

# **Trustworthy AI for Business and Society**

**Ilker Birbil**

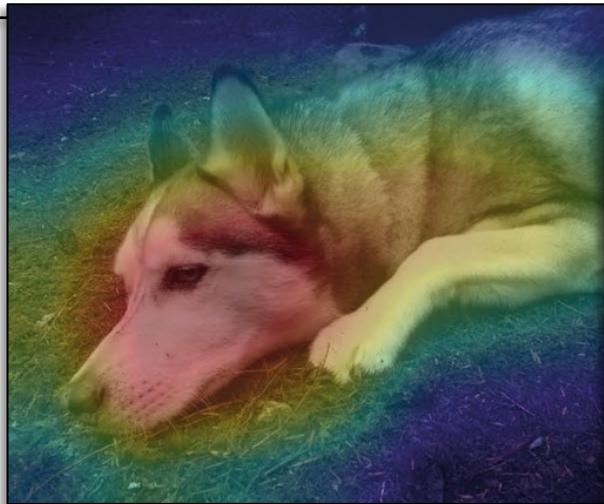
**Unboxing Methods**



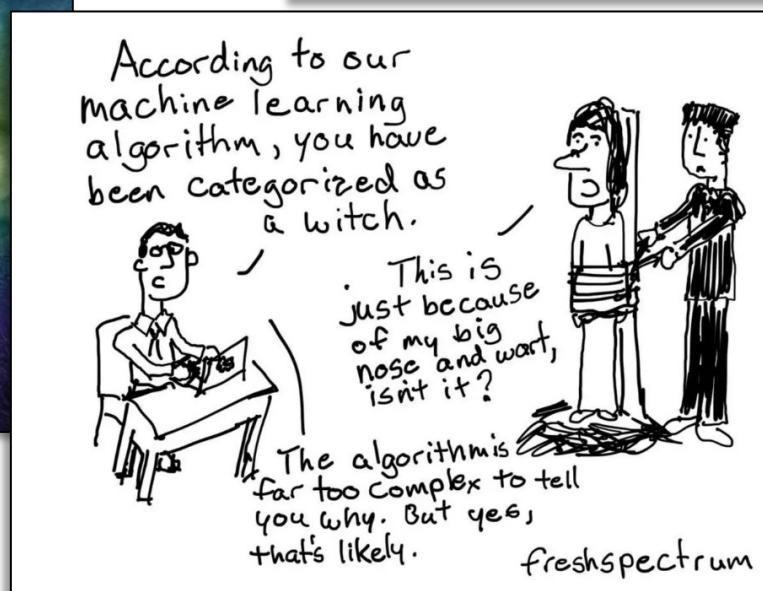
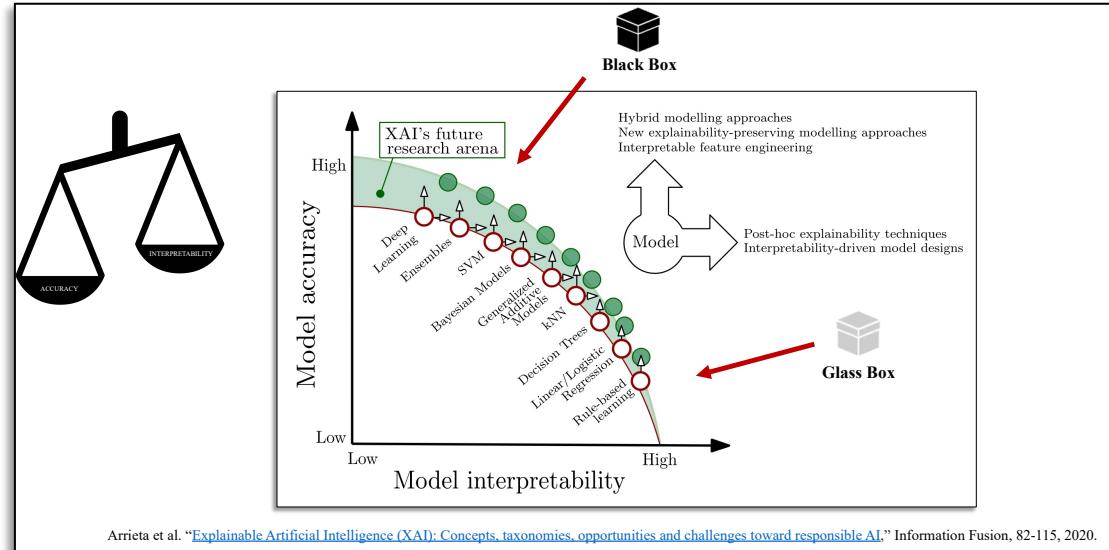
# Why Unboxing?

$$g(\mathbb{E}[Y]) = f_0 + \sum_{j=1}^p f_j(X_j)$$

contribution of each feature



[\(link\)](#)



[\(link\)](#)

**PERSPECTIVE**  
<https://doi.org/10.1038/s42256-019-0048-x>

**nature machine intelligence**

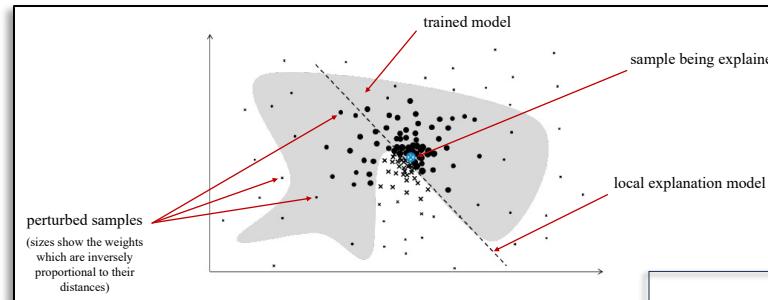
**Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead**

Cynthia Rudin

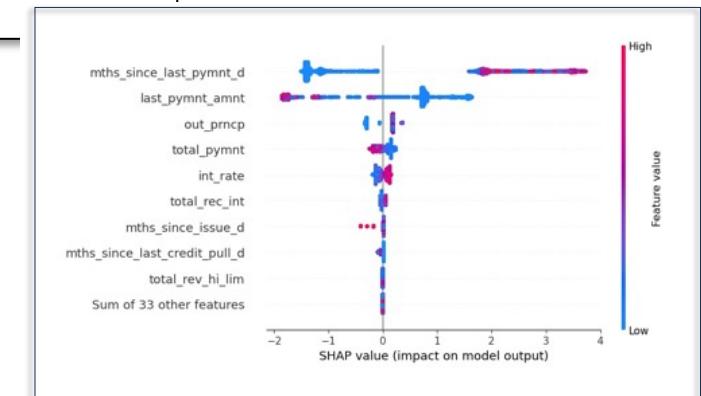
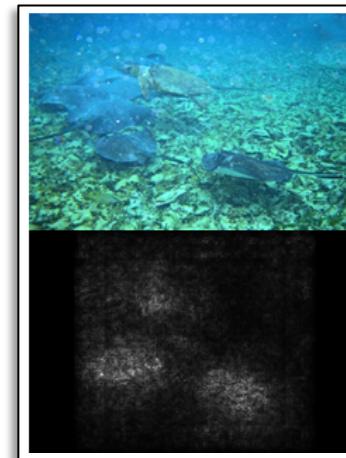


# Big Picture

- LIME



- SHAP



- Gradient Methods

- Vanilla Gradients
- SmoothGrad
- Other Methods

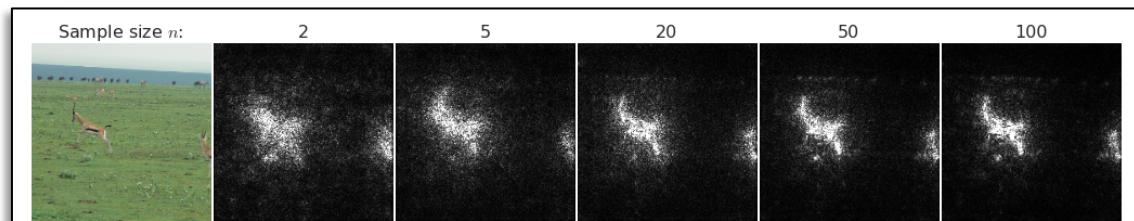
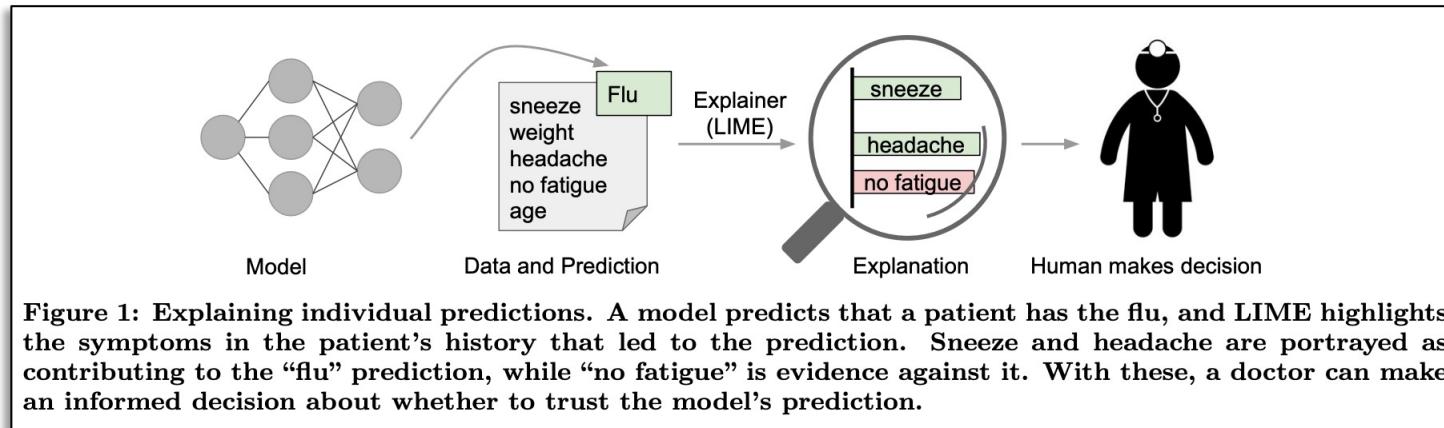


Figure 4. Effect of sample size on the estimated gradient for inception. 10% noise was applied to each image.



# LIME (Local Interpretable Model-Agnostic Explanations)



“LIME, an algorithm that can explain the predictions of any classifier or regressor in a faithful way, by approximating it locally with an interpretable model.”

“[...] an explainer should be able to explain any model, and thus be **model-agnostic** (*i.e.*, treat the original model as a black box).”

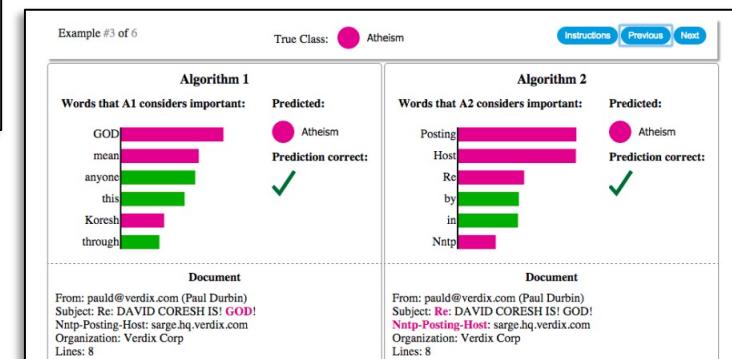
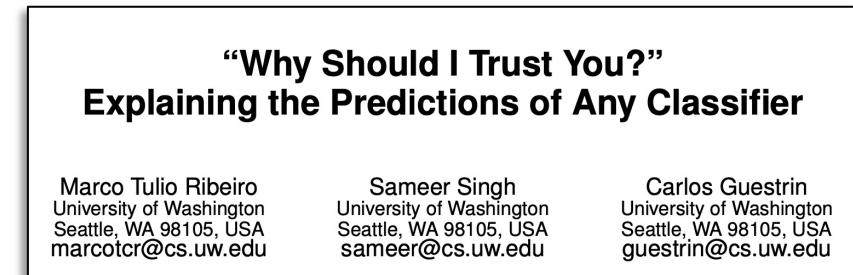
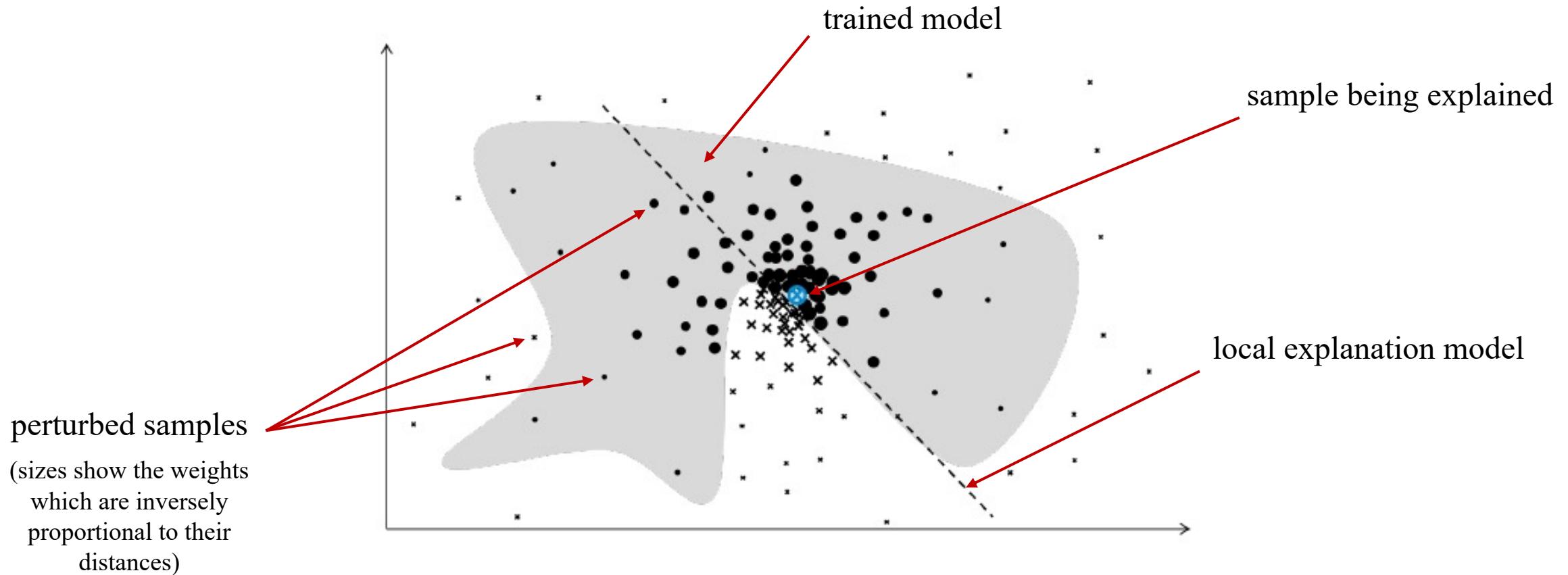


Figure 2: Explaining individual predictions of competing classifiers trying to determine if a document is about “Christianity” or “Atheism”. The bar chart represents the importance given to the most relevant words, also highlighted in the text. Color indicates which class the word contributes to (green for “Christianity”, magenta for “Atheism”).



# LIME



LIME

$\mathcal{I} = \{1, \dots, n\}$ , sample indices

$\mathcal{F} = \{1, \dots, p\}$ , feature indices

$f$ , trained model

*g*, local model

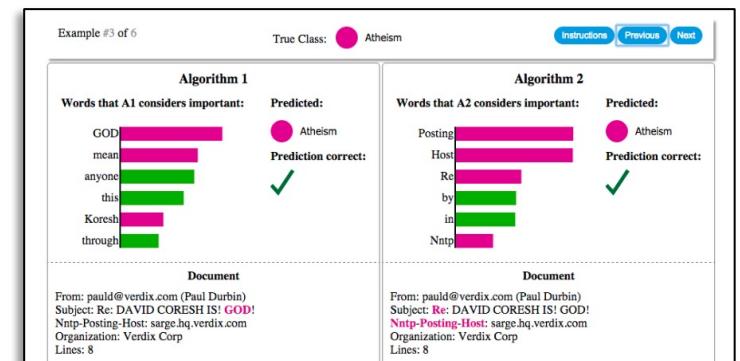
interpretable representation  
(simplified input)

$x \in \mathbb{R}^p$  →  $x' \in \{0, 1\}^{p'}$

$x = h_x(x')$

## interpretable representation (simplified input)

“a word exists (1) or not (0)”



## Local models

$$z' \approx x' \text{ and } g(z') \approx f(hx(z')) := fx(z')$$



# Optimization Problem

$f$ , trained model

$g$ , local model

$$\min_{g \in \mathcal{G}} \ell(f, g, \pi \mathbf{x}') + \Omega(g)$$

↑   ↑  
 $g(\mathbf{z}') = \mathbf{w}_g^\top \mathbf{z}'$    LIME uses  $K$ -LASSO\*

loss function

distance

model complexity

class of *potentially*  
interpretable models

$$\ell(f, g, \pi \mathbf{x}') = \sum_{\mathbf{z}' \in \mathbf{Z}} \pi \mathbf{x}'(\mathbf{z}') \left( f(\mathbf{x}') - g(\mathbf{z}') \right)^2$$

set of perturbed samples

$$\pi \mathbf{x}'(\mathbf{z}') = \exp \left( -\frac{D(\mathbf{x}', \mathbf{z}')^2}{\sigma^2} \right)$$

distance function

kernel width  $\in [0, 1]$

\* First select  $K$  features with Lasso using regularization path, then learn the weights with least squares.



# Local to Global

---

**Algorithm 2** Submodular pick (SP) algorithm
 

---

```

for  $x_i, i \in \mathcal{I}$  do
     $\mathcal{W}_i \leftarrow \text{explain}(x_i)$                                  $\triangleright$  Using Algorithm 1
end for
for  $j \in \mathcal{F}$  do
     $I_j \leftarrow \sqrt{\sum_{i \in \mathcal{I}} |\mathcal{W}_{ij}|}$      $\triangleright$  Compute feature importances
end for
 $V \leftarrow \{\}$ 
while  $|V| < B$  do       $\triangleright$  Greedy optimization of  $\max_{V, |V| \leq B} c(V, \mathcal{W}, I)$ 
     $V \leftarrow V \cup \text{argmax}_i c(V \cup \{i\}, \mathcal{W}, I)$ 
end while
return  $V$ 
  
```

---

LIME with  $K$ -Lasso

$$g(z') = \mathbf{w}_g^\top z'$$

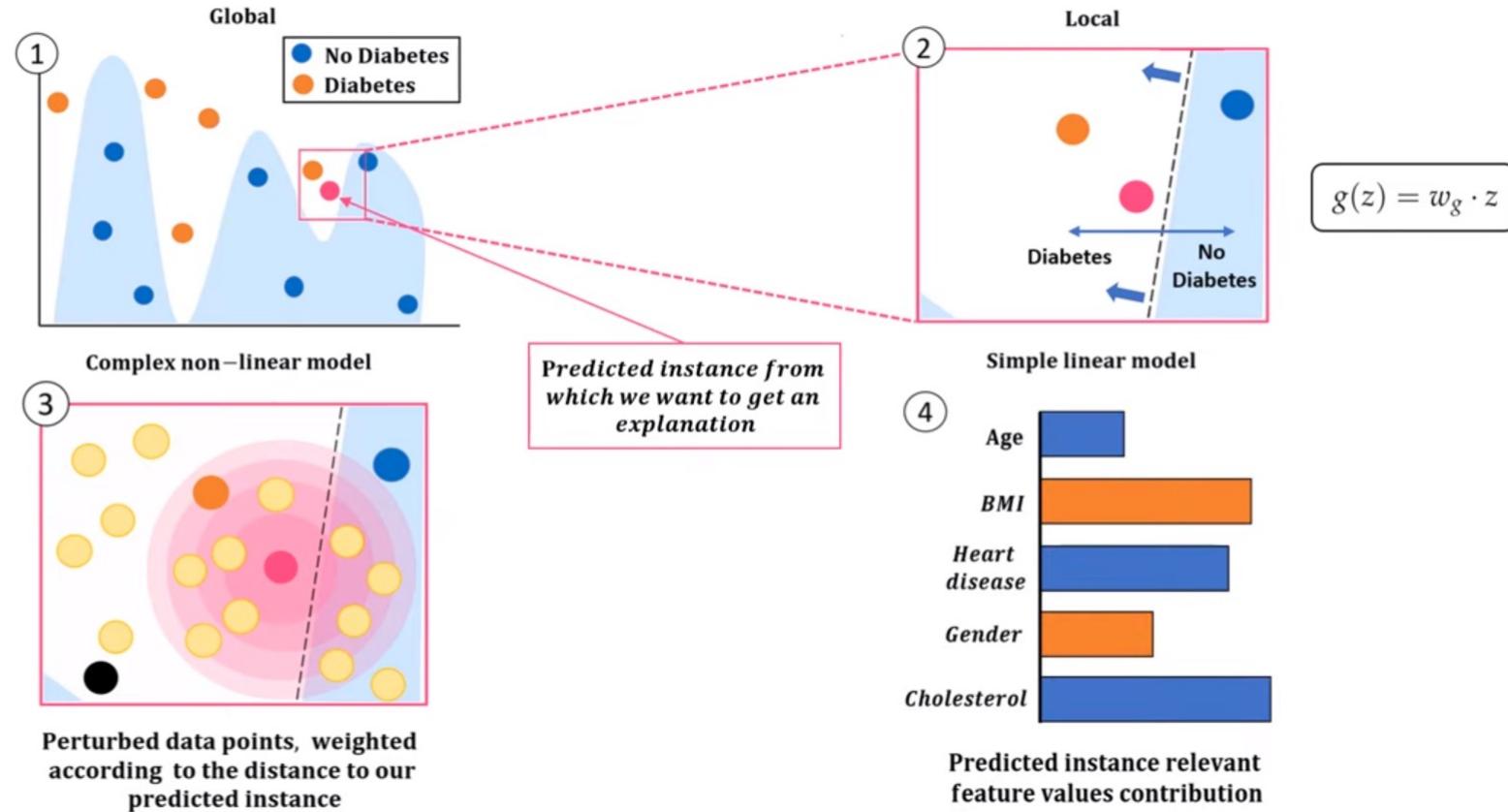
global importance

budget

coverage



# Interpretation



(source will be added soon)



# Interpretation



(a) Original Image

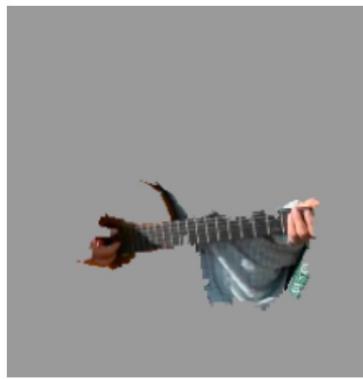
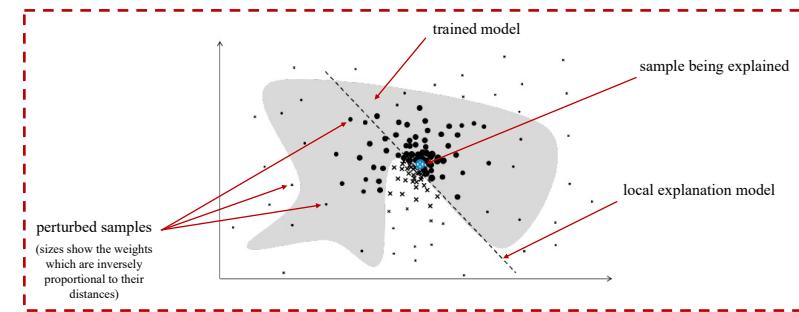
(b) Explaining *Electric guitar*(c) Explaining *Acoustic guitar*(d) Explaining *Labrador*

Figure 4: Explaining an image classification prediction made by Google's Inception neural network. The top 3 classes predicted are “Electric Guitar” ( $p = 0.32$ ), “Acoustic guitar” ( $p = 0.24$ ) and “Labrador” ( $p = 0.21$ )



# Interpretation



XGBoost model predicted an acceptance probability of 77.0%

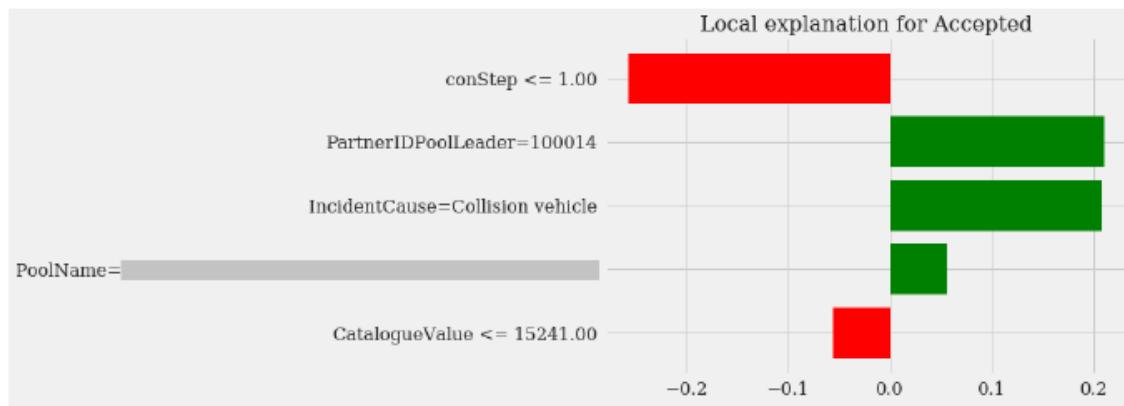


Figure 23: FP - LIME by Default (5 Features)

[Local Pred. = 0.58, Intercept = 0.42,  $R^2$  = 0.28, RMSE = 0.26, Time = 4.67 s]

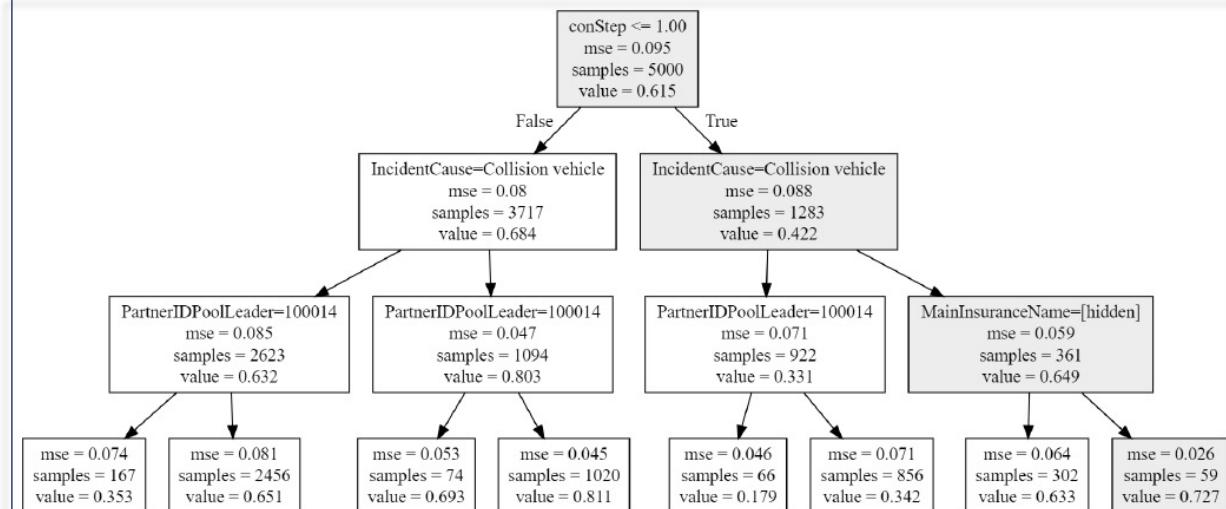


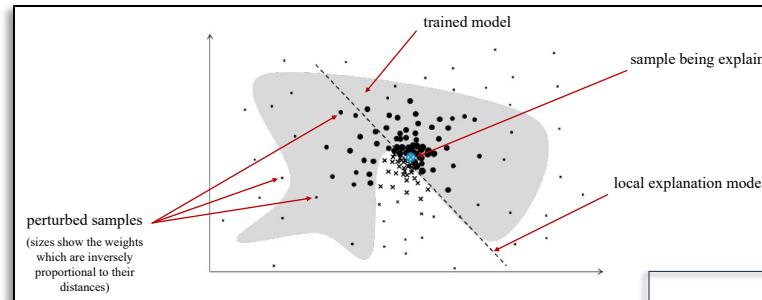
Figure 24: FP - LIME Decision Tree (3 Layers)

[Local Pred. = 0.73,  $R^2$  = 0.28, RMSE = 0.26, Time = 2.15 s]

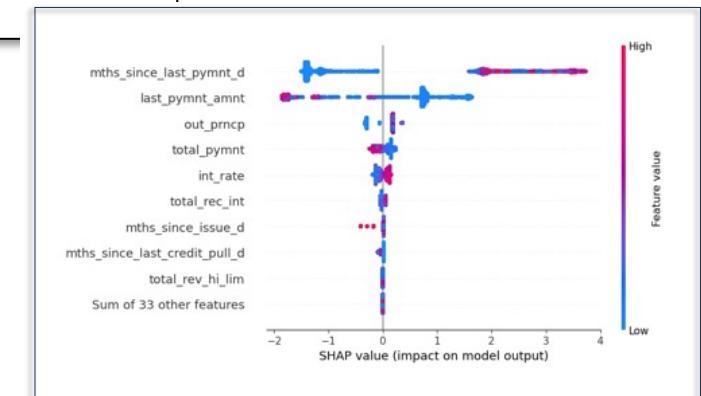
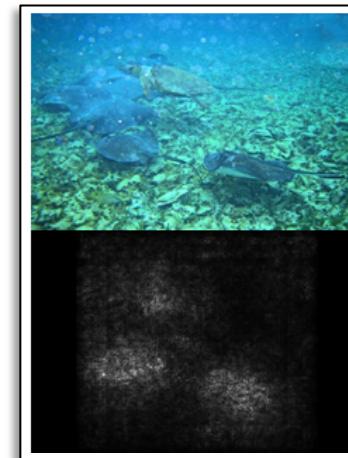


# Big Picture

- LIME



- SHAP



- Gradient Methods

- Vanilla Gradients
- SmoothGrad
- Other Methods

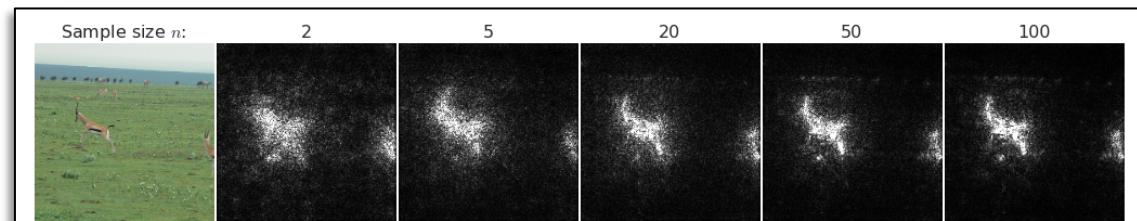


Figure 4. Effect of sample size on the estimated gradient for inception. 10% noise was applied to each image.



# Additive Feature Attribution

“SHAP assigns each feature an importance value for a particular prediction. Its novel components include: (1) the identification of a new class of **additive feature importance** measures, and (2) **theoretical results** showing there is a **unique solution** in this class with a set of desirable properties.”

$$g(\mathbf{z}') = \phi_0 + \sum_{j=1}^{p'} \phi_j z'_j$$

## A Unified Approach to Interpreting Model Predictions

**Scott M. Lundberg**

Paul G. Allen School of Computer Science  
University of Washington  
Seattle, WA 98105  
slund1@cs.washington.edu

**Su-In Lee**

Paul G. Allen School of Computer Science  
Department of Genome Sciences  
University of Washington  
Seattle, WA 98105  
suinlee@cs.washington.edu

([link](#))

presence of feature  $\in \{0, 1\}^{p'}$

attribution of each feature  $\in \mathbb{R}$



# Three Requirements

## Local Accuracy

$$f(\mathbf{x}) = g(\mathbf{x}') = \phi_0 + \sum_{j=1}^{p'} \phi_j x'_j \text{ with } \phi_0 = f_{\mathbf{x}}(\mathbf{0})$$

## Missingness

$$x'_j = 0 \implies \phi_j = 0 \text{ for all } j = 1, \dots, p'$$

## Consistency

For any two models  $f$  and  $\bar{f}$ :

$$\bar{f}_{\mathbf{x}}(\mathbf{z}') - \bar{f}_{\mathbf{x}}(\mathbf{z}'_{-j}) \geq f_{\mathbf{x}}(\mathbf{z}') - f_{\mathbf{x}}(\mathbf{z}'_{-j}) \text{ for all } \mathbf{z}' \in \mathbf{Z} \implies \phi_j(\bar{f}, \mathbf{x}) \geq \phi_j(f, \mathbf{x})$$

$\mathbf{z}'_{-j}$  means  $z'_j = 0$

$\phi_j(f, \mathbf{x})$  emphasizes the dependence of the attributions on  $f$  and  $\mathbf{x}$



# Unique Explanation

<b>Local Accuracy</b> <b>Missingness</b> <b>Consistency</b>	+	$g(\mathbf{z}') = \phi_0 + \sum_{j=1}^{p'} \phi_j z'_j$
---	---	---

$$\phi_j(f, \mathbf{x}) = \sum_{\mathbf{z}' \subseteq \mathbf{x}'} \frac{|\mathbf{z}'|!(p' - |\mathbf{z}'| - 1)!}{(p')!} \left( f_{\mathbf{x}}(\mathbf{z}') - f_{\mathbf{x}}(\mathbf{z}'_{-j}) \right)$$

$\mathcal{O}(2^p)$  subsets !

the number of non-zero entries in  $\mathbf{z}'$

all vectors  $\mathbf{z}'$  where the non-zero entries are a subset of the nonzero entries in  $\mathbf{x}'$

This follows from a classic result in game theory\* on distributing the total gain from a cooperative game.

\* Shapley, L. S. “[A value for n-person games](#),” Contributions to the Theory of Games 2.28, 307–317, 1953.



# SHAP (SHapley Additive ExPlanation) Values

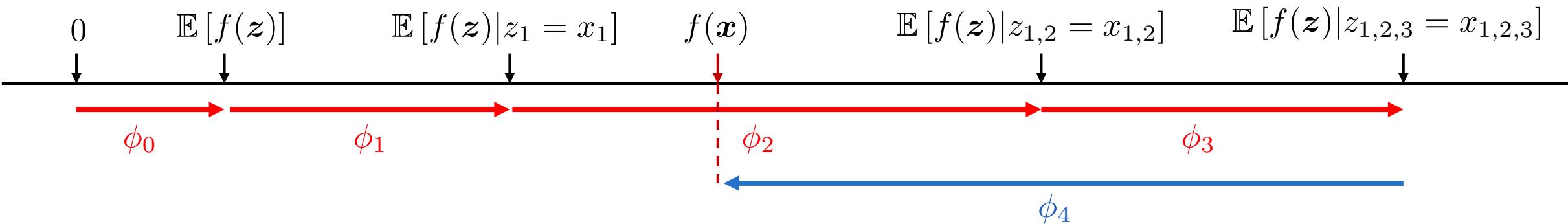
$$\phi_j(f, \mathbf{z}) = \sum_{\mathbf{z}' \subseteq \mathbf{z}} \frac{|\mathbf{z}'|!(p' - |\mathbf{z}'| - 1)!}{(p')!} \left( f_{\mathbf{x}}(\mathbf{z}') - f_{\mathbf{x}}(\mathbf{z}'_{-j}) \right)$$

$\mathcal{O}(2^p)$  subsets !

$$f_{\mathbf{x}}(\mathbf{z}') = \mathbb{E} [f(\mathbf{z}) | \mathbf{z}_{\mathcal{S}}]$$

the indices of non-zero entries in  $\mathbf{z}'$

$$h_{\mathbf{x}}(\mathbf{z}') = \mathbf{z}_{\mathcal{S}}$$



The order does matter!

$$f_{\mathbf{x}}(\mathbf{z}') \approx f([\mathbf{z}, \bar{\mathbf{z}}_{\mathcal{F} \setminus \mathcal{S}}]) \text{ with } \bar{\mathbf{z}}_{\mathcal{F} \setminus \mathcal{S}} := \mathbb{E} [\mathbf{z}_{\mathcal{F} \setminus \mathcal{S}}]$$

the set of indices not in  $\mathcal{S}$



# LIME Revisited: Kernel SHAP

**Local Accuracy ✓**

$$\min_{g \in \mathcal{G}} \ell(f, g, \pi_{\mathbf{x}'}) + \Omega(g)$$

**Missingness ✓**

**Consistency ✓**

$$\ell(f, g, \pi_{\mathbf{x}'}) = \sum_{z' \in \mathbf{Z}} \pi_{\mathbf{x}'}(z') \left( f_{\mathbf{x}}(z') - g(z') \right)^2$$

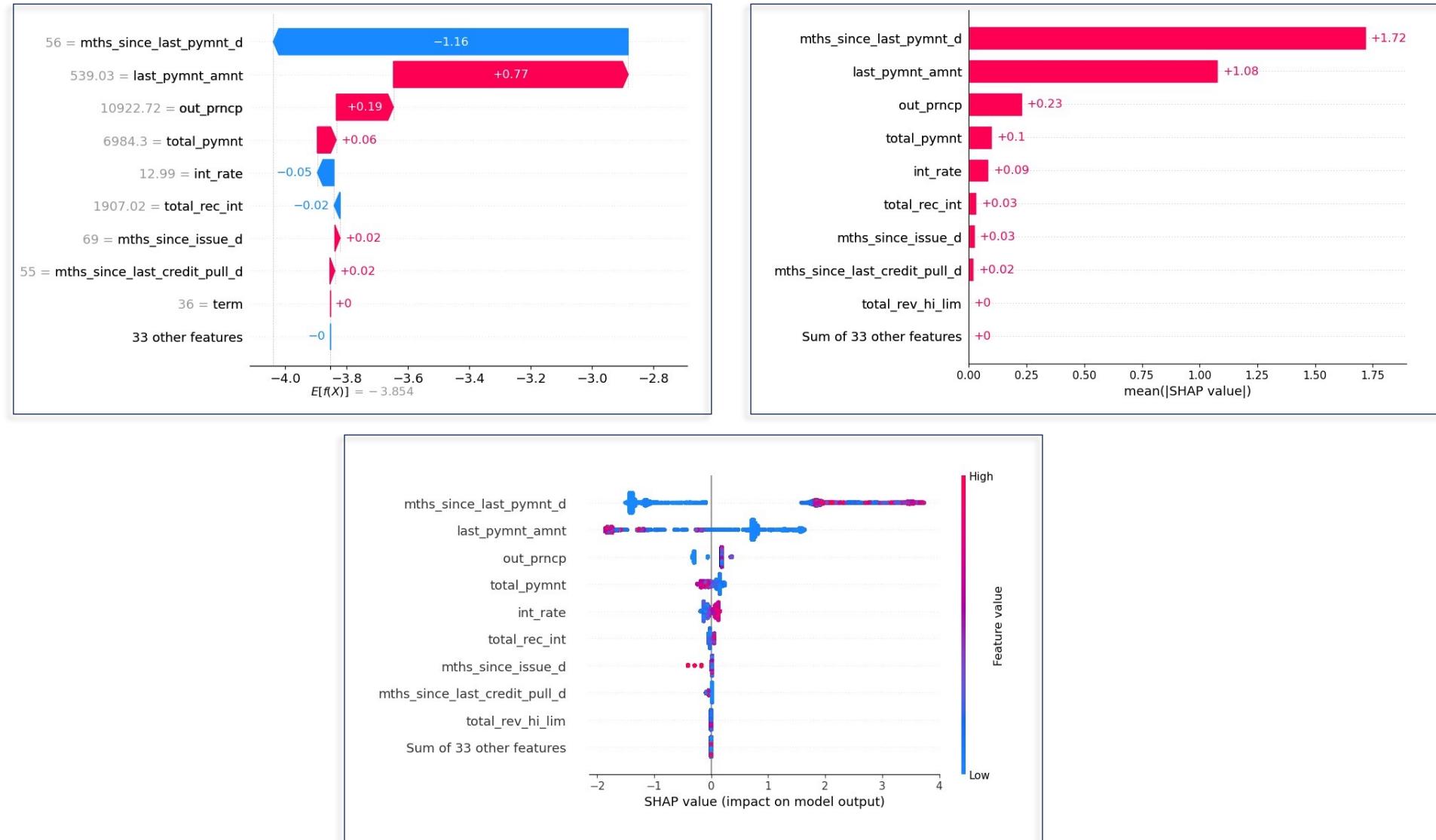
$$\Omega(g) = 0$$

$$\pi_{\mathbf{x}'}(z') = \frac{p' - 1}{\binom{p'}{|z'|} |z'| (p' - |z'|)}$$

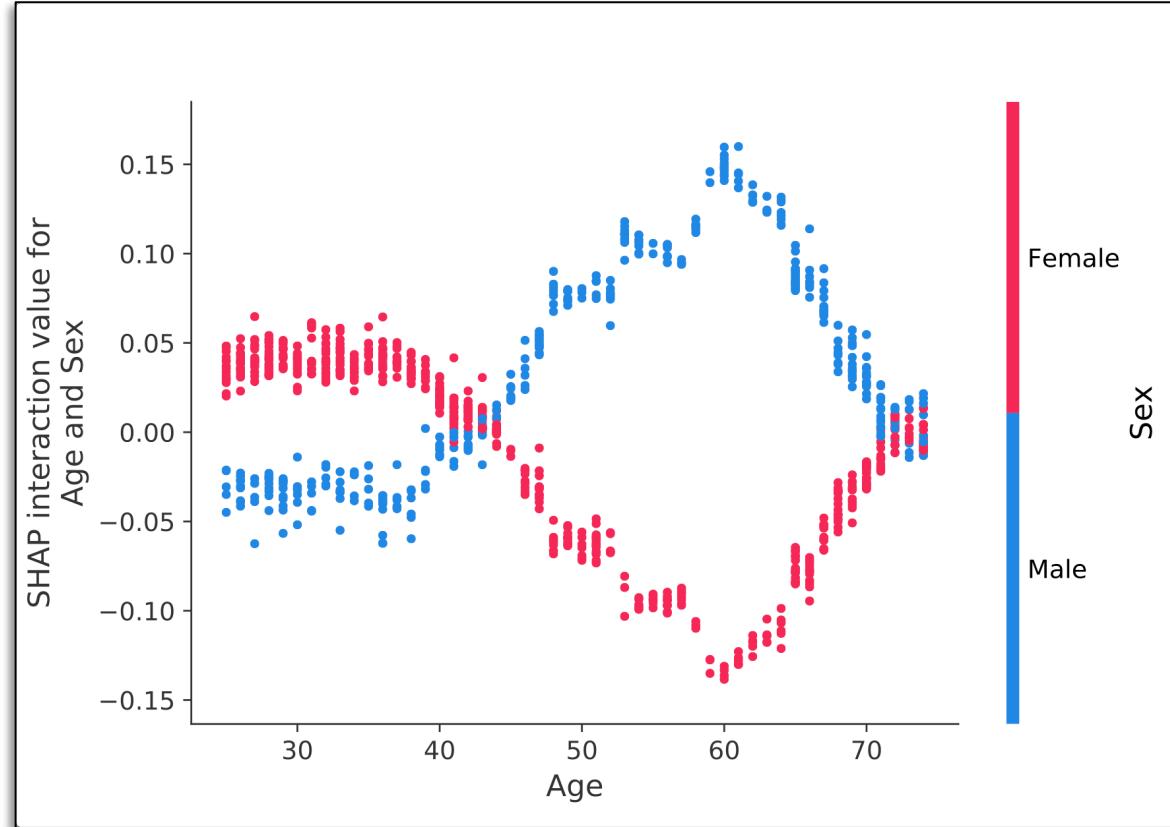
$$\pi_{\mathbf{x}'}(z') = \exp \left( -\frac{D(\mathbf{x}', z')^2}{\sigma^2} \right)$$



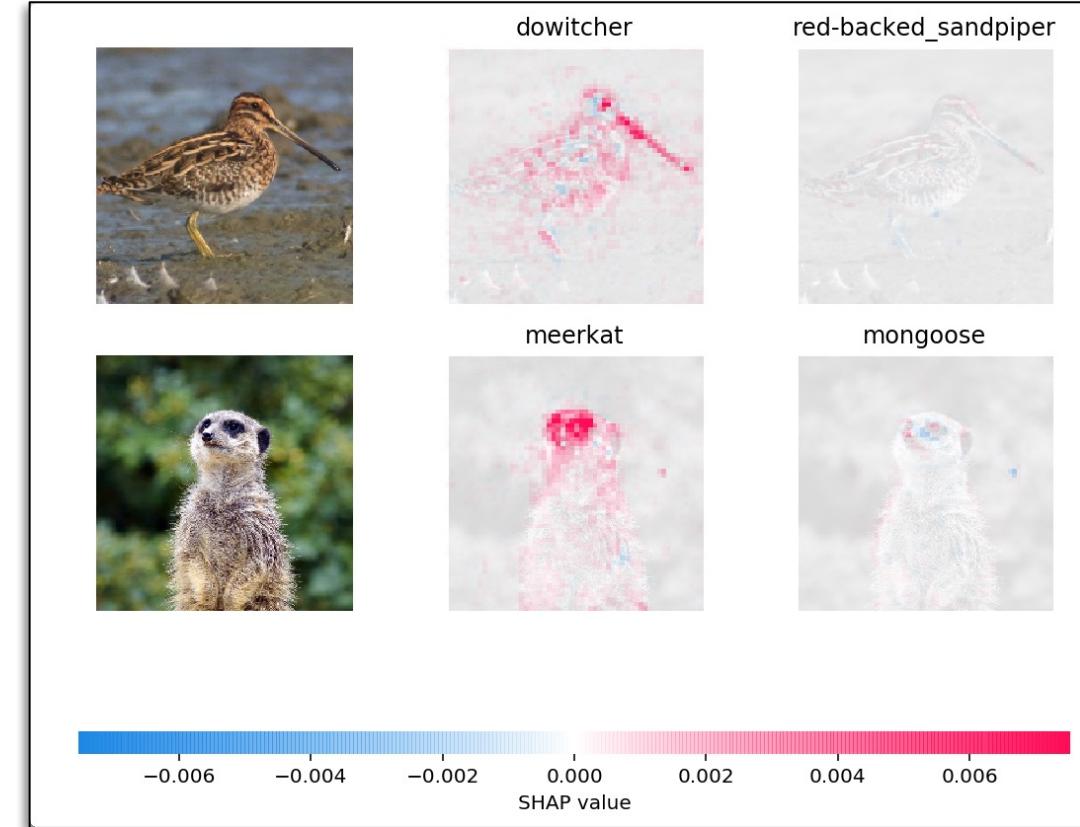
# Interpretation



# Interpretation



[\(link\)](#)

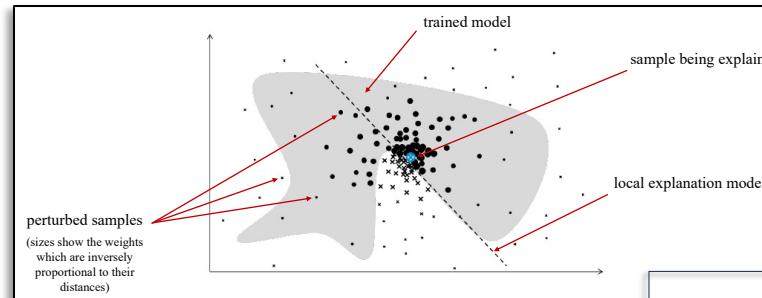


[\(link\)](#)

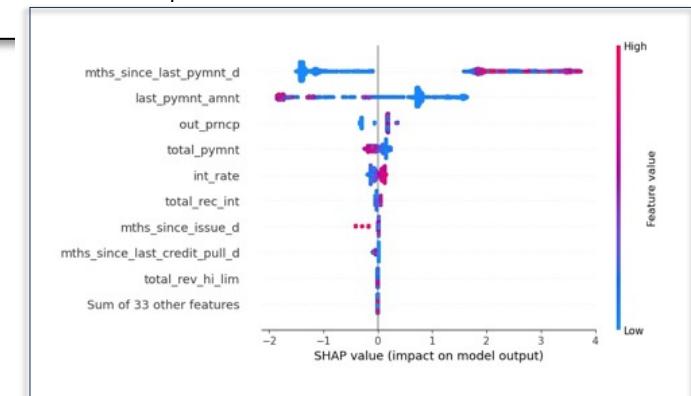
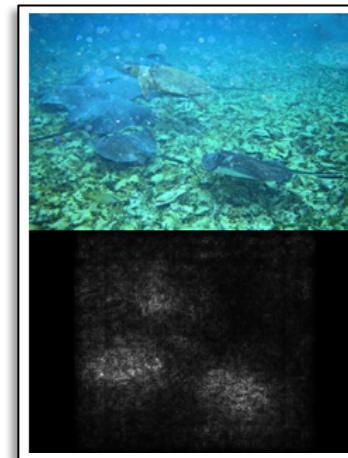


# Big Picture

## ■ LIME



## ■ SHAP



## ■ Gradient Methods

- Vanilla Gradients
- SmoothGrad
- Other Methods

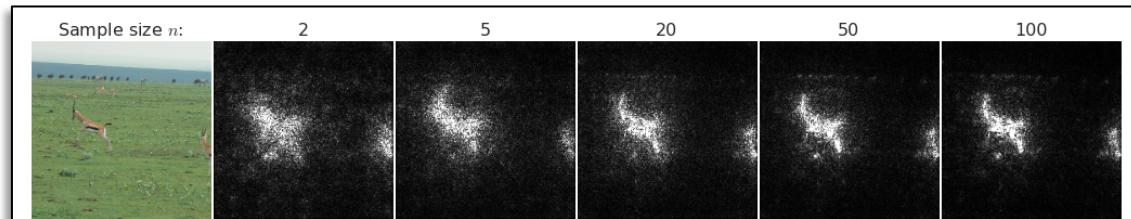


Figure 4. Effect of sample size on the estimated gradient for inception. 10% noise was applied to each image.



# Vanilla Gradient

“We consider two visualisation techniques, based on computing the gradient of the class score with respect to the input image [...] The second technique computes a class saliency map, specific to a given image and class.”

$$S_c(\mathbf{x}) \approx \mathbf{w}^\top \mathbf{x} + b$$

score of belonging to class  $c$

sample image

$$w_j = \frac{\partial S_c(\mathbf{x})}{\partial x_j}$$

## Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps

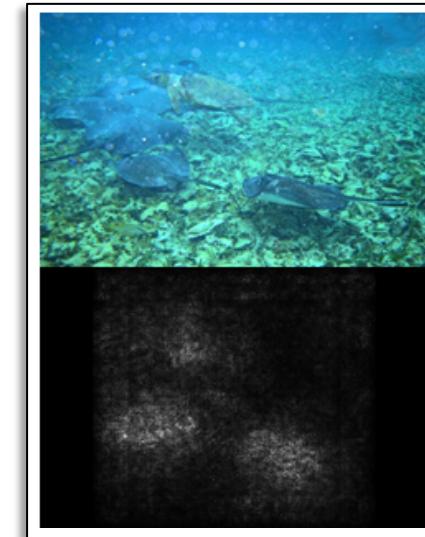
Karen Simonyan

Andrea Vedaldi

Andrew Zisserman

Visual Geometry Group, University of Oxford

{karen, vedaldi, az}@robots.ox.ac.uk



(link)



# SmoothGrad

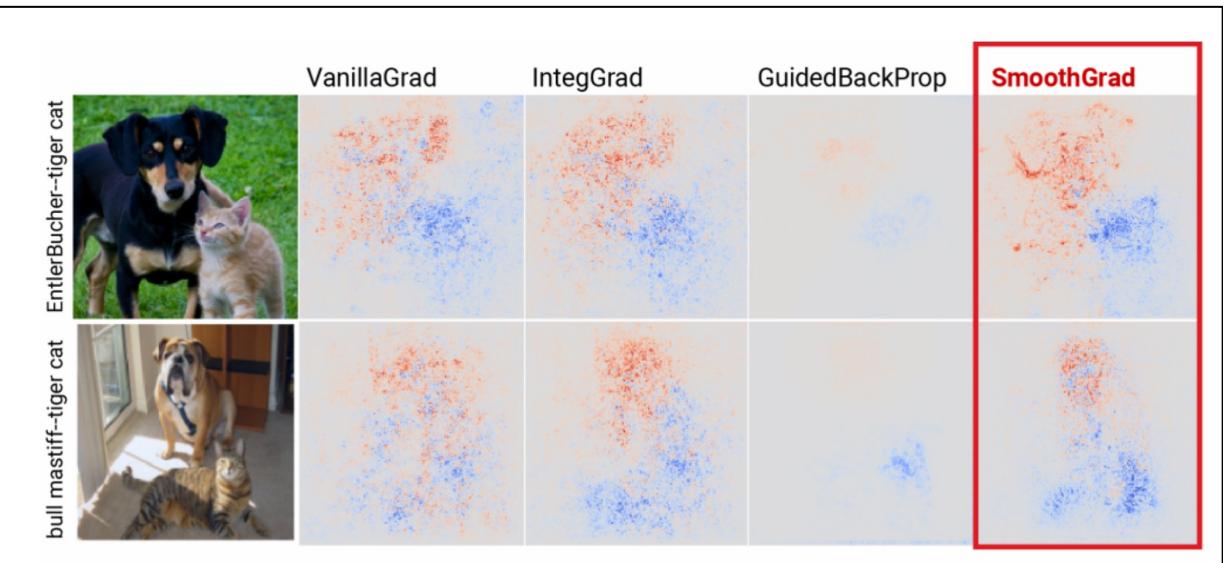
“This paper describes a very simple technique, SMOOTH-GRAD, that in practice tends to reduce visual noise, and also can be combined with other sensitivity map algorithms. The core idea is to take an image of interest, sample similar images by adding noise to the image, then take the average of the resulting sensitivity maps for each sampled image.”

“A note on terminology: although the terms ‘sensitivity map’, ‘saliency map’, and ‘pixel attribution map’ have been used in different contexts, in this paper, we will refer to these methods collectively as ‘sensitivity maps.’”

## SmoothGrad: removing noise by adding noise

Daniel Smilkov<sup>1</sup> Nikhil Thorat<sup>1</sup> Been Kim<sup>1</sup> Fernanda Viégas<sup>1</sup> Martin Wattenberg<sup>1</sup>

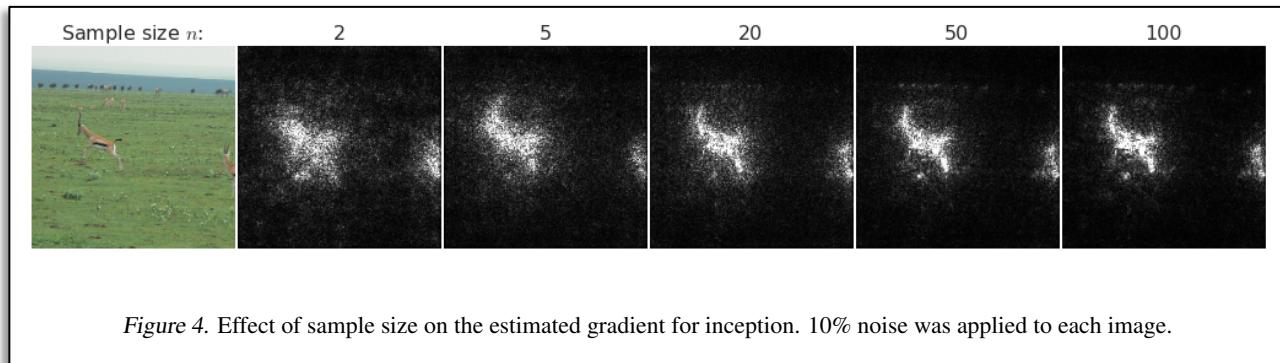
([link](#))



# SmoothGrad

$$M_c(\mathbf{x}) = \nabla S_c(\mathbf{x}) = \left[ \frac{\partial S_c(\mathbf{x})}{\partial x_j} \right]_{j \in \mathcal{F}}$$

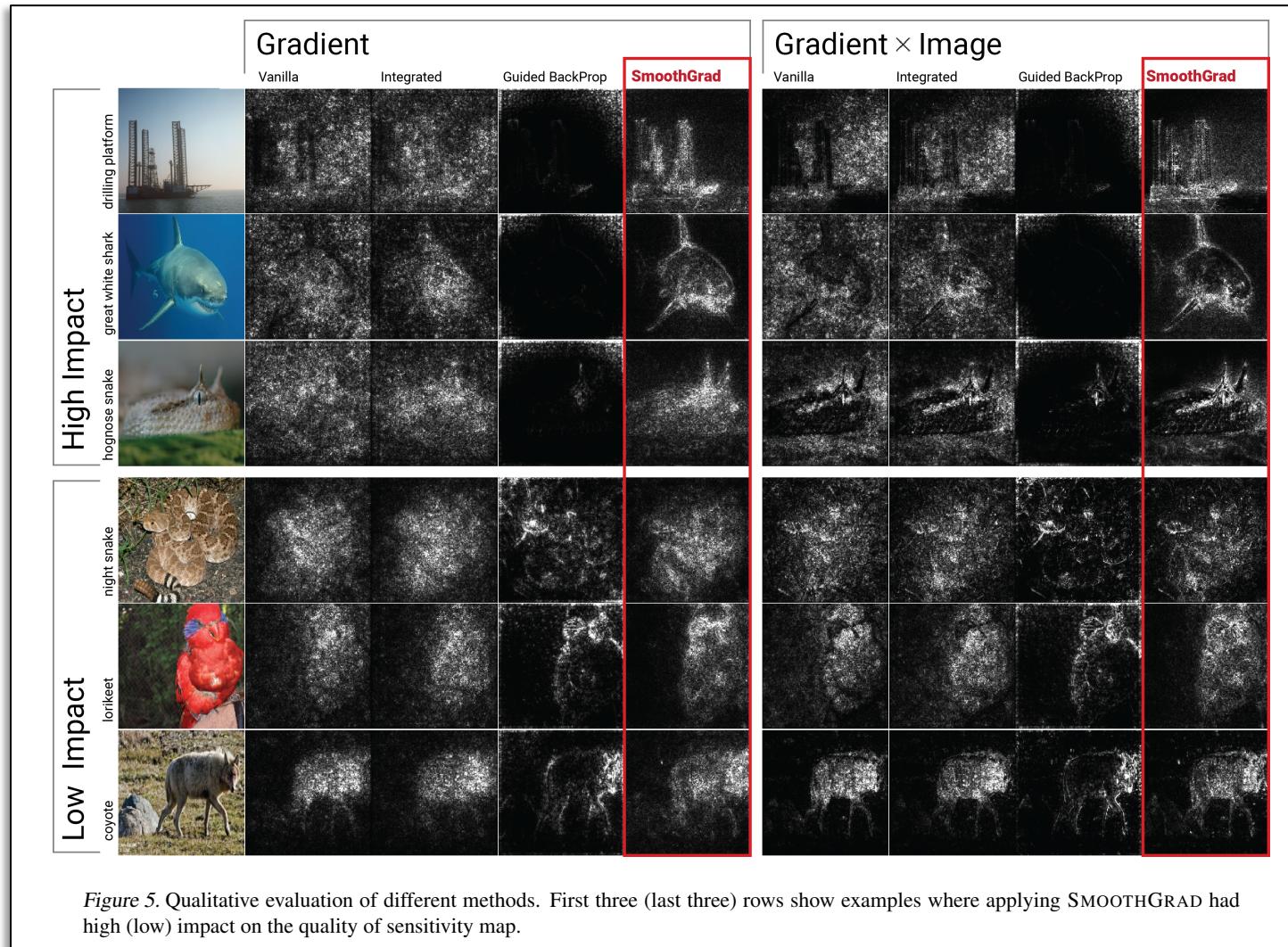
$$\hat{M}_c(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m M_c(\mathbf{x} + \boldsymbol{\varepsilon}_k) \quad \boldsymbol{\varepsilon}_k \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\sigma}^2), \quad k = 1, \dots, m$$



Smilkov, D., Thorat, N., Kim, B., Viégas, F. & Wattenberg, M. “[SmoothGrad: removing noise by adding noise](#),” arXiv:1706.03825, 2017.



# Interpretation



# Other Methods

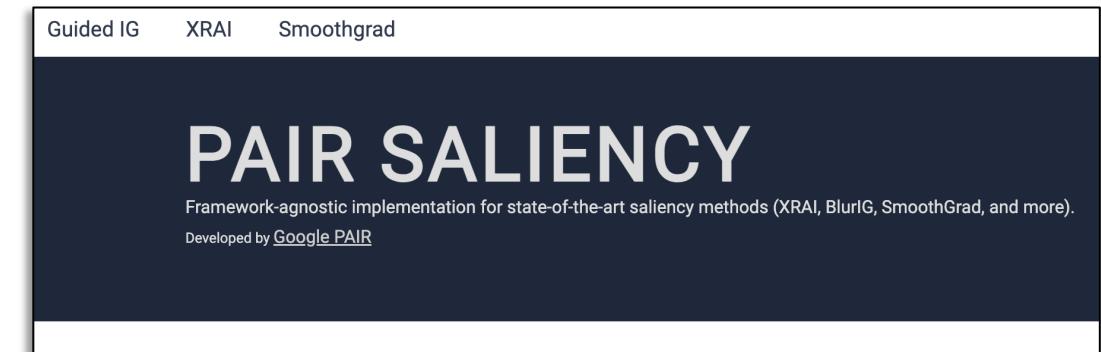
- (Guided) Integrated Gradients
- GradCAM
- FullGrad
- XRAI
- ...

## Introduction: Advanced Explainable AI for computer vision

`pip install grad-cam`

<https://github.com/jacobgil/pytorch-grad-cam>

([link](#))

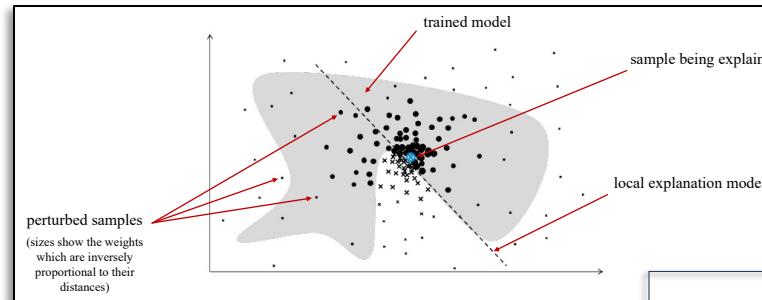


([link](#))

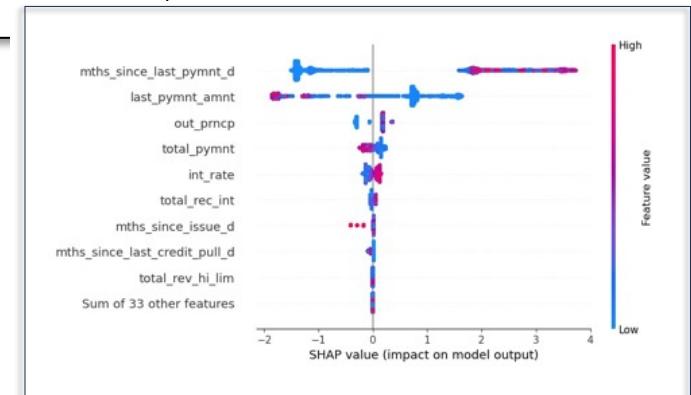
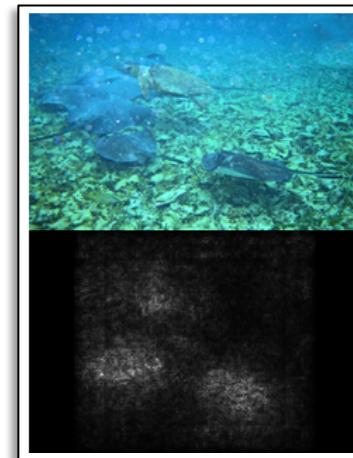


# Big Picture

## ■ LIME



## ■ SHAP



## ■ Gradient Methods

- Vanilla Gradients
- SmoothGrad
- Other Methods

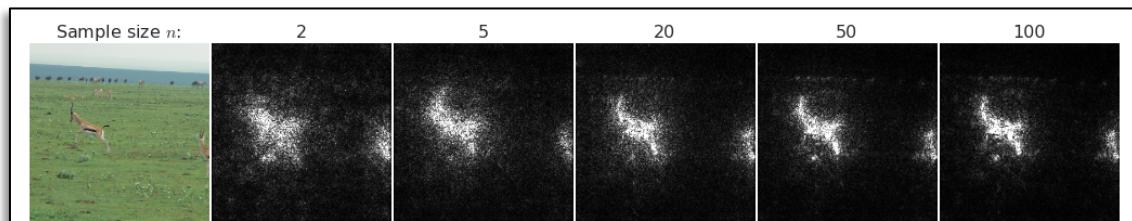


Figure 4. Effect of sample size on the estimated gradient for inception. 10% noise was applied to each image.



# Reading Material

## “Why Should I Trust You?” Explaining the Predictions of Any Classifier

Marco Tulio Ribeiro  
University of Washington  
Seattle, WA 98105, USA  
marcotcr@cs.uw.edu

Sameer Singh  
University of Washington  
Seattle, WA 98105, USA  
sameer@cs.uw.edu

Carlos Guestrin  
University of Washington  
Seattle, WA 98105, USA  
guestrin@cs.uw.edu

## A Unified Approach to Interpreting Model Predictions

Scott M. Lundberg  
Paul G. Allen School of Computer Science  
University of Washington  
Seattle, WA 98105  
slund1@cs.washington.edu

Su-In Lee  
Paul G. Allen School of Computer Science  
Department of Genome Sciences  
University of Washington  
Seattle, WA 98105  
suinlee@cs.washington.edu

## Axiomatic Attribution for Deep Networks

Mukund Sundararajan \*<sup>1</sup> Ankur Taly \*<sup>1</sup> Qiqi Yan \*<sup>1</sup>

## Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps

Karen Simonyan      Andrea Vedaldi      Andrew Zisserman  
Visual Geometry Group, University of Oxford  
{karen, vedaldi, az}@robots.ox.ac.uk

## SmoothGrad: removing noise by adding noise

Daniel Smilkov<sup>1</sup> Nikhil Thorat<sup>1</sup> Been Kim<sup>1</sup> Fernanda Viégas<sup>1</sup> Martin Wattenberg<sup>1</sup>

## Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization

Ramprasaath R. Selvaraju<sup>1\*</sup> Michael Cogswell<sup>1</sup> Abhishek Das<sup>1</sup> Ramakrishna Vedantam<sup>1\*</sup>  
Devi Parikh<sup>1,2</sup> Dhruv Batra<sup>1,2</sup>  
<sup>1</sup>Georgia Institute of Technology    <sup>2</sup>Facebook AI Research  
{ramprs, cogswell, abhshkdz, vrama, parikh, dbatra}@gatech.edu

