

Deep Learning and Explainability

University of Victoria - PHYS-555

THIS IS YOUR MACHINE LEARNING SYSTEM?

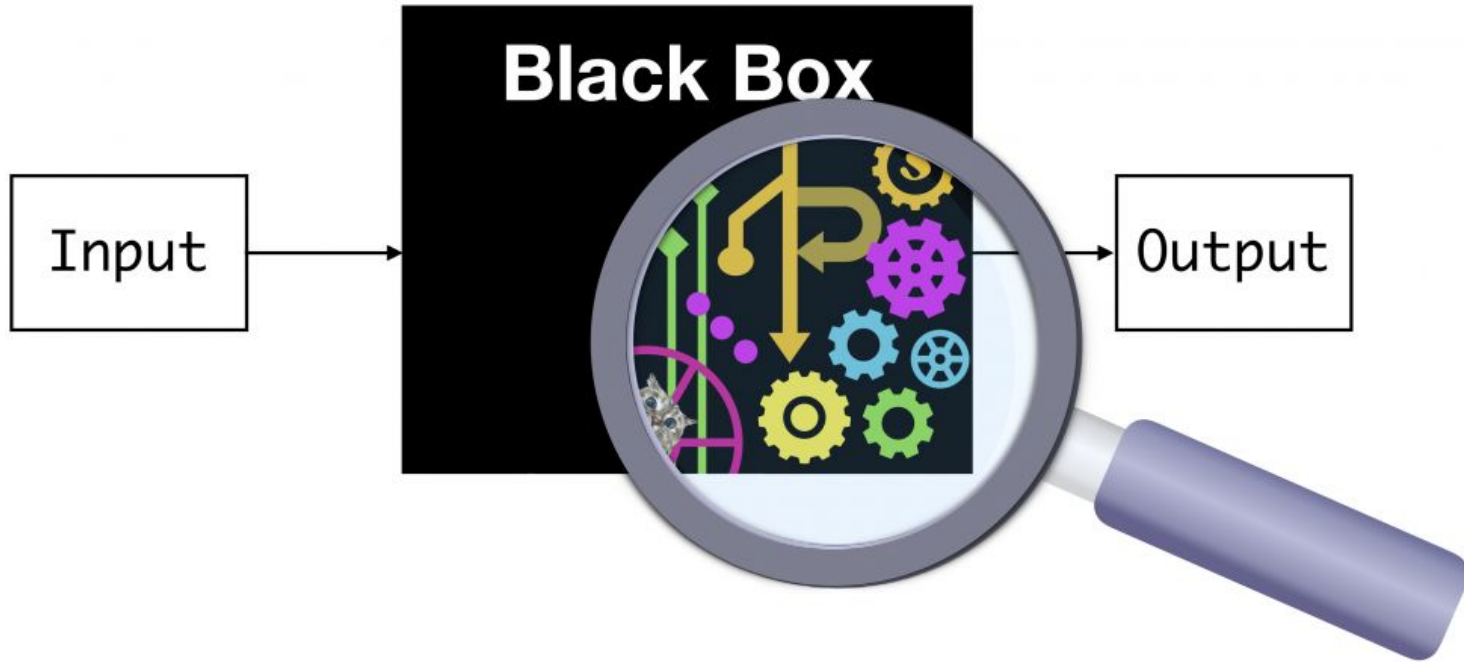
YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.



Explainable or Interpretable?



Terminology (as of today)

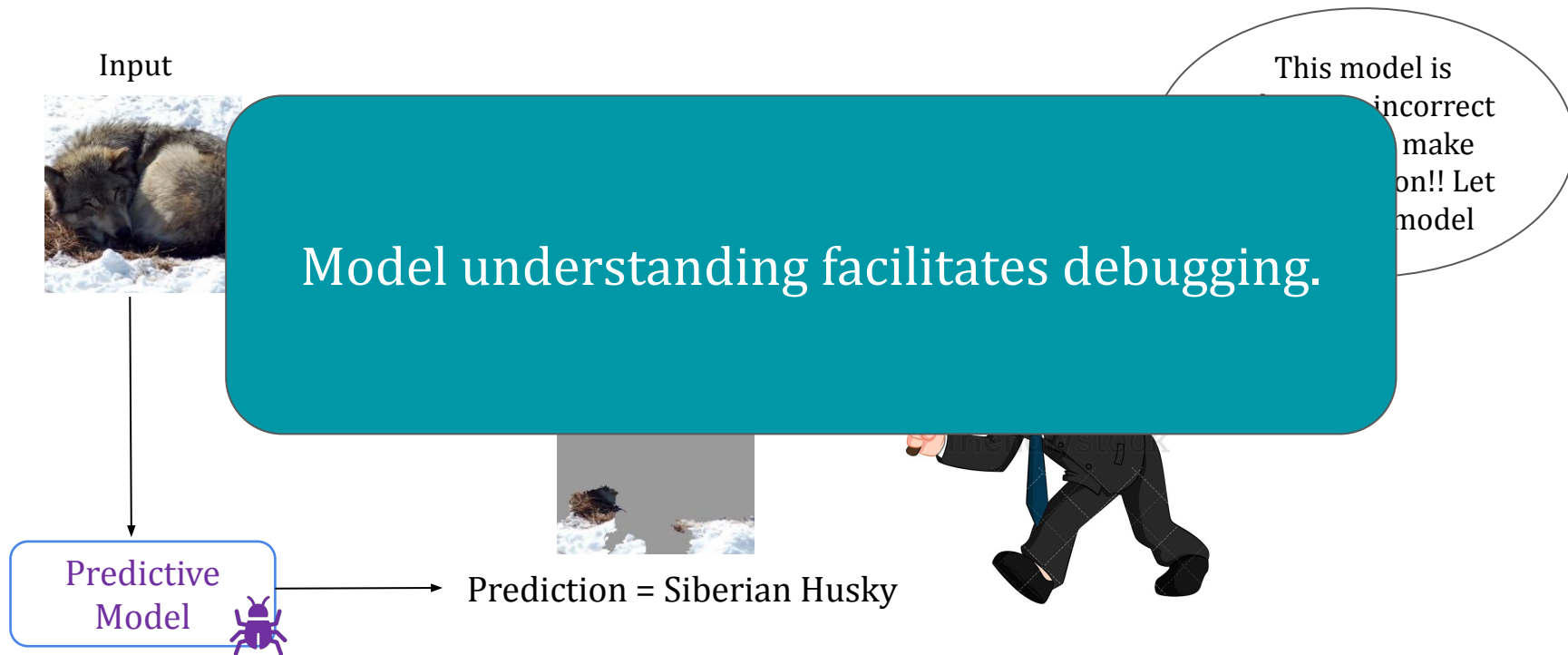
- **Interpretability**: a *passive* characteristic of a model referring to the level at which it makes sense to humans.
- **≠ Explainability**: an *active* characteristic of a model, denoting any action or procedure taken by a model with the intent of clarifying or detailing its internal functions.
- **Explainable AI (XAI)** is a set of tools and frameworks to help humans understand predictions made by AI systems.

Motivation

Model understanding is absolutely critical in several domains -- particularly those involving *high stakes decisions*



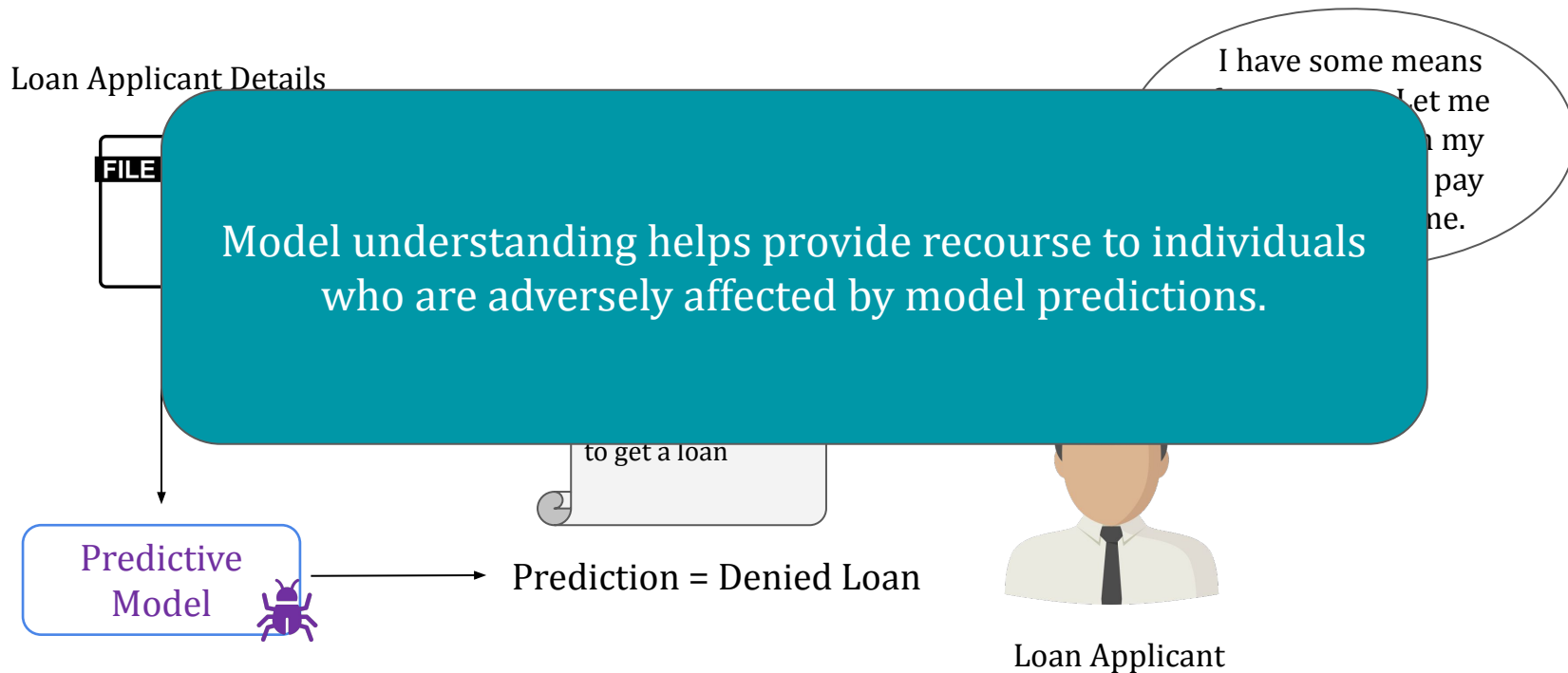
Motivation: Why Model Understanding?



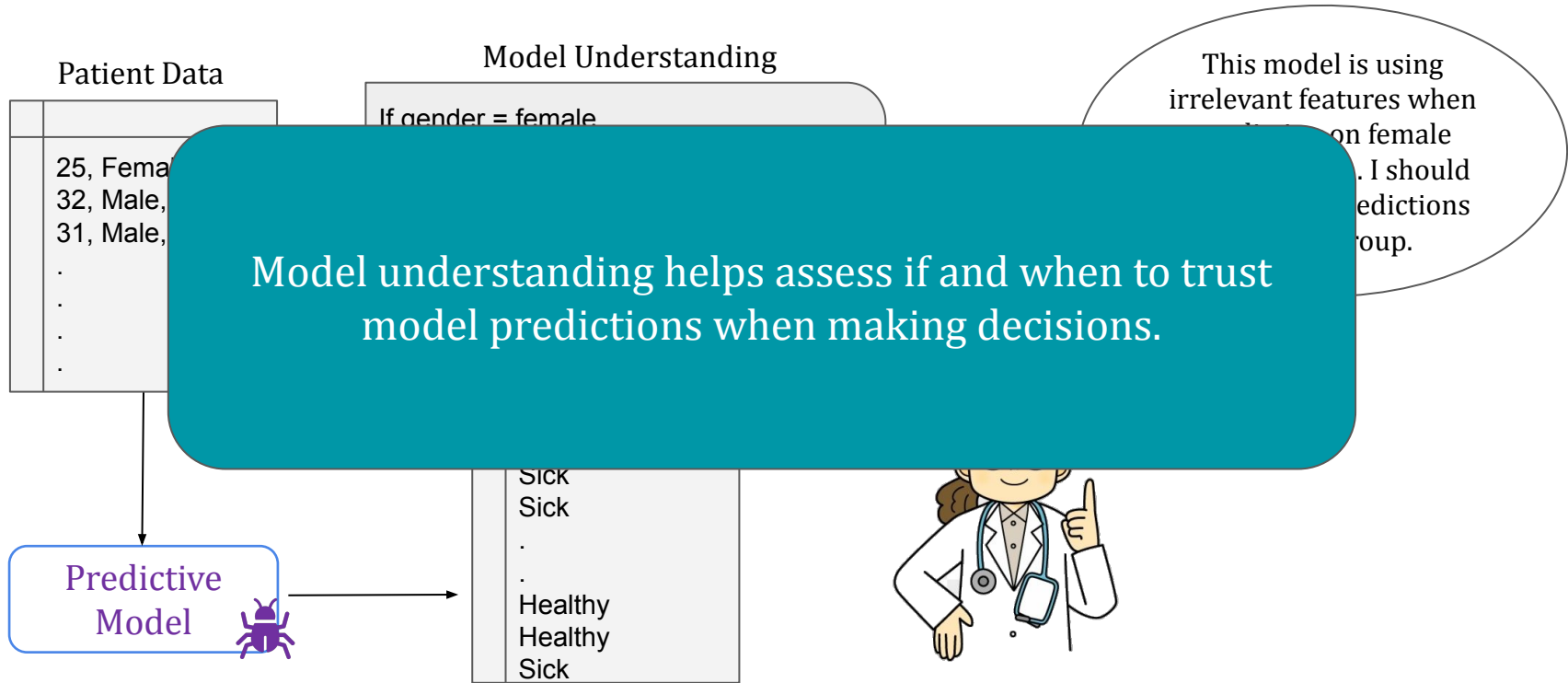
Motivation: Why Model Understanding?



Motivation: Why Model Understanding?



Motivation: Why Model Understanding?



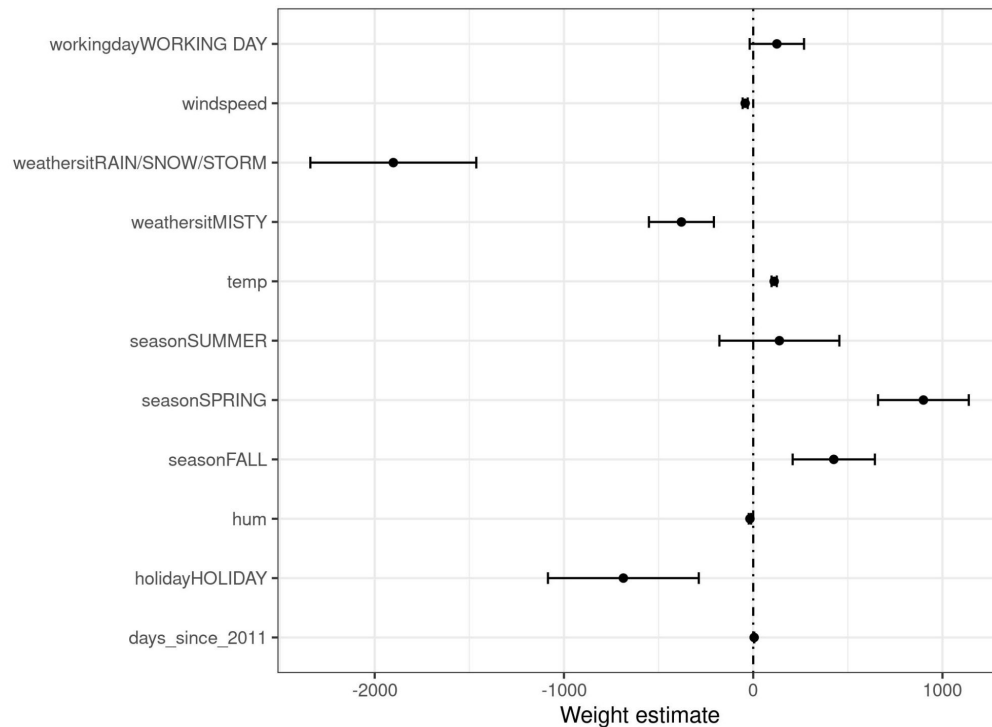
Motivation: Why Model Understanding?



Linear Models

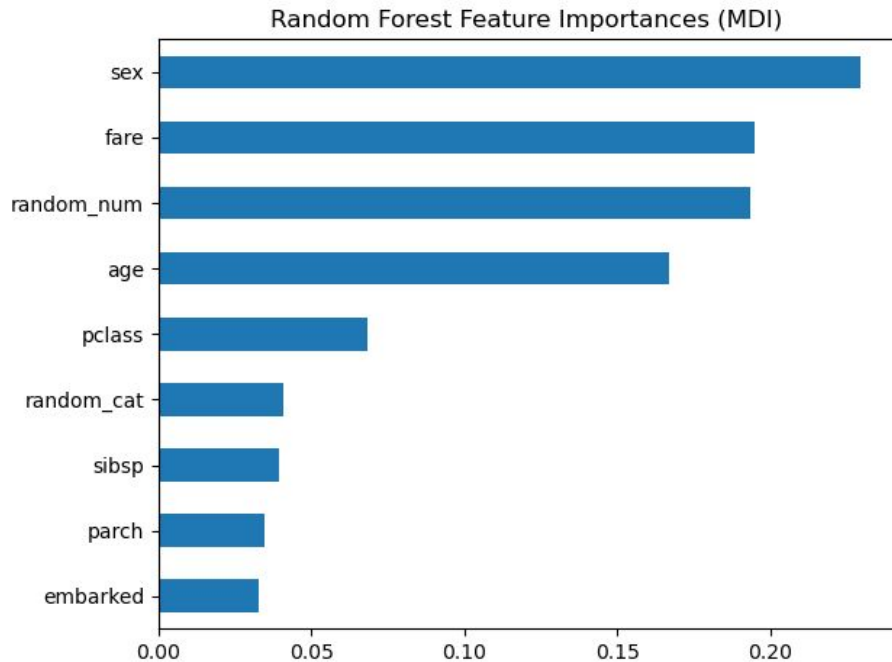
- Both Logistic regression and linear regression are easily interpretable with the weights
- Proceed with caution with the scale and correlation of features. See [scikit-learn](#) pitfall example.

Prediction of rented bikes

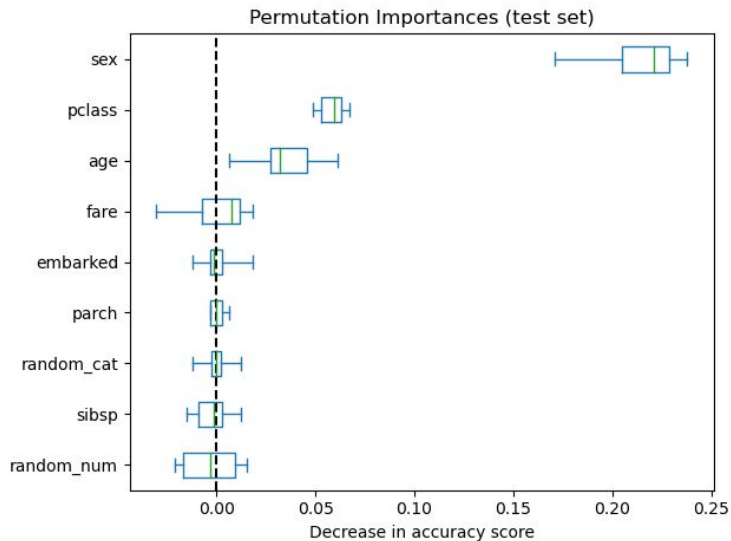


Decision Trees - Global Explanations

Features Importance with Mean Decrease Impurity

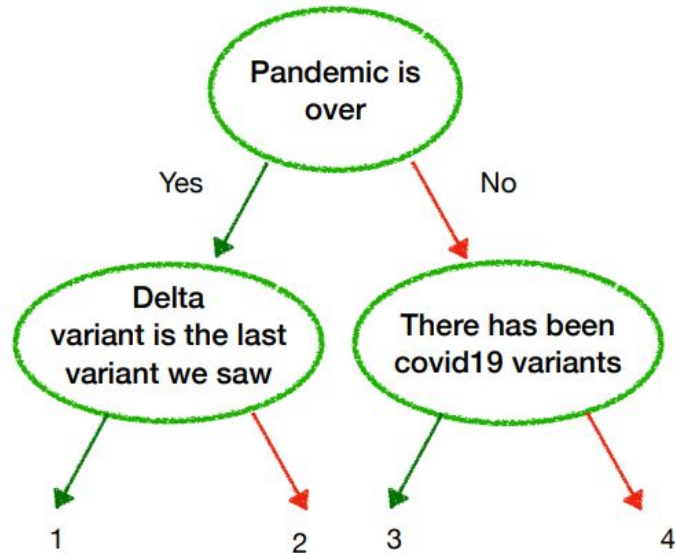


Permutation Importance

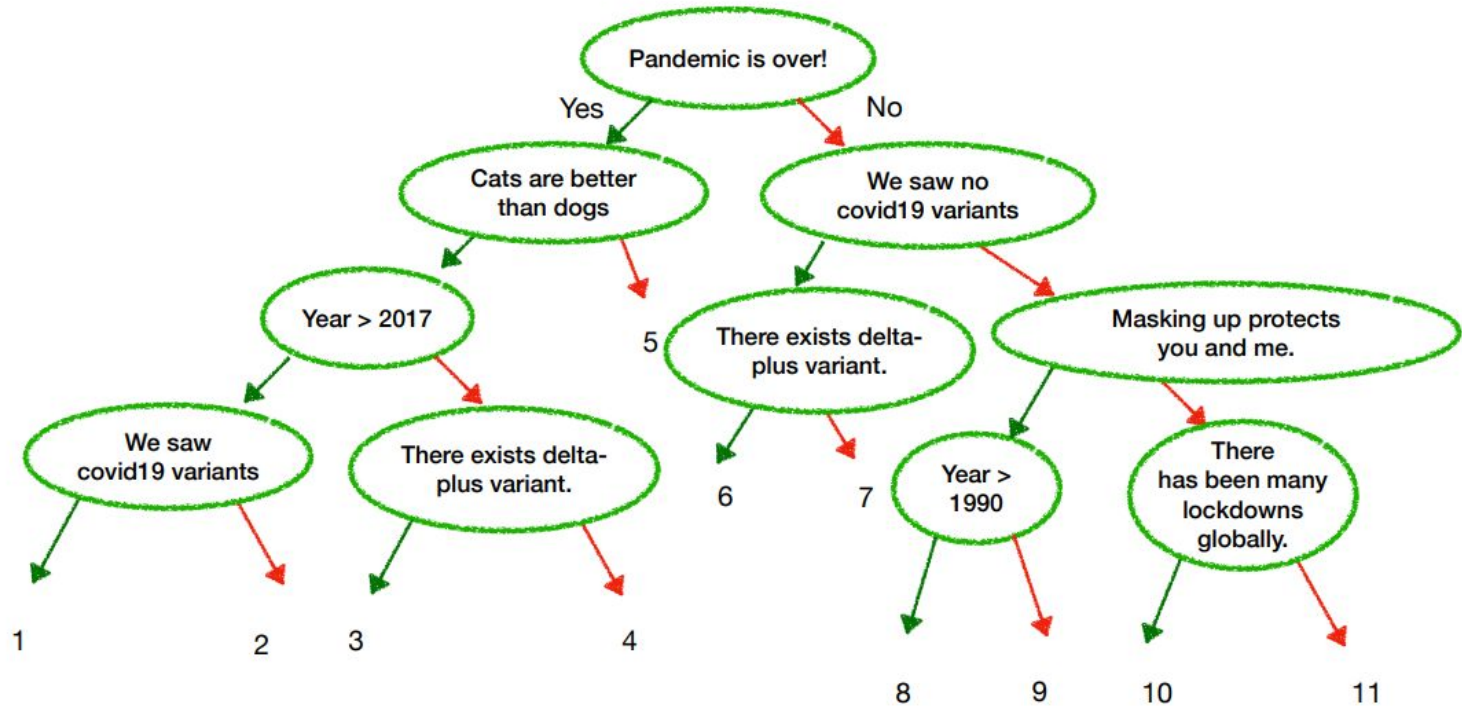


Ex: [scikit-learn RF with Titanic dataset](#)

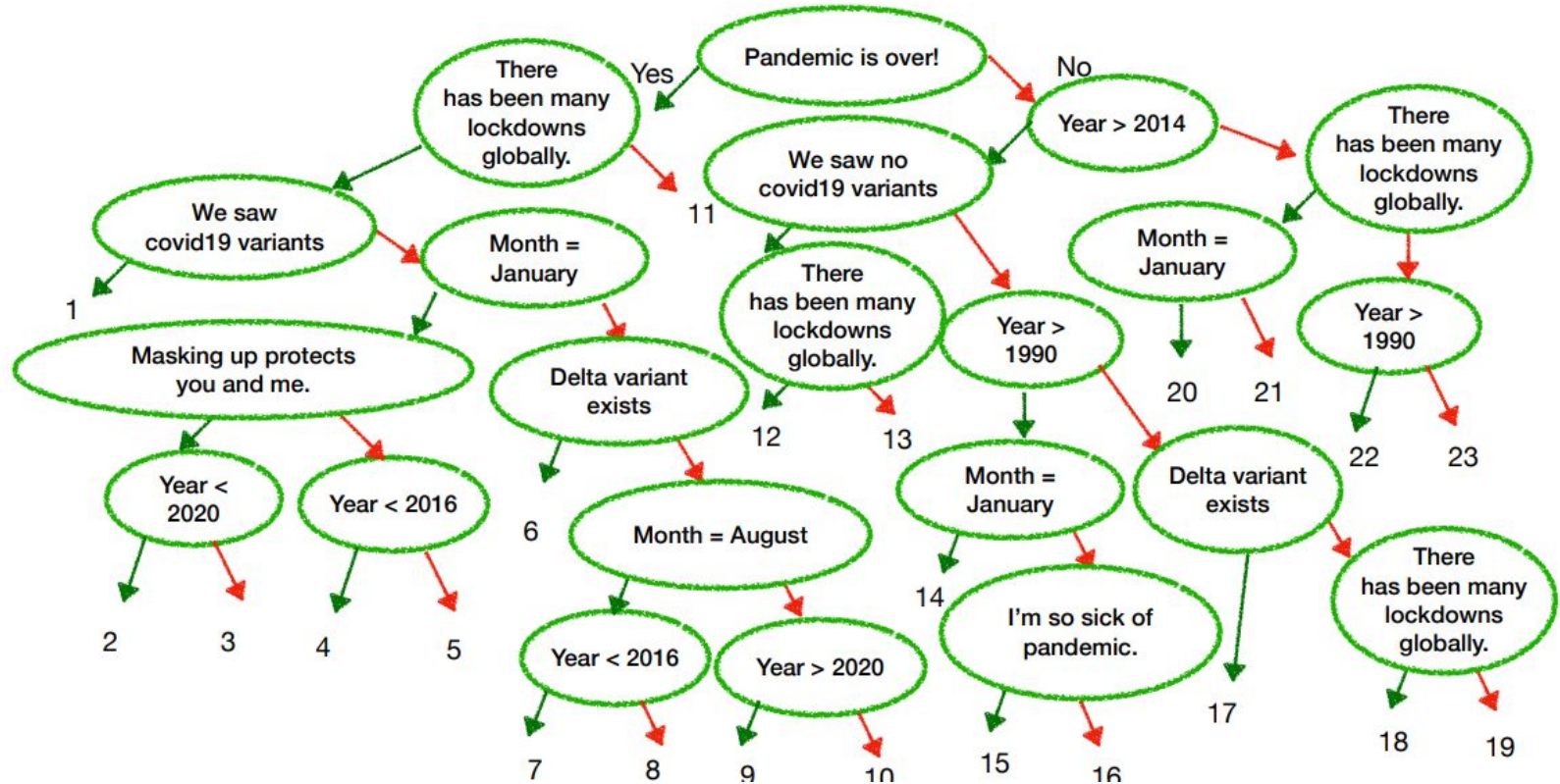
Sample Tree #1



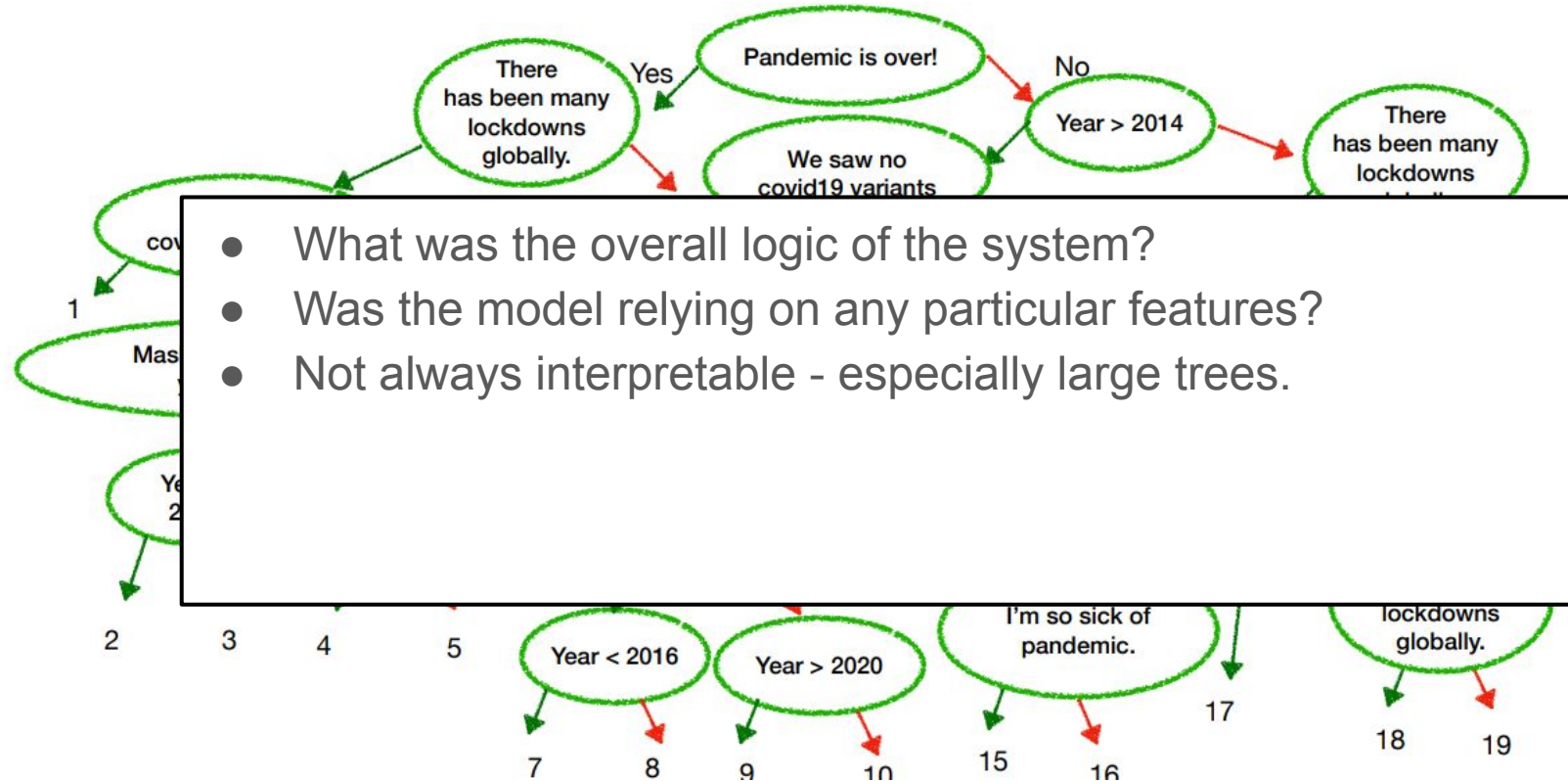
Sample Tree #2



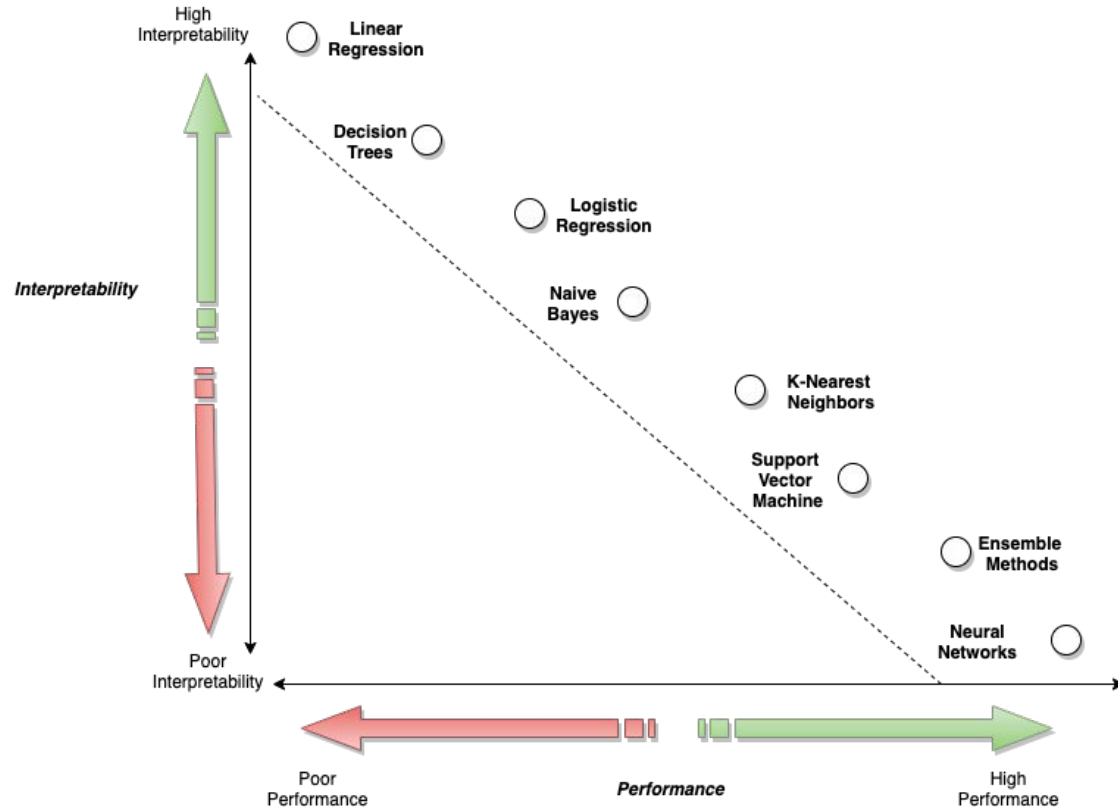
Sample Tree #3



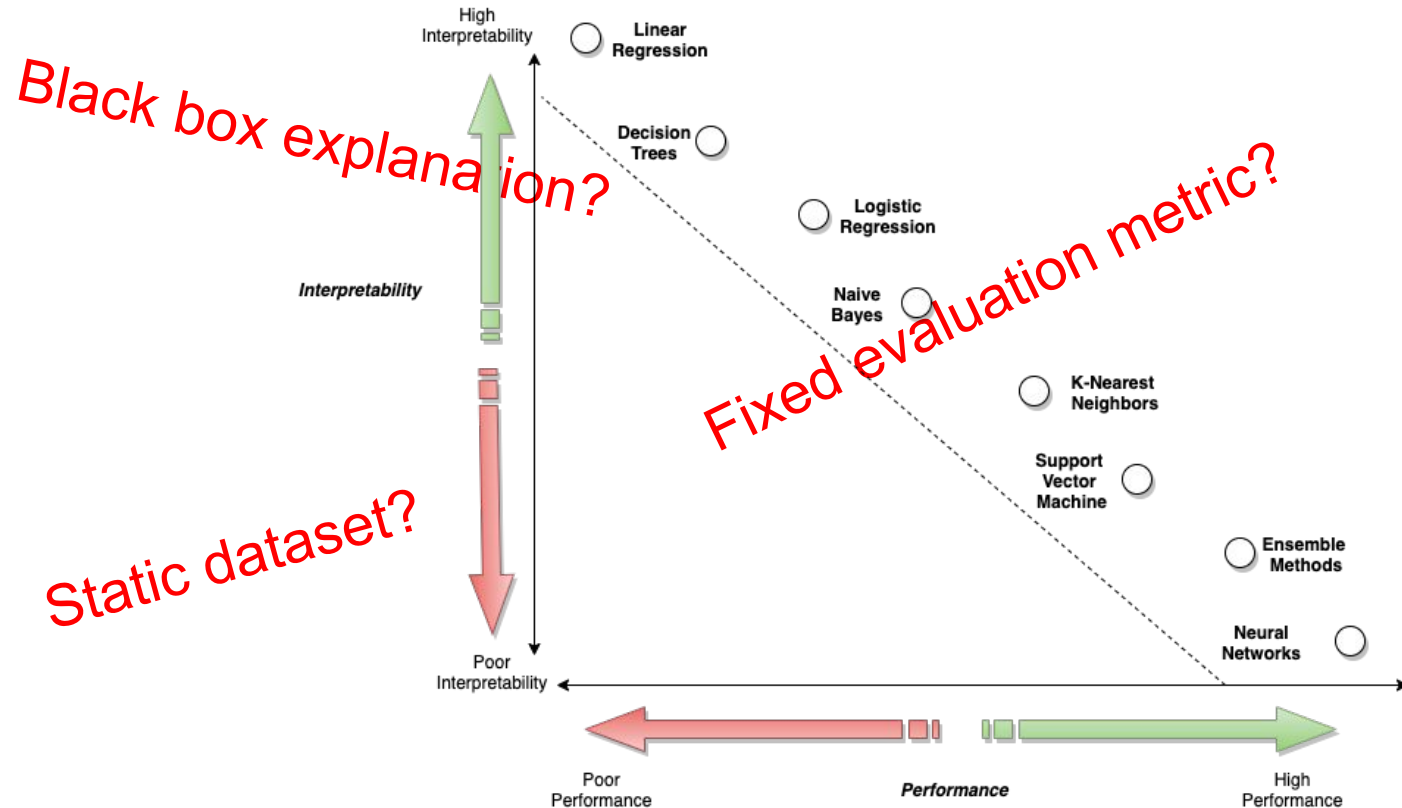
Sample Tree #3



Power / Interpretability Tradeoff ?

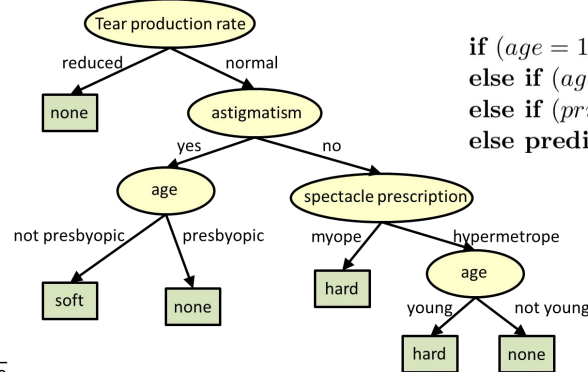
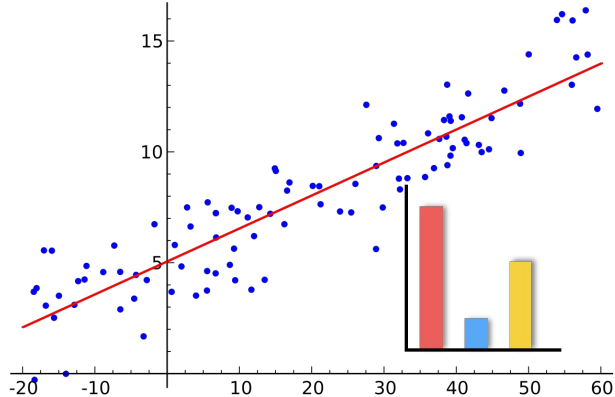


Power / Interpretability Tradeoff ?



Achieving Model Understanding

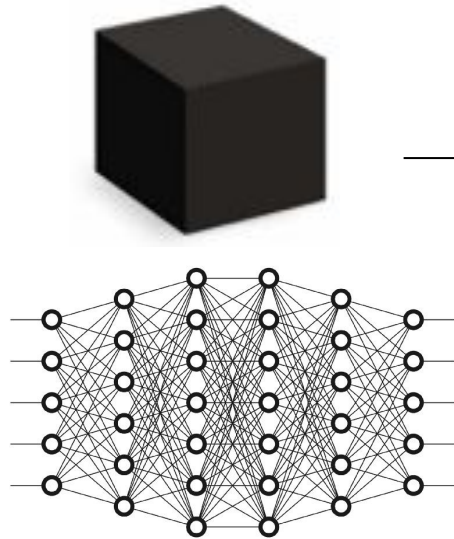
Take 1: Build *inherently interpretable* predictive models



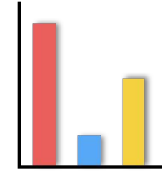
if (age = 18 – 20) and (sex = male) then predict yes
 else if (age = 21 – 23) and (priors = 2 – 3) then predict yes
 else if (priors > 3) then predict yes
 else predict no

Achieving Model Understanding

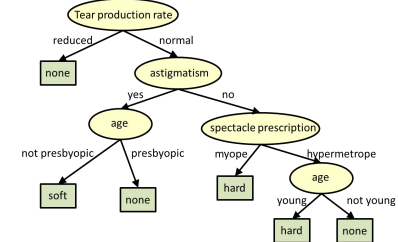
Take 2: *Explain pre-built models in a post-hoc manner*



Explainer



if ($age = 18 - 20$) and ($sex = male$) then predict *yes*
else if ($age = 21 - 23$) and ($priors = 2 - 3$) then predict *yes*
else if ($priors > 3$) then predict *yes*
else predict *no*



Inherently Interpretable Models vs. Post hoc Explanations

If you *can build* an interpretable model which is also adequately accurate for your setting, DO IT!

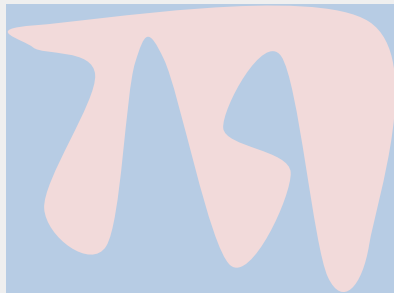
Otherwise, *post hoc explanations* come to the rescue!

What is an Explanation?

What is an Explanation?

Definition: Interpretable description of the model behavior

Classifier



Faithful

Explanation

Understandable

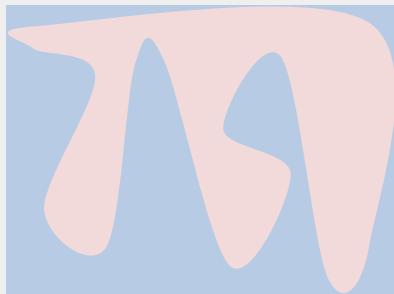
User



What is an Explanation?

Definition: Interpretable description of the model behavior

Classifier



Send all the model parameters θ ?

Send many example predictions?

Summarize with a program/rule/tree

Select most important features/points

Describe how to *flip* the model prediction

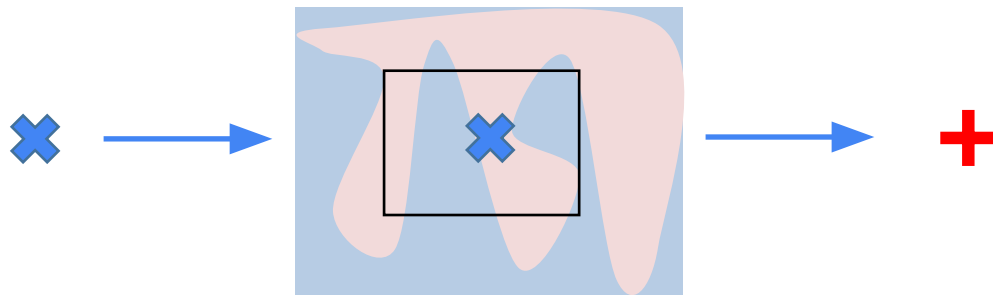
...

User



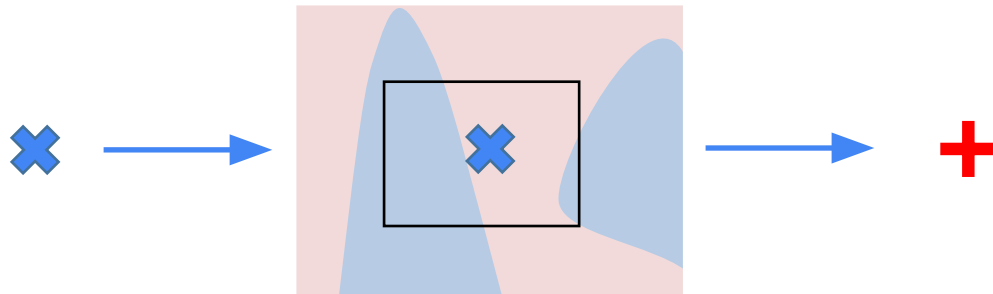
Local vs. Global Explanations

Global explanation may be too complicated



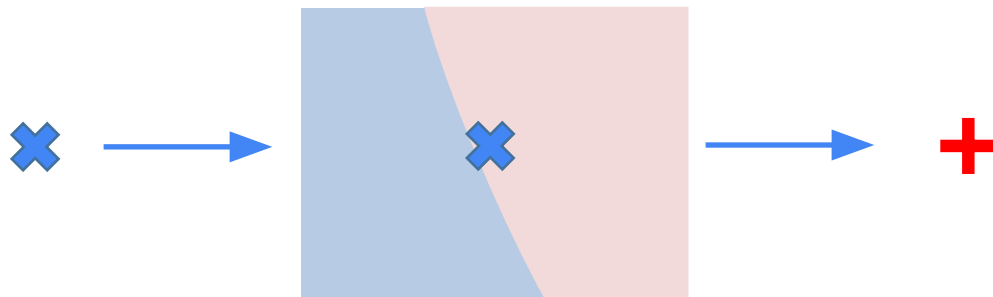
Local vs. Global Explanations

Global explanation may be too complicated



Local vs. Global Explanations

Global explanation may be too complicated

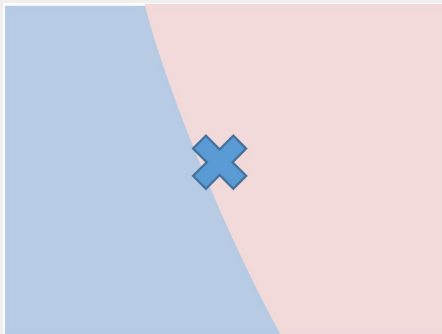


Definition: Interpretable description of the model behavior
in a target neighborhood.

Local Explanations

Definition: Interpretable description of the model behavior
in a target neighborhood.

Classifier



Send many example predictions?

Summarize with a program/rule/tree

Select most important features/points

Describe how to *flip* the model prediction

...

User



Local Explanations vs. Global Explanations

Explain individual predictions

Help unearth biases in the *local neighborhood* of a given instance

Help vet if individual predictions are being made for the right reasons

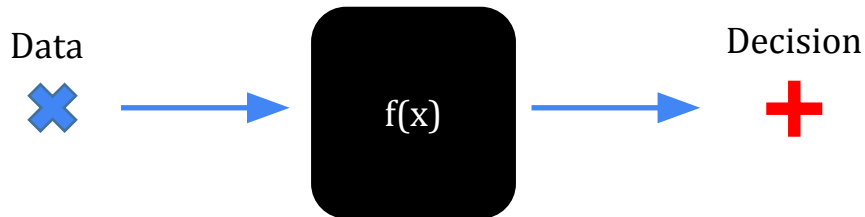
Explain complete behavior of the model

Help shed light on *big picture biases* affecting larger subgroups

Help vet if the model, at a high level, is suitable for deployment

Being Model-Agnostic...

No access to the internal structure...

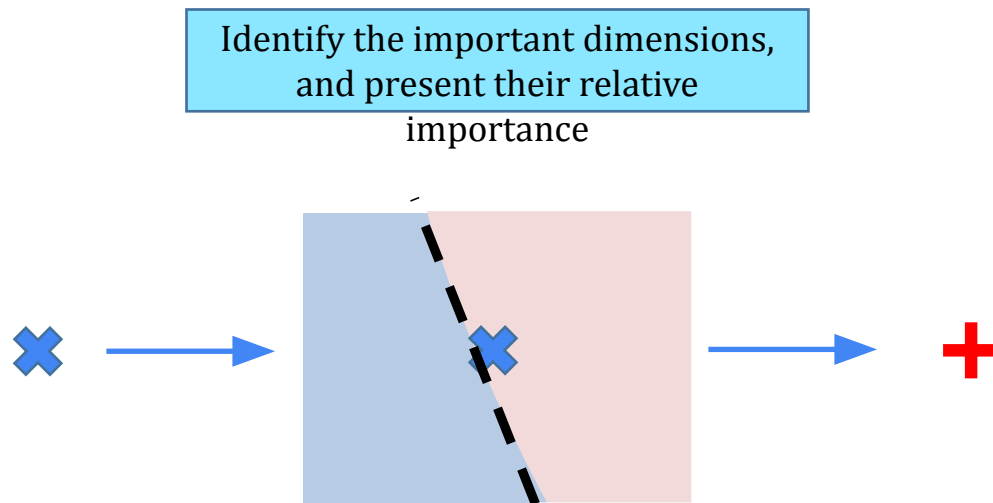


Not restricted to specific models

Practically easy: not tied to PyTorch, Tflow, etc.

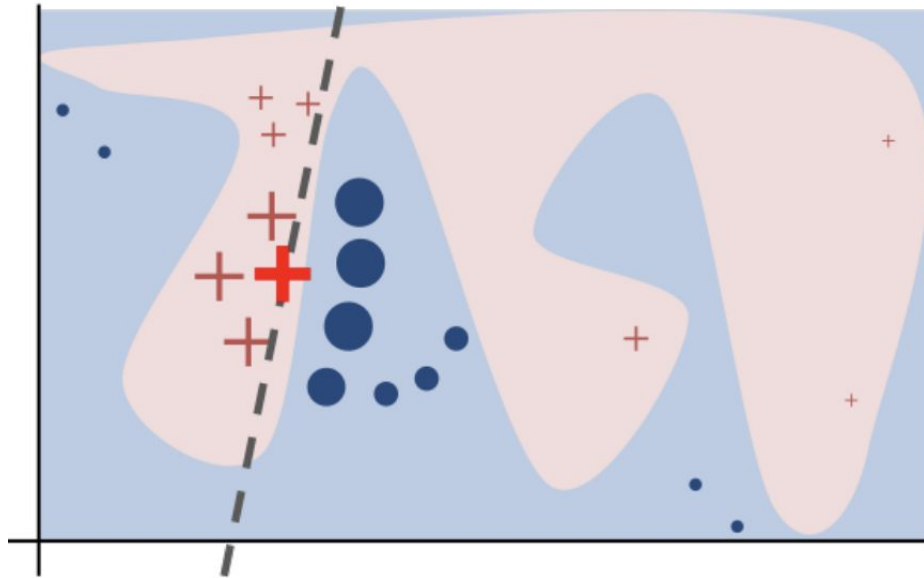
Study models that you don't have access to!

LIME: Sparse, Approximate Linear Explanations



LIME

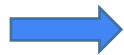
- Approximate model with something explainable (e.g. linear model)
- The approximation only needs to hold locally, i.e. on similar inputs.

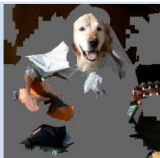




LIME Example - Images

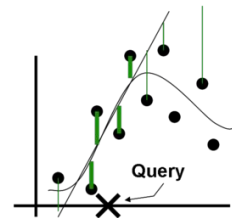


Original Image
 $P(\text{labrador}) = 0.21$

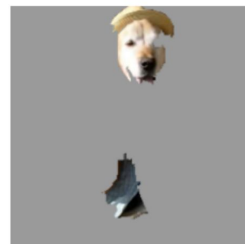


Perturbed Instances	$P(\text{Labrador})$
	<div><div></div>0.92</div>
	<div><div></div>0.001</div>
	<div><div></div>0.34</div>

Maybe to a fault?



Locally weighted
regression



Explanation

LIME is quite customizable:

- How to perturb?
- Distance/similarity?
- How *local* you want it to be?
- How to express explanation

Predict Wolf vs Husky

Only 1 mistake!



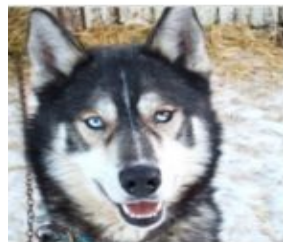
Predicted: **wolf**
True: **wolf**



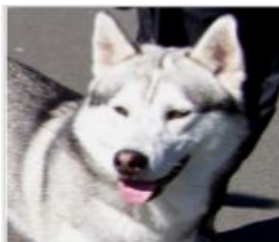
Predicted: **husky**
True: **husky**



Predicted: **wolf**
True: **wolf**



Predicted: **wolf**
True: **husky**

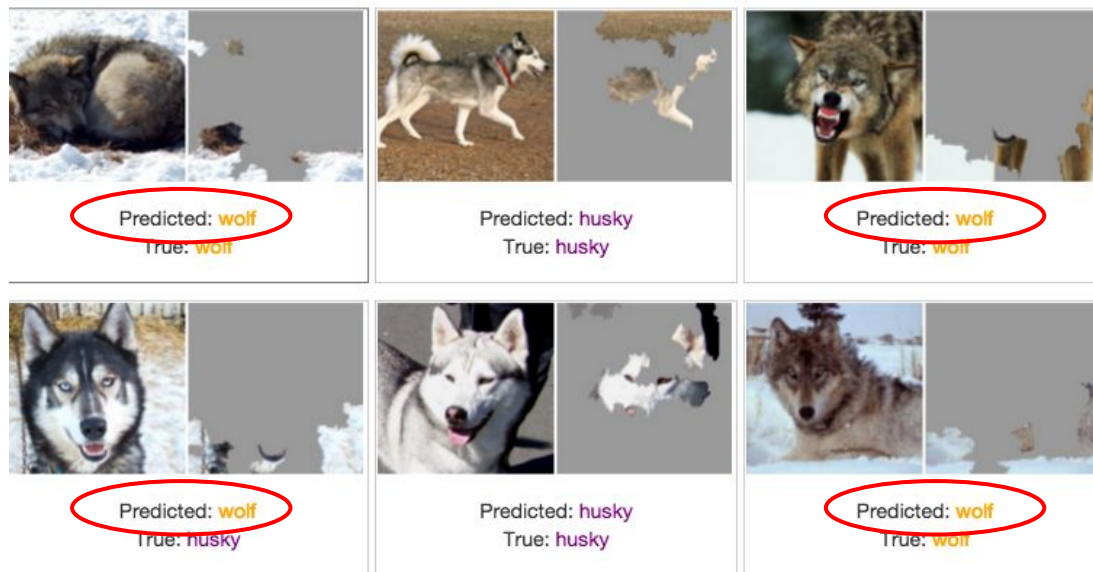


Predicted: **husky**
True: **husky**



Predicted: **wolf**
True: **wolf**

Predict Wolf vs Husky



We've built a great snow detector...

Shapley Values (from Game Theory)

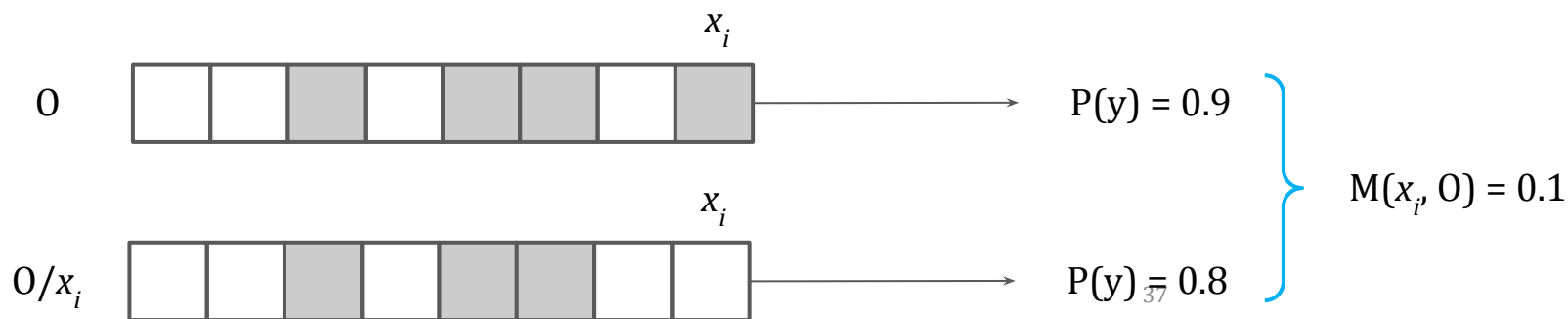
Main Analogy:

Features are “players” who play the cooperative game of making a prediction

- Prediction can be explained by assuming that each feature value of the instance is a “player” in a game.
- The contribution of each player is measured by adding and removing the player from all subsets of the rest of the players.
- The Shapley Value for one player is the weighted sum of all its contributions.

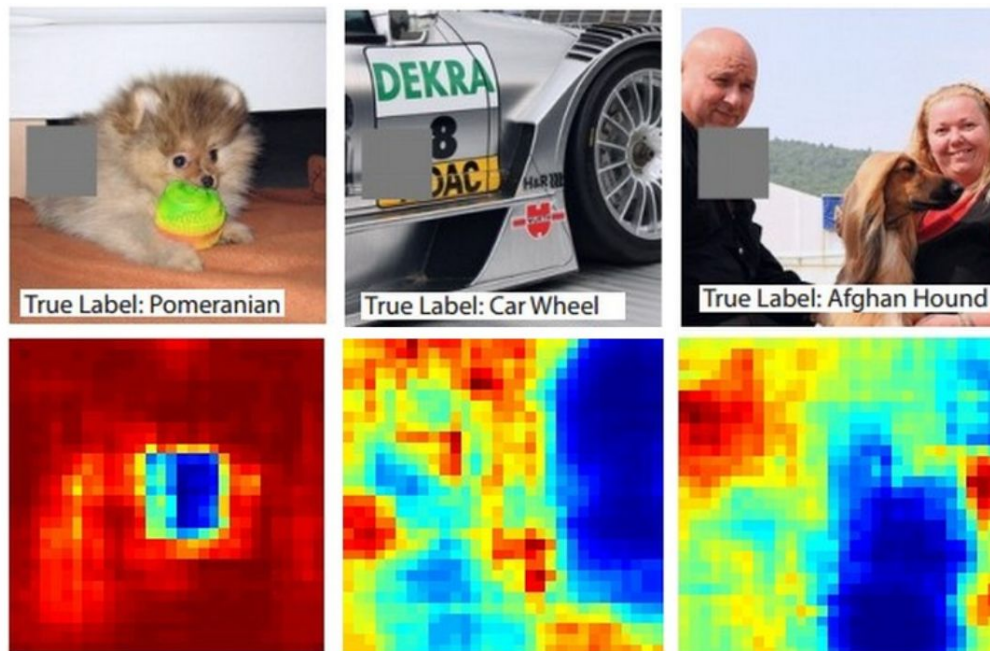
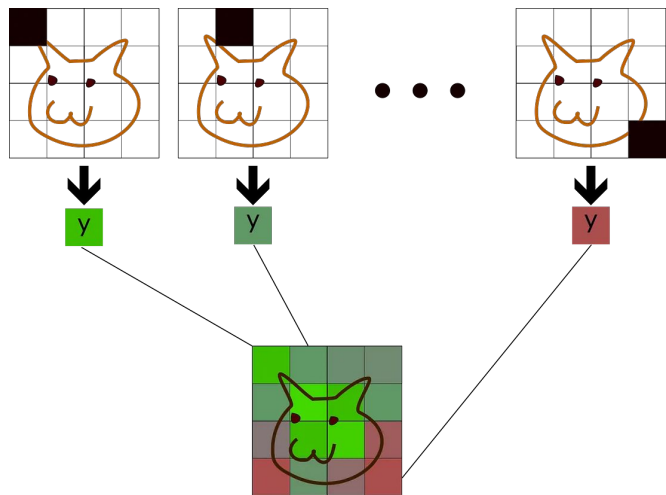
SHAP: Shapley Values as Importance

Marginal contribution of each feature towards the prediction, averaged over all possible permutations.



Fairly attributes the prediction to all the features.

Explanation by Occlusion



[Zeiler+ \(2013\)](#)

Feature/Pixel Attribution

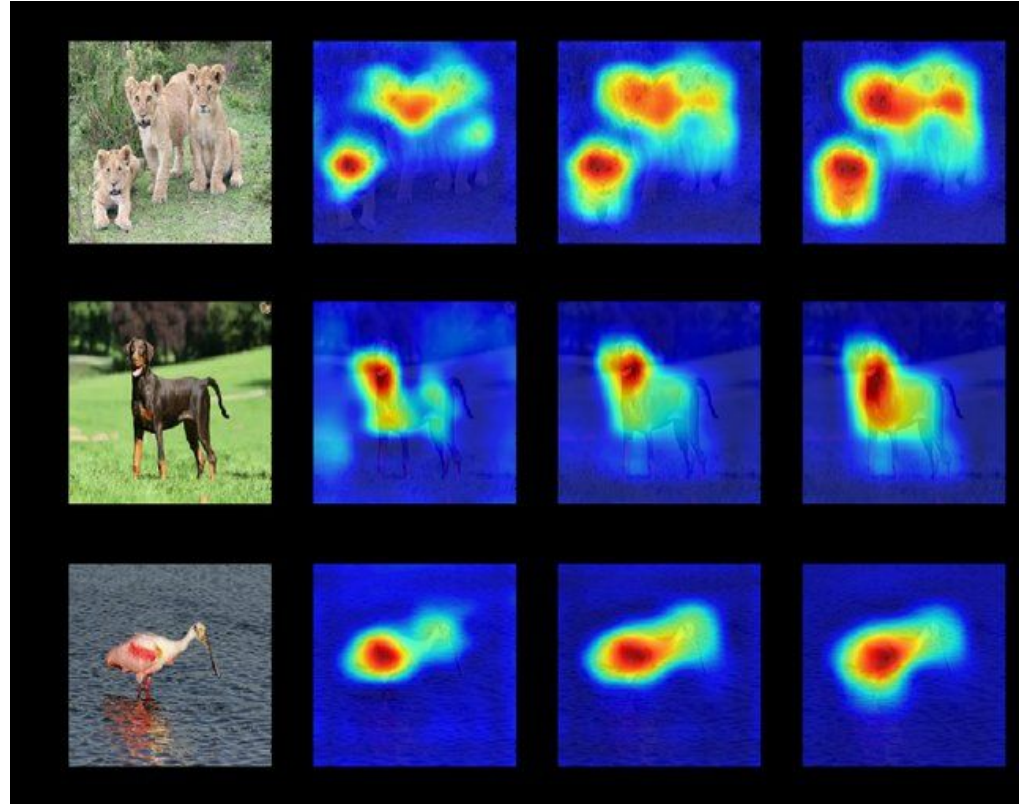
A standard neural network implicitly “focus” its attention on relevant input features to obtain the target outputs.

Can we check where it focused?

See the [CNN Explainer demo](#)

Saliency Maps / Feature Attribution

Task: Classification
Network: CNN



Neural Network Jacobian

\mathbf{x} = input vector, size k

\mathbf{y} = output vector, size m

\mathbf{J} = $m \times k$ Jacobian

Compute with back-prop:

set output 'errors' = output activations

$$J_{ij} = \frac{\partial y_i}{\partial x_j}$$

$$\mathbf{J} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_k} \end{bmatrix}$$

Saliency Map Overview

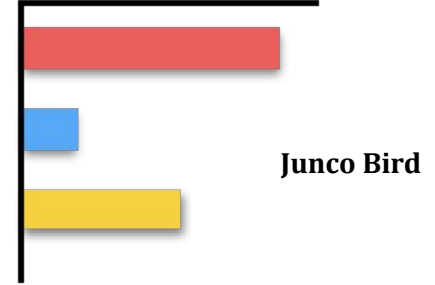
Input



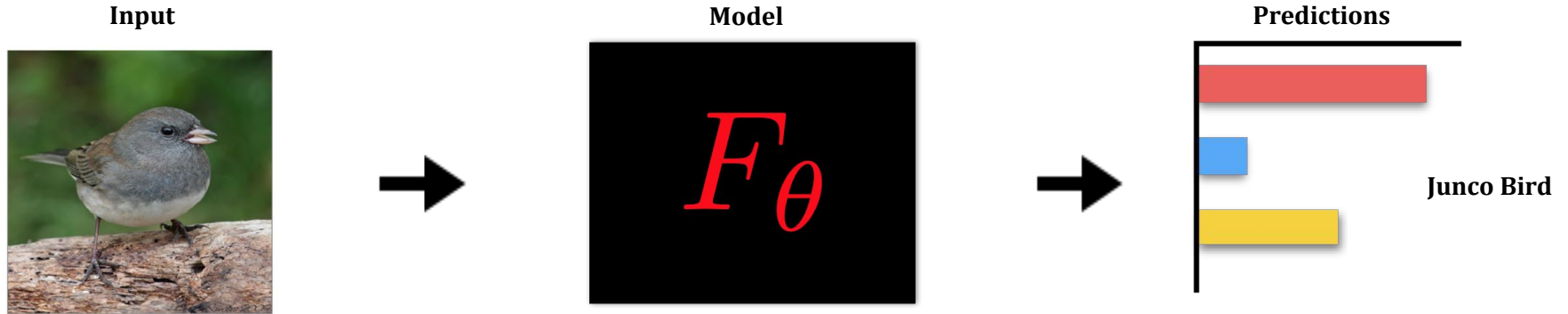
Model



Predictions

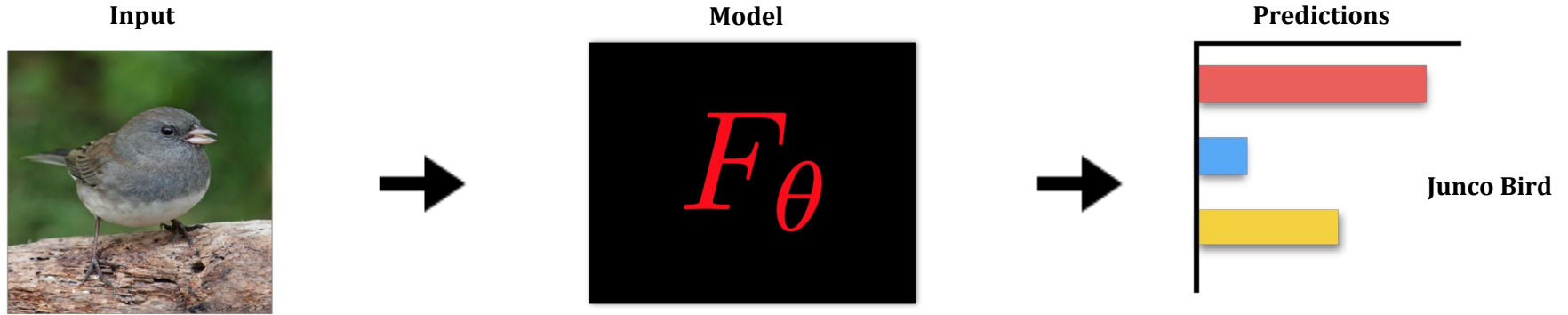


Saliency Map Overview

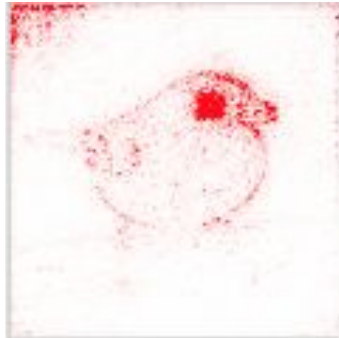


What parts of the input are most relevant for the model's prediction: **'Junco Bird'**?

Saliency Map Overview

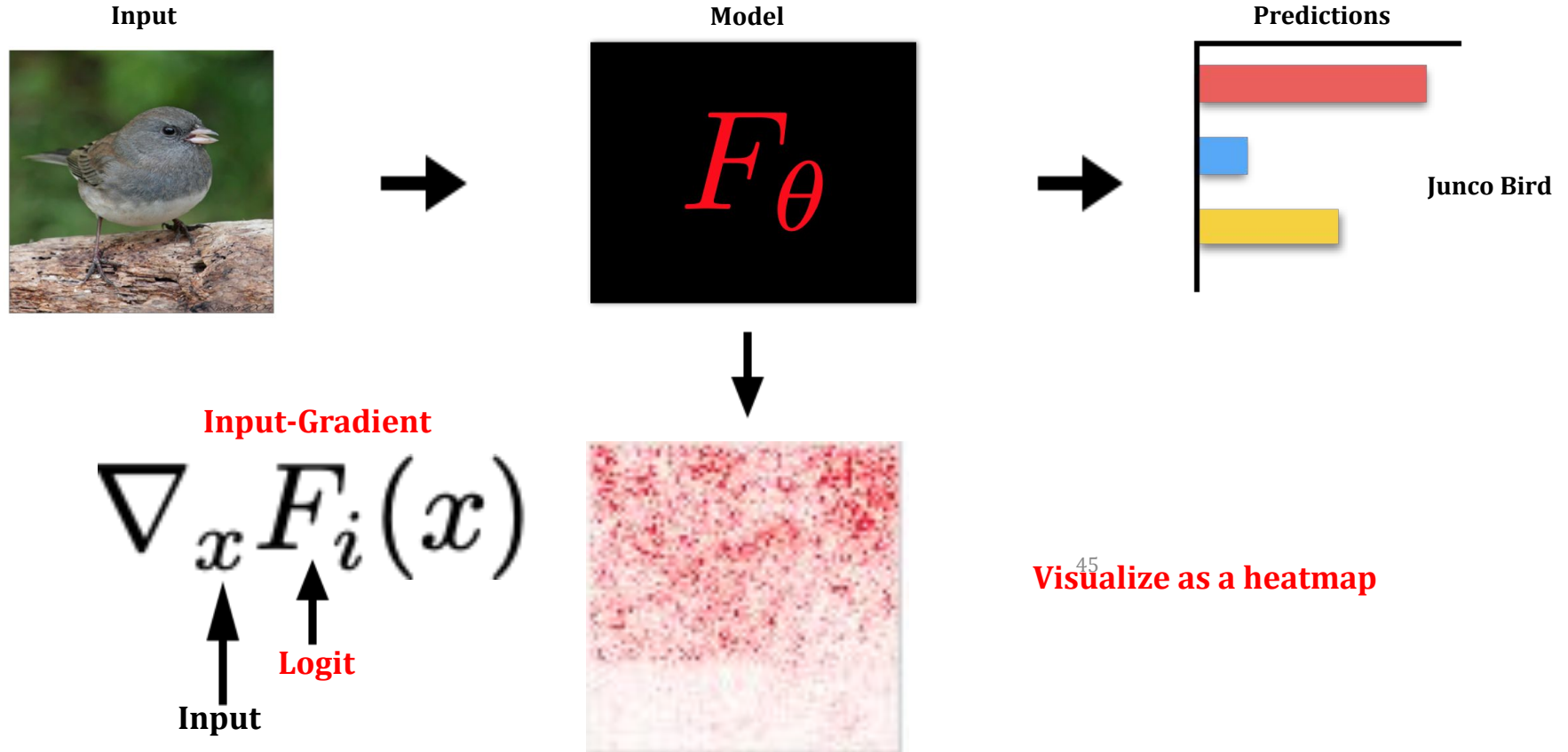


What parts of the input are most relevant for the model's prediction: **'Junco Bird'**?

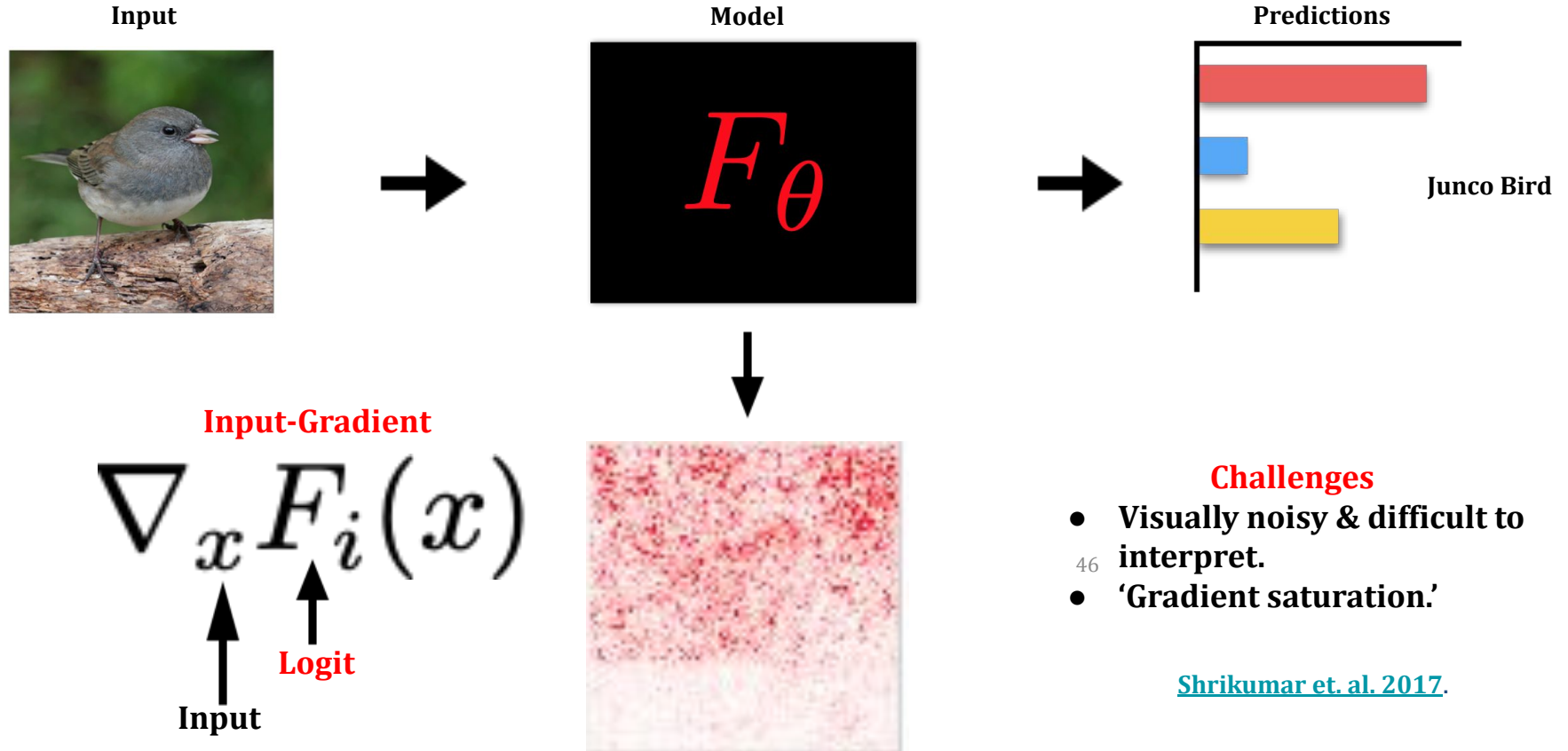


- Feature Attribution
- ^{4.4} 'Saliency Map'
- Heatmap

Input-Gradient



Input-Gradient



SmoothGrad

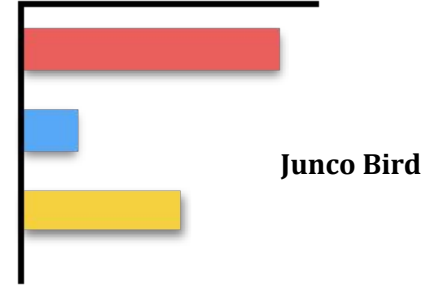
Input



Model



Predictions

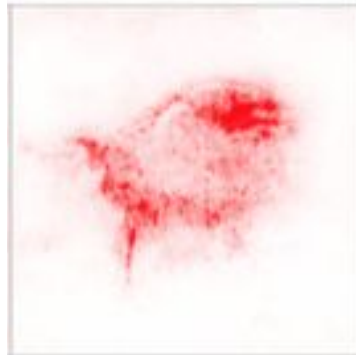


SmoothGrad

$$\frac{1}{N} \sum_i^N \nabla_{(x+\epsilon)} F_i(x + \epsilon)$$



Gaussian noise



Average Input-gradient of
'noisy' inputs.

Integrated Gradients

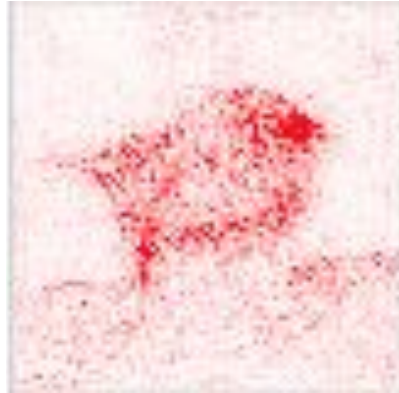
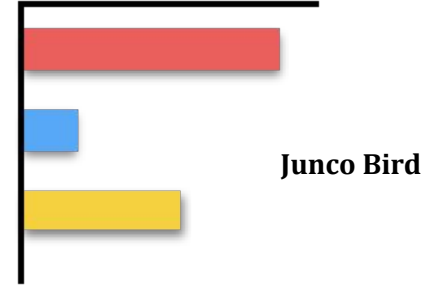
Input




Model



Predictions



Path integral: 'sum' of
interpolated gradients

$$(x - \tilde{x}) \times \int_{\alpha=0}^1 \frac{\partial F(\tilde{x} + \alpha \times (x - \tilde{x}))}{\partial x}$$


Baseline input

Recap

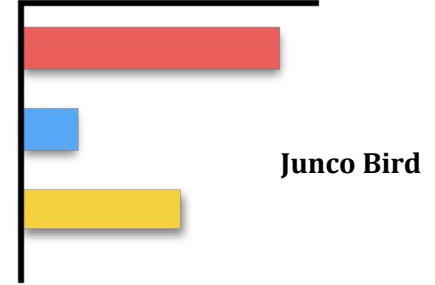
Input



Model



Predictions



LIME



SHAP



Recap

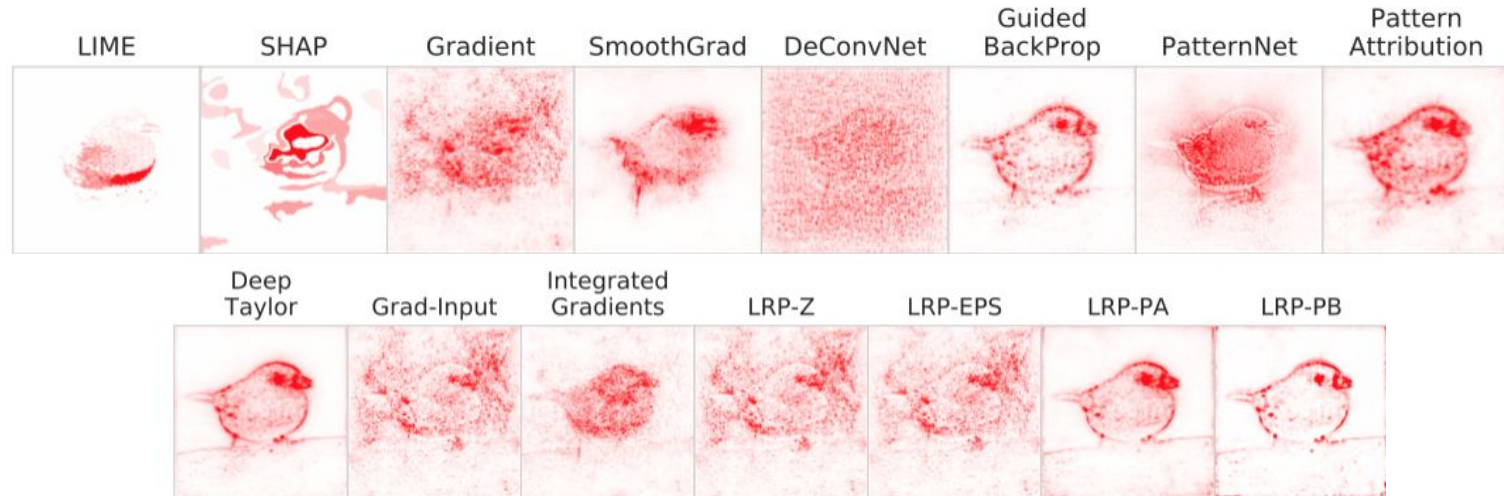
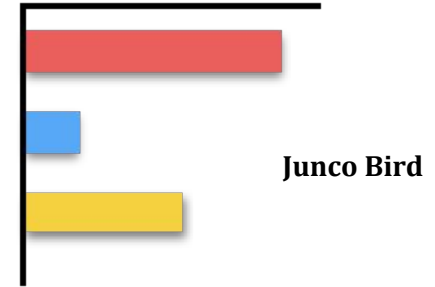
Input



Model



Predictions



Limitations to post-hoc explainability

- **Faithfulness/Fidelity**

Some explanation methods do not '*reflect*' the underlying model.

- **Fragility**

Post-hoc explanations can be easily manipulated.

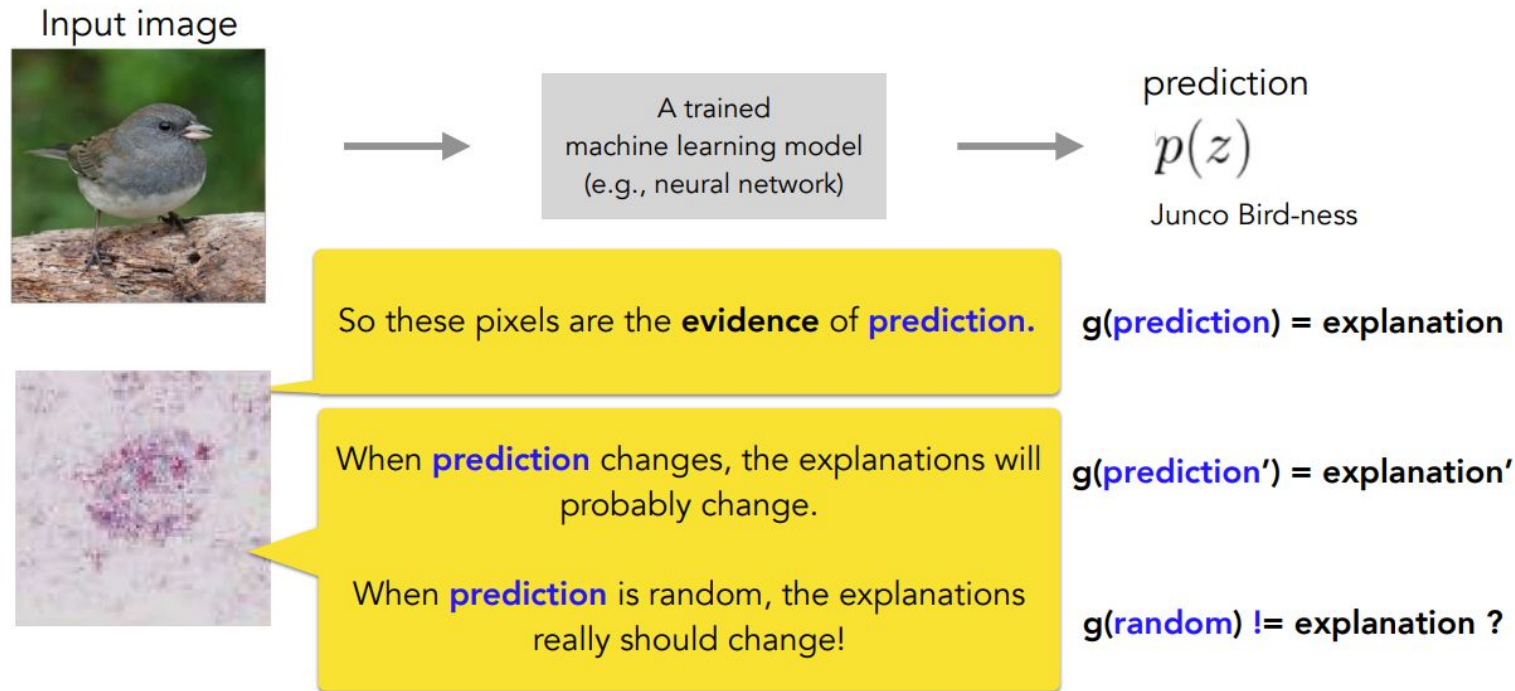
- **Stability**

Slight changes to inputs can cause large changes in explanations.

- **Useful in practice?**

Unclear if a data scientist (ML engineer)/end-user can use explanations to isolate errors, improve 'trust' or simulate the model.

Sanity Checks



Sanity Checks for Saliency Map

Model parameter randomization test

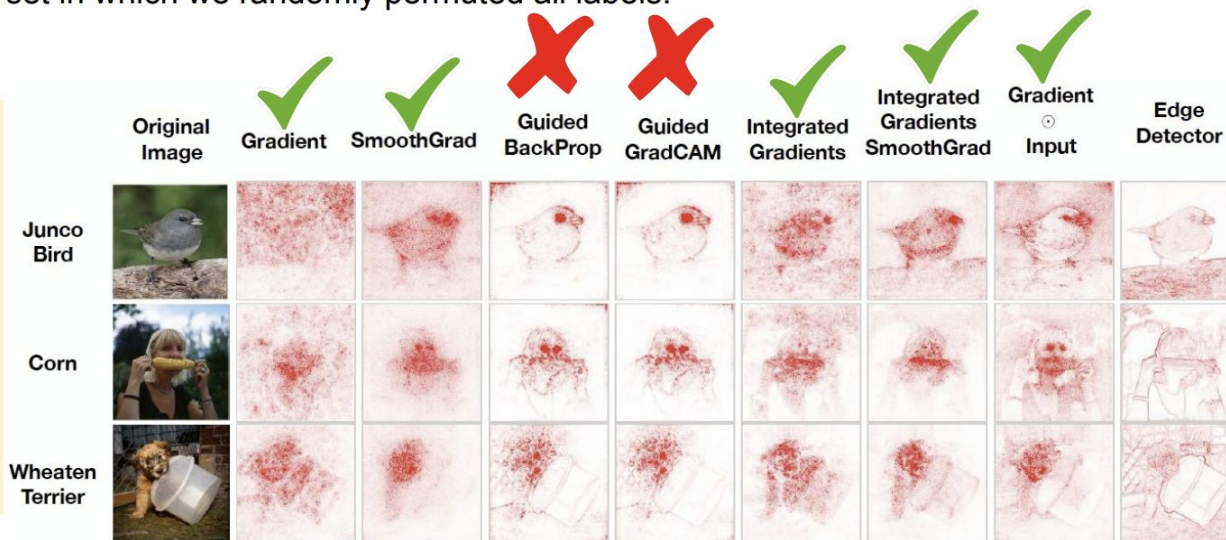
compares the output of a saliency method on a trained model with the output of the saliency method on a randomly initialized untrained network of the same architecture.

Data randomization test

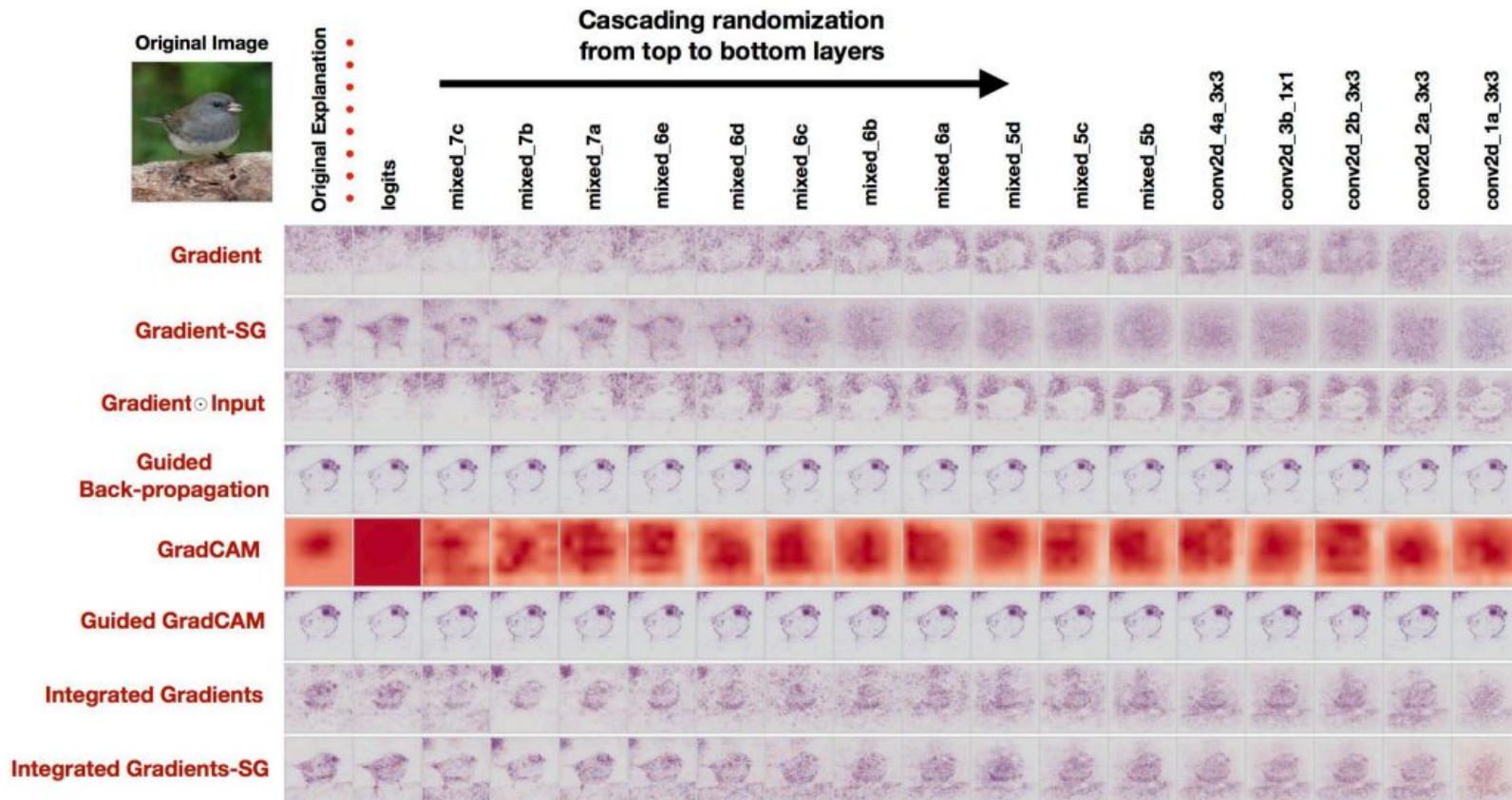
compares a given saliency method applied to a model trained on a labeled data set with the method applied to the same model architecture but trained on a copy of the data set in which we randomly permuted all labels.

“We find that reliance, solely, on visual assessment can be misleading.

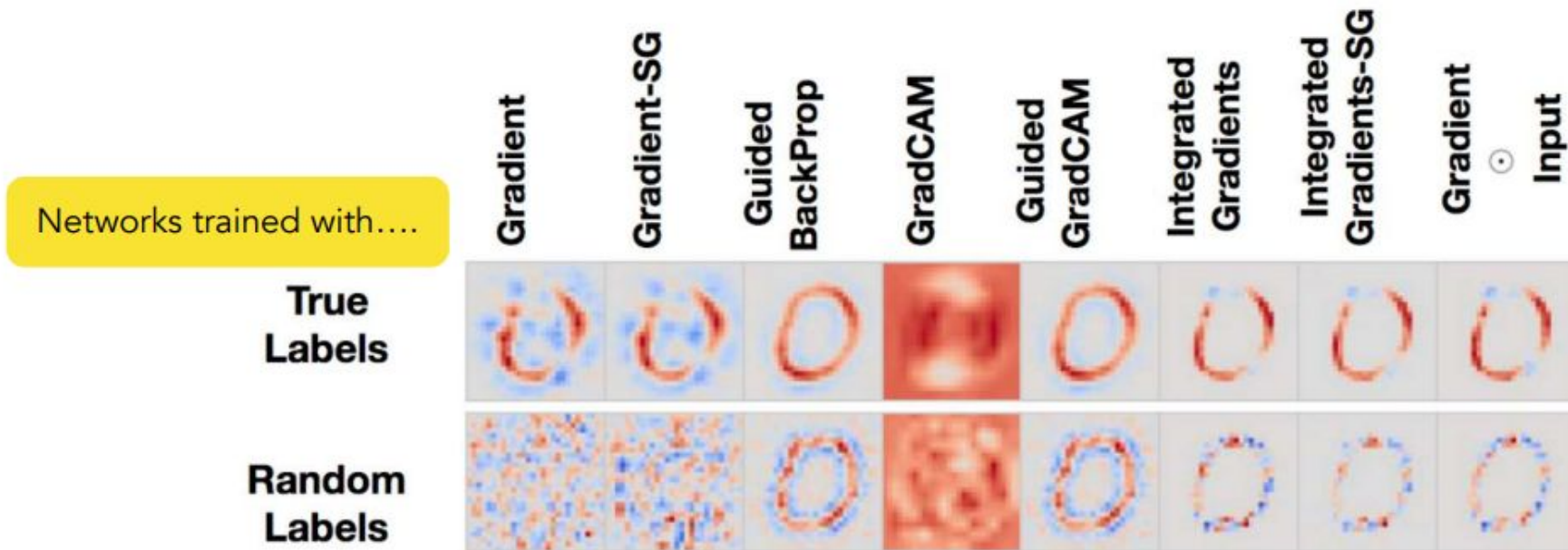
Through extensive experiments we show that some existing saliency methods are independent both of the model and of the data generating process.”



Sanity Check 1: Replace Weight with Random Layer



Sanity Check 2: Replace the Trained Model with Bogus



See [Adebayo+ \(2018\)](#)

Baked-in interpretable CNN

Why is this bird classified as a clay-colored sparrow?



Because this part of the bird

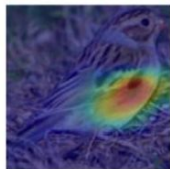


looks like

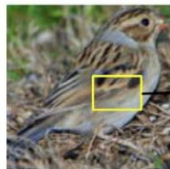
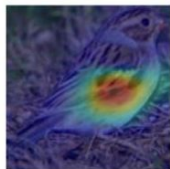


that part

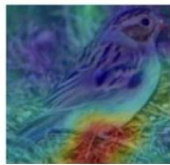
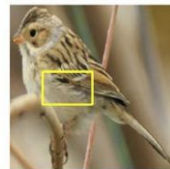
of a prototypical clay-colored sparrow



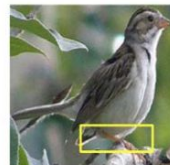
looks like



looks like



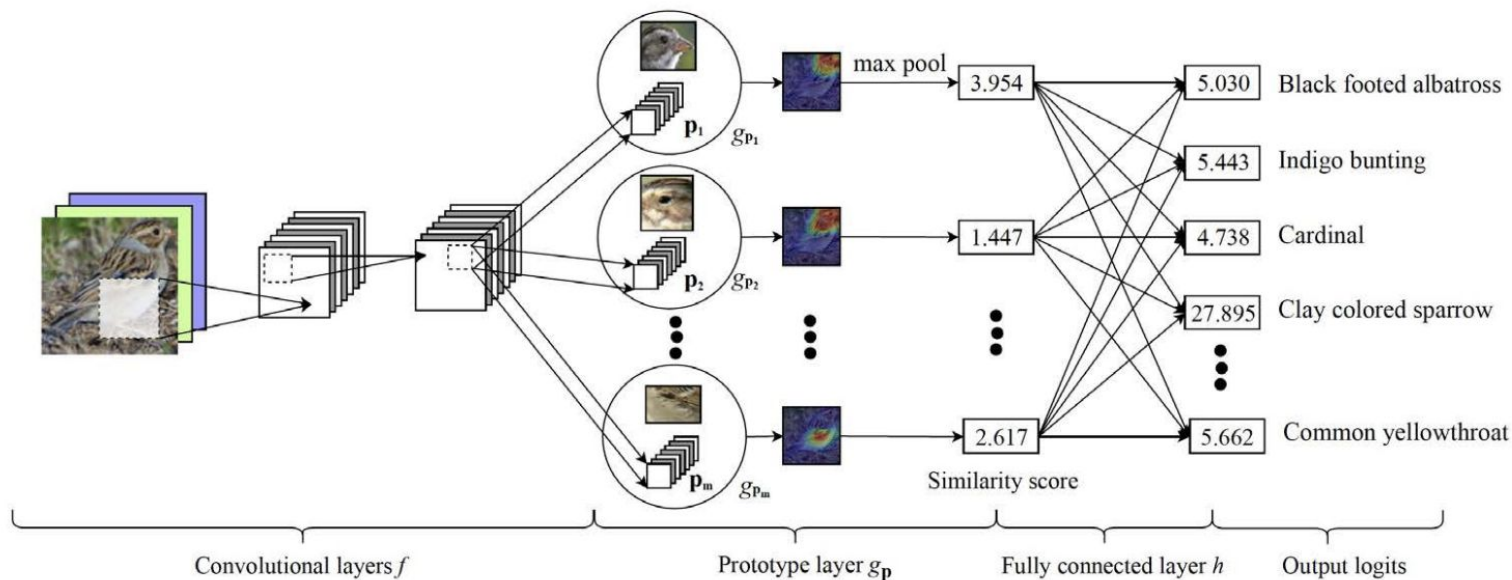
looks like



Baked-in interpretable CNN

[Chen+ \(2019\)](#)

Take any “standard” black box CNN...
And transform it to be interpretable



Influence Functions



- Remove specific data points or features, and measure their influence on a performance metric
- Largest change in performance indicates the most influential data points or features

New Journal of Physics

The open access journal at the forefront of physics

PAPER • OPEN ACCESS

Phase detection with neural networks: interpreting the black box

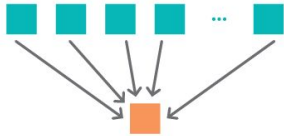
Anna Dawid^{4,1,2} , Patrick Huembeli² , Michal Tomza¹ , Maciej Lewenstein^{2,3}  and Alexandre Dauphin² 

Published 12 November 2020 • © 2020 The Author(s). Published by IOP Publishing Ltd on behalf of the Institute of Physics and Deutsche Physikalische Gesellschaft

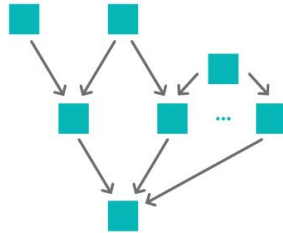
[Dawid+ \(2020\)](#)

Structural Methods

Machine Learning

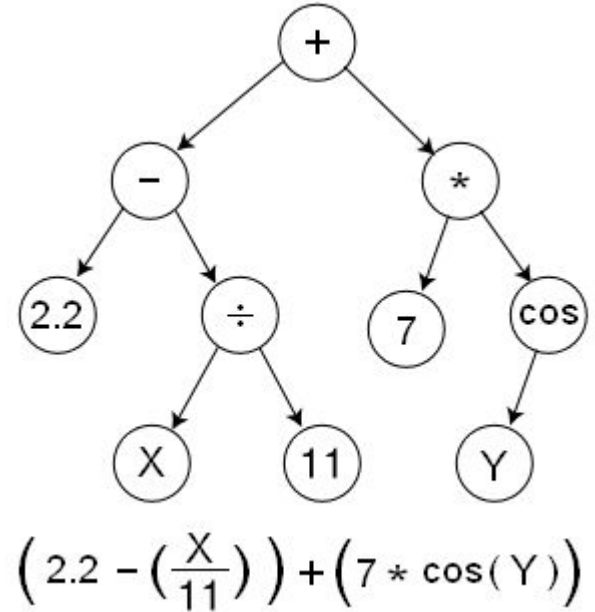


Casuality



Causal Machine Learning
tries to make sense on the
data generation process

Symbolic Regression tries to
find the best mathematical
expression to match a
dataset



Model Transparency

- Specify training set, metrics, limitations

[Mitchell+ \(2018\)](#)

[Model Card creation tool.](#)

Model Card for Breast Cancer Wisconsin (Diagnostic) Dataset

Model Details

Overview

This model predicts whether breast cancer is benign or malignant based on image measurements.

Version

name: bba6bec9-9d5c-4f0e-a291-72125c2c534a
date: 2020-09-25

Owners

- Model Cards Team, model-cards@google.com

References

- [https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))
- <https://minds.wisconsin.edu/bitstream/handle/1793/59692/TR1131.pdf>

Considerations

Intended Users

- Medical professionals
- ML researchers

Use Cases

- Breast cancer diagnosis

Limitations

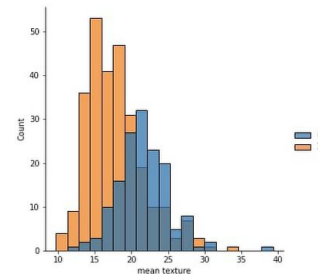
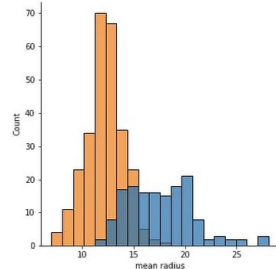
- Breast cancer diagnosis

Ethical Considerations

- Risk: Manual selection of image sections to digitize could create selection bias
Mitigation Strategy: Automate the selection process

Train Set

426 rows with 30 features



Resources

- Tutorial: [Interpretable Machine Learning MLSS](#), Kim (2021)
- Book: [Interpretable Machine Learning](#), Molnar (2023)
- Software: [shap](#), [interpret.ml](#), [captum](#), [PySR](#), [causalml](#)