



An Introduction to Explainable ML

BOSTON UNIVERSITY
MACHINE INTELLIGENCE
COMMUNITY

Devin de Hueck, Jianqi Ma
Spring, 2019

What we will cover

- Motivation for Interpretability
- Overview of Interpretable ML Methods
- Counterfactual Explanations
- A method to generate saliency maps for black-box models
- Future Work

The problem of explainable ML

Why does a ML model make the prediction that it does?





Why do we want interpretable models?

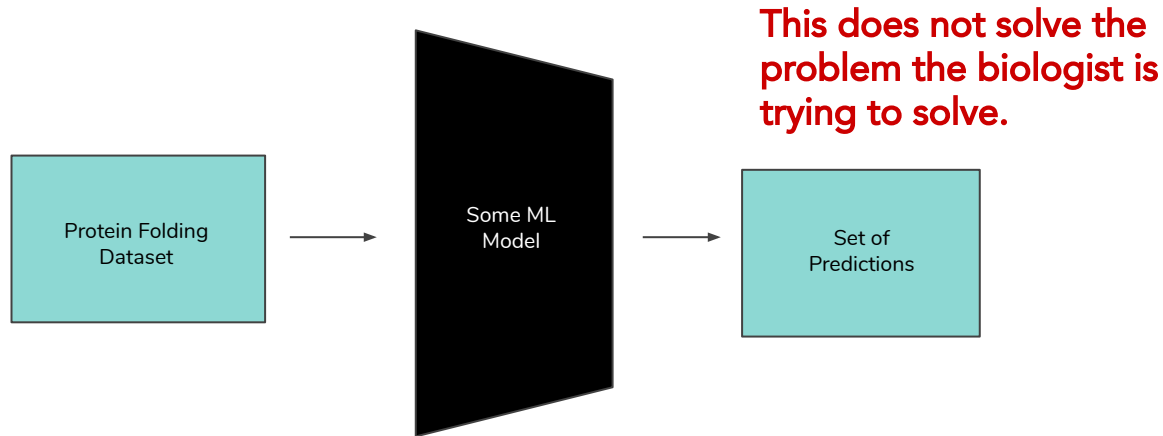
What is your problem?

Do you want to just know what was predicted or do you want to know why it was predicted?

Problems that requires interpretability

Consider the case of a biologist studying protein behavior

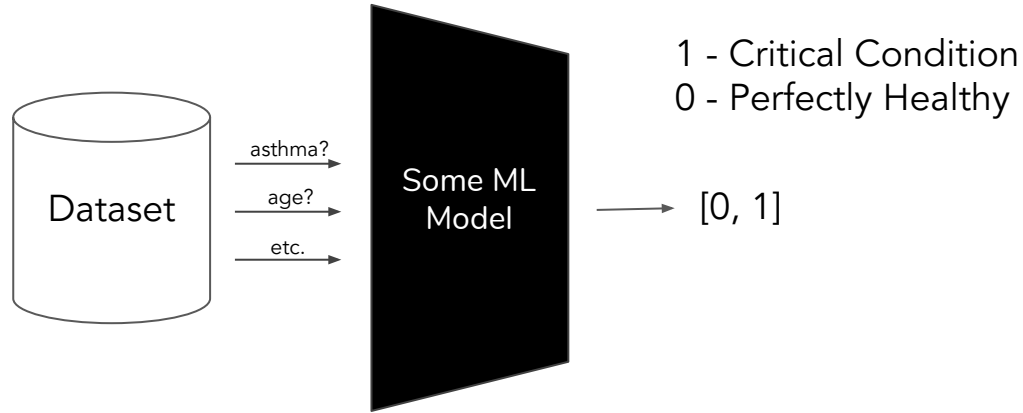
This biologist has the goal of understanding **why/how** proteins behave in the way they do



Problems that requires interpretability cont.

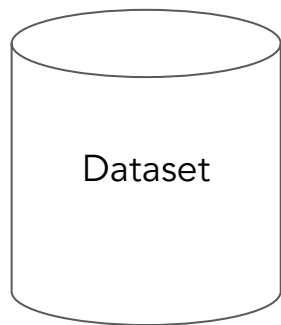
A hospital wants to predict the likelihood a patient has pneumonia.

It trains a model to predict a score between 0 and 1



Problems that requires interpretability cont.

A hospital wants to predict the likelihood a patient has pneumonia.



Consider the case that the dataset showed that patients with asthma actually **had a better prognosis than those without.**

Imagine if the hospital was prioritizing treatment with this model. Yikes!



Overview of Explainable ML Methods

3 Categories

Surrogate Models

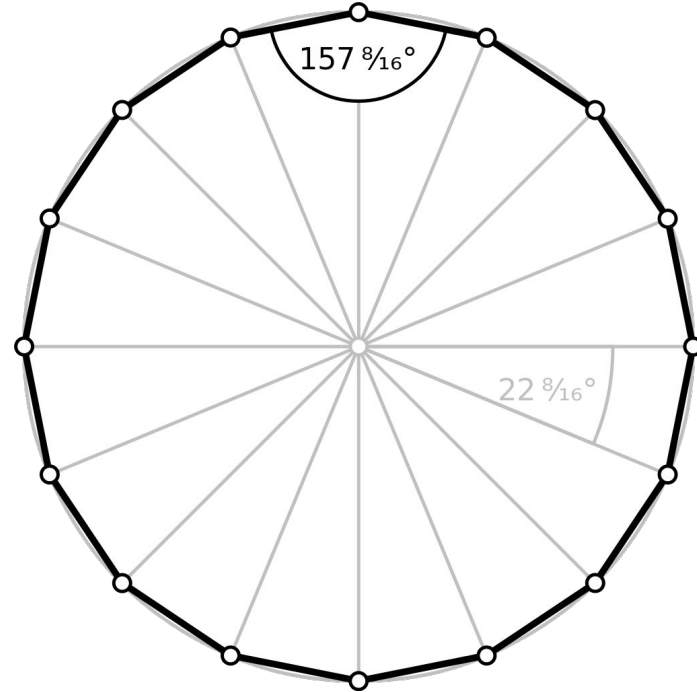
Feature Importance

Feature Effects

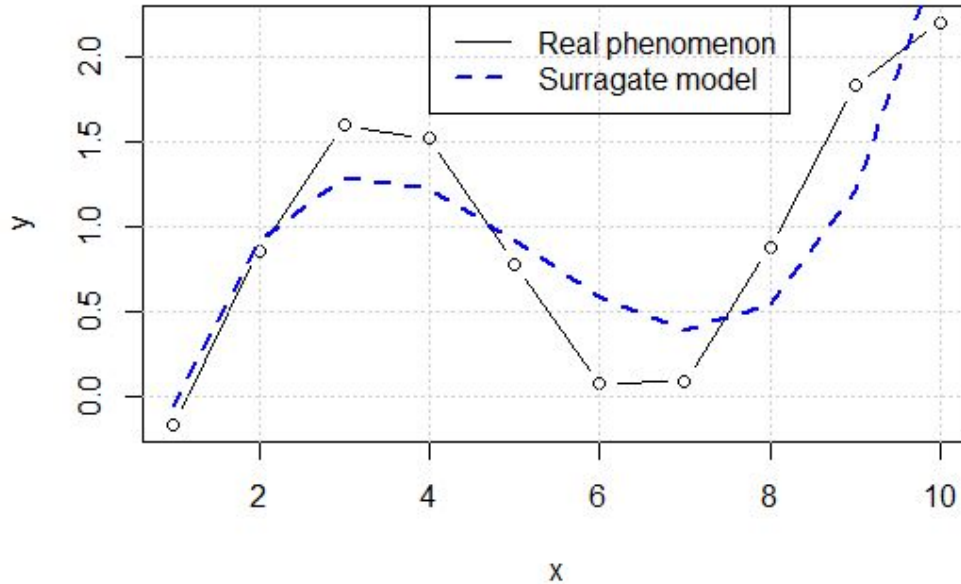
Surrogate Models

Why using Surrogate Models:

If a model is expensive, time-consuming or otherwise difficult to measure, a **cheap** and **fast** surrogate model of the outcome can be used instead.



Surrogate Models



The purpose of surrogate models is to approximate the predictions of the ML model as **accurately** as possible and to be **interpretable** at the same time.

Feature Importance

Feature Importance = Increase in Prediction Error

- A feature is “important” if shuffling its values increases the model error, because in this case the model relied on the feature for the prediction.
- A feature is “unimportant” if shuffling its values leaves the model error unchanged, because in this case the model ignored the feature for the prediction.

Feature Effects

Main Idea: summarize the average **effect** a feature has on the **prediction**

- **Partial Dependence Plot (PDP):**

The method considers all instances and gives a statement about the global relationship of a feature with the predicted outcome.

- **Individual Conditional Expectation (ICE):**

ICE plot visualizes the dependence of the prediction on a feature for each instance separately, resulting in one line per instance.

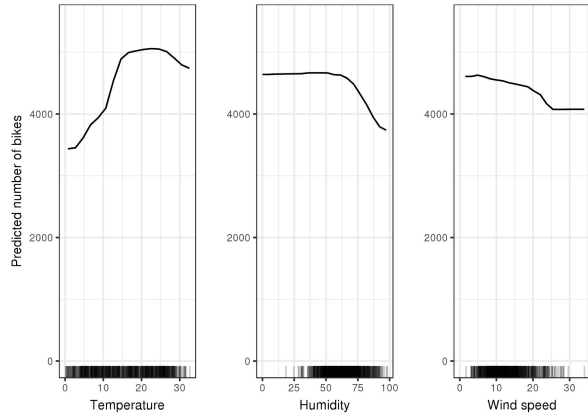
- **Accumulated Local Effects (ALE):**

Accumulated local effects describe how features influence the prediction of a machine learning model on average.

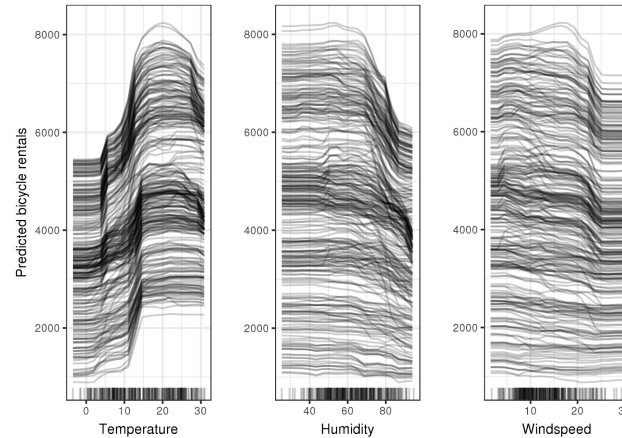


Visualizing Feature Effects Methods

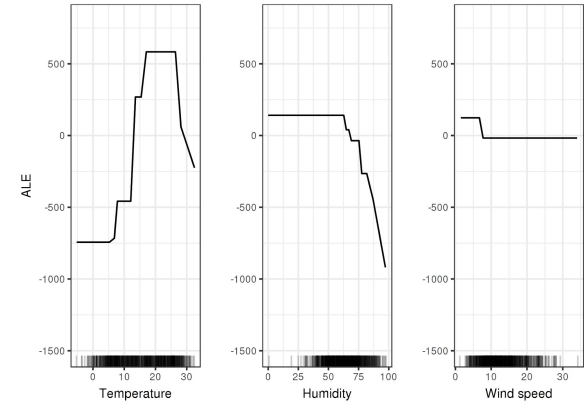
Partial Dependence Plot(PDP):



Individual Conditional Expectation(ICE):



Accumulated Local Effects(ALE):



Properties of Explainable Models

- Global vs. Local

Is the interpretation for the entire model? Or maybe just for an individual prediction?

- Model-Specific vs. Model-Agnostic

Whitebox vs. Blackbox

- Type of Data
 - Tabular
 - Text
 - Image



Counterfactual Explanations

What is a Counterfactual Explanation?

If X had not occurred, Y would not have occurred

We contradict observed facts to “imagine” a reality in which these facts did not occur

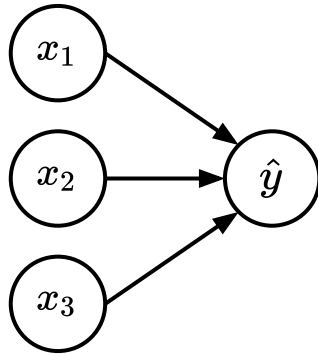
E.g. If I hadn't taken a sip of this hot coffee, I wouldn't have burned my tongue

What is a Counterfactual Explanation in the Context of ML?

Our inputs **cause** our outputs
so an input **explains** an output.

We explain one output with examples

- Example Based/Local



Some Scenarios for Counterfactual Explanations

Given a model which outputs a decision of denial or acceptance for a loan based on factors such as income, assets, age, etc.

We ask the question: What would have to change for an individual who has been rejected to be granted the loan?

Plausible Counterfactual:

“You were denied a loan because your annual income was \$30,000. If your income had been \$45,000, you would have been offered a loan.”



Some Scenarios for Counterfactual Explanations

Given a model which outputs a diabetes risk score for a woman of Pima heritage based on factors such as age, BMI, insulin level, etc.

We ask the question: What would have to change for an individual to have a risk score of 0.5?

Plausible Counterfactual:

“If the individuals 2-Hour serum insulin level was 169.5, the risk score would be 0.51”

This is already very similar to how doctors communicate risks. E.g. “A BMI greater than X indicates obesity and an increase in health risks.”



Generating Counterfactuals - Notation

We asked the question: *What would have to change for an individual to have a risk score of 0.5?*

General Case: What in the input \mathbf{x} would have to change such that the output has the desired outcome \mathbf{y}' ?

\mathbf{x}' - The counterfactual that we are generating

\mathbf{x} - An input value

\mathbf{y}' - The desired outcome

\mathbf{x}' is the smallest change in \mathbf{x} such that the prediction is the desired outcome \mathbf{y}' .

Generating Counterfactuals - The Loss Function

Our assertion that \mathbf{x}' is the smallest change in \mathbf{x} such that the prediction is the desired outcome \mathbf{y}' is characterized by the loss function:

$$L(x, x', y', \lambda)$$

Generating Counterfactuals - The Loss Function

Our assertion that \mathbf{x}' is the smallest change in \mathbf{x} such that the prediction is the desired outcome \mathbf{y}' is characterized by the loss function:

$$L(x, x', y', \lambda) = \lambda \cdot (\hat{f}(x') - y')^2$$

The hyperparameter
lambda

Distance between the
prediction from the
current counterfactual
and desired outcome

Generating Counterfactuals - The Loss Function

Our assertion that \mathbf{x}' is the smallest change in \mathbf{x} such that the prediction is the desired outcome \mathbf{y}' is characterized by the loss function:

$$L(x, x', y', \lambda) = \lambda \cdot (\hat{f}(x') - y')^2 + d(x, x')$$

The hyperparameter
lambda

Distance between the
prediction from the
current counterfactual
and desired outcome

Distance between the
counterfactual
example and the
original example

Generating Counterfactuals - The Loss Function

Our assertion that \mathbf{x}' is the smallest change in \mathbf{x} such that the prediction is the desired outcome \mathbf{y}' is characterized by the loss function:

$$L(x, x', y', \lambda) = \lambda \cdot (\hat{f}(x') - y')^2 + d(x, x')$$

The hyperparameter
lambda

Distance between the
prediction from the
current counterfactual
and desired outcome

Distance between the
counterfactual
example and the
original example

Find the values of \mathbf{x}' and λ such that: $\arg \min_{x'} L(x, x', y', \lambda)$



Generating Counterfactuals - The Algorithm

1. Initialize:

\mathcal{X} - You input that is to be explained

y' - The desired outcome

ϵ - The tolerance (an acceptable distance for $f(x')$ to be from y')

λ - Hyperparameter to balance the two parts of the loss function (usually set to a low value, > 1)

x' - An random input instance



Generating Counterfactuals - The Algorithm

2. Minimize the loss function (find an x' really close to x):

$$L(x, x', y', \lambda) = \lambda \cdot (\hat{f}(x') - y')^2 + d(x, x')$$

Generating Counterfactuals - The Algorithm

2. Minimize the loss function (find an x' really close to x):

$$L(x, x', y', \lambda) = \lambda \cdot (\hat{f}(x') - y')^2 + d(x, x')$$

Until satisfied:

Use an optimization algorithm like gradient descent to update x' to minimize the loss:

$$x' = x' - \alpha \nabla_{x'} L(x, x', y', \lambda)$$

Generating Counterfactuals - The Algorithm

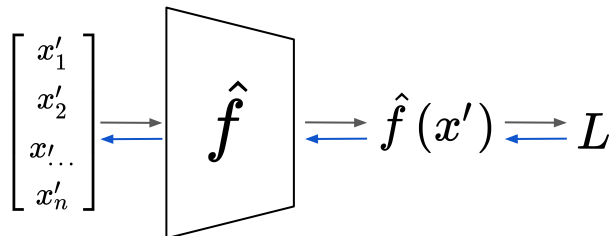
2. Minimize the loss function (find an x' really close to x):

$$L(x, x', y', \lambda) = \lambda \cdot (\hat{f}(x') - y')^2 + d(x, x')$$

Until satisfied:

Use an optimization algorithm like gradient descent to update x' to minimize the loss:

$$x' = x' - \alpha \nabla_{x'} L(x, x', y', \lambda)$$



Generating Counterfactuals - The Algorithm

3. Maximize λ until we have our desired outcome while maintaining a x' similar to x

While $|f(x') - y'| > \epsilon$

$\lambda = \lambda + \alpha$ - Increase lambda by a small step size α



Generating Counterfactuals - The Algorithm

3. Maximize λ until we have our desired outcome while maintaining a x' similar to x

While $|f(x') - y'| > \epsilon$

$\lambda = \lambda + \alpha$ - Increase lambda by a small step size α

Use an optimization algorithm like gradient descent to update x' to minimize the loss:

$$x' = x' - \alpha \nabla_{x'} L(x, x', y', \lambda)$$

Generating Counterfactuals - The Algorithm

3. Maximize λ until we have our desired outcome while maintaining a x' similar to x

While $|f(x') - y'| > \epsilon$

$\lambda = \lambda + \alpha$ - Increase lambda by a small step size α

Use an optimization algorithm like gradient descent to update x' to minimize the loss:

$$x' = x' - \alpha \nabla_{x'} L(x, x', y', \lambda)$$

Return x' and Repeat

What is a Good Counterfactual Explanation?

\mathbf{x}' is the smallest change in \mathbf{x} such that the prediction is the desired outcome \mathbf{y}' .

All features must be of reasonable scale

- The counterfactual: "If the apartment had 300 bathrooms the rent prediction would 1500/month" is not a meaningful counterfactual

Limit features to be within the max and min values in the dataset

Advantages

- Very clear explanations
 - Humans understand these type of explanations very well
- Does not require much domain knowledge to understand
 - Don't need to be a doctor to understand a counterfactual concerning a medical model

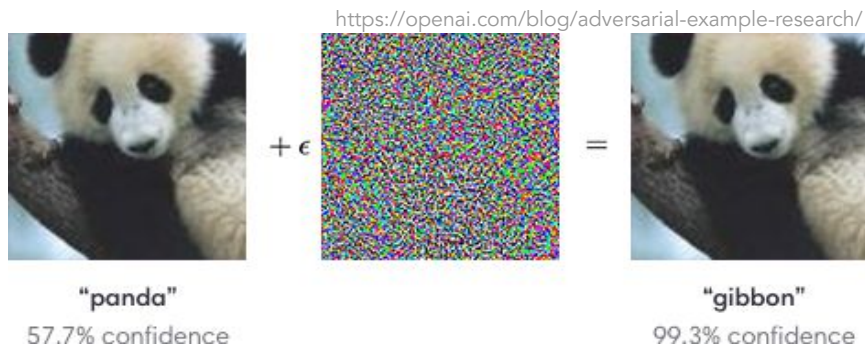
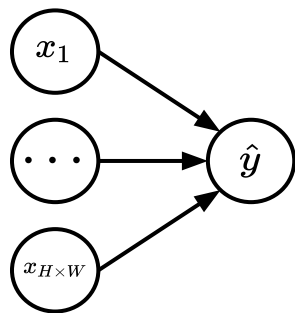
Disadvantages

- Rashmon Effect
 - There can be multiple *valid* counterfactual examples for each input
- No guarantee of finding a counterfactual within the set tolerance (ϵ)

Adversarial Examples vs Counterfactual Explanations

x' is the smallest change in x such that the prediction is the desired outcome y' .

What happens when our input is an image? Finding a counterfactual would give us an adversarial example.



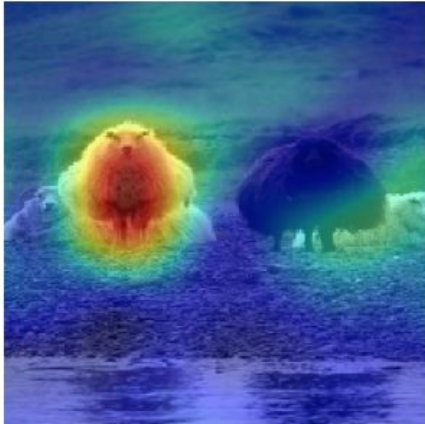


Saliency Map Generation

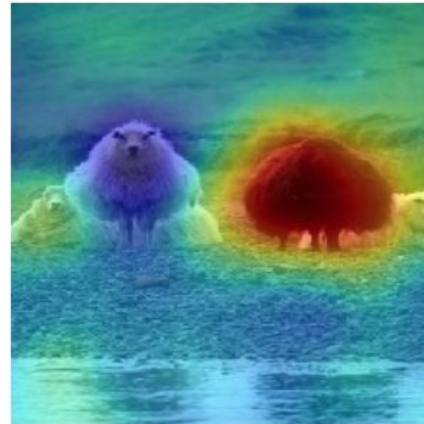
What Are Saliency Maps?

Saliency maps tell us how important each pixel of an image is for a prediction

Saliency map for prediction of class **sheep**

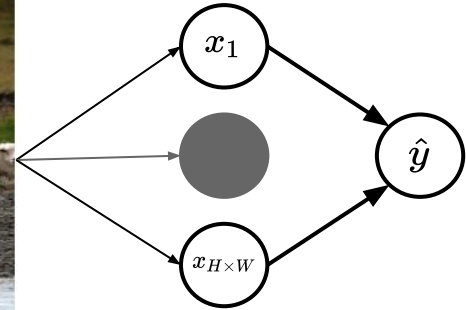
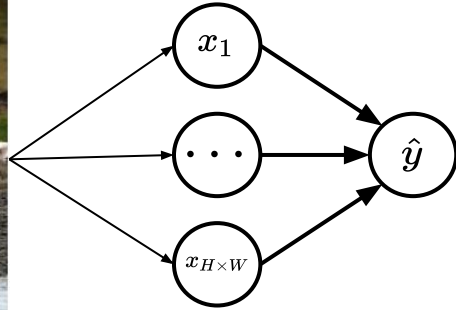


Saliency map for prediction of class **cow**



Causality and Saliency Maps

What will happen if we delete some inputs/pixels? How will the predicted class score change?

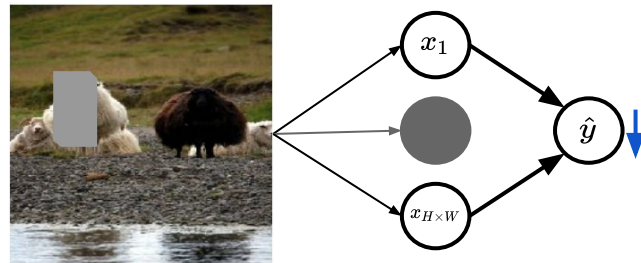
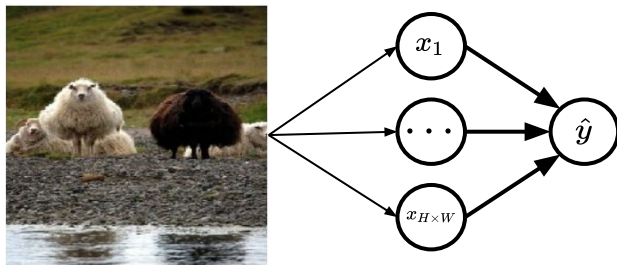


Causality and Saliency Maps

What will happen if we delete some inputs/pixels? How will the predicted class score change? **It depends!**

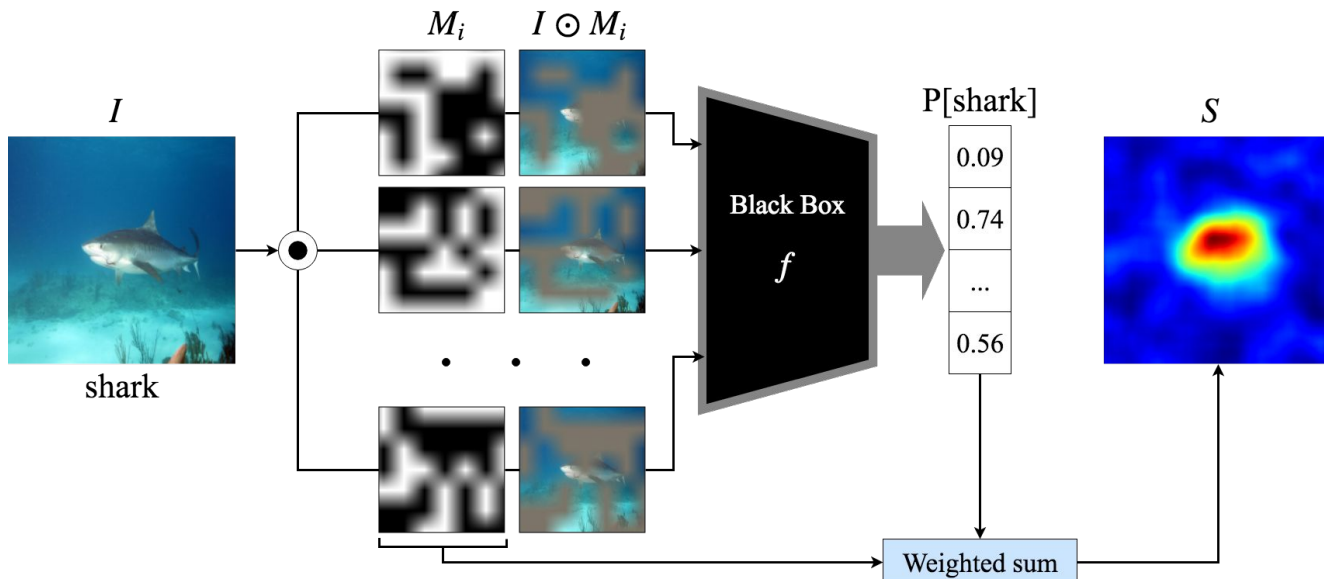
If the inputs/pixel that were removed **were important** for making the prediction then the predicted class score would **decrease** based on the importance.

If the inputs/pixel that were removed **were not important** for making the prediction then the predicted class score would **stay about the same**.



RISE: Randomized Input Sampling for Explanation of Black-box Models

If we randomly remove pixels from different parts of an image and make a prediction on these modified images then we will eventually learn the importance of each pixel.



RISE - Notation

I - The image to be explained which is of size $H \times W \times 3$

M - A binary mask which is of size $H \times W$

$I \odot M$ - A masked image which is of size $H \times W \times 3$

$f(I \odot M)$ - The prediction from a masked image



RISE - The importance of a pixel

“We define importance of pixel λ as the expected score over all possible masks \mathbf{M} conditioned on the event that pixel λ is observed, i.e., $\mathbf{M}(\lambda) = 1$ ”

$$S_{I,f}(\lambda)$$

↑
Saliency/Importance
of pixel λ



RISE - The importance of a pixel

“We define importance of pixel λ as the expected score over all possible masks M conditioned on the event that pixel λ is observed, i.e., $M(\lambda) = 1$ ”

$$S_{I,f}(\lambda) = \mathbb{E}_M[f(I \odot M)]$$

Saliency/Importance
of pixel λ

Black-box model's
prediction of a
masked image



RISE - The importance of a pixel

“We define importance of pixel λ as the expected score over all possible masks M conditioned on the event that pixel λ is observed, i.e., $M(\lambda) = 1$ ”

$$S_{I,f}(\lambda) = \mathbb{E}_M[f(I \odot M) | M(\lambda) = 1]$$

Saliency/Importance
of pixel λ

Black-box model's
prediction of a
masked image

Visibility of pixel λ - 0
if hidden by mask, 1 if
visible

RISE - The importance of a pixel

“We define importance of pixel λ as the expected score over all possible masks M conditioned on the event that pixel λ is observed, i.e., $M(\lambda) = 1$ ”

$$S_{I,f}(\lambda) = \mathbb{E}_M [f(I \odot M) | M(\lambda) = 1]$$

Saliency/Importance
of pixel λ

Black-box model's
prediction of a
masked image

Visibility of pixel λ - 0
if hidden by mask, 1 if
visible

To get the complete saliency map we solve this equation for every pixel in the image



RISE - The importance of a pixel - In practice

In practice we perform a Monte-Carlo simulation.

The importance of one pixel is the average predicted score of N masked images where that pixel is visible normalized by the expected value of a mask (0.5).

$$S_{I,f}(\lambda) \overset{\text{MC}}{\approx}$$

RISE - The importance of a pixel - In practice

In practice we perform a Monte-Carlo simulation.

The importance of one pixel is the average predicted score of N masked images where that pixel is visible normalized by the expected value of a mask (0.5).

$$S_{I,f}(\lambda) \stackrel{\text{MC}}{\approx} \sum_i f(I \odot M_i)$$

Sum of masked images
prediction scores

RISE - The importance of a pixel - In practice

In practice we perform a Monte-Carlo simulation.

The importance of one pixel is the average predicted score of N masked images where that pixel is visible normalized by the expected value of a mask (0.5).

$$S_{I,f}(\lambda) \stackrel{\text{MC}}{\approx} \sum_i f(I \odot M_i) \cdot \frac{M_i(\lambda)}{0.5}$$

Only want to sum if visible

RISE - The importance of a pixel - In practice

In practice we perform a Monte-Carlo simulation.

The importance of one pixel is the average predicted score of N masked images where that pixel is visible normalized by the expected value of a mask (0.5).

$$S_{I,f}(\lambda) \stackrel{\text{MC}}{\approx} \frac{1}{N} \sum_{i=1}^N f(I \odot M_i) \cdot M_i(\lambda)$$

Get the average
predicted score



RISE - The importance of a pixel - In practice

In practice we perform a Monte-Carlo simulation.

The importance of one pixel is the average predicted score of N masked images where that pixel is visible normalized by the expected value of a mask (0.5).

$$S_{I,f}(\lambda) \stackrel{\text{MC}}{\approx} \frac{1}{\mathbb{E}[M]} \cdot \frac{1}{N} \sum_{i=1}^N f(I \odot M_i) \cdot M_i(\lambda)$$

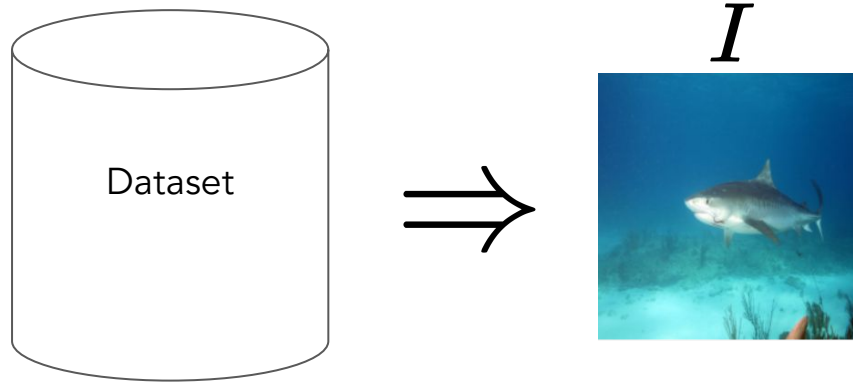
Normalize by the
expected value of a
mask (0.5)

The average predicted
score of N masked
images



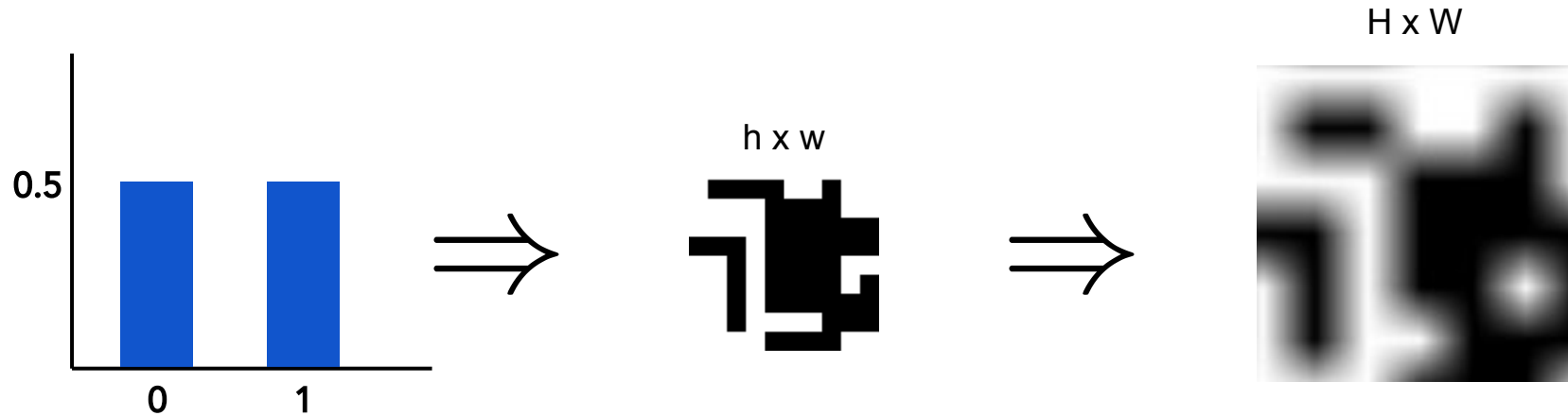
RISE - The algorithm

1. Select an image to have explained



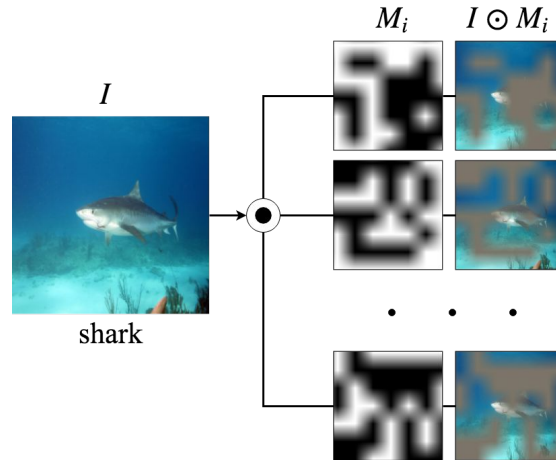
RISE - The algorithm

2. Generate a set of N masks



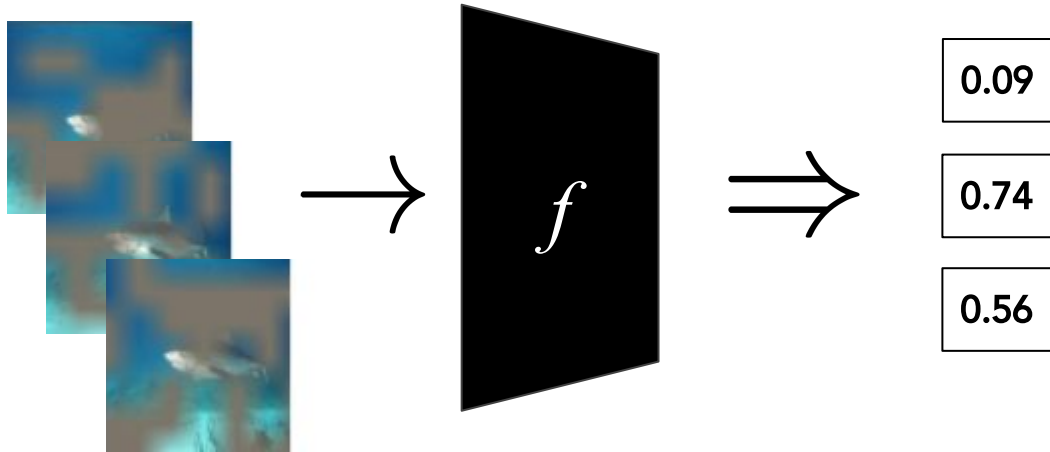
RISE - The algorithm

3. Generate N masked images through element-wise multiplication of a mask and the image.



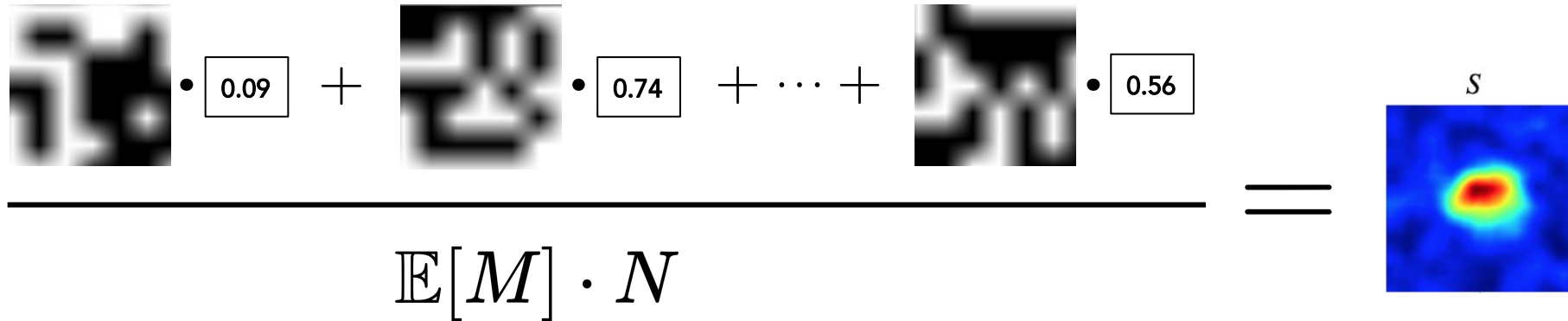
RISE - The algorithm

4. Get N predictions from your black-box model on the each masked image



RISE - The algorithm

5. Take the dot product between the N masks and N predictions then normalize to get one saliency map for each possible prediction


$$\frac{\text{Mask}_1 \cdot 0.09 + \text{Mask}_2 \cdot 0.74 + \dots + \text{Mask}_N \cdot 0.56}{\mathbb{E}[M] \cdot N} = S$$

The diagram shows a sequence of grayscale masks, each followed by a dot product with a numerical value (0.09, 0.74, ..., 0.56). These are summed and then divided by the expression $\mathbb{E}[M] \cdot N$ to produce a saliency map S , which is a heatmap with a central red/yellow region on a blue background.

And that's it!



What does the future hold?

The Future of Explainable ML

Focus will be on Model-Agnostic methods

Explainable ML add-ons to currently existing libraries

References & Further Reading

1. <https://christophm.github.io/interpretable-ml-book/counterfactual.html>
2. Petsiuk, Vitali, Abir Das, and Kate Saenko. "[RISE: Randomized Input Sampling for Explanation of Black-box Models](#)." arXiv preprint arXiv:1806.07421 (2018).
3. Doshi-Velez, Finale, and Been Kim. "[Towards a rigorous science of interpretable machine learning](#)." arXiv preprint arXiv:1702.08608 (2017).
4. Hall, Patrick, and Navdeep Gill. [An Introduction to Machine Learning Interpretability](#). O'Reilly Media Inc., 2018.



RISE - The algorithm

1. Select an image to have explained
2. Generate N masks
3. Generate N masked images through element-wise multiplication of a mask and the image.
4. Get N predictions from your black-box model on the each masked image
5. Take the dot product between the N masks and N predictions then normalize to get one saliency map for each possible prediction