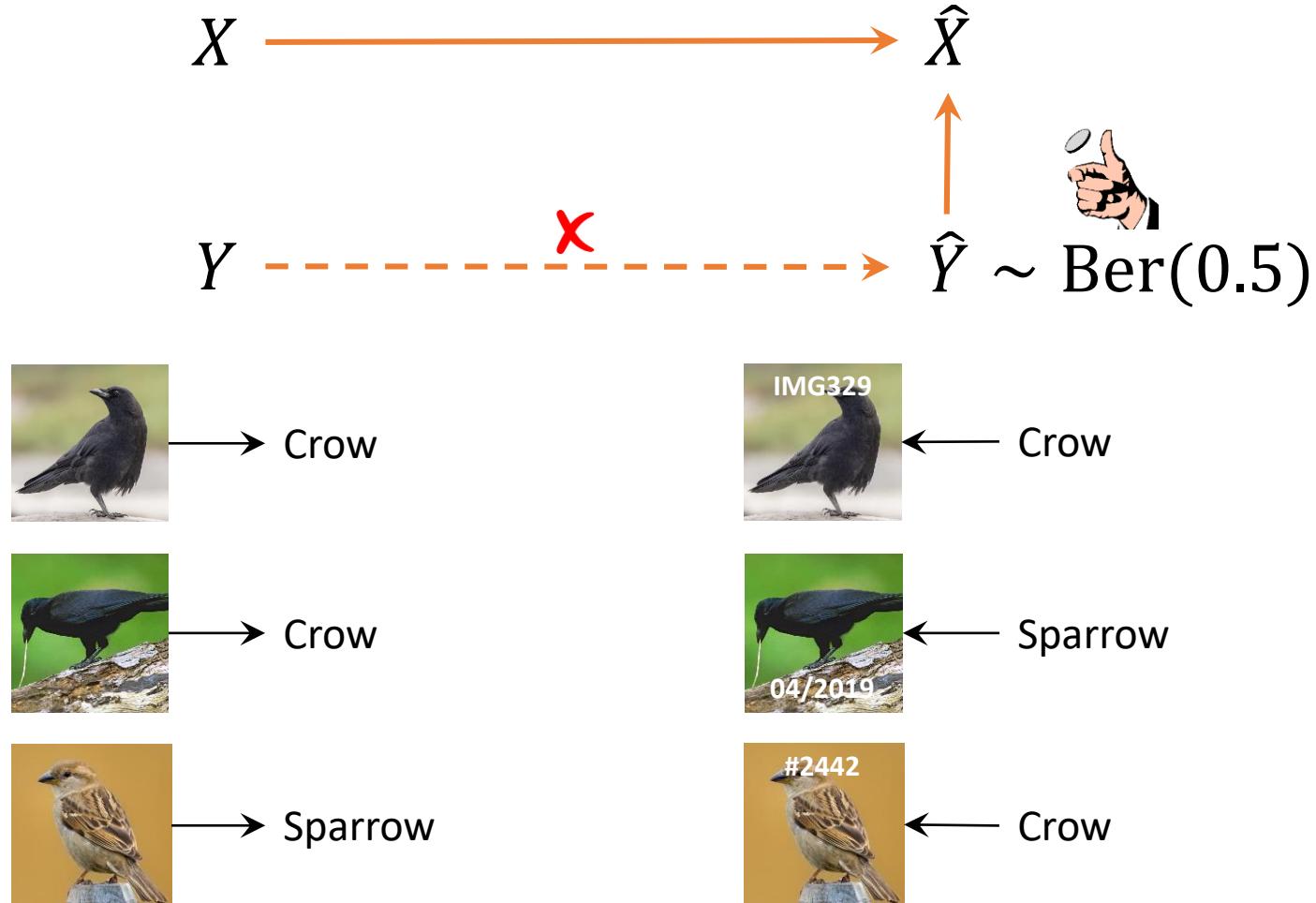
\hat{Y} : modified output label

\hat{X} : modified input image



Slides and Resources

Dataset Modification





Slides and
Resources

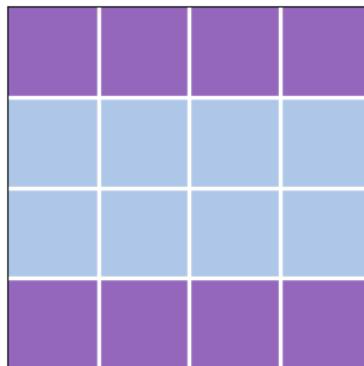
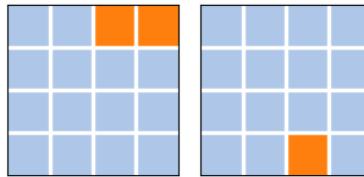
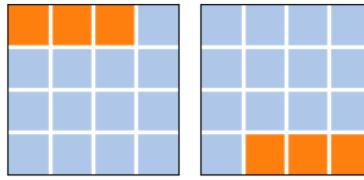
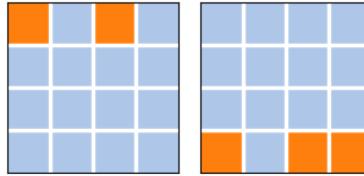
Dataset Modification





Slides and
Resources

Dataset Modification



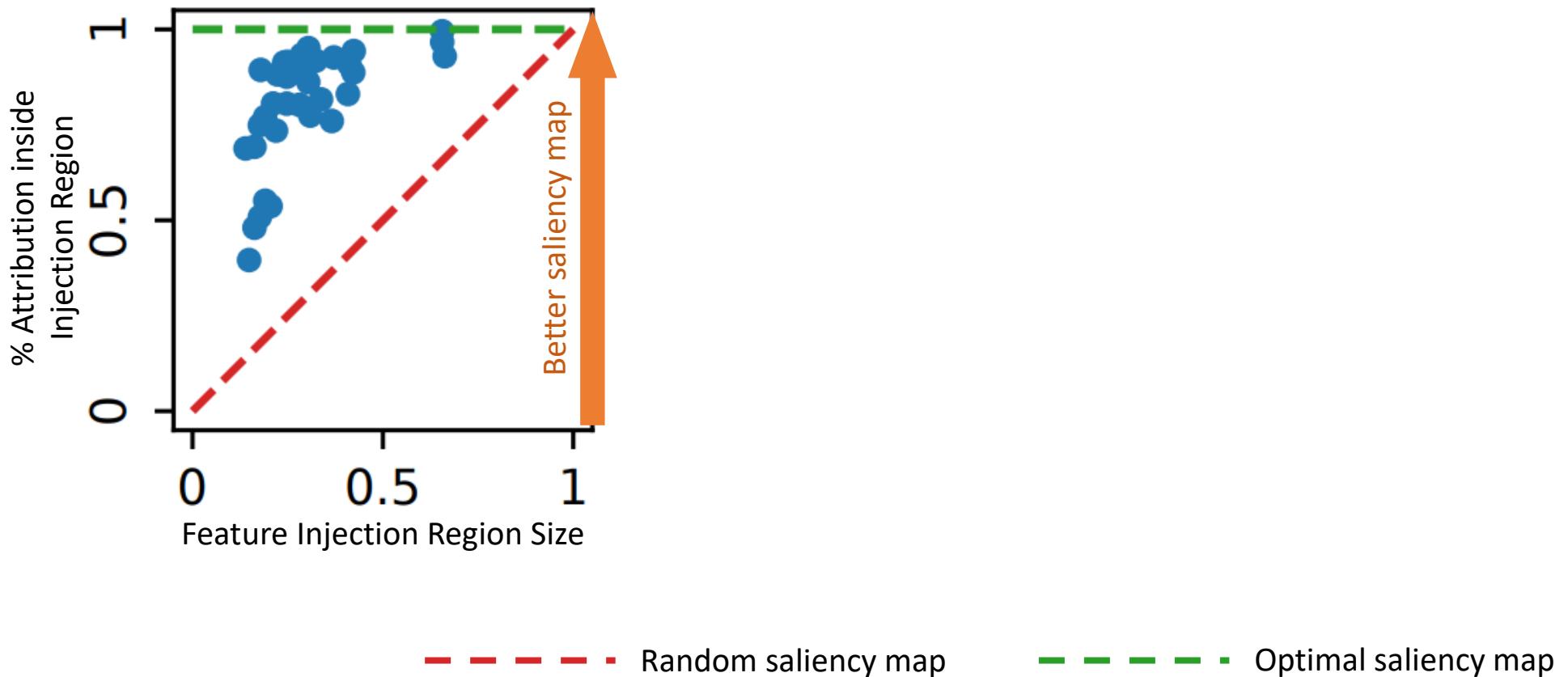
$$\% \text{ Attribution} = \frac{\sum \text{purple } A}{\sum \text{light blue } A + \sum \text{purple } A}$$



Slides and
Resources

Evaluating Image Saliency Maps

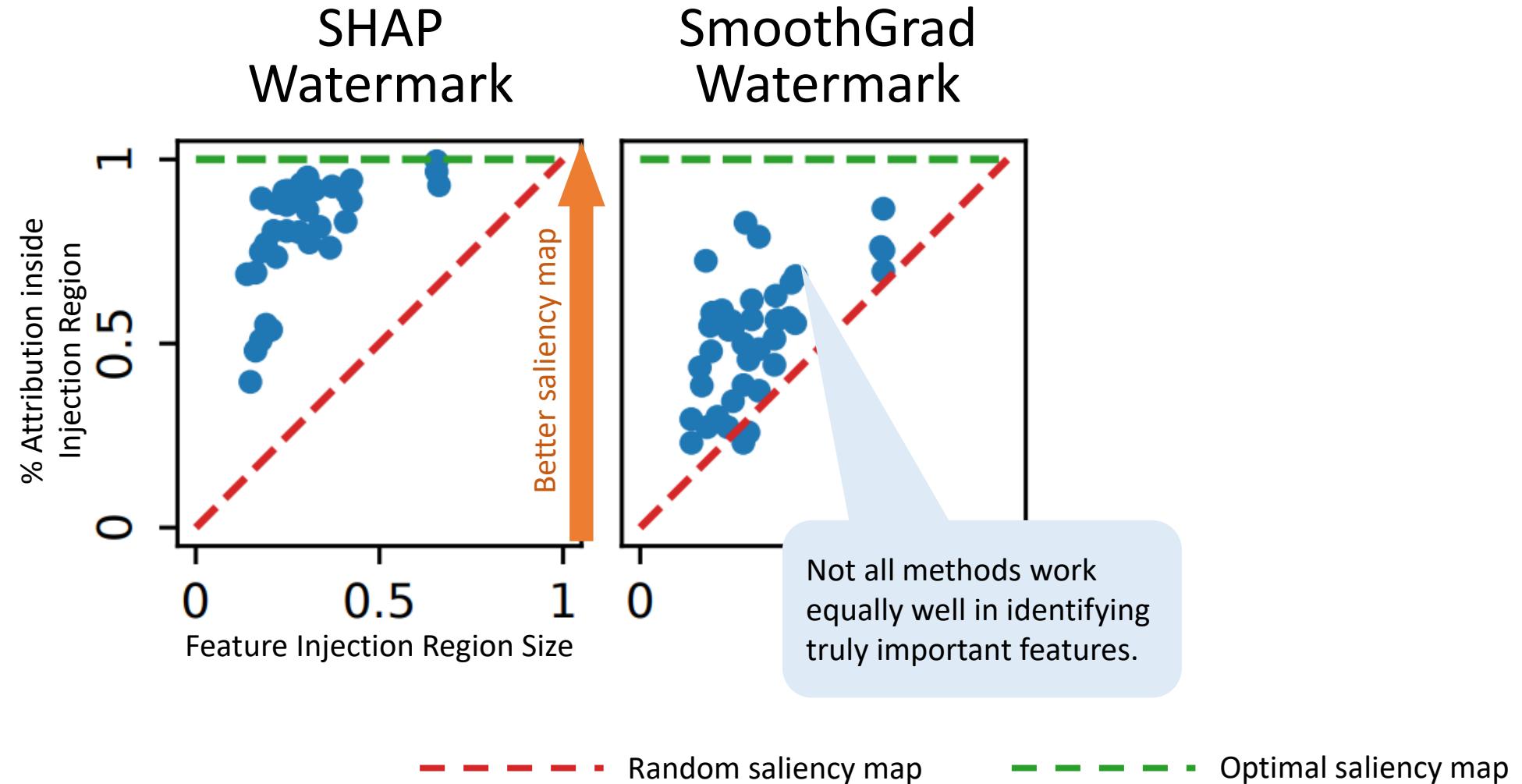
SHAP
Watermark





Slides and
Resources

Evaluating Image Saliency Maps



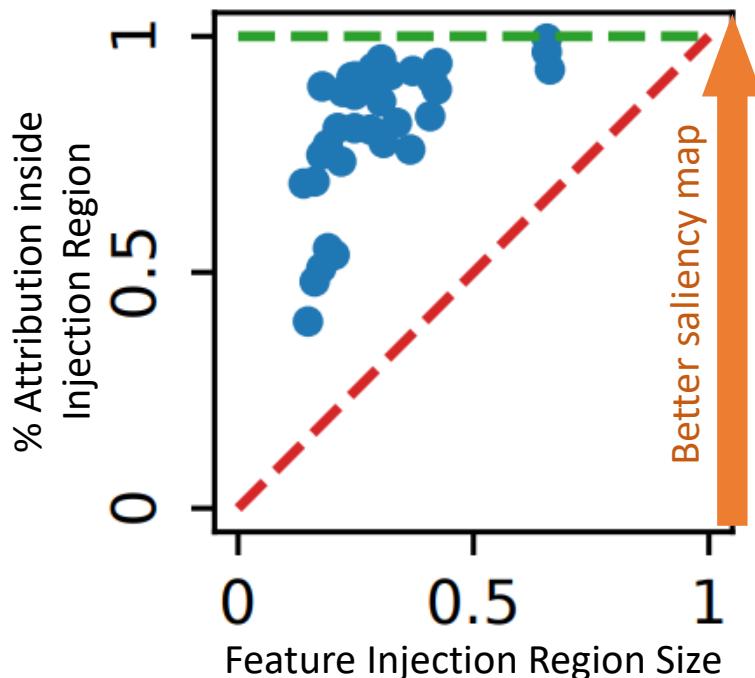


Slides and Resources

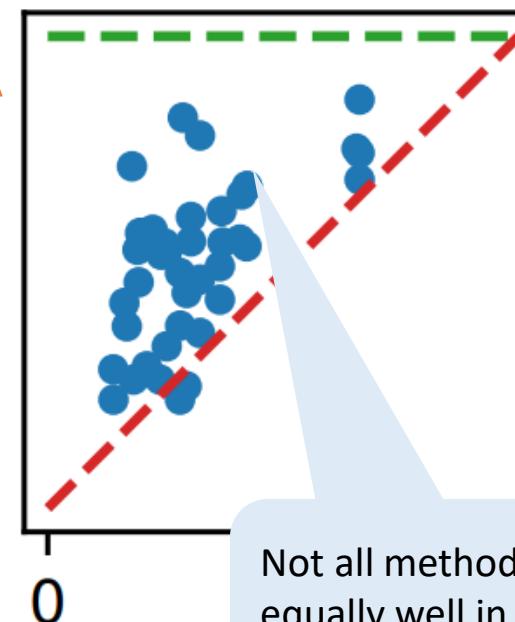
Evaluating Image Saliency Maps

Some methods even struggle against the random baseline on certain feature types.

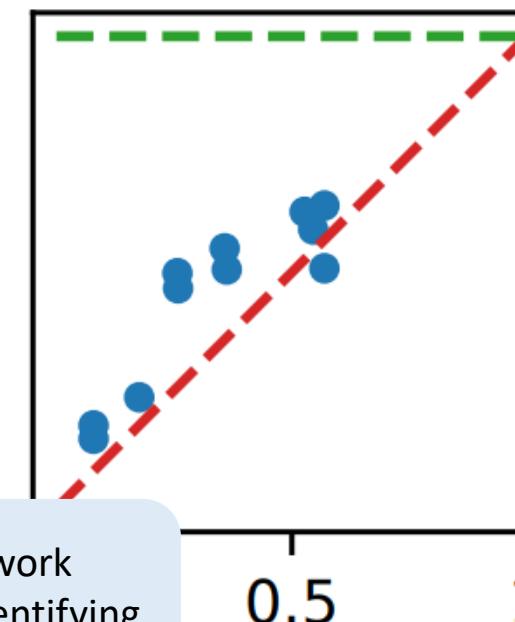
SHAP
Watermark



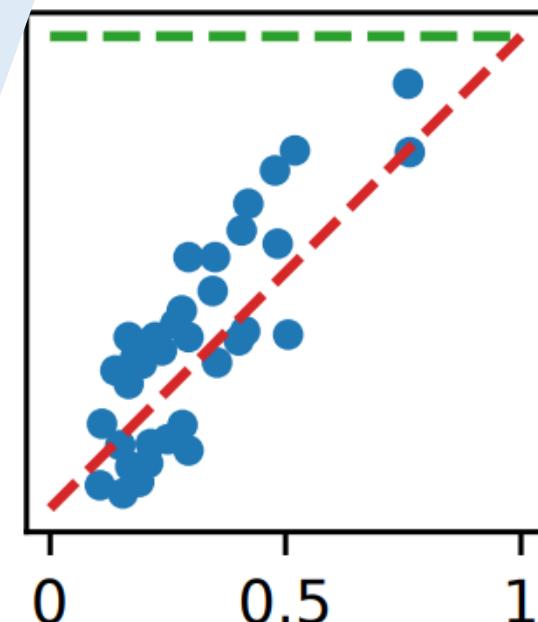
SmoothGrad
Watermark



SmoothGrad
Blurring



GradCAM
Brightness



Not all methods work
equally well in identifying
truly important features.

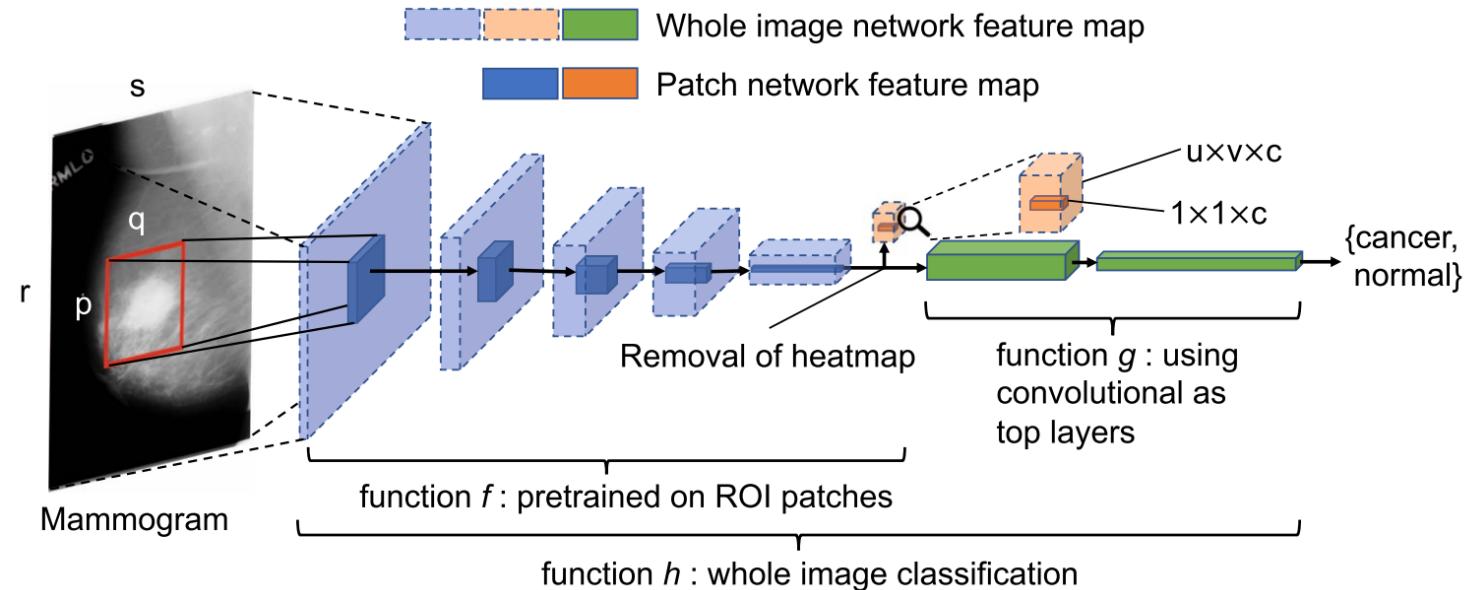
— Random saliency map

— Optimal saliency map



Slides and Resources

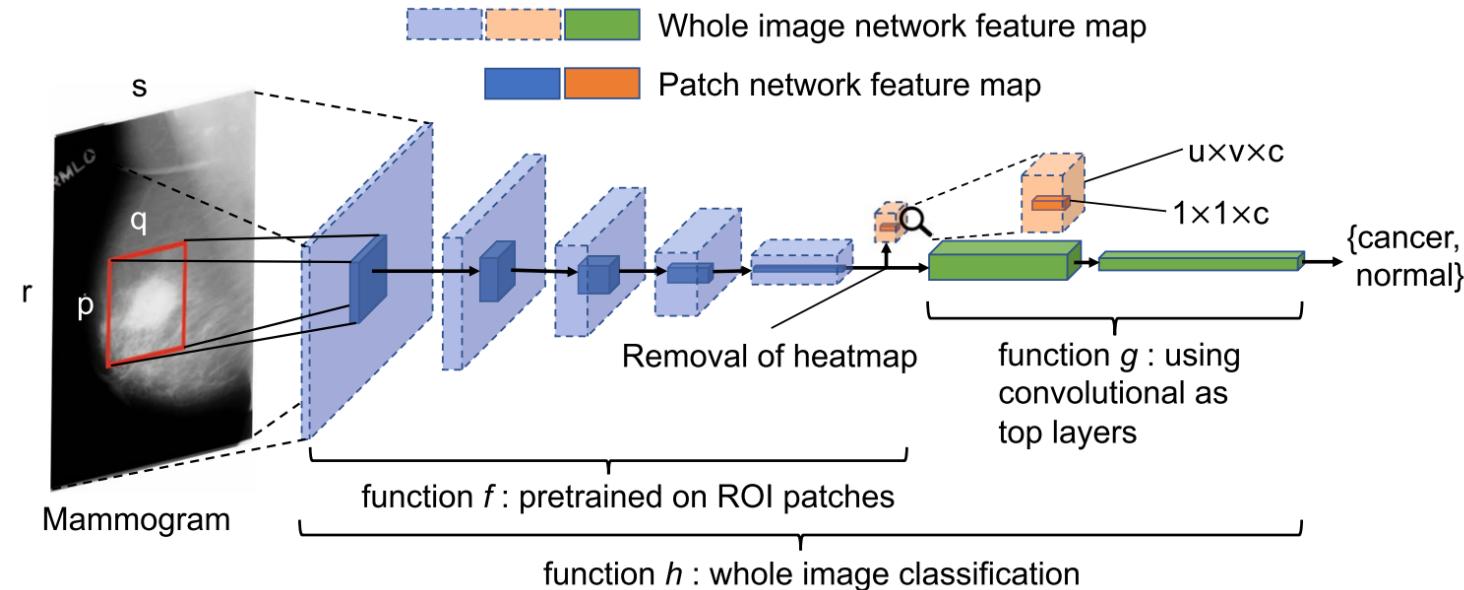
Practical Implications





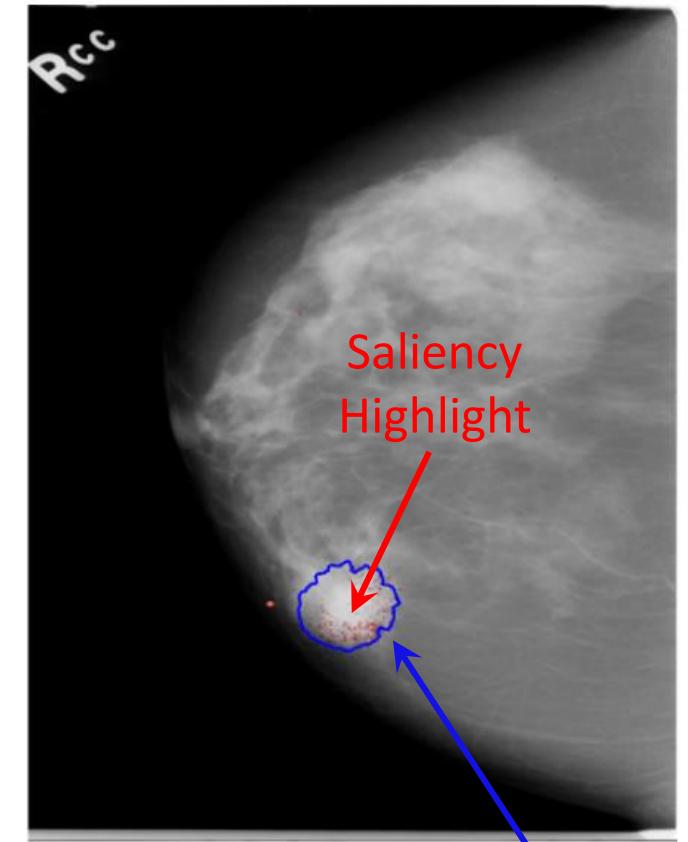
Slides and Resources

Practical Implications



"A saliency map illustrates which area of the input image is considered to be responsible for the cancer [...] Figure 4a [...] shows that the image classifier was able to **correctly locate the cancerous region** on which its decision was based."

(Shen et al., *Scientific Reports*, 2019)



Tumor
Segmentation



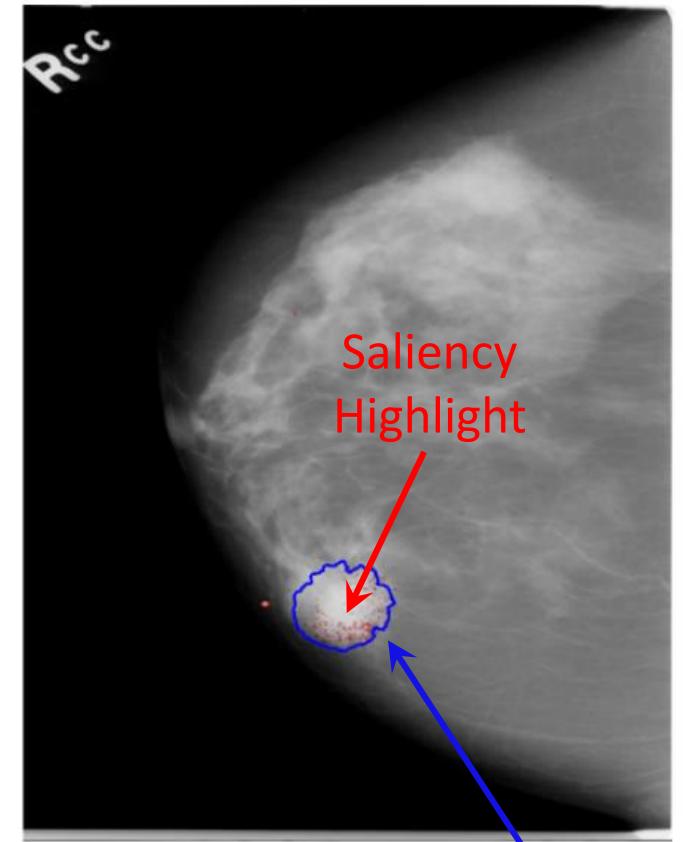
Slides and Resources

Practical Implications



"A saliency area of the input image is considered to be the cancer [...] Figure 4a [...] shows that the image classifier was able to correctly locate the cancerous region on which its decision was based."

(Shen et al., *Scientific Reports*, 2019)



(a)
Tumor
Segmentation



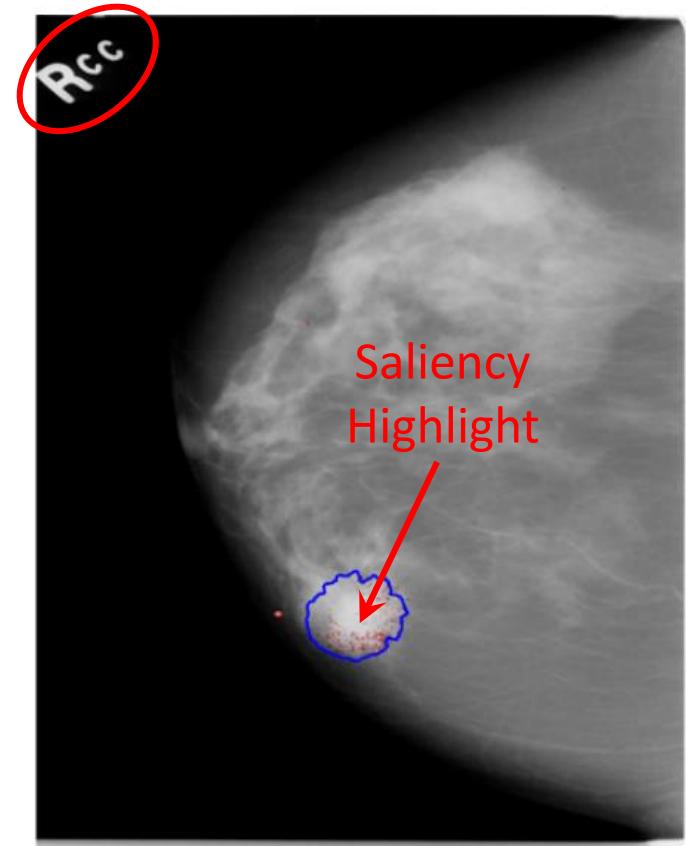
Slides and Resources

Practical Implications



"A saliency area of the input image is considered to be the cancer [...] Figure 4a [...] shows that the image classifier was able to correctly locate the cancerous region on which its decision was based."

(Shen et al., *Scientific Reports*, 2019)



Unknown spurious correlation
Can we trust these explainers?



Slides and
Resources

Explanation Understandability

Do the explanations correctly explain the
model prediction logic?



Slides and
Resources

Explanation Understandability

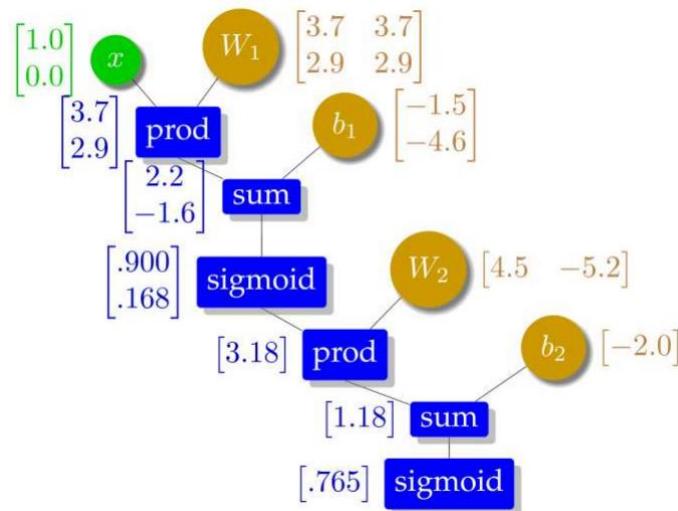
Do the explanations correctly explain the model prediction logic?

Do people correctly understand the model prediction logic from the explanations?



Slides and
Resources

Correct but Not Understandable Explanations

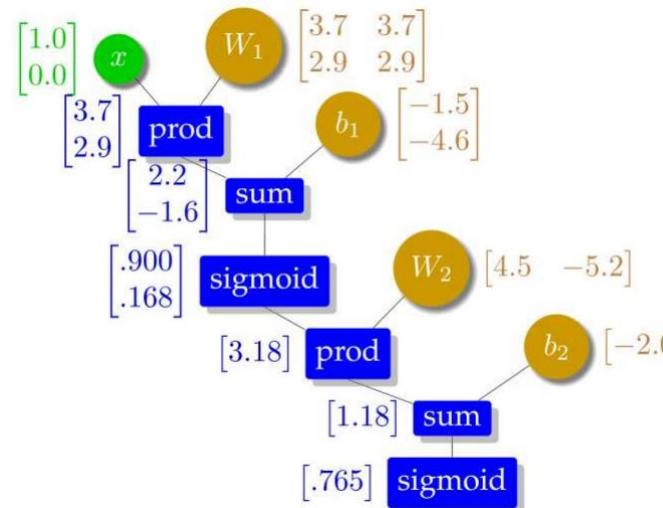


Computation trace: a fully correct explanation for the prediction

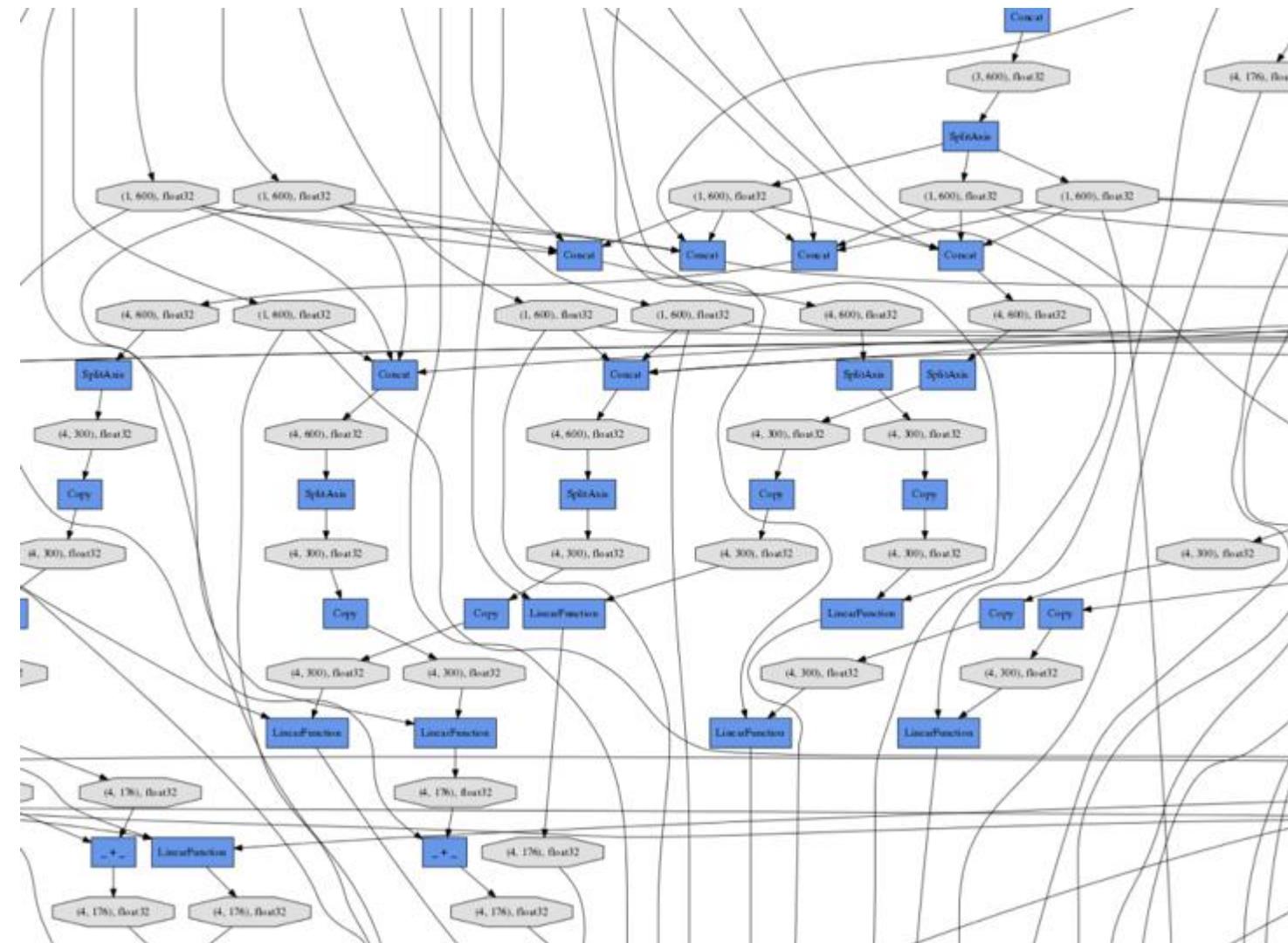


Slides and Resources

Correct but Not Understandable Explanations



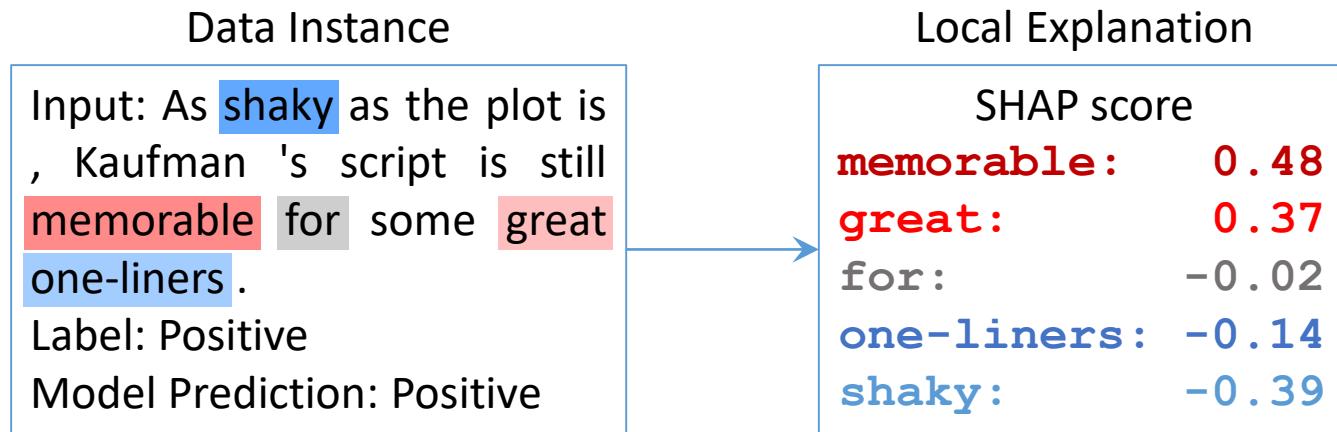
Computation trace: a fully correct
but totally not understandable
explanation for the prediction





Slides and
Resources

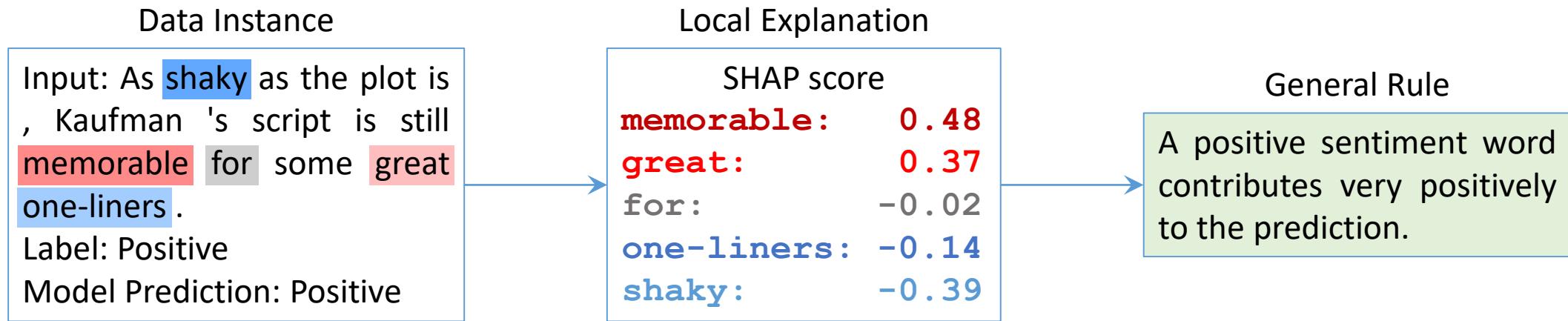
Understanding a Sentiment Classifier





Slides and
Resources

Understanding a Sentiment Classifier





Slides and
Resources

Understanding a Sentiment Classifier

As shaky as the plot is, Kaufman's script is still **memorable** for some **great** one-liners.

Label = Positive
Prediction = Positive

memorable: 0.48
great: 0.37

A positive sentiment word contributes very positively to the prediction



Slides and
Resources

Understanding a Sentiment Classifier

As shaky as the plot is, Kaufman's script is still **memorable** for some **great** one-liners.

memorable: 0.48
great: 0.37

Label = Positive
Prediction = Positive

Just because it really happened to you, honey, does n't mean that it's **attractive** to anyone else.

attractive: 0.00

Label = Negative
Prediction = Negative

A positive sentiment word contributes very positively to the prediction

- Supporting instance
- Opposing instance



Slides and
Resources

Understanding a Sentiment Classifier

As shaky as the plot is, Kaufman's script is still **memorable** for some **great** one-liners.

memorable: 0.48
great: 0.37

Label = Positive
Prediction = Positive

Just because it really happened to you, honey, does n't mean that it's **attractive** to anyone else.

attractive: 0.00

Label = Negative
Prediction = Negative

A positive sentiment word contributes very positively to the prediction...

unless it is negated

- Supporting instance
- Opposing instance



Slides and
Resources

Understanding a Sentiment Classifier

As shaky as the plot is, Kaufman's script is still **memorable** for some **great** one-liners.

memorable: 0.48
great: 0.37

Label = Positive
Prediction = Positive

Just because it really happened to you, honey, does n't mean that it's **attractive** to anyone else.

attractive: 0.00

Label = Negative
Prediction = Negative

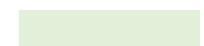
A positive sentiment word contributes very positively to the prediction...

unless it is negated

Daring, mesmerizing and exceedingly hard to forget.

Daring: 0.13
mesmerizing: 0.06

Label = Positive
Prediction = Positive

 Supporting instance

 Opposing instance



Slides and
Resources

Understanding a Sentiment Classifier

As shaky as the plot is, Kaufman's script is still **memorable** for some **great** one-liners.

memorable: 0.48
great: 0.37

Label = Positive
Prediction = Positive

Just because it really happened to you, honey, does n't mean that it's **attractive** to anyone else.

attractive: 0.00

Label = Negative
Prediction = Negative

A positive sentiment word contributes very positively to the prediction...

unless it is negated, or near another positive word

Daring, mesmerizing and exceedingly hard to forget.

Daring: 0.13
mesmerizing: 0.06

Label = Positive
Prediction = Positive





Slides and
Resources

Understanding a Sentiment Classifier

As shaky as the plot is, Kaufman's script is still **memorable** for some **great** one-liners.

memorable: 0.48
great: 0.37

Label = Positive
Prediction = Positive

Just because it really happened to you, honey, does n't mean that it's **attractive** to anyone else.

attractive: 0.00

Label = Negative
Prediction = Negative

A positive sentiment word contributes very positively to the prediction...

unless it is negated, or near another positive word

Daring, mesmerizing and exceedingly hard to forget.

Daring: 0.13
mesmerizing: 0.06

Label = Positive
Prediction = Positive

Ranks among Williams' **best** screen work.

best: 0.05

Label = Positive
Prediction = Positive

Supporting instance

Opposing instance



Slides and
Resources

Understanding a Sentiment Classifier

As shaky as the plot is, Kaufman's script is still **memorable** for some **great** one-liners.

memorable: 0.48
great: 0.37

Label = Positive
Prediction = Positive

Just because it really happened to you, honey, does n't mean that it's **attractive** to anyone else.

attractive: 0.00

Label = Negative
Prediction = Negative

A positive sentiment word **sometimes** contributes very positively to the prediction...
unless it is negated, or near another positive word

Daring, mesmerizing and exceedingly hard to forget.

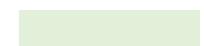
Daring: 0.13
mesmerizing: 0.06

Label = Positive
Prediction = Positive

Ranks among Williams' **best** screen work.

best: 0.05

Label = Positive
Prediction = Positive

 Supporting instance

 Opposing instance



Slides and
Resources

Understanding a Sentiment Classifier

As shaky as the plot is, Kaufman's script is still **memorable** for some **great** one-liners.

memorable: 0.48
great: 0.37

Label = Positive
Prediction = Positive

Just because it really happened to you, honey, does n't mean that it's **attractive** to anyone else.

attractive: 0.00

Label = Negative
Prediction = Negative

A positive sentiment word **sometimes** contributes very positively to the prediction...
unless it is negated, or near another positive word

Daring, mesmerizing and exceedingly hard to forget.

Daring: 0.13
mesmerizing: 0.06

Label = Positive
Prediction = Positive

Ranks among Williams' **best** screen work.

best: 0.05

Label = Positive
Prediction = Positive



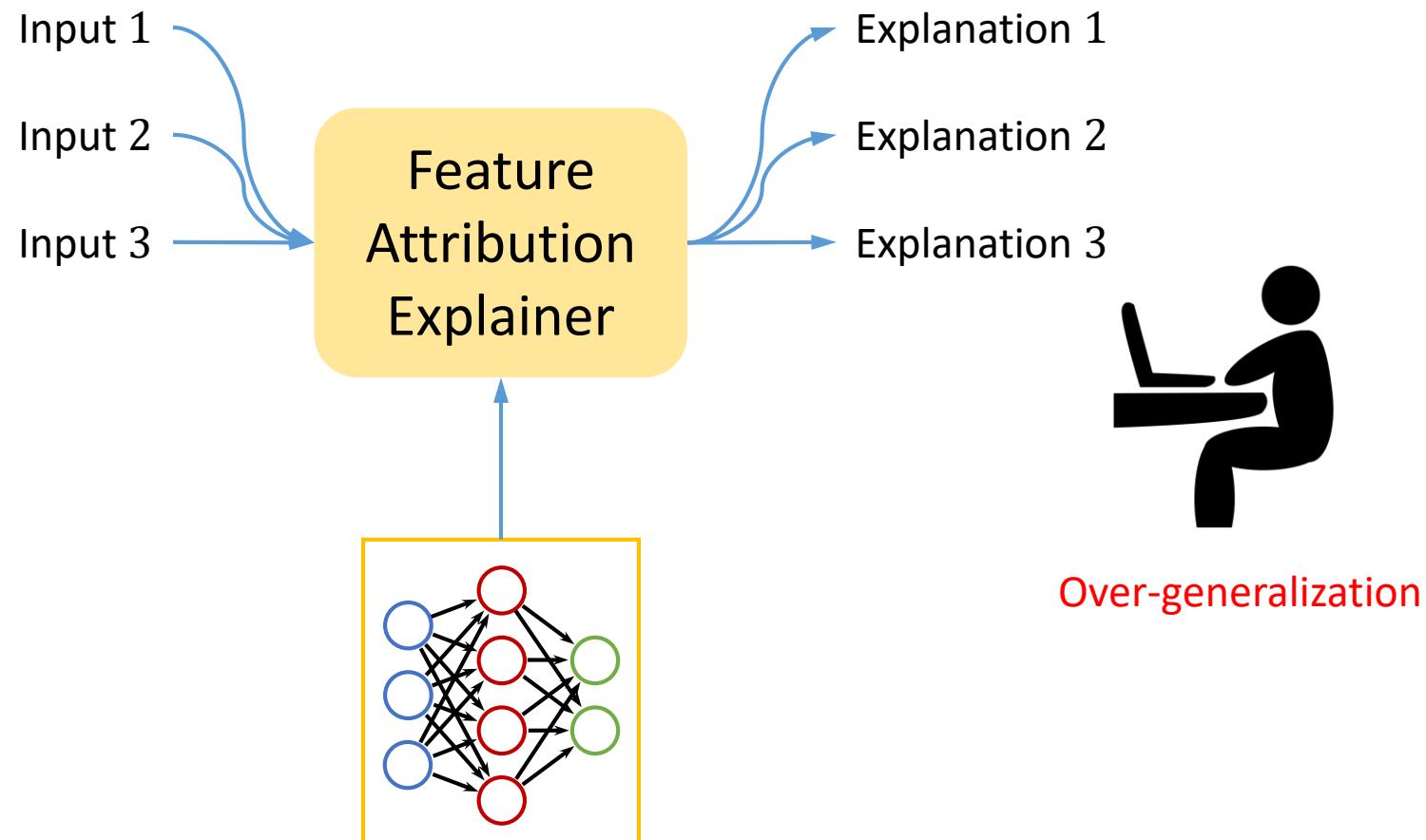
(many more counter-examples not shown)

- Supporting instance
- Opposing instance



Slides and
Resources

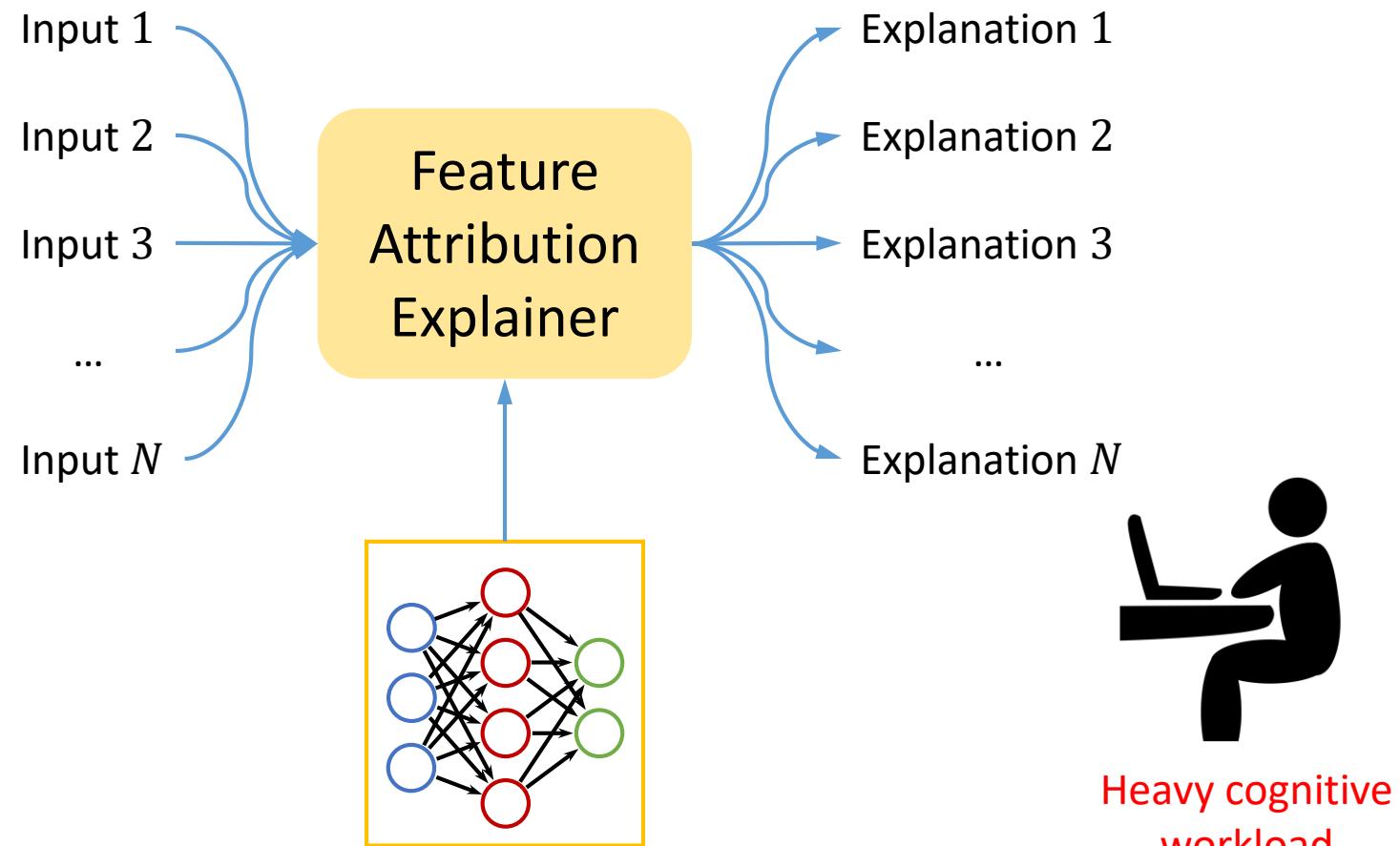
Explanation Summary (ExSum)





Slides and
Resources

Explanation Summary (ExSum)

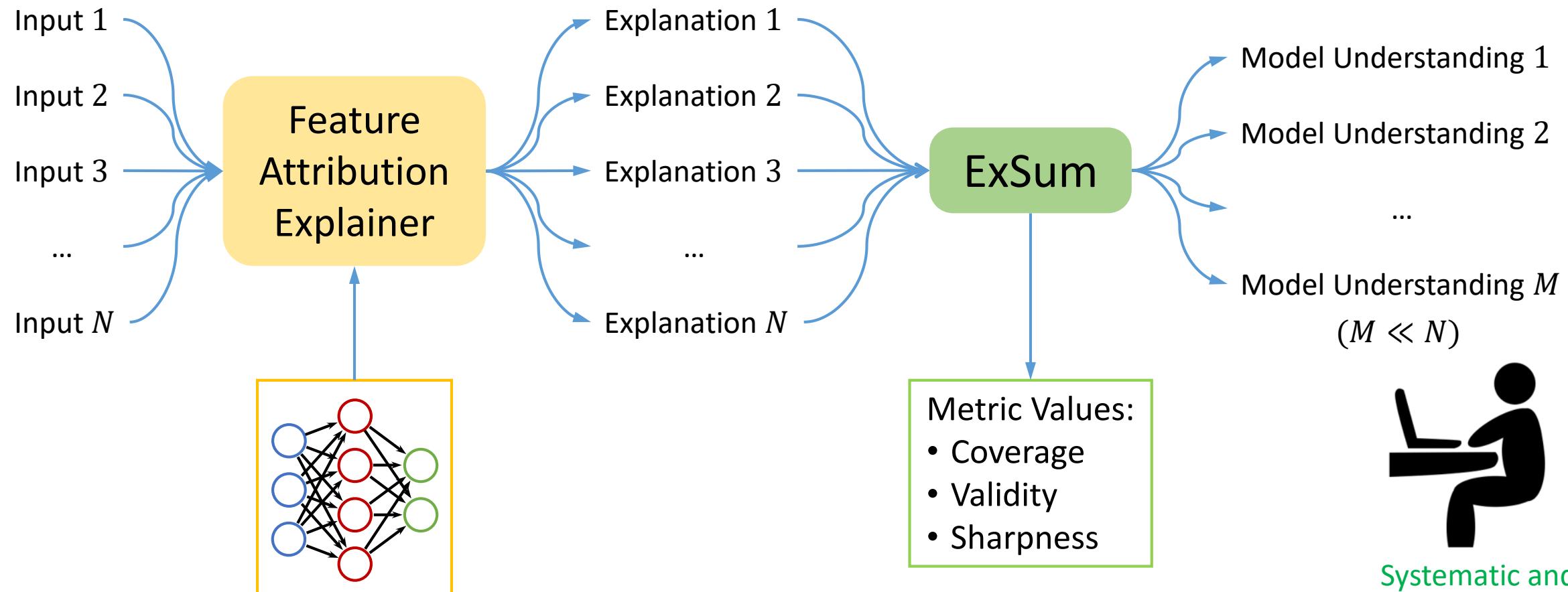


Heavy cognitive
workload



Slides and Resources

Explanation Summary (ExSum)



Systematic and quantitative model understanding



Slides and Resources

GUI for Developing ExSum Rules

ExSum Inspection

Rule Union: (((R1 > R4) > R3) > R5) > R6) > R7
CF Without Rule 7: (((R1 > R4) > R3) > R5) > R6

Rule Selection

- Rule 1: negation
- Rule 2: highly positive adjectives have positive saliency
- Rule 3: highly negative adjectives have negative saliency
- Rule 4: highly positive words have positive saliency
- Rule 5: highly negative words have negative saliency
- Rule 6: Person names have small saliency
- Rule 7: Stop words have small saliency

Reset **Save**

Metric Values

Metrics	CF → Full	Selected
Coverage	0.130 → 0.605	0.475
Validity	0.911 → 0.915	0.916
Sharpness	0.539 → 0.269	0.195

Parameter Values

Applicability Function Parameters	Behavior Function Parameters
(None)	saliency lower range AutoTune -0.1
	saliency upper range AutoTune 0.13

Example Visualization

Rule Union → Selected Rule Sentence → FEU All → Invalid New Examples

y=0 : 0.00 If you collected all the moments of coherent dialogue , they still would n't add up to the time required to boil a four - minute egg .

y=1 : 1.00 Ranks among Williams ' best screen work .

y=1 : 1.00 Like the film 's almost anthropologically detailed realization of early - '80s suburbia , it 's significant without being overstated .

y=0 : 0.00 It is a comedy that 's not very funny and an action movie that is not very thrilling (and an uneasy alliance , at that) .

y=1 : 1.00 A triumph , relentless and beautiful in its downbeat darkness .

y=0 : 0.00 They felt like the same movie to me .

y=1 : 1.00 The story feels more like a serious read , filled with heavy doses of always enticing Sayles dialogue .

y=1 : 1.00 Behind the snow games and lovable Siberian huskies (plus one sheep dog) , the picture hosts a parka-wrapped dose of heart .

y=1 : 1.00 It is a film that will have people walking out halfway through , will encourage others to stand up and applaud , and will , undoubtedly , leave both camps engaged in a ferocious debate for years to come .

y=0 : 0.05 De Niro may enjoy the same free ride from critics afforded to Clint Eastwood in the lazy Bloodwork .

y=0 : 0.01 The film might have been more satisfying if it had , in fact , been fleshed out a little more instead of going for easy smiles .



Slides and
Resources

GUI for Developing ExSum Rules

ExSum Inspection

Rule Union: (((R1 > R4) > R3) > R5) > R6) > R7
CF Without Rule 7: (((R1 > R4) > R3) > R5) > R6

Rule Selection Metric Values Example Visualization

Rule 1: negation Metrics CF → Full Selected Rule Union → Selected Rule Sentence → FEU All → Invalid New Examples

\$ pip install exsum
\$ git clone https://github.com/YilunZhou/exsum-demos
\$ cd exsum-demos
\$ exsum sst_rule_union.py

Open up a browser to localhost:5000 to interact with the GUI

More information at <https://yilunzhou.github.io/exsum/>

(None) saliency lower range AutoTune -0.1 saliency upper range AutoTune 0.13

picture hosts a parka-wrapped dose of heart .
y=1 : 1.00 It is a film that will have people walking out halfway through , will encourage others to stand up and applaud , and will , undoubtedly , leave both camps engaged in a ferocious debate for years to come .

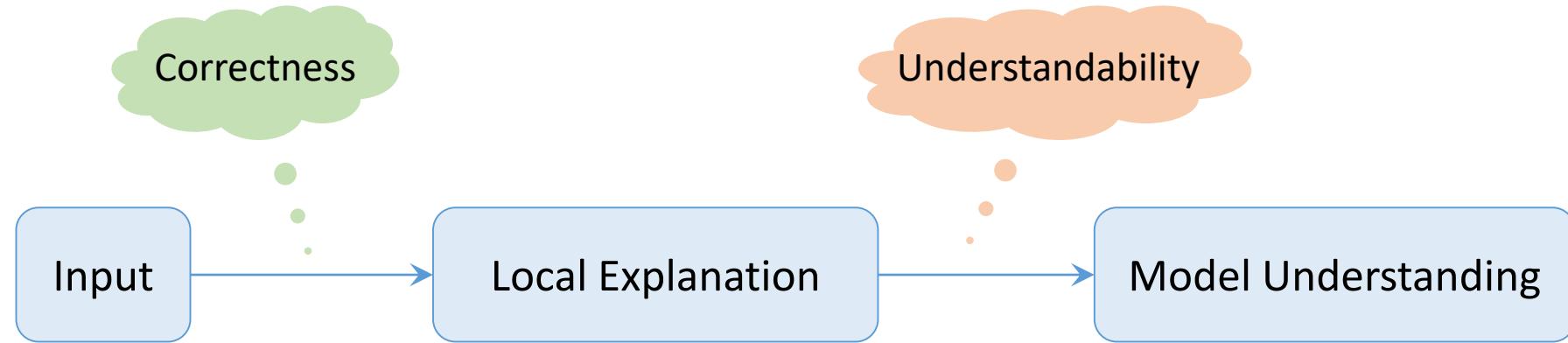
y=0 : 0.05 De Niro may enjoy the same free ride from critics afforded to Clint Eastwood in the lazy Bloodwork .

y=0 : 0.01 The film might have been more satisfying if it had , in fact , been fleshed out a little more instead of going for easy smiles .



Slides and
Resources

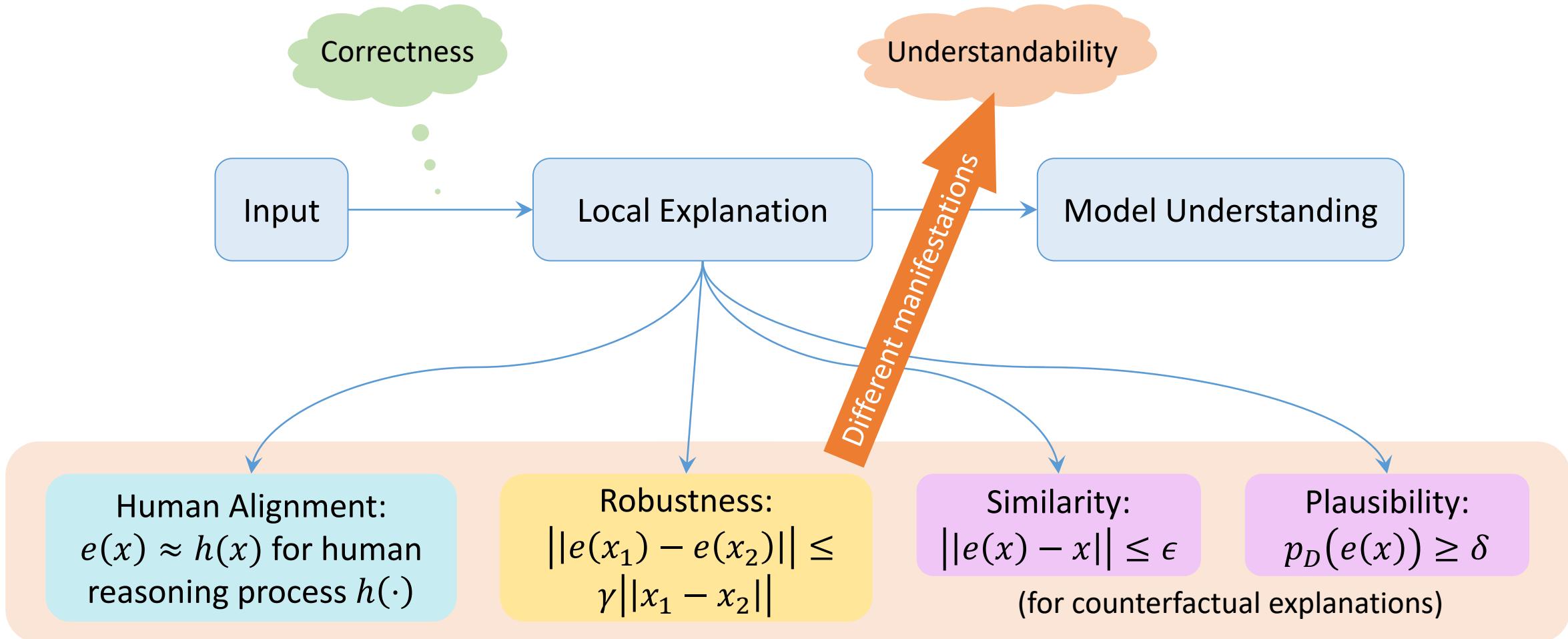
The Many Faces of Understandability





Slides and
Resources

The Many Faces of Understandability





Slides and
Resources

Definition vs. Evaluation

Definition

Gradient

$$D_g(x) = \nabla_x f(x)$$

Evaluation

Comprehensiveness

$$E_\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$



Slides and
Resources

Definition vs. Evaluation

Definition

Gradient

$$D_g(x) = \nabla_x f(x)$$

Evaluation

Comprehensiveness

$$E_\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Gradient

$$E_g(x, e) = -\|\nabla_x f(x) - e\|$$



Slides and
Resources

Definition vs. Evaluation

Definition

Gradient

$$D_g(x) = \nabla_x f(x)$$

Comprehensiveness

$$D_\kappa(x) = \operatorname{argmax}_e \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Evaluation

Comprehensiveness

$$E_\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Gradient

$$E_g(x, e) = -\|\nabla_x f(x) - e\|$$



Slides and
Resources

Definition-Evaluation Duality

Definition

Duality

Evaluation

Gradient

$$D_g(x) = \nabla_x f(x)$$

Comprehensiveness

$$E_\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Comprehensiveness

$$D_\kappa(x) = \operatorname{argmax}_e \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Gradient

$$E_g(x, e) = -\|\nabla_x f(x) - e\|$$



Slides and
Resources

Definition-Evaluation Duality

Definition

Duality

Evaluation

Gradient

$$D_g(x) = \nabla_x f(x)$$

Comprehensiveness

$$E_\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Comprehensiveness

$$D_\kappa(x) = \operatorname{argmax}_e \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Gradient

$$E_g(x, e) = -\|\nabla_x f(x) - e\|$$



Slides and
Resources

Definition-Evaluation Duality

Definition

Duality

Evaluation

Gradient

$$D_g(x) = \nabla_x f(x)$$

$$D_g \gg \text{"ease"} D_\kappa$$

Comprehensiveness

$$E_\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Comprehensiveness

$$D_\kappa(x) = \operatorname{argmax}_e \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Gradient

$$E_g(x, e) = -\|\nabla_x f(x) - e\|$$



Slides and
Resources

Definition-Evaluation Duality

Definition

Duality

Evaluation

Gradient

$$D_g(x) = \nabla_x f(x)$$

$$D_g \gg \text{"ease"} D_\kappa ?$$

Comprehensiveness

$$E_\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Comprehensiveness

$$D_\kappa(x) = \operatorname{argmax}_e \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

Gradient

$$E_g(x, e) = -\|\nabla_x f(x) - e\|$$



Slides and
Resources

Beam Search Results

$$\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

$$e^* = \operatorname{argmax}_e \kappa(x, e)$$

Algorithm 1: Beam search for finding e^* .

```
1 Input: beam size  $B$ , metric  $m$ , sentence  $x$ 
   of length  $L$ ;
2 Let  $e^{(0)}$  be an empty length- $L$  explanation;
3  $\text{beams} \leftarrow \{e^{(0)}\}$ ;
4 for  $l = 1, \dots, L$  do
5    $\text{beams} \leftarrow \bigcup_{e \in \text{beams}} \text{ext}(e, L - l + 1)$ ;
6    $\text{beams} \leftarrow \text{choose\_best}(\text{beams}, B)$ ;
7 end
8  $e^* \leftarrow \text{choose\_best}(\text{beams}, 1)$ ;
9  $e^* \leftarrow \text{shift}(e^*)$ ;
10 return  $e^*$ ;
```



Beam Search Results

$$\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

$$e^* = \operatorname{argmax}_e \kappa(x, e)$$

Algorithm 1: Beam search for finding e^* .

```
1 Input: beam size  $B$ , metric  $m$ , sentence  $x$ 
   of length  $L$ ;
2 Let  $e^{(0)}$  be an empty length- $L$  explanation;
3  $\text{beams} \leftarrow \{e^{(0)}\}$ ;
4 for  $l = 1, \dots, L$  do
5    $\text{beams} \leftarrow \bigcup_{e \in \text{beams}} \text{ext}(e, L - l + 1)$ ;
6    $\text{beams} \leftarrow \text{choose\_best}(\text{beams}, B)$ ;
7 end
8  $e^* \leftarrow \text{choose\_best}(\text{beams}, 1)$ ;
9  $e^* \leftarrow \text{shift}(e^*)$ ;
10 return  $e^*$ ;
```

A worthy tribute to a great humanitarian and
her vibrant ‘co-stars.’

So stupid, so ill-conceived, so badly drawn,
it created whole new levels of ugly.



Beam Search Results

$$\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^L f(x) - f(\bar{x}_e^{(l)})$$

$$e^* = \operatorname{argmax}_e \kappa(x, e)$$

Algorithm 1: Beam search for finding e^* .

```

1 Input: beam size  $B$ , metric  $m$ , sentence  $x$ 
   of length  $L$ ;
2 Let  $e^{(0)}$  be an empty length- $L$  explanation;
3  $\text{beams} \leftarrow \{e^{(0)}\}$ ;
4 for  $l = 1, \dots, L$  do
5    $\text{beams} \leftarrow \bigcup_{e \in \text{beams}} \text{ext}(e, L - l + 1)$ ;
6    $\text{beams} \leftarrow \text{choose\_best}(\text{beams}, B)$ ;
7 end
8  $e^* \leftarrow \text{choose\_best}(\text{beams}, 1)$ ;
9  $e^* \leftarrow \text{shift}(e^*)$ ;
10 return  $e^*$ ;

```

A worthy tribute to a great humanitarian and
her vibrant ‘co-stars.’

So stupid, so ill-conceived, so badly drawn,
it created whole new levels of ugly.

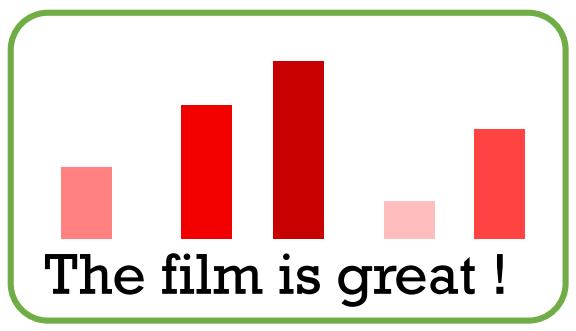
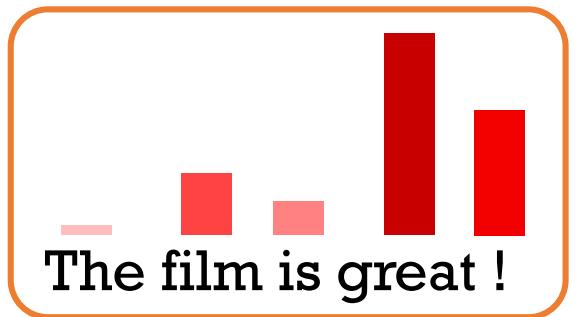
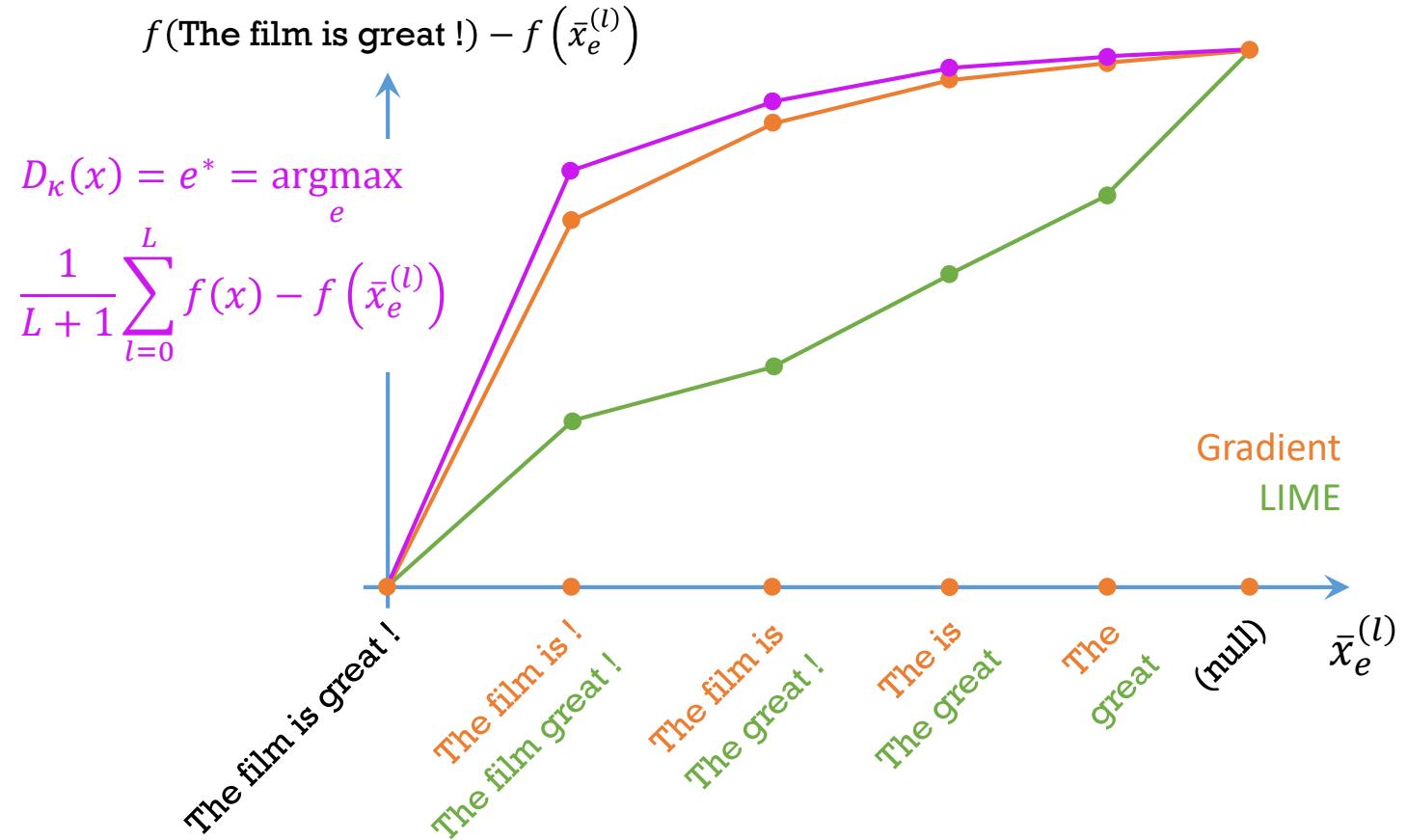
Explainer	Comp $\kappa \uparrow$
Grad	0.327
IntG	0.525
LIME	0.682
SHAP	0.612
Occl	0.509
E*	0.740
Random	0.218



Slides and Resources

Proxy Metrics for Explanation Quality

- Feature importance \Leftrightarrow model prediction change with feature removal





Slides and
Resources

Evaluation on Other Aspects

- E^* is competitive on other metrics

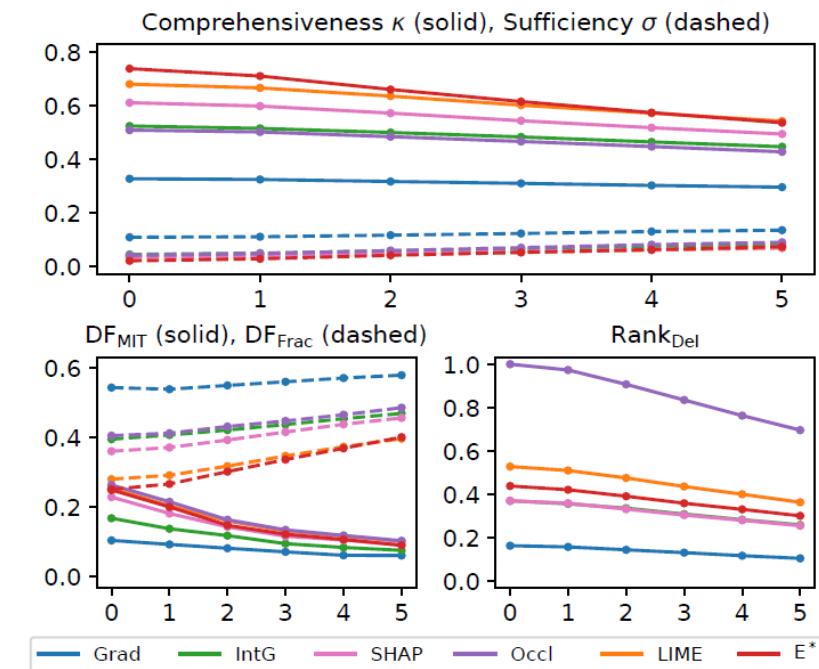
Explainer	$DF_{MIT} \uparrow$	$DF_{Frac} \downarrow$	$Rank_{Del} \uparrow$
Grad	10.5%	54.5%	0.162
IntG	16.9%	39.6%	0.369
LIME	25.5%	28.1%	0.527
SHAP	23.0%	36.1%	0.369
Occl	26.4%	40.6%	1.000
<hr/>			
E^*	25.0%	25.2%	0.438
Random	3.4%	72.3%	0.004



Evaluation on Other Aspects

- E^* is competitive on other metrics
- E^* is robust to perturbations

Explainer	$DF_{MIT} \uparrow$	$DF_{Frac} \downarrow$	$Rank_{Del} \uparrow$
Grad	10.5%	54.5%	0.162
IntG	16.9%	39.6%	0.369
LIME	25.5%	28.1%	0.527
SHAP	23.0%	36.1%	0.369
Occl	26.4%	40.6%	1.000
E^*	25.0%	25.2%	0.438
Random	3.4%	72.3%	0.004





Slides and
Resources

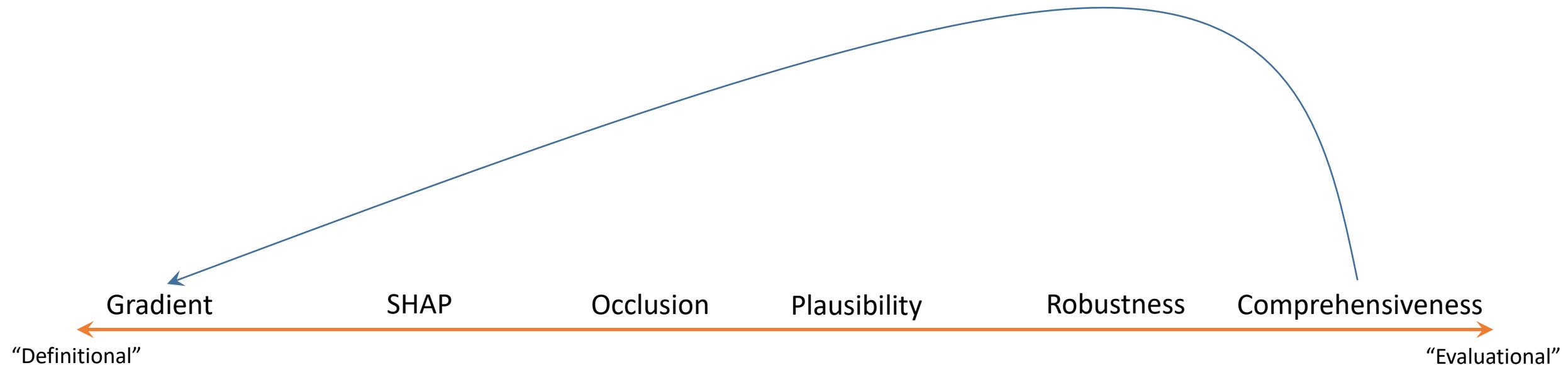
A New Paradigm of Developing Explainers?





Slides and
Resources

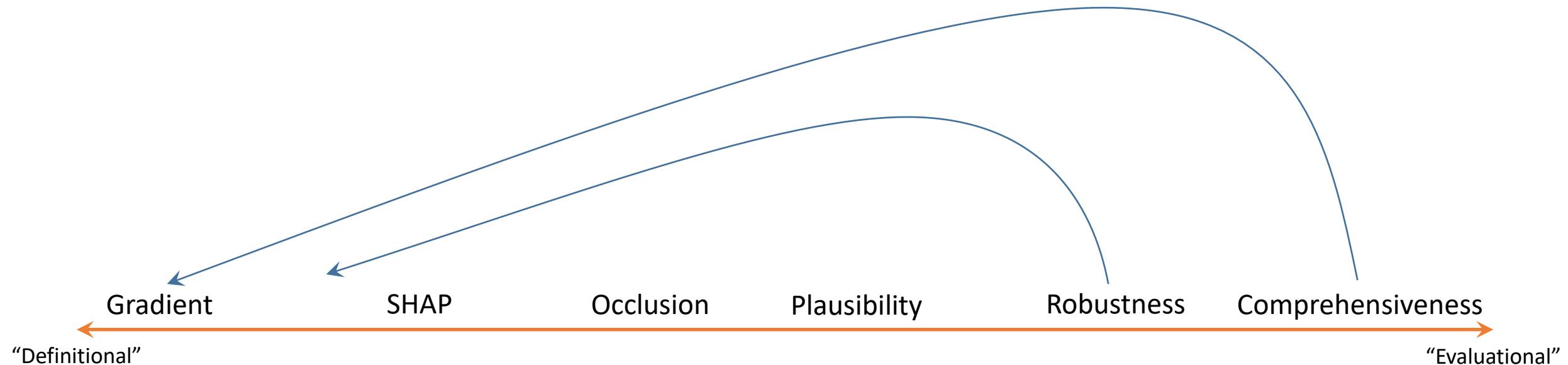
A New Paradigm of Developing Explainers?





Slides and
Resources

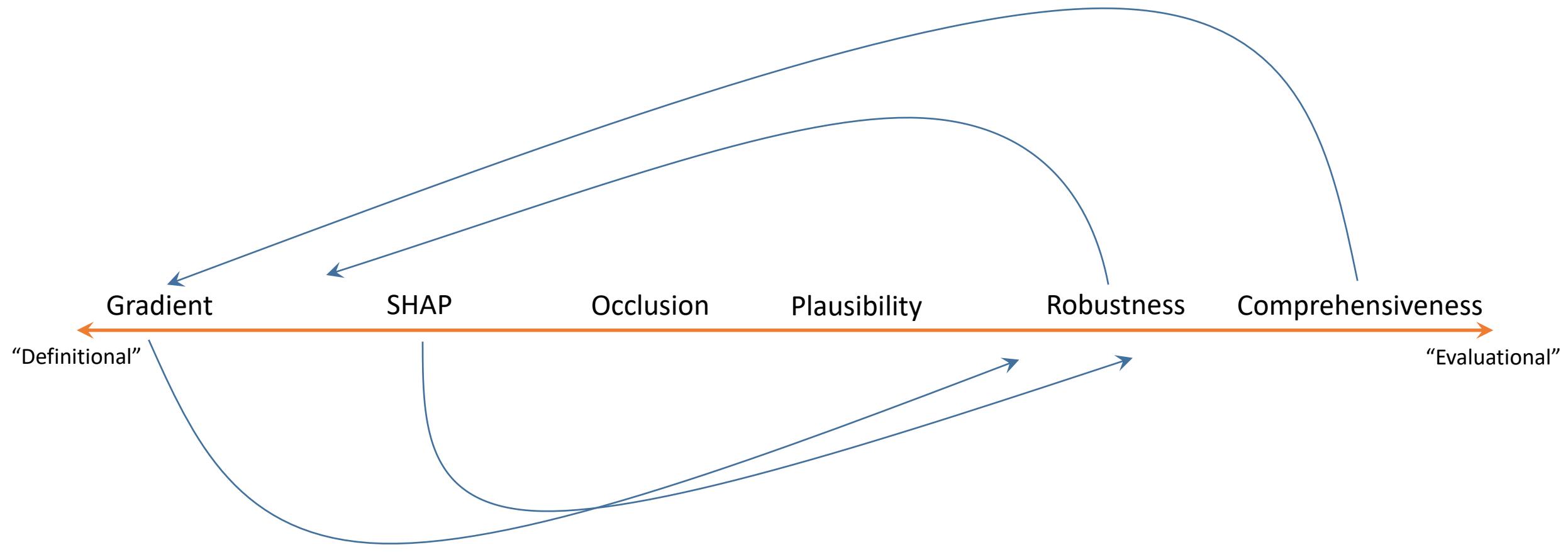
A New Paradigm of Developing Explainers?





Slides and
Resources

A New Paradigm of Developing Explainers?





Slides and
Resources

Solvability-Based Explainer

`pip install solvex`

<https://yilunzhou.github.io/solvability/>

Summary

NLP - Word

NLP - Sentence

CV - Grid Superpixel

CV - Custom Superpixel

Tabular

Read this tutorial as a Jupyter notebook [here](#).

In this demo, we compute word-level explanations for the Huggingface [textattack/roberta-base-SST-2](#) model, also the setup presented in the paper.

We first load required packages and the RoBERTa model. Two classes are needed to compute the explanations. `BeamSearchExplainer` implements the beam search algorithm, and `*Masker` implements the feature masking. In this demo, we use `TextWordMasker` since we need to mask out individual words from a text input. The other demos showcase other `*Masker`s.

```
from solvex import BeamSearchExplainer, TextWordMasker
import torch
from transformers import AutoTokenizer, AutoModelForSequenceClassification
```



Slides and
Resources

Solvability-Based Explainer

`pip install solvex`

Explained label: 1

Function value for label 1: 1.000

Contrary to other reviews, I have zero complaints about the service or the prices. I have been getting tire service here for the past 5 years now, and compared to my experience with places like Pep Boys, these guys are experienced and know what they're doing. Also, this is one place that I do not feel like I am being taken advantage of, just because of my gender. Other auto mechanics have been notorious for capitalizing on my ignorance of cars, and have sucked my bank account dry. But here, my service and road coverage has all been well explained - and let up to me to decide. And they just renovated the waiting room. It looks a lot better than it did in previous years.



Slides and
Resources

Solvability-Based Explainer

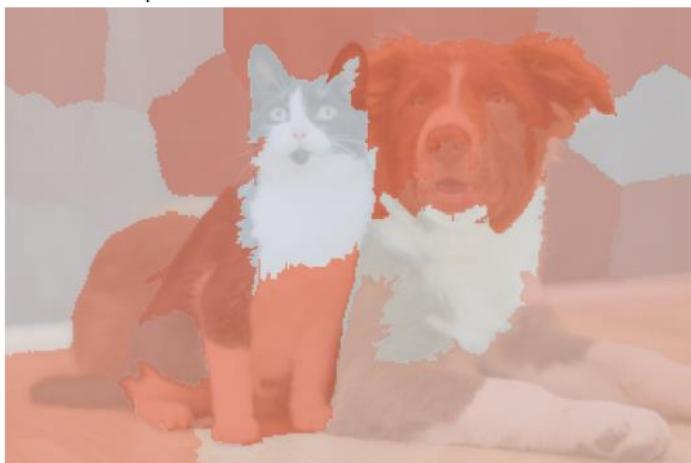
`pip install solvex`

Explained label: 1

Function value for label 1: 1.000

Contrary to other reviews, I have zero complaints about the service or the prices. I have been getting tire service here for the past 5 years now, and compared to my experience with places like Pep Boys, these guys are experienced and know what they're doing. Also, this is one place that I do not feel like I am being taken advantage of, just because of my gender. Other auto mechanics have been notorious for capitalizing on my ignorance of cars, and have sucked my bank account dry. But here, my service and road coverage has all been well explained - and let up to me to decide. And they just renovated the waiting room. It looks a lot better than it did in previous years.

Explained label: 232. Function value: 0.159





Slides and
Resources

Solvability-Based Explainer

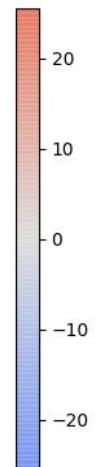
`pip install solvex`

Explained label: 1

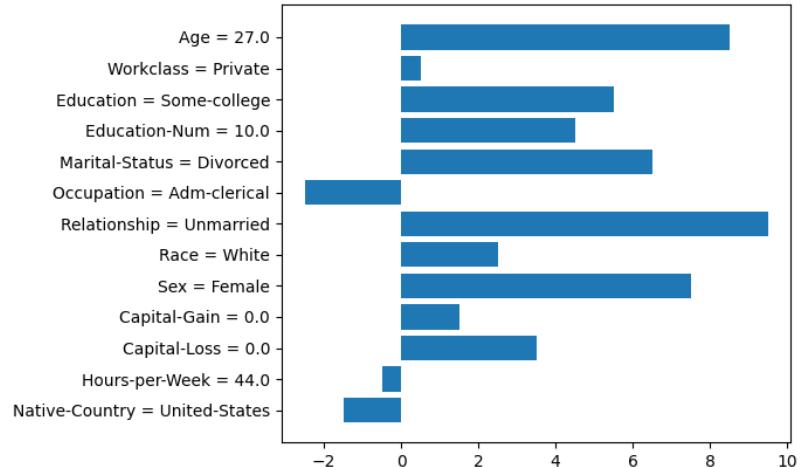
Function value for label 1: 1.000

Contrary to other reviews, I have zero complaints about the service or the prices. I have been getting tire service here for the past 5 years now, and compared to my experience with places like Pep Boys, these guys are experienced and know what they're doing. Also, this is one place that I do not feel like I am being taken advantage of, just because of my gender. Other auto mechanics have been notorious for capitalizing on my ignorance of cars, and have sucked my bank account dry. But here, my service and road coverage has all been well explained - and let up to me to decide. And they just renovated the waiting room. It looks a lot better than it did in previous years.

Explained label: 232. Function value: 0.159



Explained label: 0. Function value: 0.980





Slides and
Resources

References

- Definitions
 - Simonyan et al. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. arXiv 2014
 - Zeiler and Fergus. Visualizing and Understanding Convolutional Networks. ECCV 2014
 - Ribeiro et al. “Why Should I Trust You?” Explaining the Predictions of Any Classifier. KDD 2016
 - Smilkov et al. SmoothGrad: Removing Noise by Adding Noise. arXiv 2017
 - Sundararajan et al. Axiomatic Attribution for Deep Networks. ICML 2017
 - Lundberg et al. A Unified Approach to Interpreting Model Predictions. NIPS 2017
- Evaluations
 - Samek et al. Evaluating the Visualization of What a Deep Neural Network Has Learned. T-NNLS 2016
 - Petsiuk et al. RISE: Randomized Input Sampling for Explanation of Black-box Models. BMVC 2018
 - Ghorbani et al. Interpretation of Neural Networks is Fragile. AAAI 2019
 - DeYoung et al. ERASER: A Benchmark to Evaluate Rationalized NLP Models. ACL 2020
 - Ross et al. Explaining NLP Models via Minimal Contrastive Editing (MiCE). ACL 2021 (Findings)
 - Zhou et al. Do Feature Attribution Methods Correctly Attribute Features? AAAI 2022
 - Zhou et al. ExSum: From Local Explanations to Model Understanding. NAACL 2022
- Duality
 - Zhou and Shah. The Solvability of Interpretability Evaluation Metrics. EACL 2023 (Findings)