

6.874 Spring 2021

Lecture 05

Interpretable Deep Learning

Prof. Manolis Kellis



Slides by Beomsu Kim, KAIST

Interpretable Deep Learning

1. Intro to Interpretability

1a. Interpretability definition: Convert implicit NN information to human-interpretable information

1b. Motivation: Verify model works as intended; debug classifier; make discoveries; Right to explanation

1c. Ante-hoc (train interpretable model) vs. **Post-hoc** (interpret complex model; degree of “locality”)

2. Interpreting Deep Neural Networks

2a. Interpreting Models (macroscopic, understand internals) vs. **decisions** (microscopic, practical applications)

2b. Interpreting Models: Weight visualization, Surrogate model, Activation maximization, Example-based

2c. Interpreting Decisions:

- Example-based
- Attribution Methods: why are gradients noisy?
- Gradient-based Attribution: SmoothGrad, Interior Gradient
- Backprop-based Attribution: Deconvolution, Guided Backpropagation

3. Evaluating Attribution Methods

3a. Qualitative: Coherence: Attributions should highlight discriminative features / objects of interest

3b. Qualitative: Class Sensitivity: Attributions should be sensitive to class labels

3c. Quantitative: Sensitivity: Removing feature with high attribution → large decrease in class probability

3d. Quantitative: ROAR & KAR. Low class prob cuz image unseen → remove pixels, retrain, measure acc. drop

What is Interpretability?

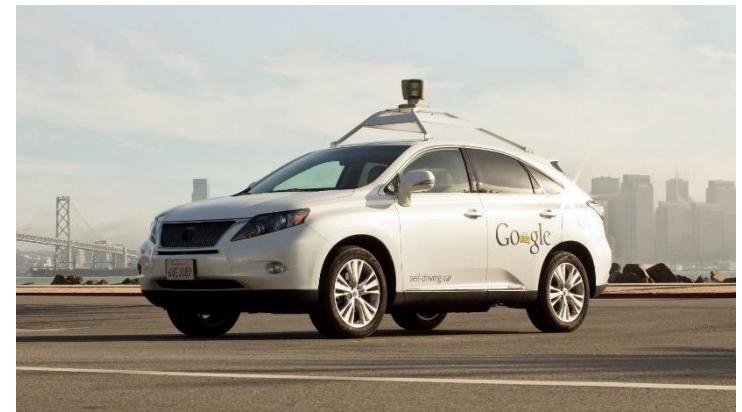
AlphaGo vs. Lee Sedol



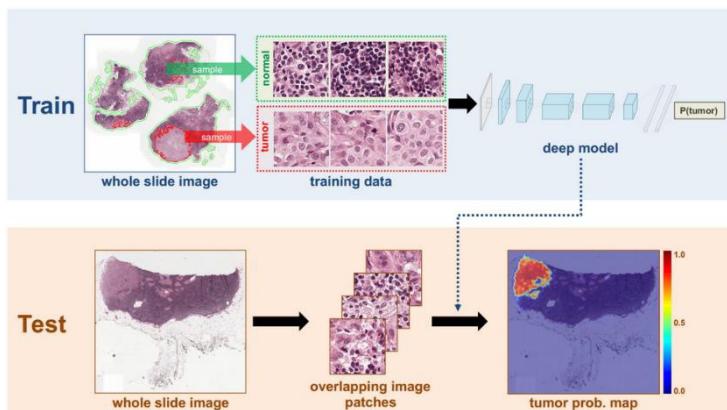
ImageNet Challenge



Self-driving Cars



Disease Diagnosis



Neural Machine Translation



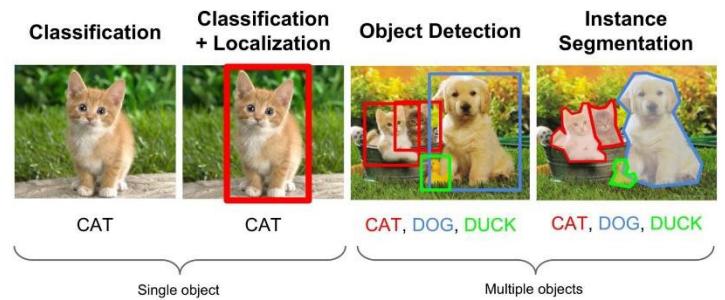
& More to Come!

What is Interpretability?

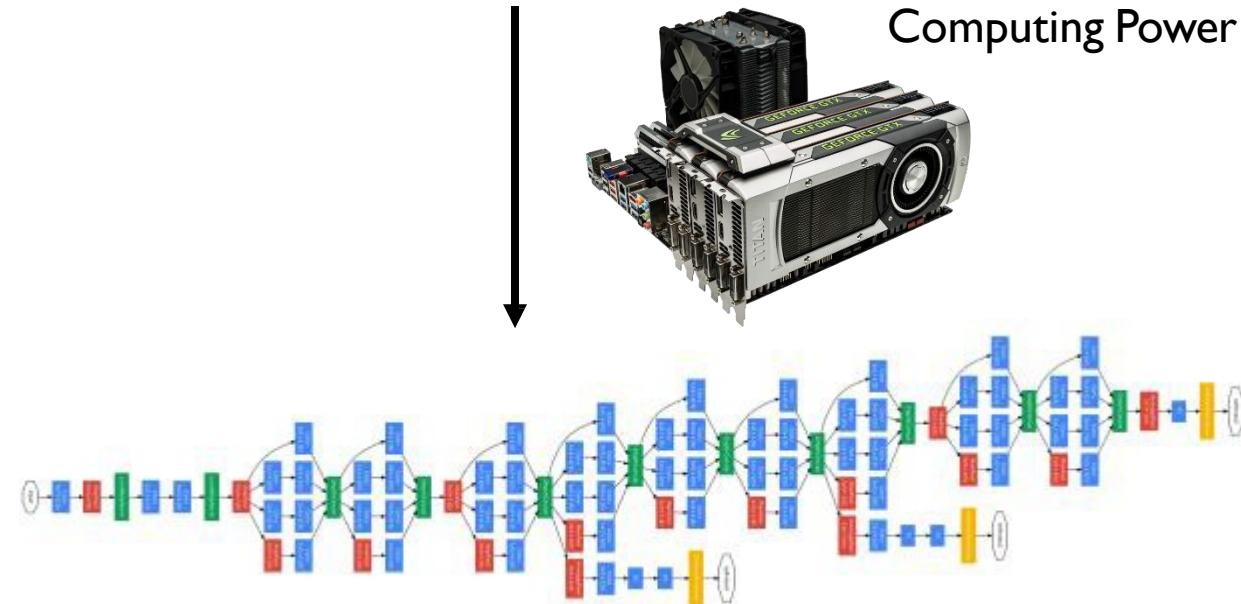
Large Dataset



Task Solving



Computing Power



Deep Neural Networks

Implicit Information

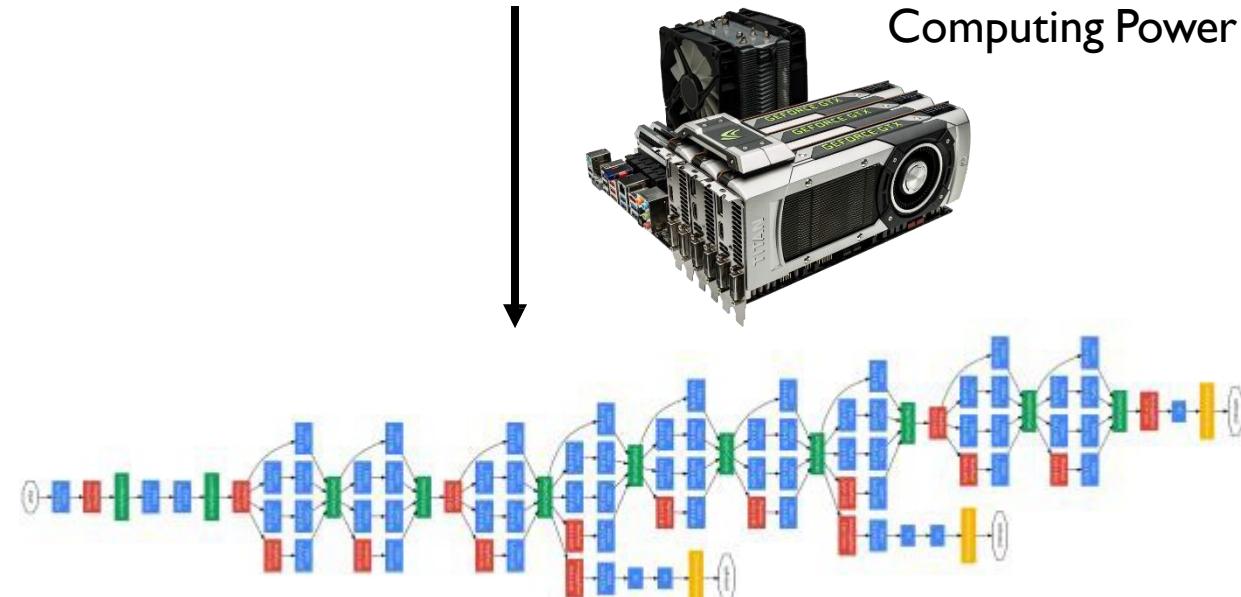
What is Interpretability?

Large Dataset



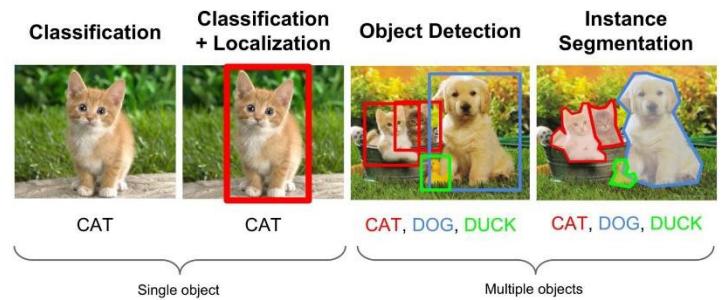
Interpretable Information

Computing Power



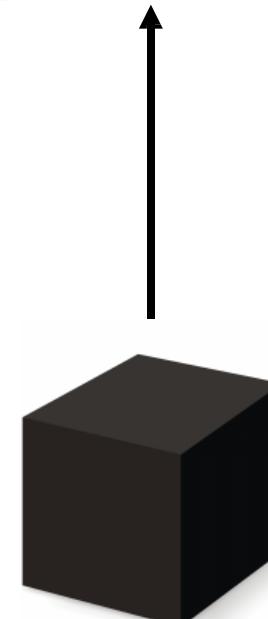
Deep Neural Networks

Task Solving



Single object

Multiple objects



Implicit Information

Why Interpretability?

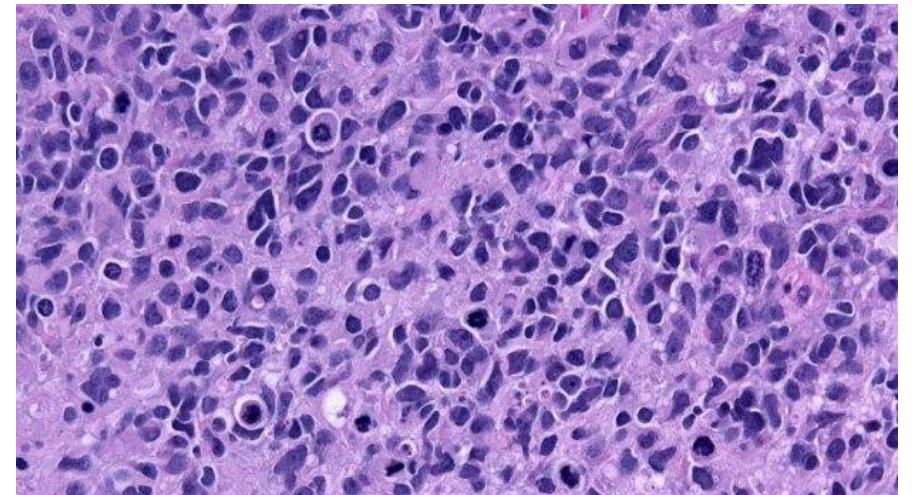
I. Verify that model works as expected

Wrong decisions can be costly and dangerous

Self-driving Uber kills Arizona woman in first fatal crash involving pedestrian

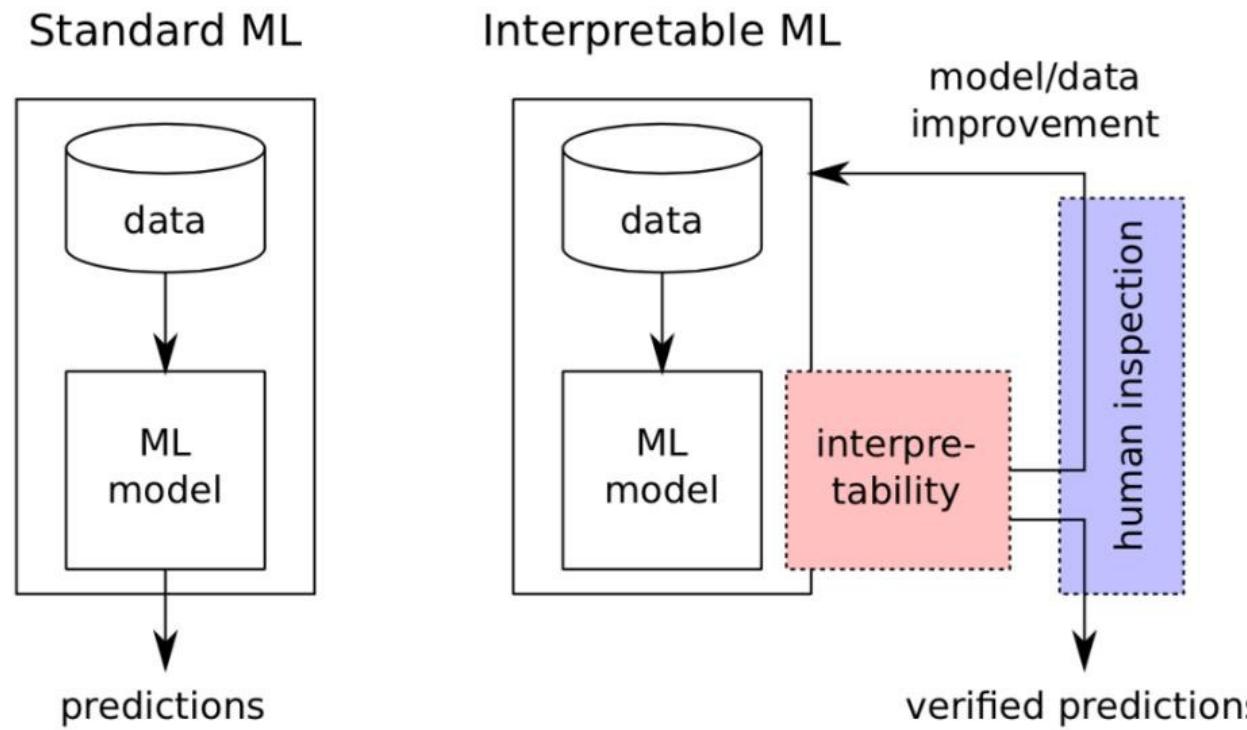


Disease Misclassification



Why Interpretability?

2. Improve / Debug classifier



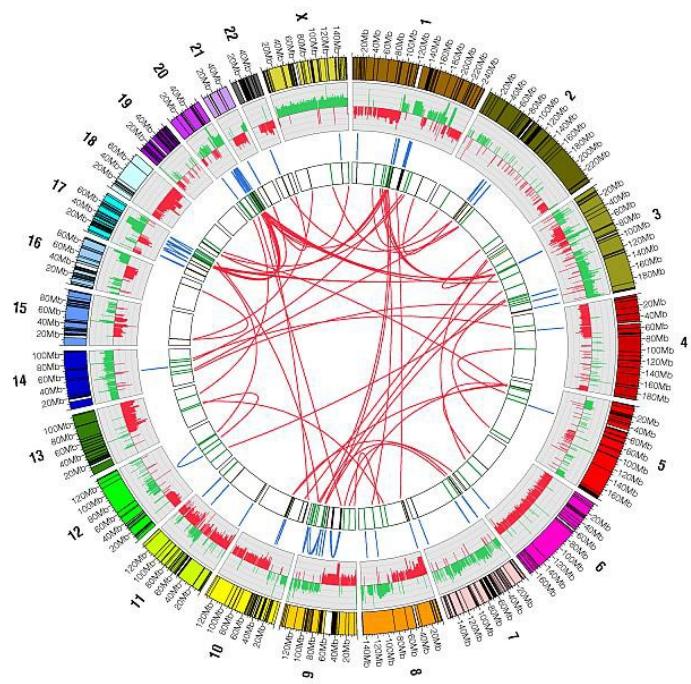
Generalization error

Generalization error + human experience

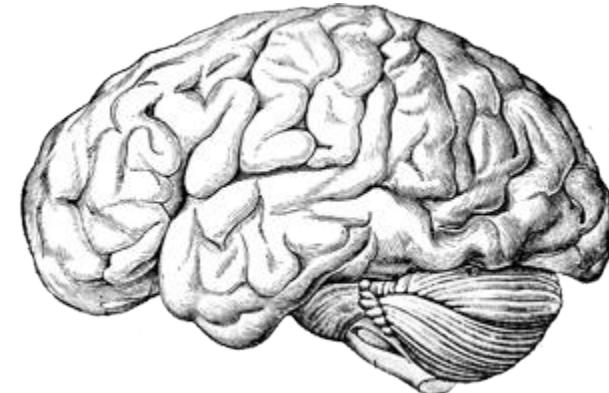
Why Interpretability?

3. Make new discoveries

Learn about the physical / biological / chemical mechanisms



Learn about the human brain



Why Interpretability?

4. Right to explanation

“Right to be given an explanation for an output of the algorithm”

Ex.

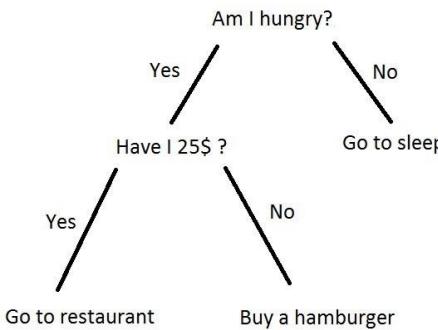
- US Equal Credit Opportunity Act
- The European Union General Data Protection Regulation
- France Digital Republic Act

Types of Interpretability in ML

Ante-hoc Interpretability

Choose an **interpretable model** and train it.

Ex.



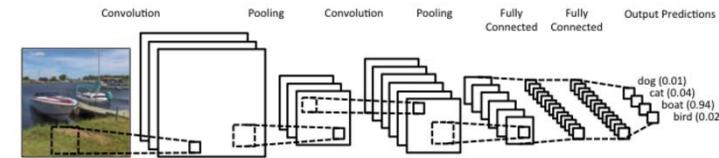
Decision Tree

Problem. Is the model expressive enough to predict the data?

Post-hoc Interpretability

Choose a **complex model** and develop a special technique to interpret it.

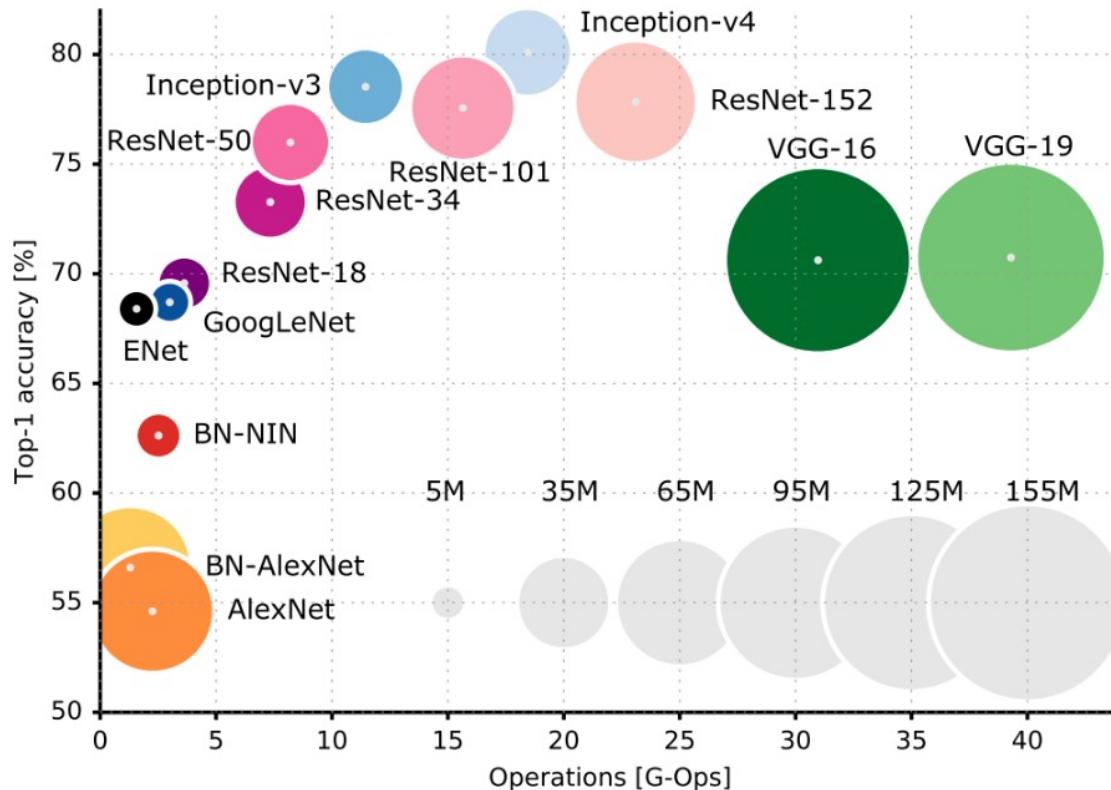
Ex.



Deep Neural Networks

Problem. How to interpret millions of parameters?

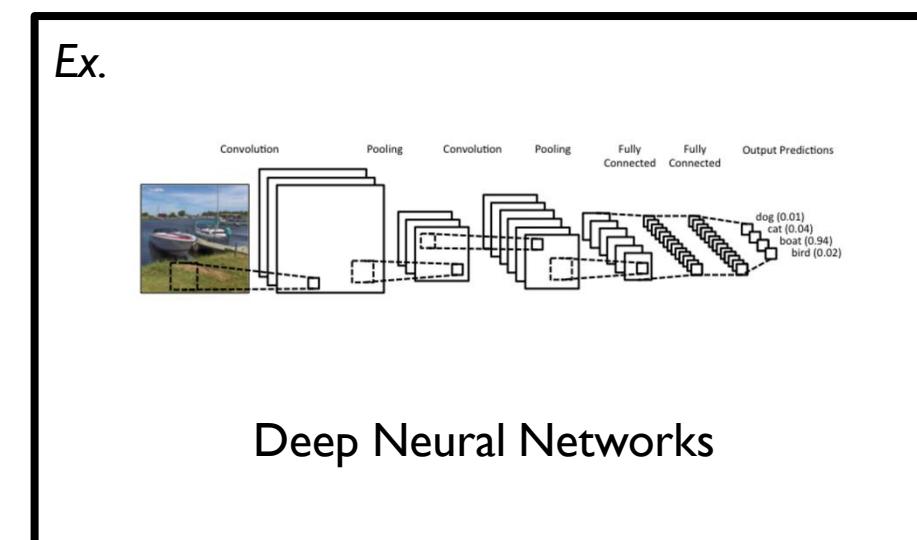
Types of Interpretability in ML



At least **5 million** parameters!

Post-hoc Interpretability

Choose a complex model and develop a special technique to interpret it.



Problem. How to interpret millions of parameters?

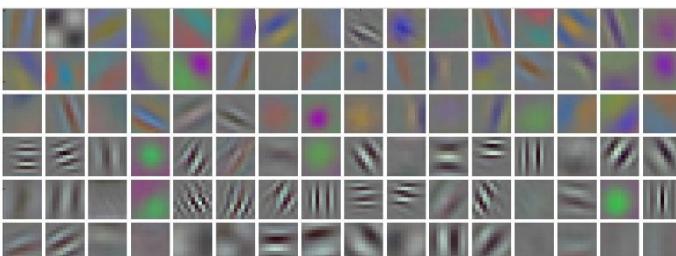
Types of Post-hoc Interpretability

Post-hoc interpretability techniques
can be classified by degree of “locality”

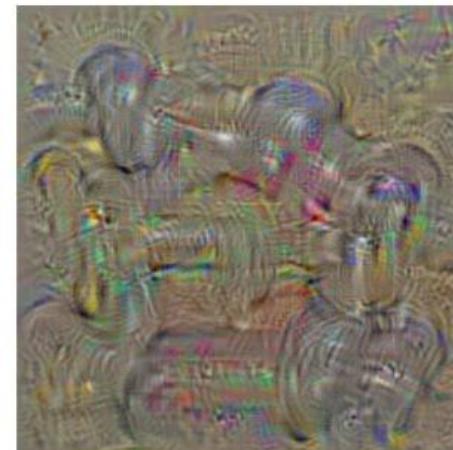
Model



What representations have
the DNN learned?



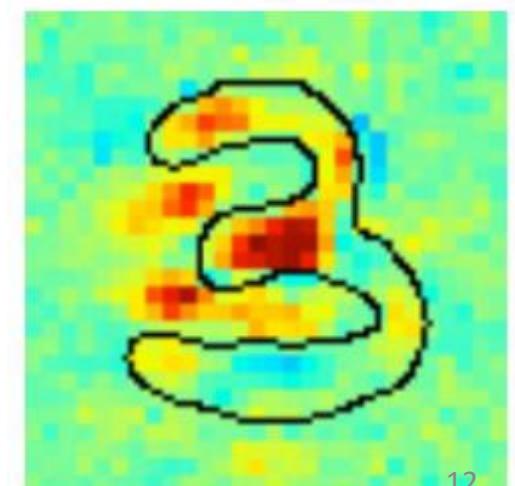
What pattern / image maximally
activates a particular neuron?



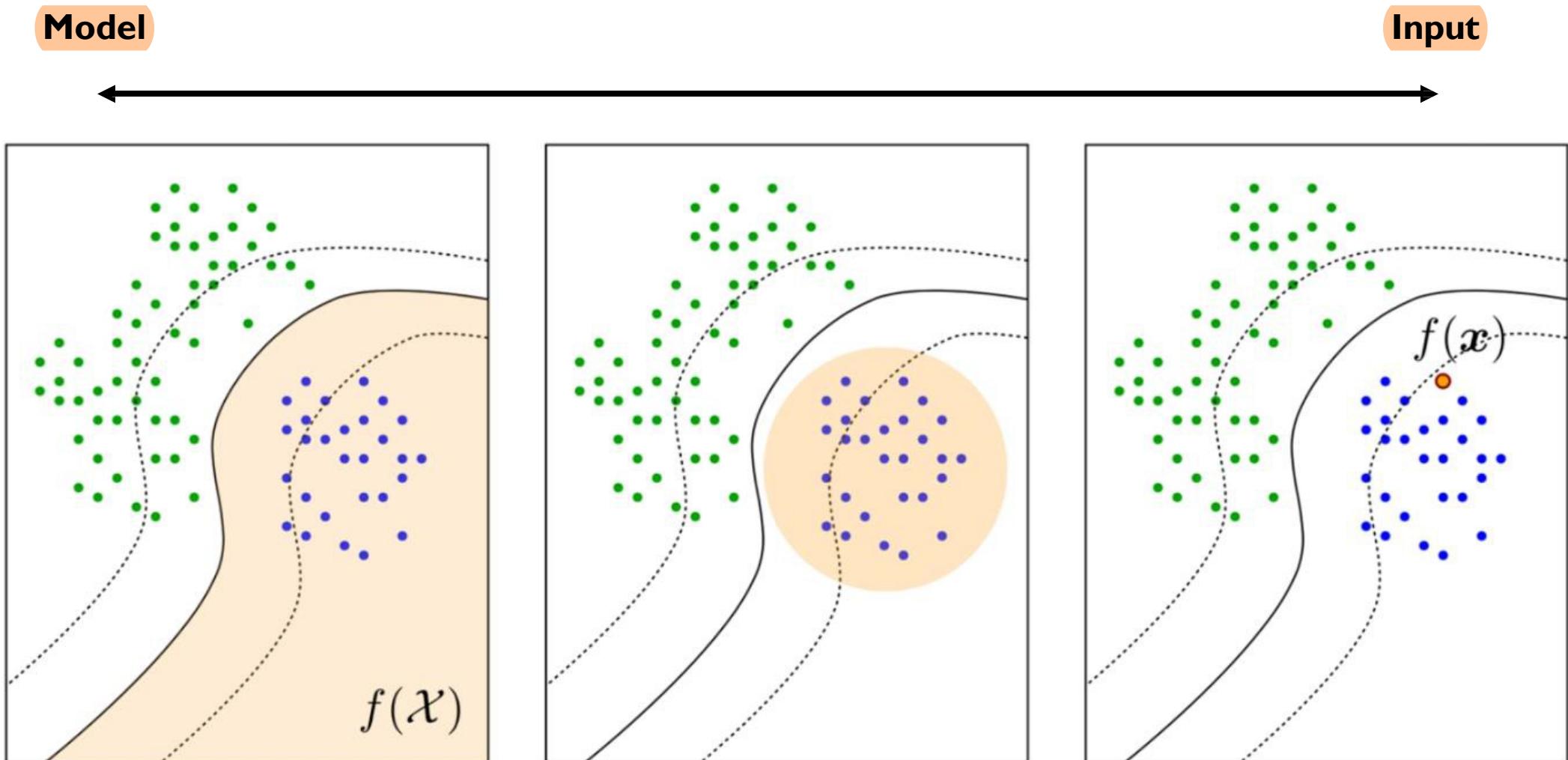
dumbbell

Input

Explain why input x has
been classified as $f(x)$.



Types of Post-hoc Interpretability



Interpretable Deep Learning

1. Intro to Interpretability

- 1a. **Interpretability definition:** Convert implicit NN information to human-interpretable information
- 1b. **Motivation:** Verify model works as intended; debug classifier; make discoveries; Right to explanation
- 1c. **Ante-hoc** (train interpretable model) vs. **Post-hoc** (interpret complex model; degree of “locality”)

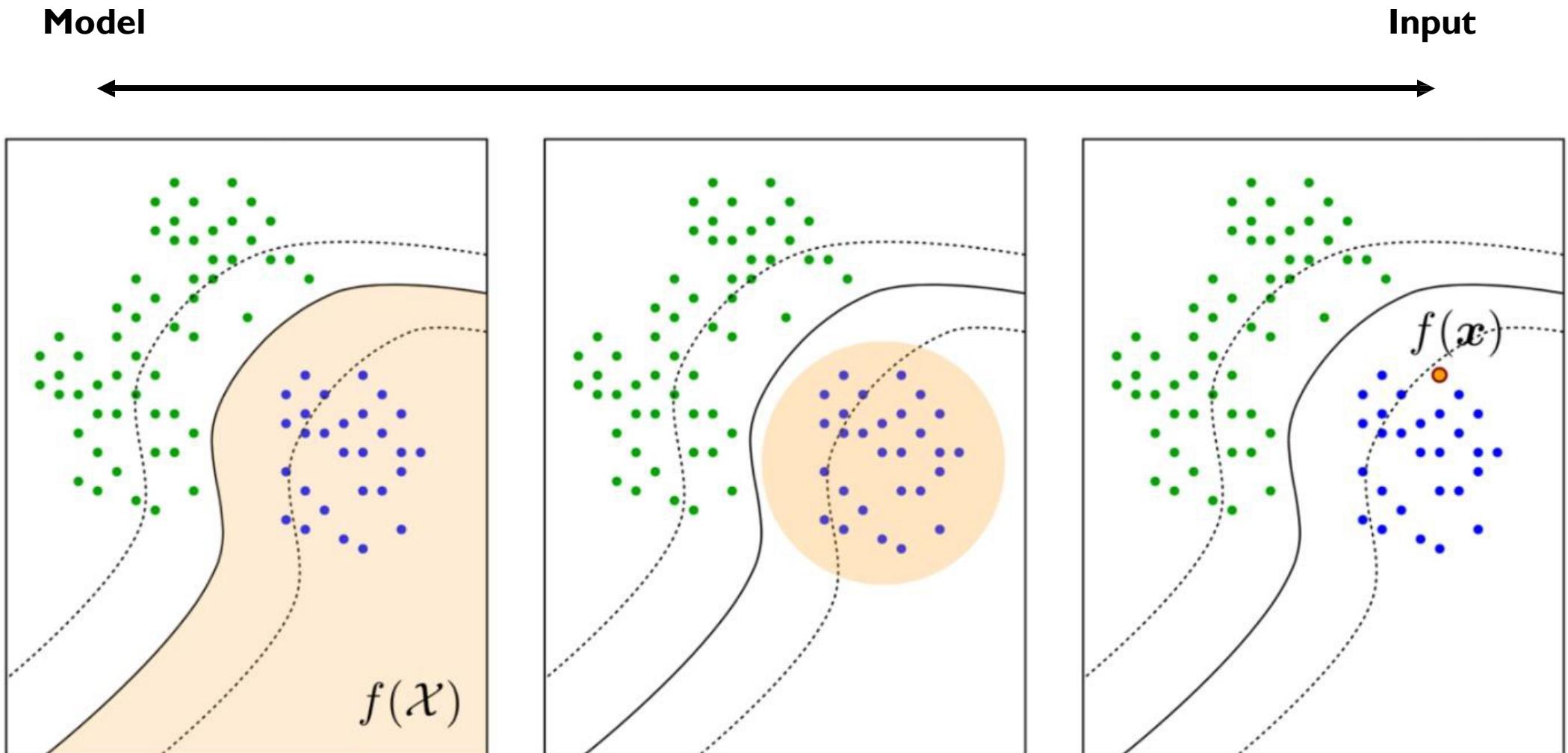
2. Interpreting Deep Neural Networks

- 2a. **Interpreting Models** (macroscopic, understand internals) vs. **decisions** (microscopic, practical applications)
- 2b. **Interpreting Models:** Weight visualization, Surrogate model, Activation maximization, Example-based
- 2c. **Interpreting Decisions:**
 - Example-based
 - Attribution Methods: why are gradients noisy?
 - Gradient-based Attribution: SmoothGrad, Interior Gradient
 - Backprop-based Attribution: Deconvolution, Guided Backpropagation

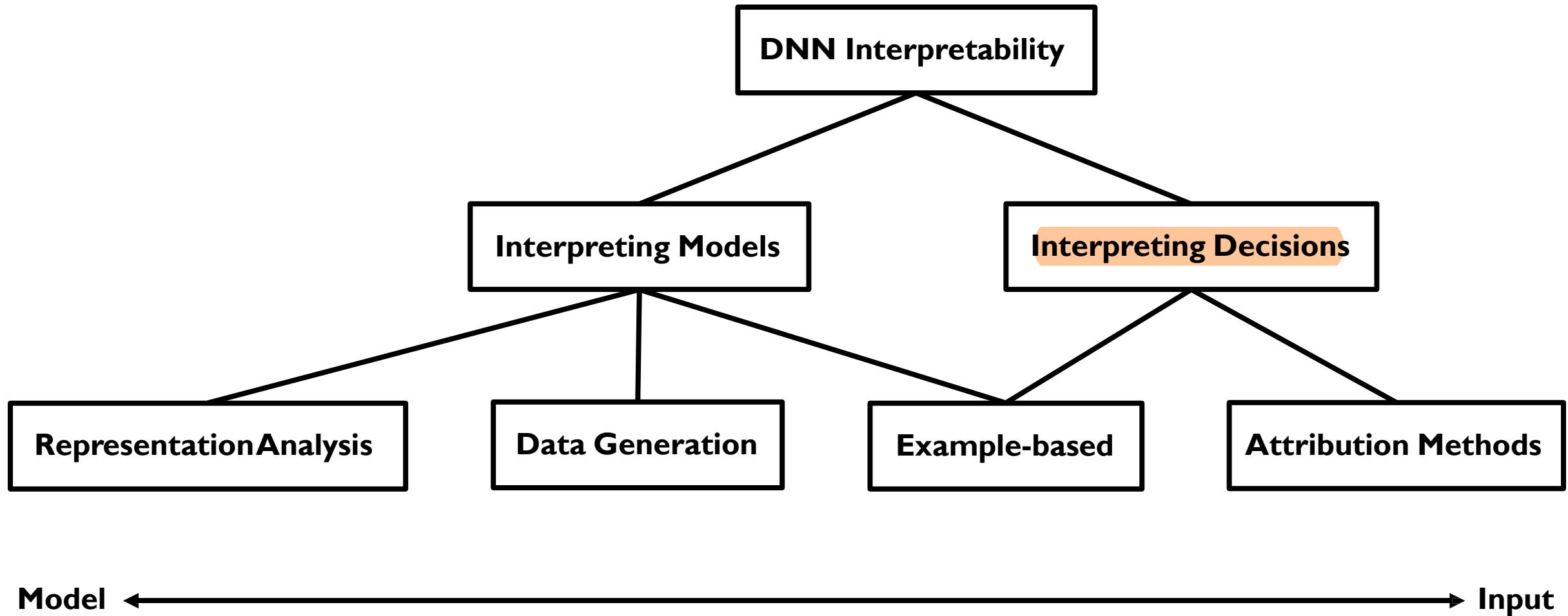
3. Evaluating Attribution Methods

- 3a. **Qualitative: Coherence:** Attributions should highlight discriminative features / objects of interest
- 3b. **Qualitative: Class Sensitivity:** Attributions should be sensitive to class labels
- 3c. **Quantitative: Sensitivity:** Removing feature with high attribution → large decrease in class probability
- 3d. **Quantitative: ROAR & KAR.** Low class prob cuz image unseen → remove pixels, retrain, measure acc. drop

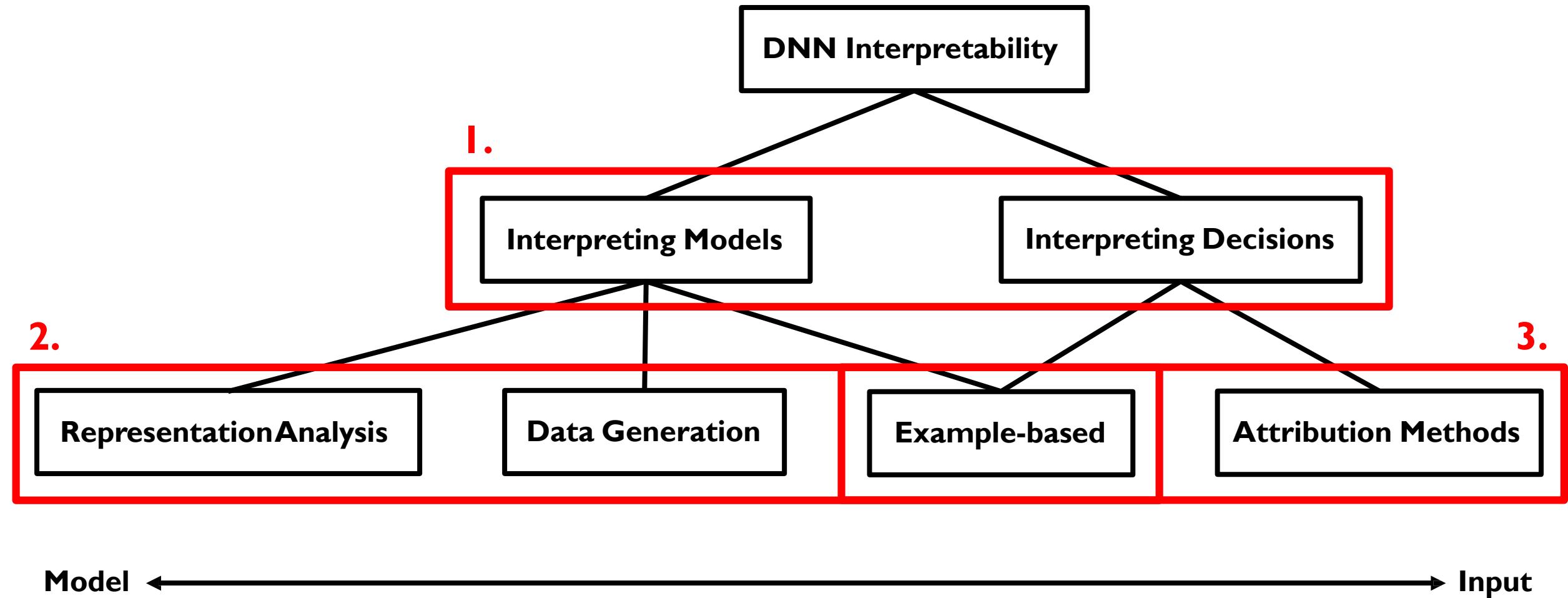
Types of Post-hoc Interpretability



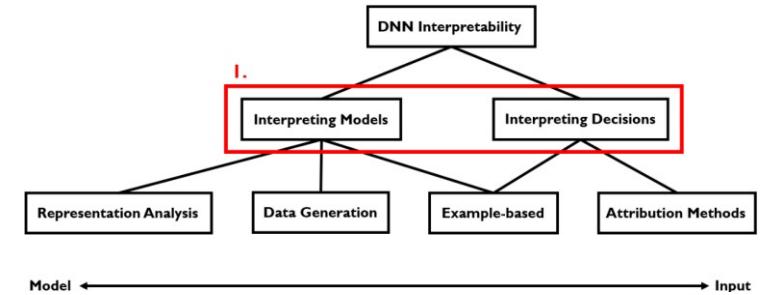
Types of DNN Interpretability



Types of DNN Interpretability



Types of DNN Interpretability



Interpreting Models (Macroscopic)

- “Summarize” DNN with a simpler model (e.g. decision tree)
- Find prototypical example of a category
- Find pattern maximizing activation of a neuron

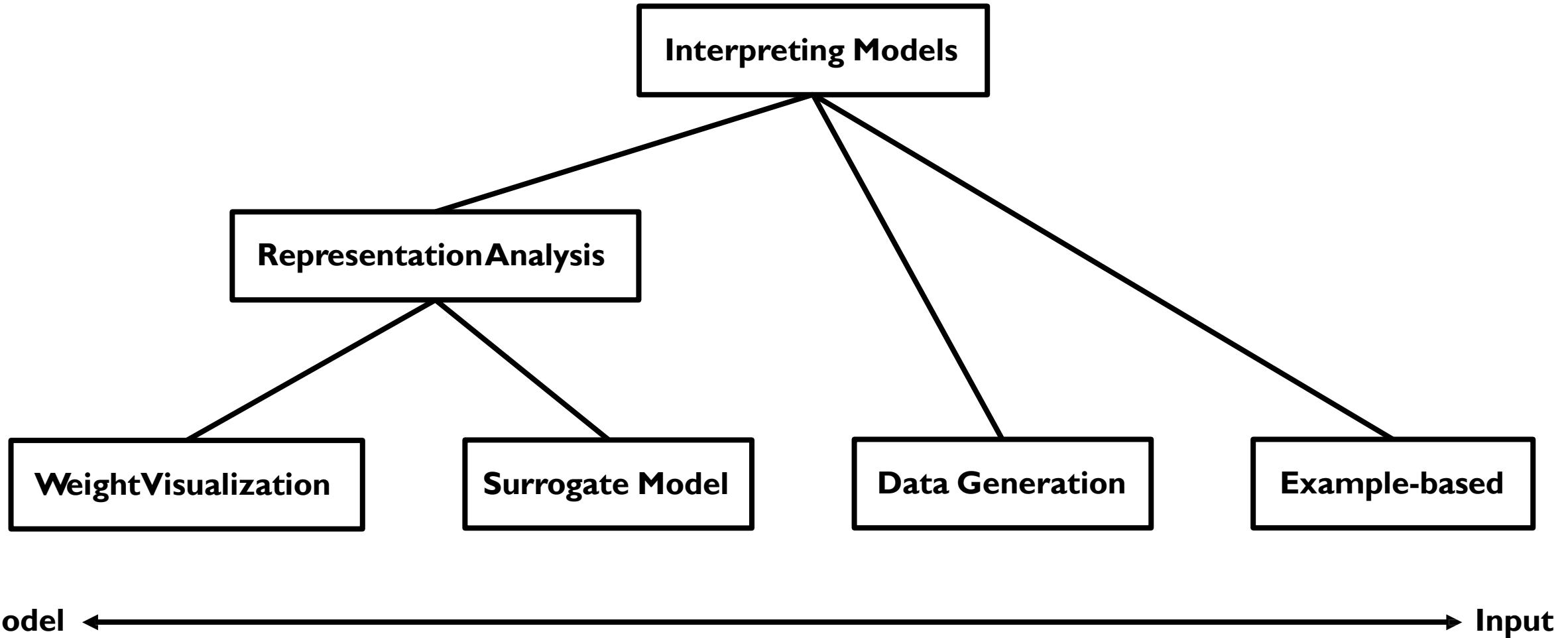
Better understand internal representations

Interpreting Decisions (Microscopic)

- Why did DNN make this decision
- Verify that model behaves as expected
- Find evidence for decision

Important for practical applications

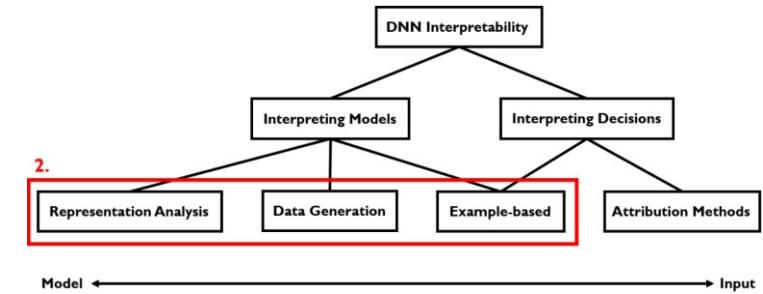
Types of DNN Interpretability



2b – Interpreting models:

- (i) Representation analysis: Weight Visualization
- (ii) Representation analysis: Surrogate Model
- (iii) Data Generation / Activation Maximization
- (iv) Example based

Types of DNN Interpretability



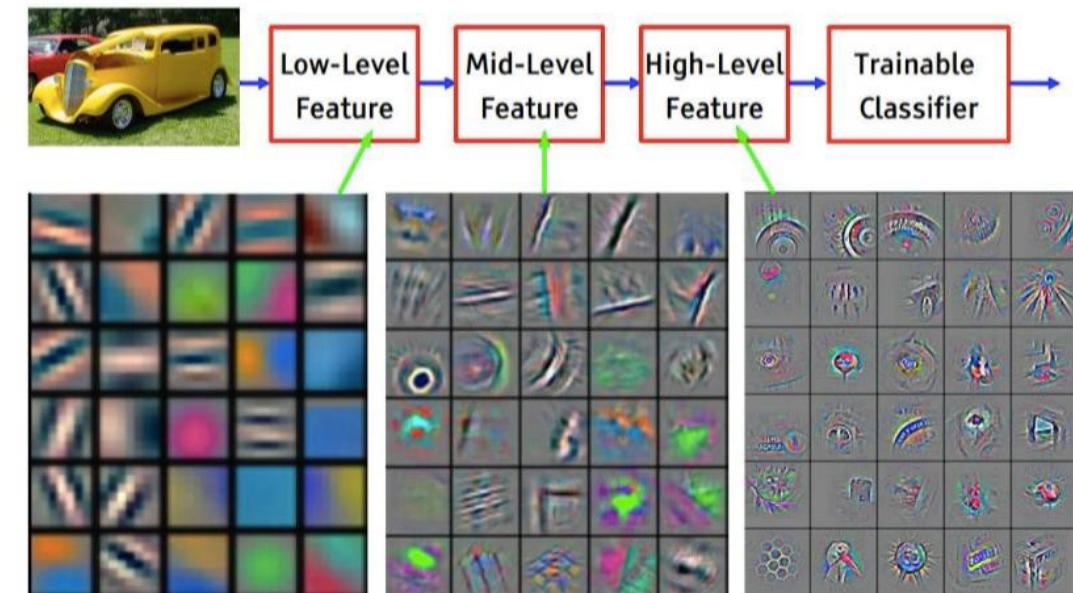
Weight Visualization

- Filter visualization in Convolutional Neural Networks
- Can understand what kind of features CNN has learned
- Still too many filters!

Surrogate Model

Data Generation

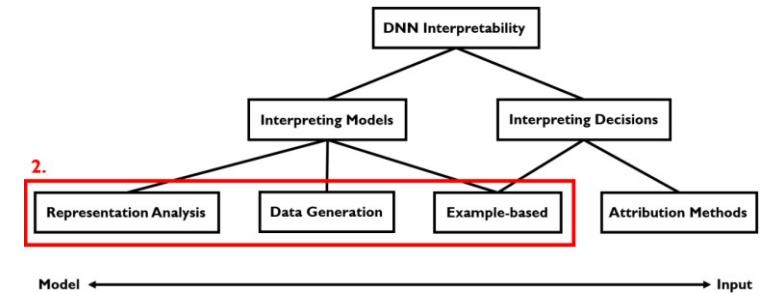
Example-based



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

2b – Interpreting models:

- (i) Weight Visualization
- (ii) Surrogate Model
- (iii) Data Generation / Activation Maximization
- (iv) Example based



Types of DNN Interpretability

Weight Visualization

- “Summarize” DNN with a simpler model
- E.g. Decision trees, graphs or linear models
(aka. Meta-model, approximation model, response surface model, emulator...)

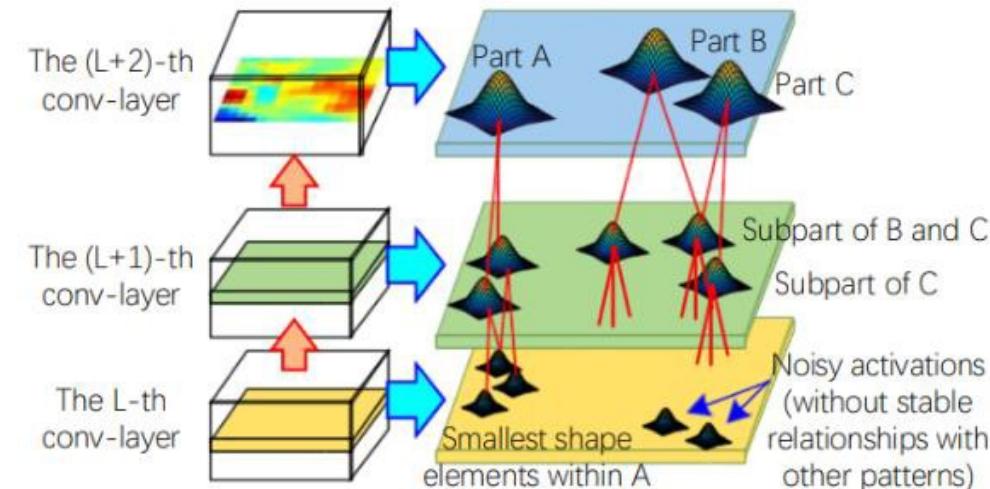
Idea: Train an *Interpretable* Machine Learning model on the Outputs of our “Black Box” model with the specific goal of interpreting it. Not exact, but close and interpretable.

Model agnostic. Approximate the predictions, not the actual real world.

Surrogate Model

Data Generation

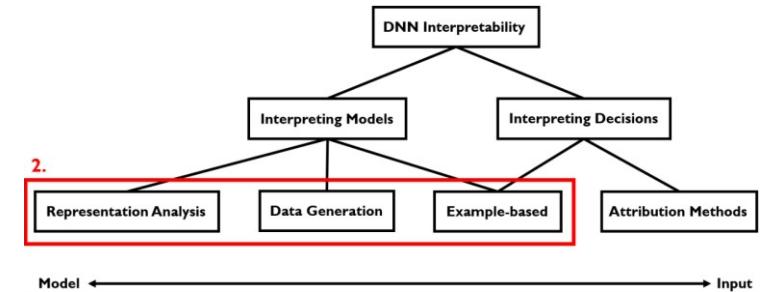
Example-based



2b – Interpreting models:

- (i) Weight Visualization
- (ii) Surrogate Model
- (iii) Data Generation / Activation Maximization
- (iv) Example based

Types of DNN Interpretability



Weight Visualization

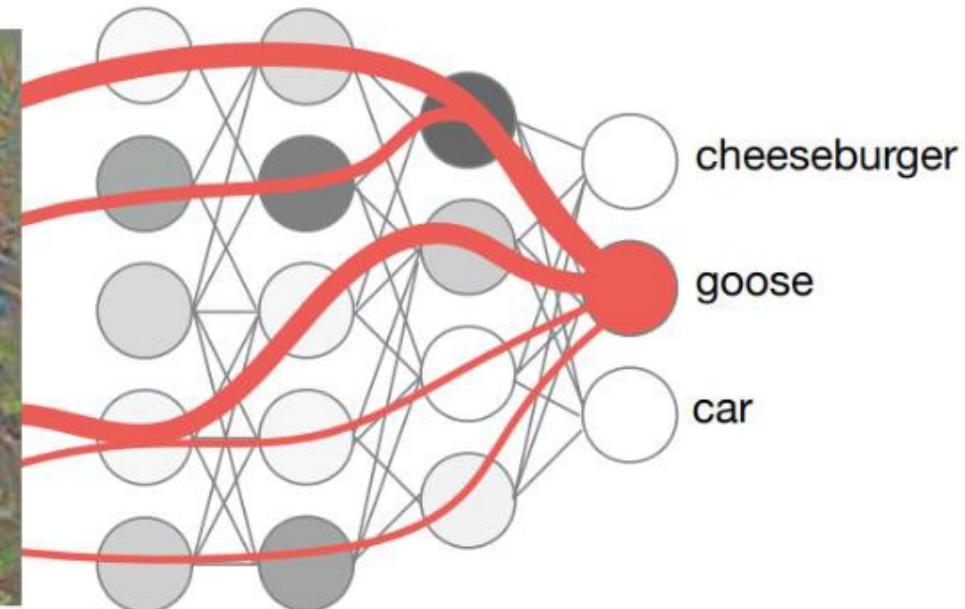
Surrogate Model

Data Generation

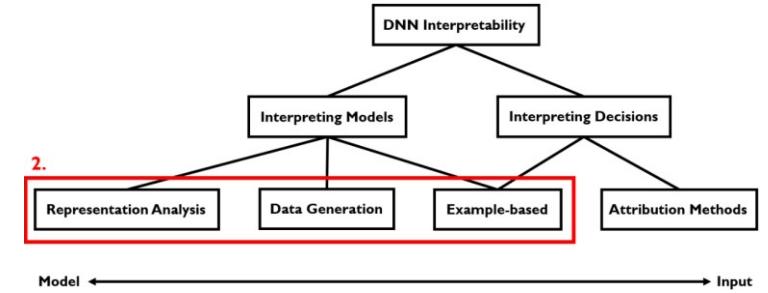
Example-based

Activation Maximization

- Find pattern maximizing activation of a neuron



Types of DNN Interpretability



Weight Visualization

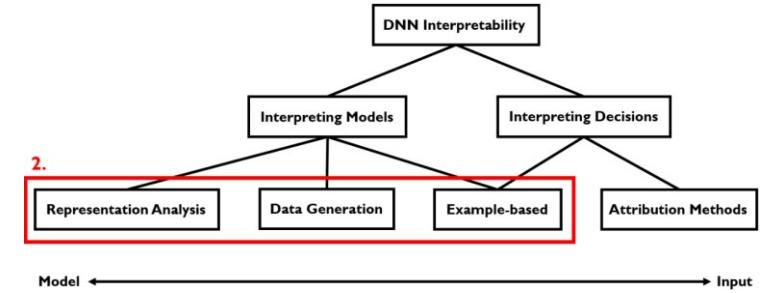
Surrogate Model

Data Generation

Example-based

$$\max_{x \in \mathcal{X}} p_{\theta}(\omega_c | x) + \lambda \Omega(x)$$

Class Probability Regularization Term



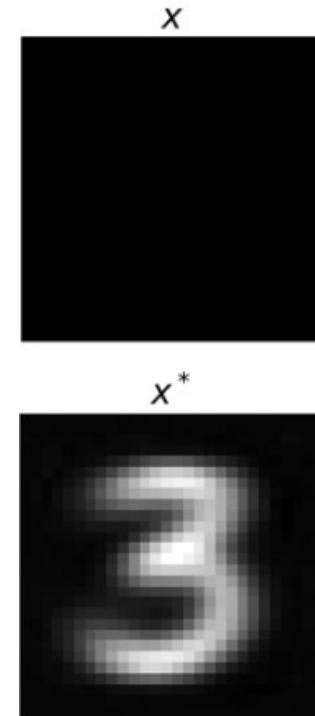
Types of DNN Interpretability

Weight Visualization

Surrogate Model

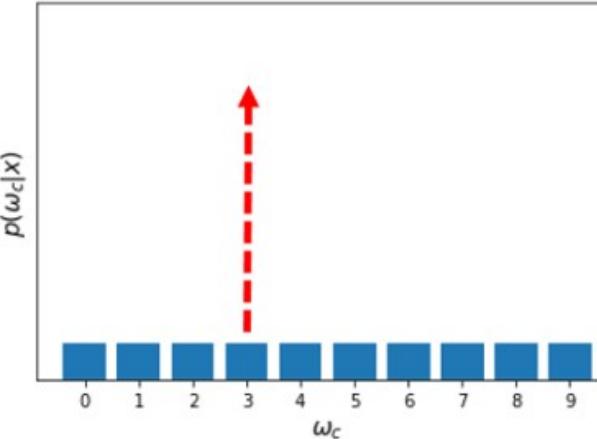
Data Generation

Example-based



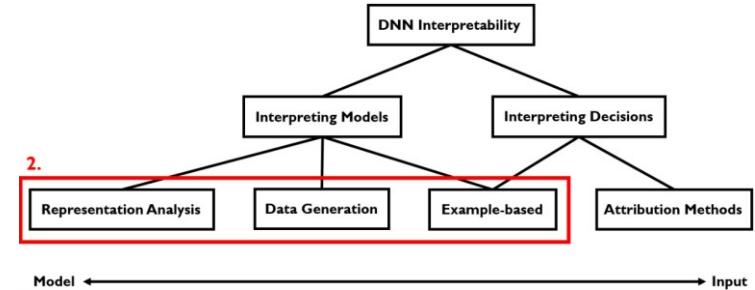
DNN
 $p(\omega_c | x)$

$$\max_{x \in \mathcal{X}} p_\theta(\omega_c | x)$$



Find x that maximizes class posterior probability

Types of DNN Interpretability

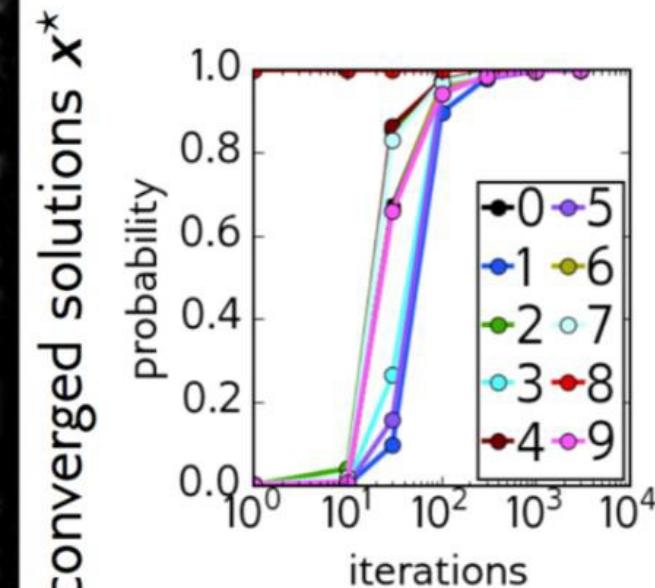
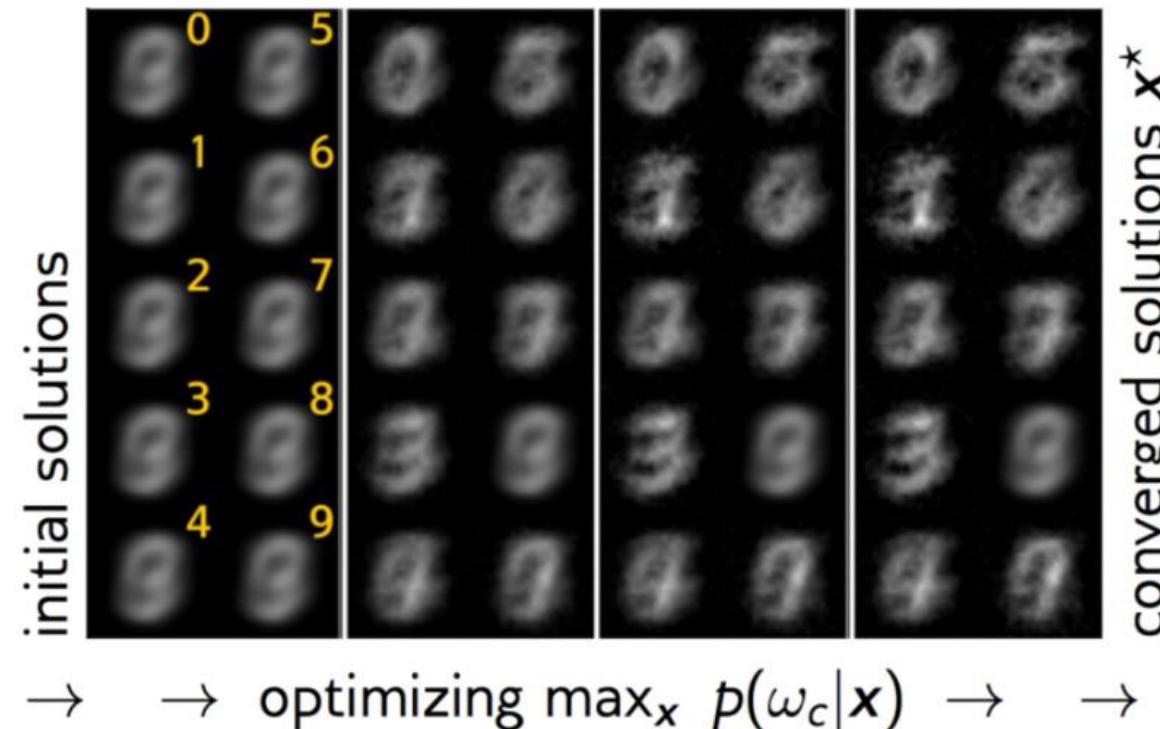


Weight Visualization

Surrogate Model

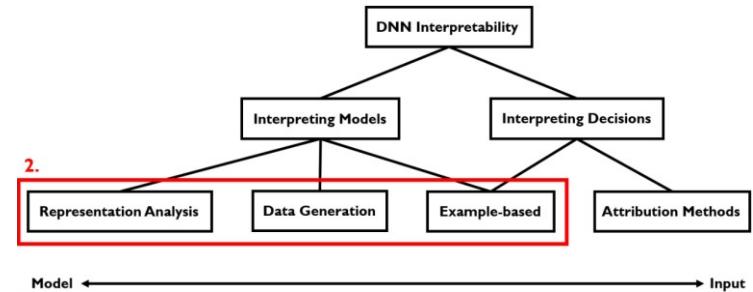
Data Generation

Example-based



Find x that maximizes class posterior probability: search

Types of DNN Interpretability



Weight Visualization

Surrogate Model

Data Generation

Example-based

goose

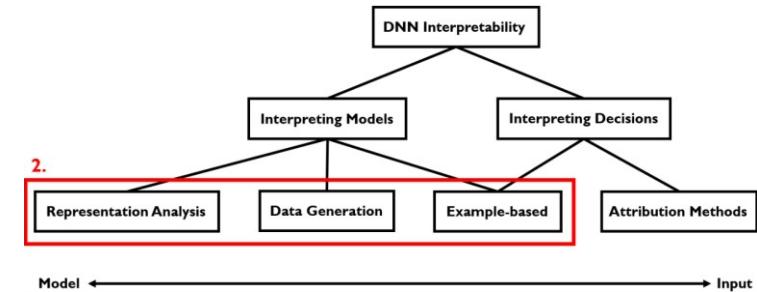


ostrich



Images from Simonyan et al. 2013 "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps"

Types of DNN Interpretability



Weight Visualization

Surrogate Model

Data Generation

Example-based

Advantages

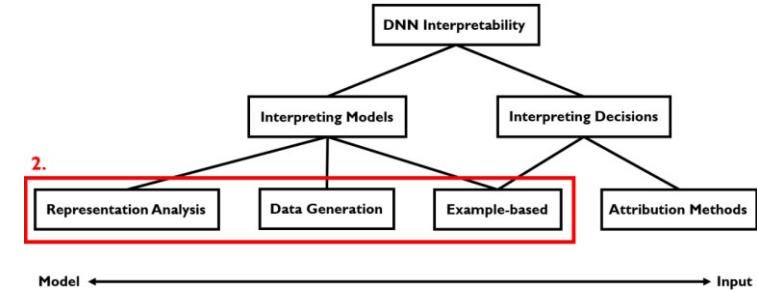
- Activation maximization (AM) builds typical patterns for given classes (e.g. beaks, legs)
- Unrelated background objects are not present in the image

Disadvantages

- Does not resemble class-related patterns
- Lowers the quality of the interpretation for given classes

Redefine optimization problem!

Types of DNN Interpretability



Weight Visualization

Surrogate Model

Data Generation

Example-based

- Does not resemble class-related patterns
- Lowers the quality of the interpretation for given classes

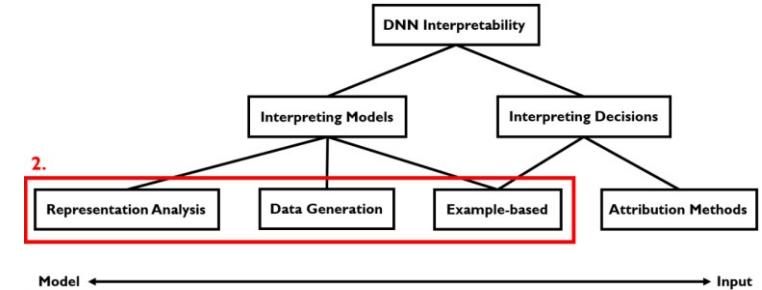
Redefine optimization problem!

Force the generated data x^* to match the data more closely

Find the input pattern that maximizes class probability

Find the most likely input pattern for a given class

Types of DNN Interpretability



Weight Visualization

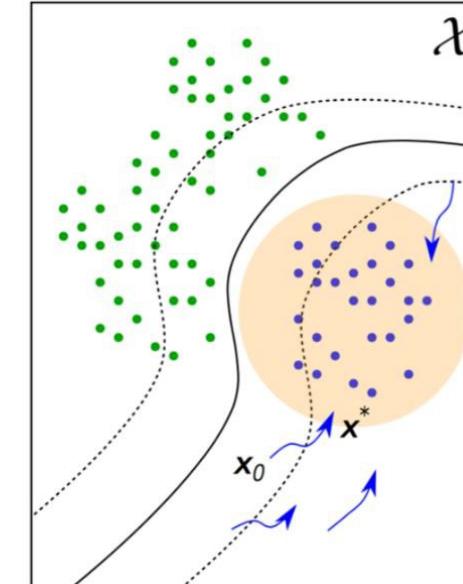
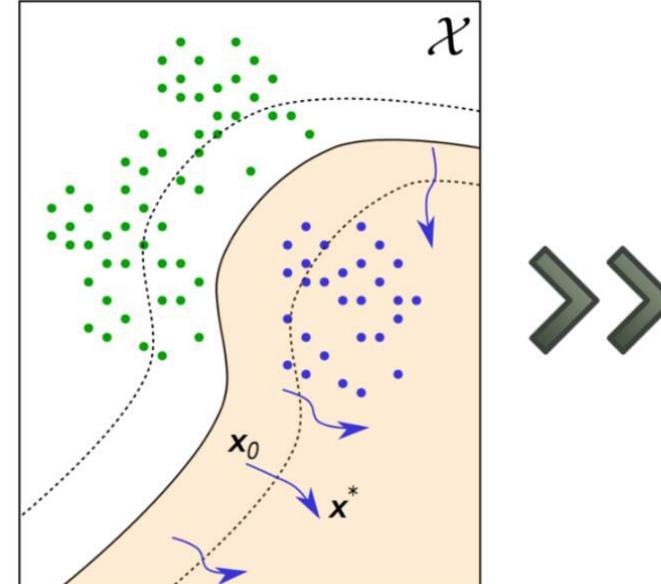
Surrogate Model

Data Generation

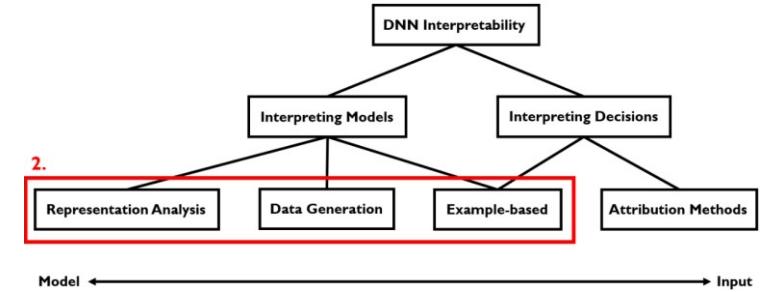
Example-based

Find the input pattern that maximizes class probability

Find the most likely input pattern for a given class



Types of DNN Interpretability



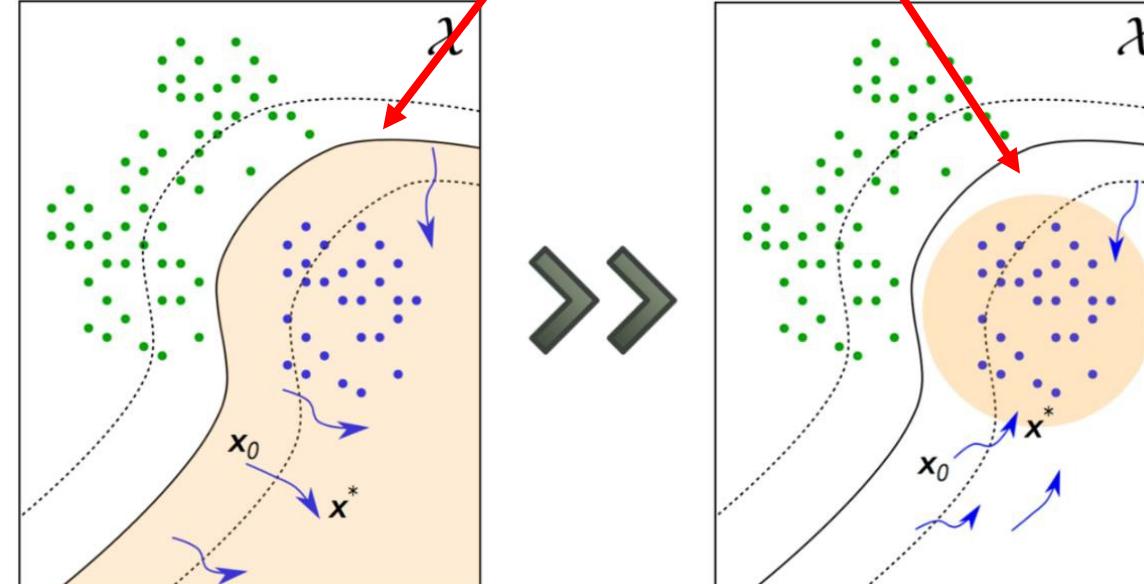
Weight Visualization

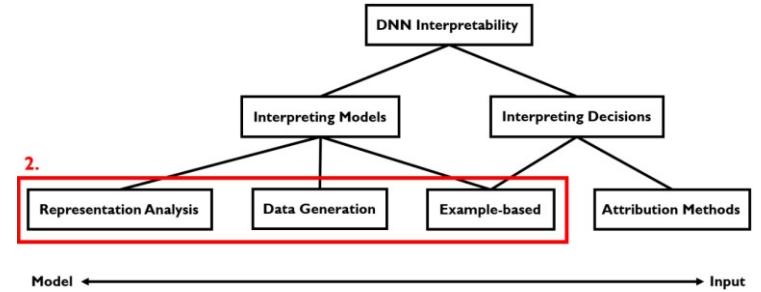
Surrogate Model

Data Generation

Example-based

$$\max_{x \in \mathcal{X}} p_\theta(\omega_c | x) + \lambda \Omega(x)$$





Types of DNN Interpretability

Weight Visualization

Surrogate Model

Data Generation

Example-based

Find the input pattern that maximizes class probability

Find the most likely input pattern for a given class

Activation Maximization with Expert

$$p(x|\omega_c) \propto p(\omega_c|x) \cdot p(x)$$

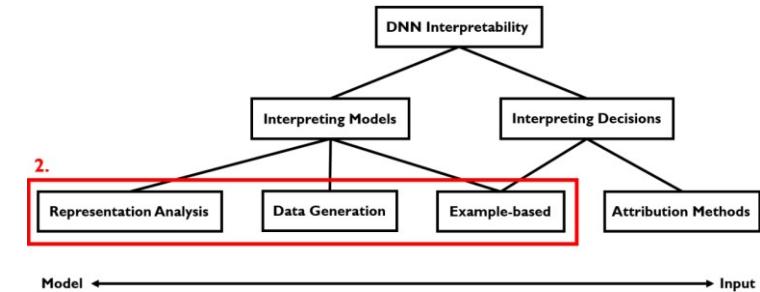
original

Activation Maximization in Code Space

$$\max_{z \in Z} p(\omega_c | \underbrace{g(z)}_x) + \lambda \|z\|^2 \quad x^* = g(z^*)$$

These two techniques require an **unsupervised model of the data**, either a density model $p(x)$ or a generator $g(z)$

Types of DNN Interpretability



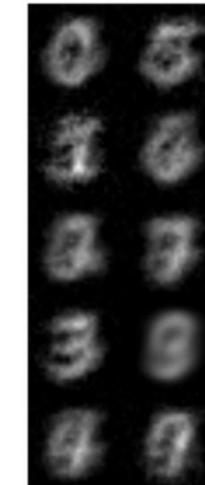
Weight Visualization

Surrogate Model

Data Generation

Example-based

simple AM
(initialized
to mean)



*Original
Activation
Maximization*

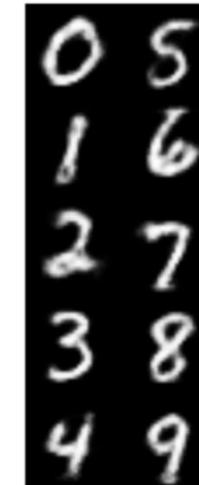
simple AM
(init. to
class
means)



AM-density
(init. to
class
means)



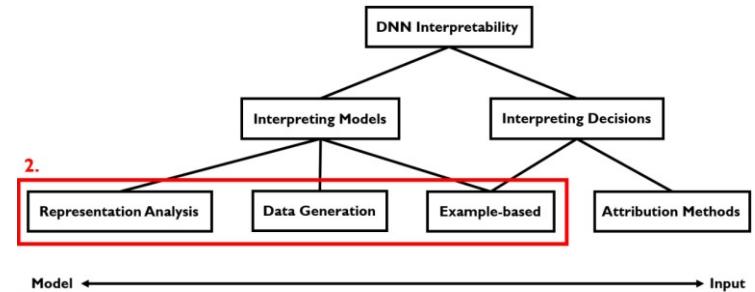
AM-gen
(init. to
class
means)



*Constrained
Activation
Maximization*

Observation: Connecting to the **data** leads to **sharper** visualizations.

Types of DNN Interpretability



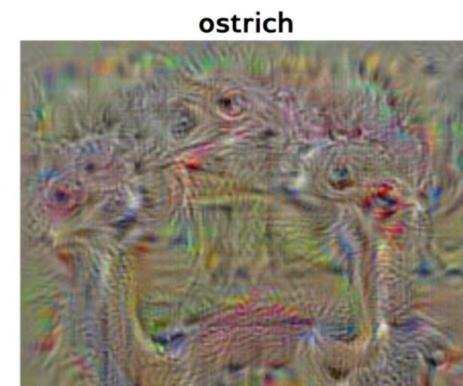
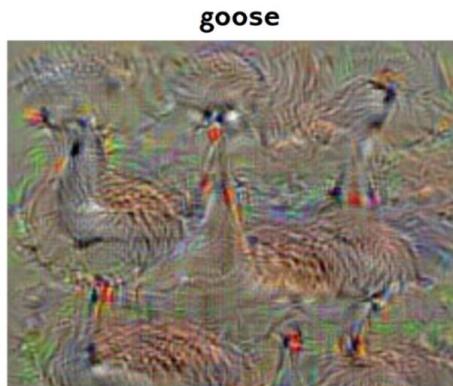
Weight Visualization

Surrogate Model

Data Generation

Example-based

Activation Maximization



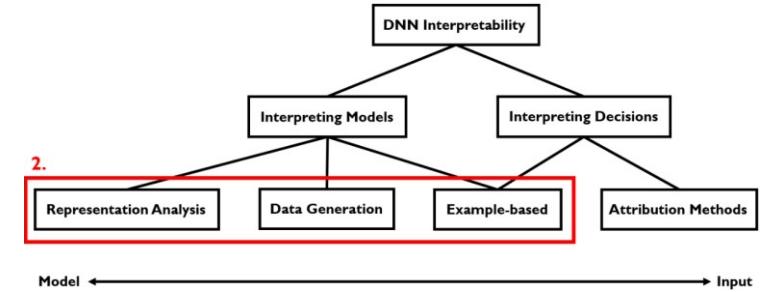
Images from Simonyan et al. 2013 "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps"

Activation Maximization in Code Space

Images from Nguyen et al. 2016. "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks"



Observation: Connecting to the **data** leads to **sharper** visualizations.



Types of DNN Interpretability

Weight Visualization

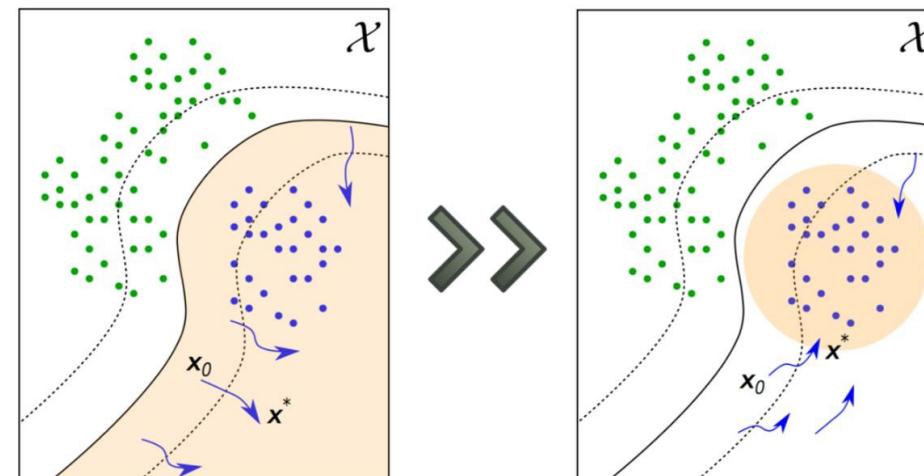
Surrogate Model

Data Generation

Example-based

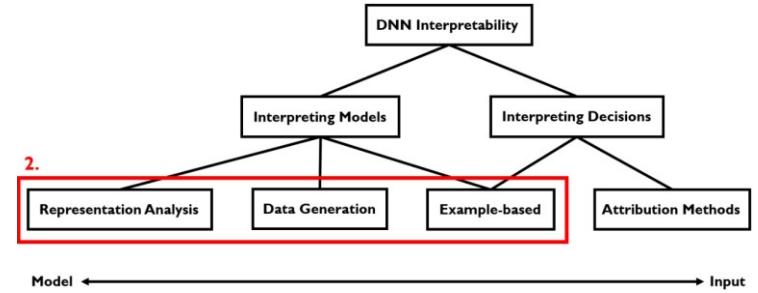
Summary

- DNNs can be interpreted by finding input patterns that maximize a certain output quantity.
- Connecting to the data improves the interpretability of the visualization.



2b – Interpreting models:

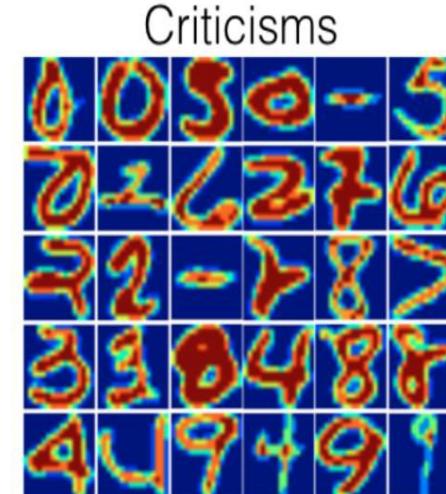
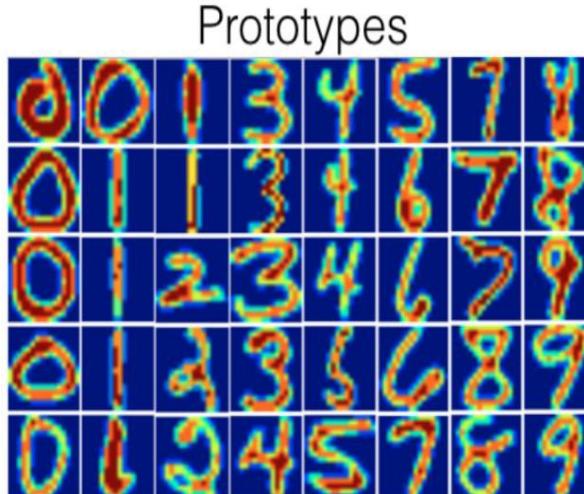
- (i) Weight Visualization
- (ii) Surrogate Model
- (iii) Data Generation / Activation Maximization
- (iv) Example based



Types of DNN Interpretability

Weight Visualization

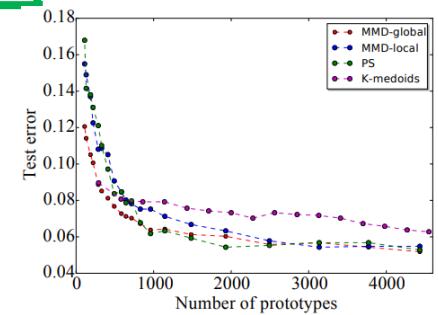
- Find image instances that represent / do not represent the image class
- Maximum Mean Discrepancy – MMD-critic, efficient prototype selection
- Nearest prototype classifier
- Example: digits; Classifying proximal dog breeds; dogs in costumes misclassified



Surrogate Model

Data Generation

Example-based



Interpretable Deep Learning

1. Intro to Interpretability

- 1a. **Interpretability definition:** Convert implicit NN information to human-interpretable information
- 1b. **Motivation:** Verify model works as intended; debug classifier; make discoveries; Right to explanation
- 1c. **Ante-hoc** (train interpretable model) vs. **Post-hoc** (interpret complex model; degree of “locality”)

2. Interpreting Deep Neural Networks

- 2a. **Interpreting Models** (macroscopic, understand internals) vs. **decisions** (microscopic, practical applications)
- 2b. **Interpreting Models:** Weight visualization, Surrogate model, Activation maximization, Example-based

2c. Interpreting Decisions:

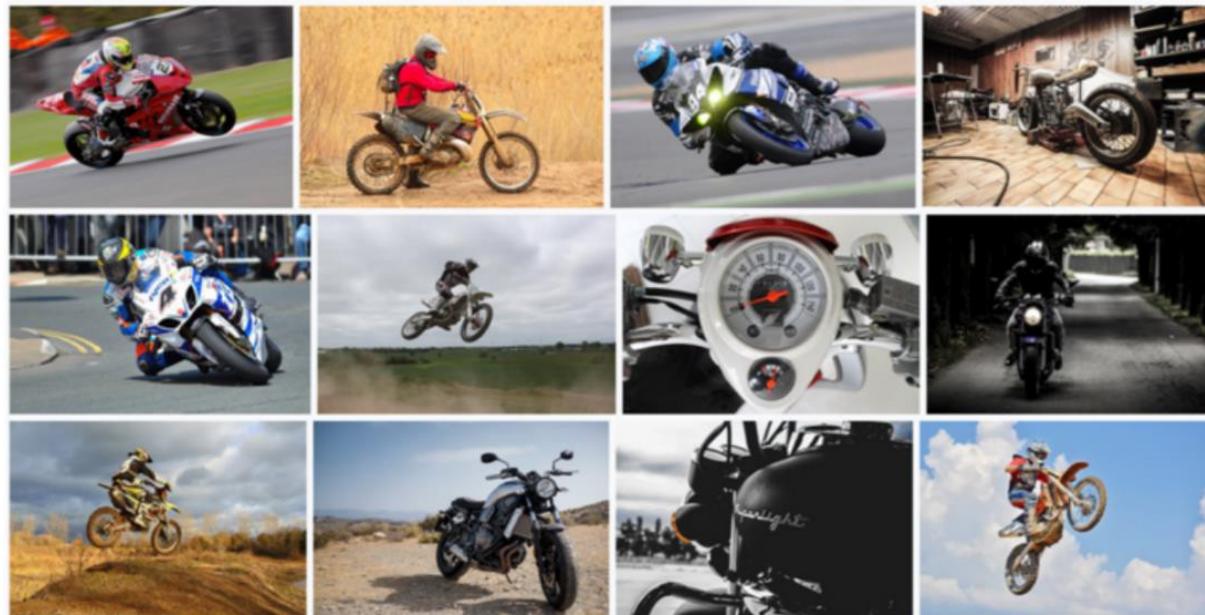
- Example-based
- Attribution Methods: why are gradients noisy?
- Gradient-based Attribution: SmoothGrad, Interior Gradient
- Backprop-based Attribution: Deconvolution, Guided Backpropagation

3. Evaluating Attribution Methods

- 3a. **Qualitative: Coherence:** Attributions should highlight discriminative features / objects of interest
- 3b. **Qualitative: Class Sensitivity:** Attributions should be sensitive to class labels
- 3c. **Quantitative: Sensitivity:** Removing feature with high attribution → large decrease in class probability
- 3d. **Quantitative: ROAR & KAR.** Low class prob cuz image unseen → remove pixels, retrain, measure acc. drop

Limitation of Model Interpretations

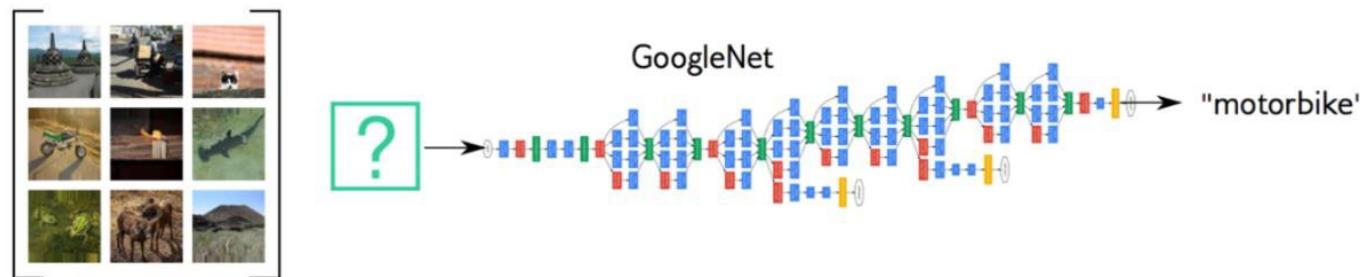
Question: What would be the best image to interpret the class “motorcycle”?



- Summarizing a concept or a category like “motorcycle” into a single image is difficult.
- A good interpretation would grow as large as the diversity of the concept to interpret.

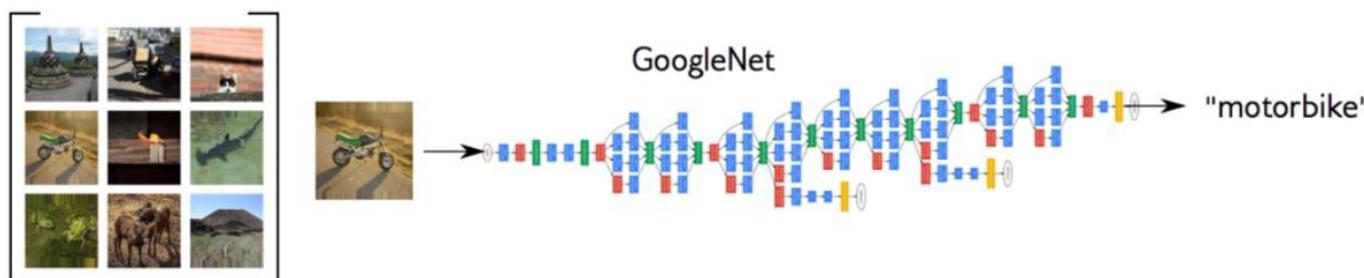
Limitation of Model Interpretations

Finding a prototype:



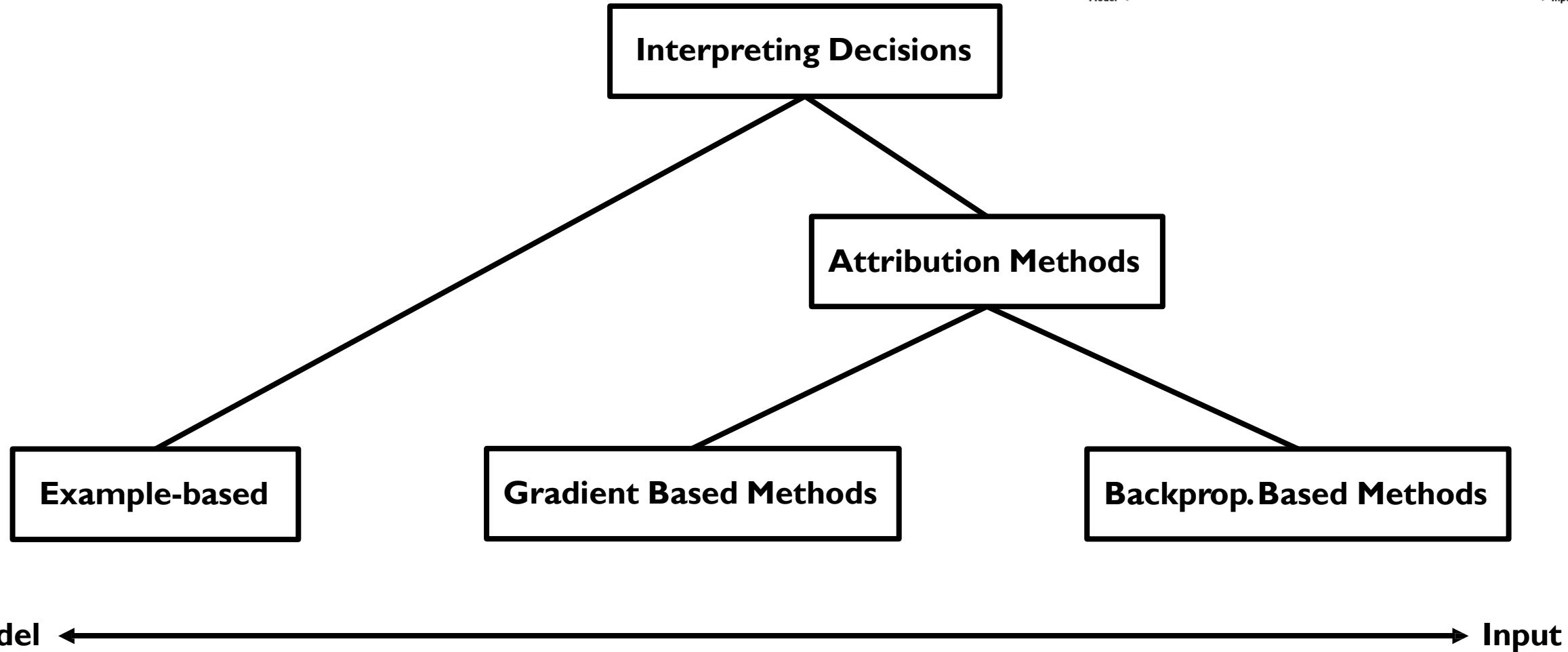
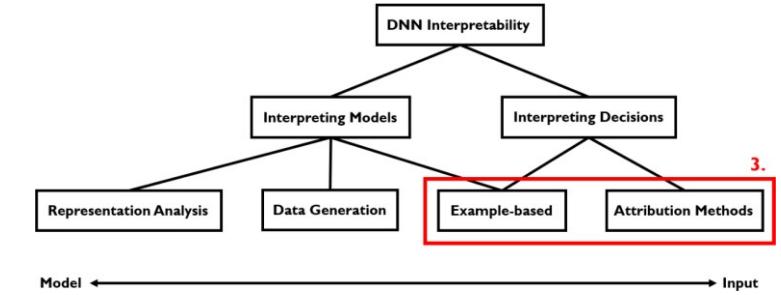
Question: How does a “motorbike” **typically** look like?

Decision explanation:



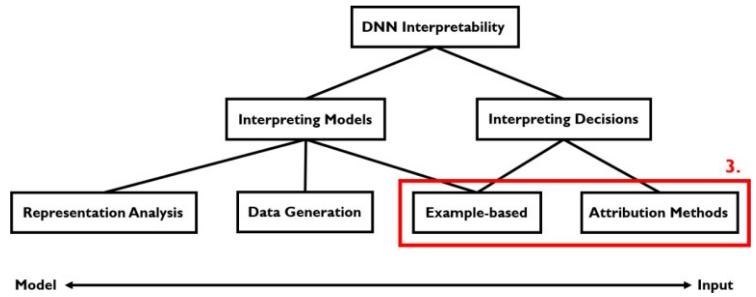
Question: Why is this example classified as a motorbike?

Types of DNN Interpretability



2c – Interpreting decisions:

- (i) Example based
- (ii) Attribution Methods: why are gradients noisy?
- (iii) Gradient-based Attribution: SmoothGrad, Interior Gradient
- (iv) Backprop-based Attribution: Deconvolution, Guided Backprop



Types of DNN Interpretability

Example-based

Attribution Methods

Gradient Based

Backprop. Based

- Which training instance influenced the decision most?
- Still does not specifically highlight which features were important
- Influence functions for interpreting black-box methods. Fragility of NN model interpretation.

'Sunflower': 59.2% conf.



Original

Influence: 0.09



Influence: 0.14



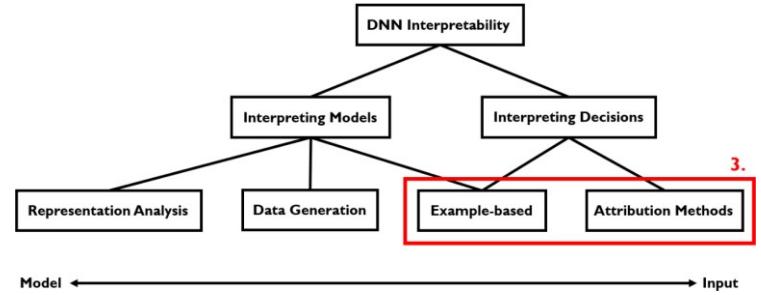
Influence: 0.42



2c – Interpreting decisions:

- (i) Example based
- (ii) Attribution Methods: why are gradients noisy?
- (iii) Gradient-based Attribution: SmoothGrad, Interior Gradient
- (iv) Backprop-based Attribution: Deconvolution, Guided Backprop

Types of DNN Interpretability



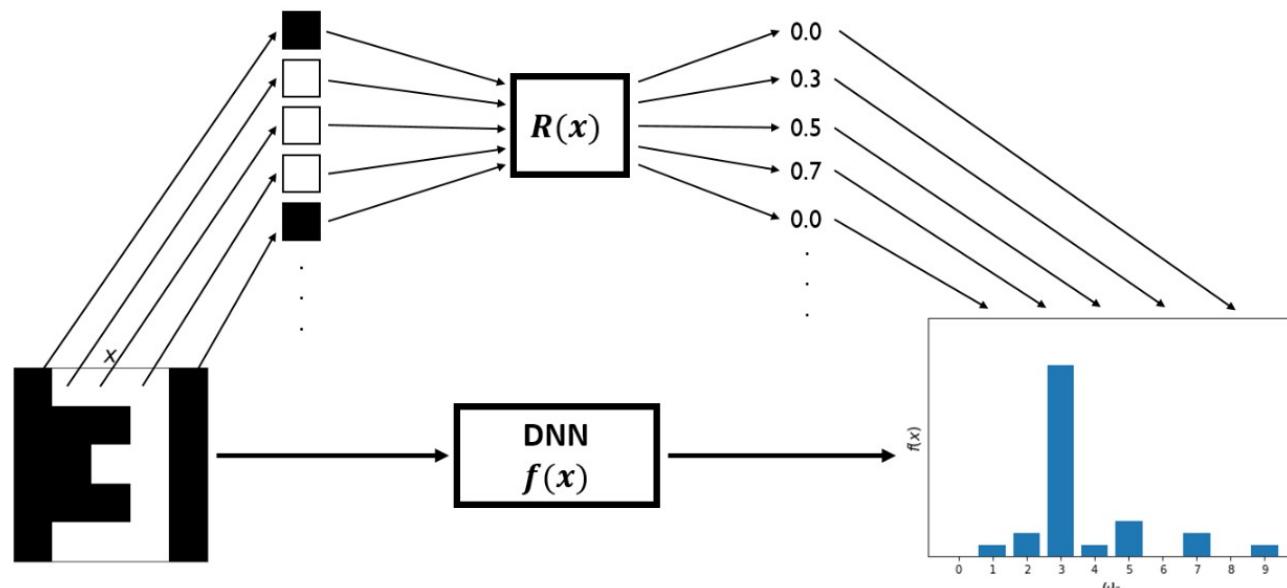
Example-based

Attribution Methods

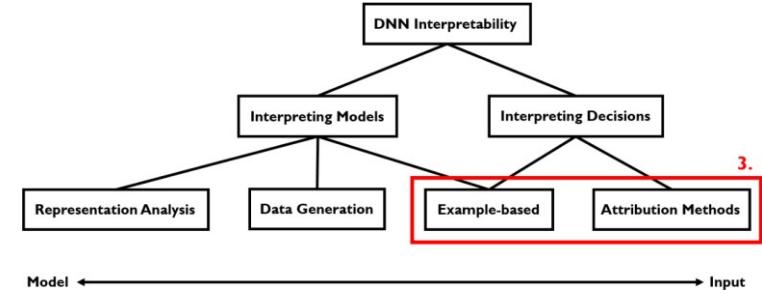
Gradient Based

Backprop. Based

Given an image $x \in \mathbb{R}^n$ and a decision $f(x)$,
assign to each pixel x_1, x_2, \dots, x_n attribution values $R_1(x), R_2(x), \dots, R_n(x)$.



Types of DNN Interpretability



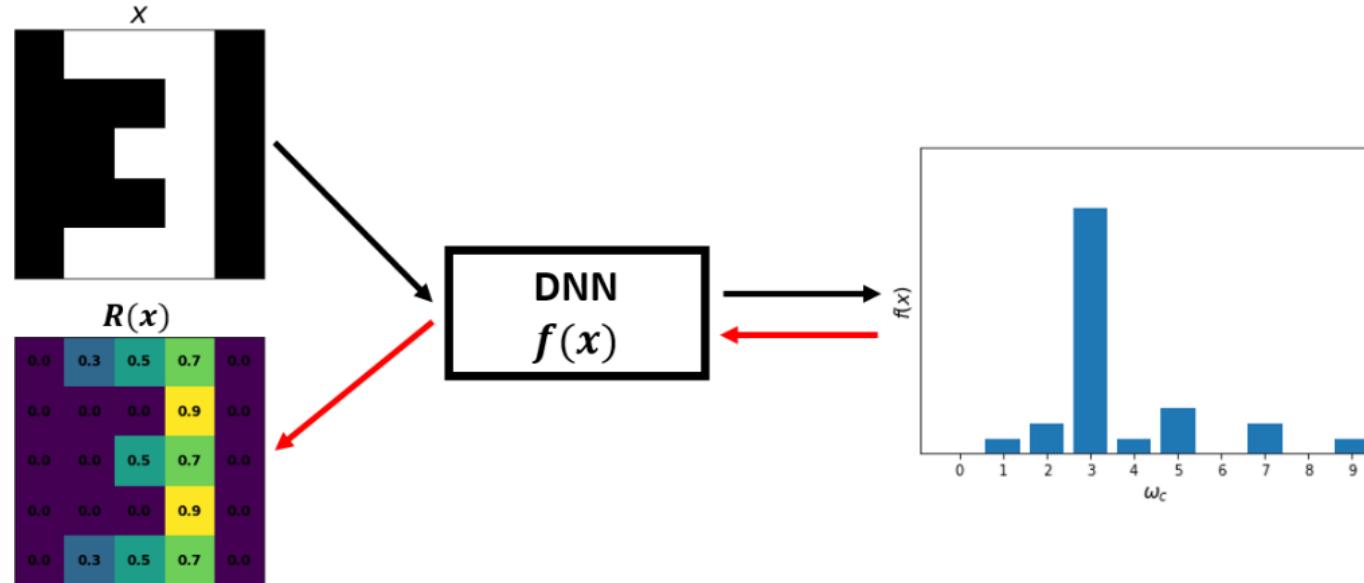
Example-based

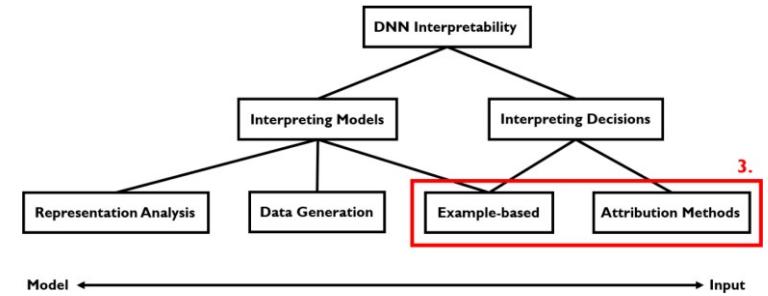
Attribution Methods

Gradient Based

Backprop. Based

Attributions visualized as heatmaps





Types of DNN Interpretability

Example-based

Attribution Methods

Gradient Based

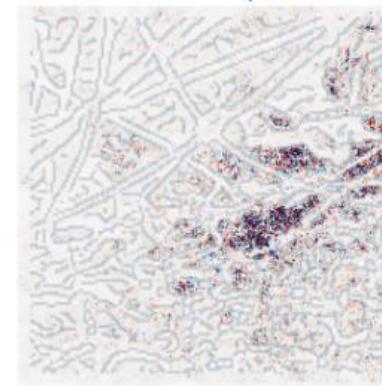
Backprop. Based

Attributions visualized as **heatmaps**

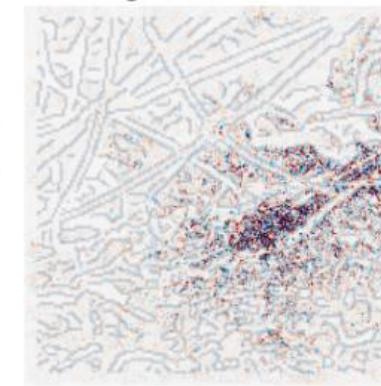
Original (label: "garter snake")



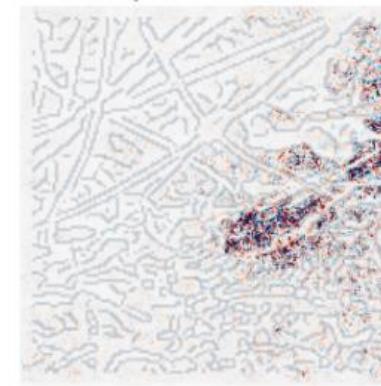
Grad * Input



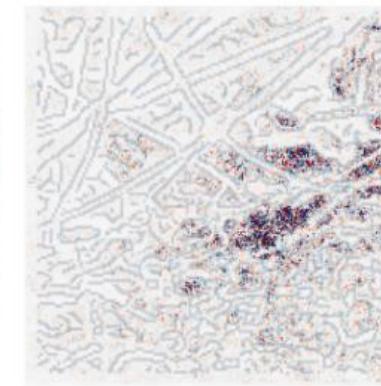
Integrated Gradients

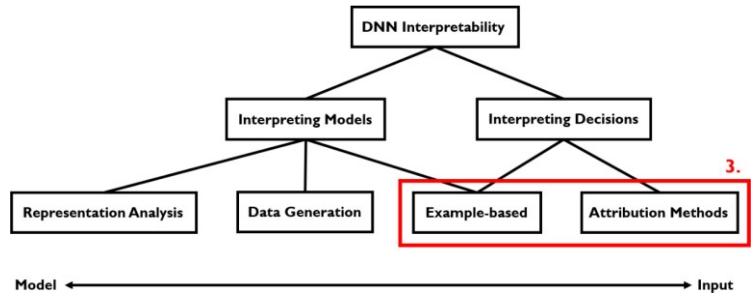


DeepLIFT (Rescale)



ϵ -LRP





Types of DNN Interpretability

Example-based

Attribution Methods

Gradient Based

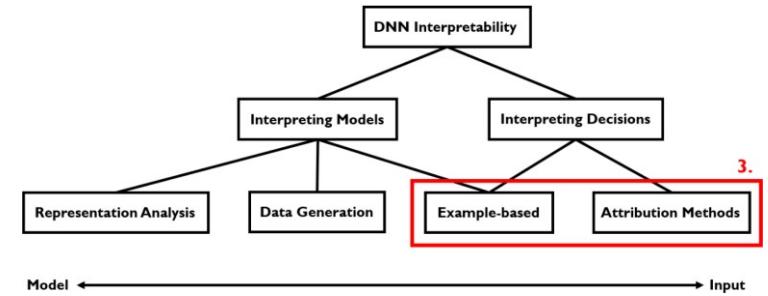
Backprop. Based

The Baseline Attribution Method Saliency Map

- Gradient of the decision $f(x)$ with respect to each input image x :

$$\text{Saliency}(x) := \nabla_x f(x) = \frac{\partial f(x)}{\partial x}$$

- Can be calculated through backpropagation.



Types of DNN Interpretability

Example-based

Attribution Methods

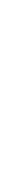
Gradient Based

Backprop. Based

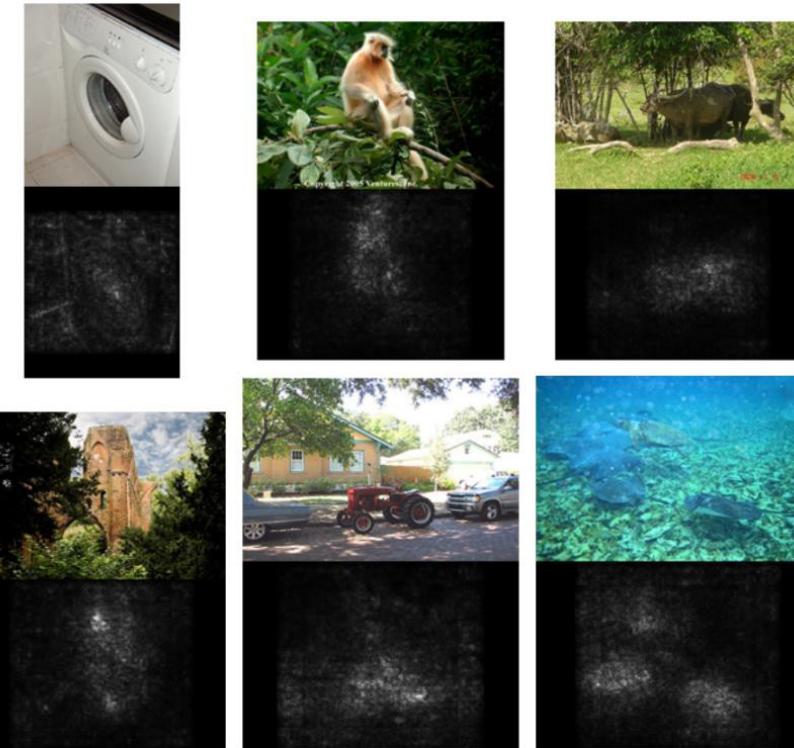
The Baseline Attribution Method Saliency Map

- Saliency maps are very **noisy!**
- Only roughly correlated with the object(s) of interest.

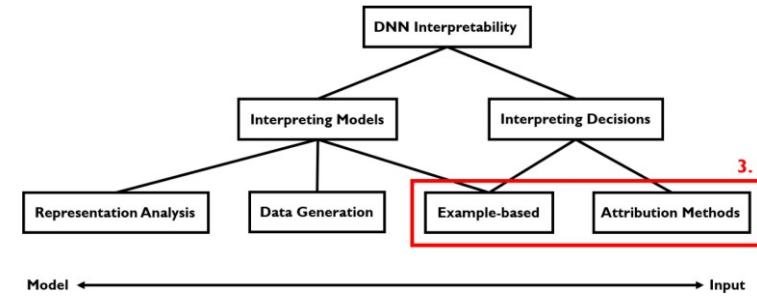
Question: How to improve saliency maps?



Question: Why are saliency maps noisy?



Types of DNN Interpretability



Example-based

Attribution Methods

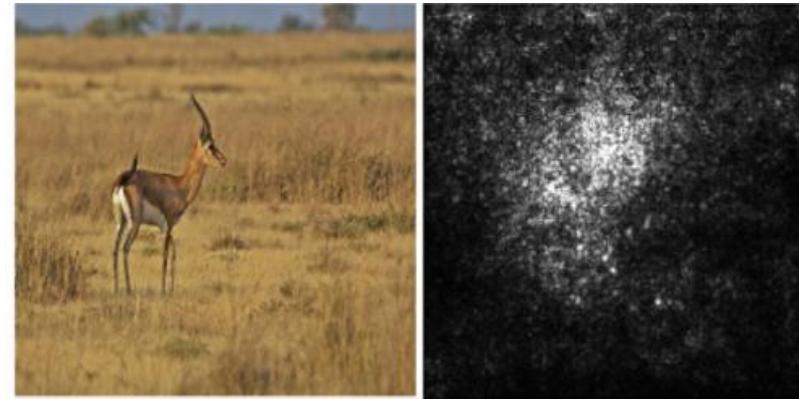
Gradient Based

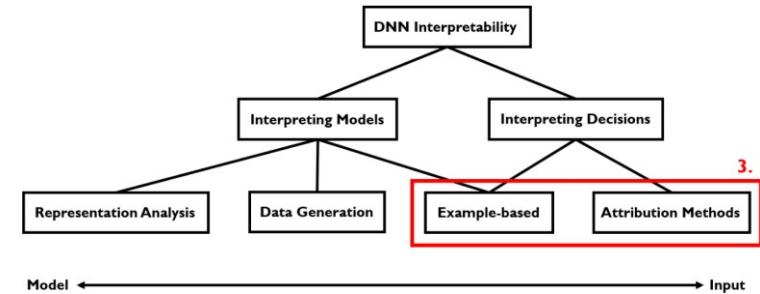
Backprop. Based

Question: Why are saliency maps noisy?

Hypothesis I – Saliency maps are truthful

- Certain pixels scattered randomly across the image are central to how the network is making a decision.
 - Noise is important!





Types of DNN Interpretability

Example-based

Attribution Methods

Gradient Based

Backprop. Based

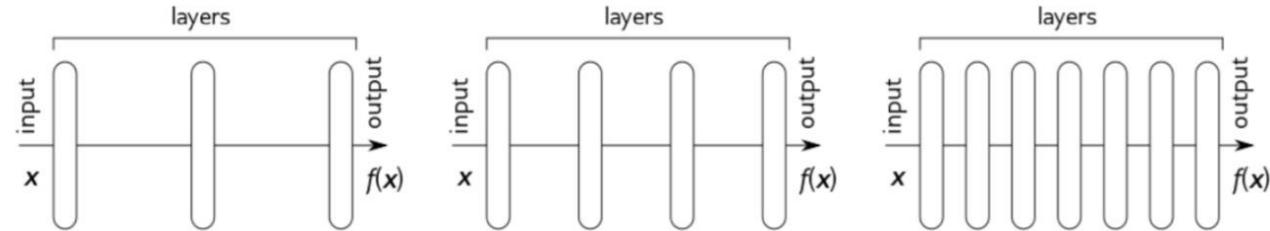
Question: Why are saliency maps noisy?

Hypothesis 2 – Gradients are discontinuous

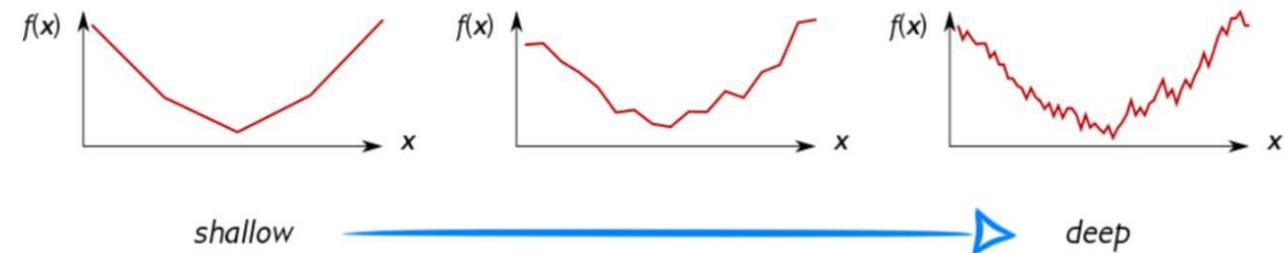
- DNN uses piecewise-linear functions (ReLU activation, max-pooling, etc.).

- Sudden jumps in the importance score over infinitesimal changes in the input.

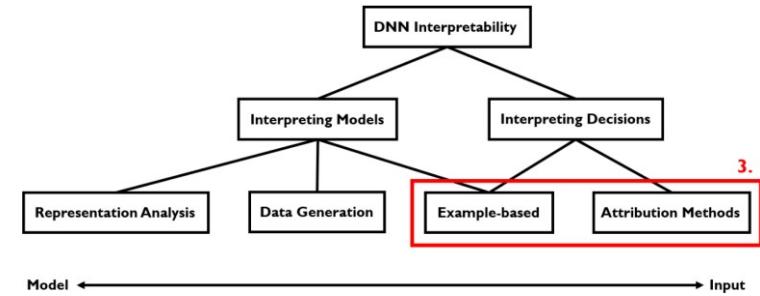
Structure's view



Function's view (cartoon)



Types of DNN Interpretability



Example-based

Attribution Methods

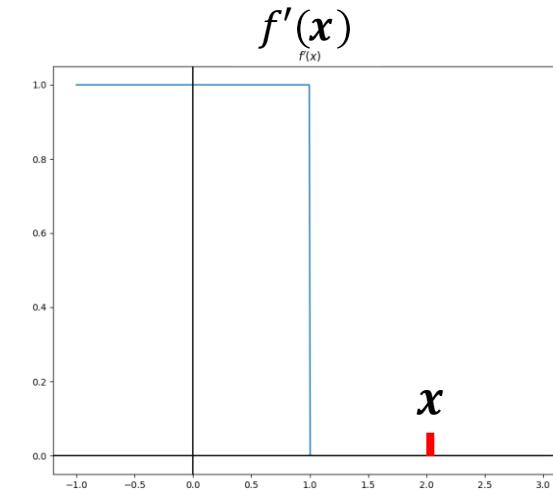
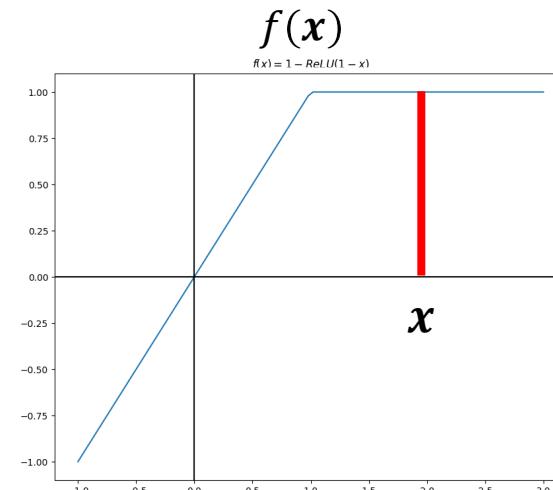
Gradient Based

Backprop. Based

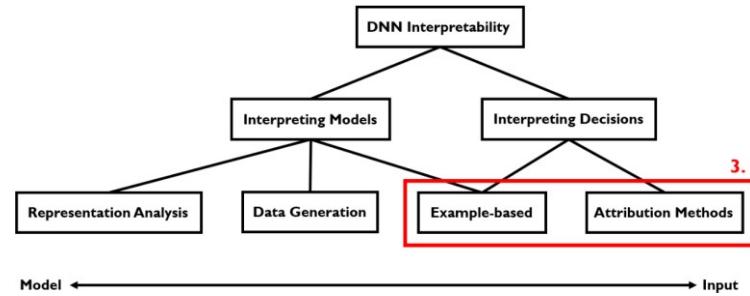
Question: Why are saliency maps noisy?

Hypothesis 3 – $f(x)$ saturates

-A feature may have a strong effect globally, but with a small derivative locally.



Types of DNN Interpretability



Example-based

Attribution Methods

Gradient Based

Backprop. Based

Question: How to improve saliency maps?

$$\text{Saliency}(\mathbf{x}) := \nabla_{\mathbf{x}} f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$$

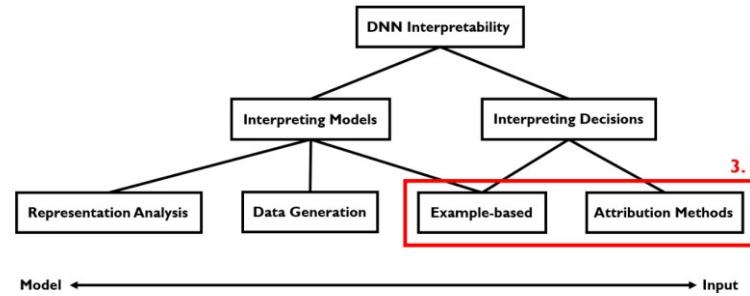
Gradient-based Methods

- Add noise: Perturb the input \mathbf{x} to \mathbf{x}^* and use $\nabla_{\mathbf{x}^*} f(\mathbf{x}^*)$.
- Some methods take the average over the perturbation set $\{\mathbf{x}_1^*, \mathbf{x}_2^*, \dots, \mathbf{x}_n^*\}$.

Backprop-based Methods

- Modify the backpropagation algorithm.

Types of DNN Interpretability



Example-based

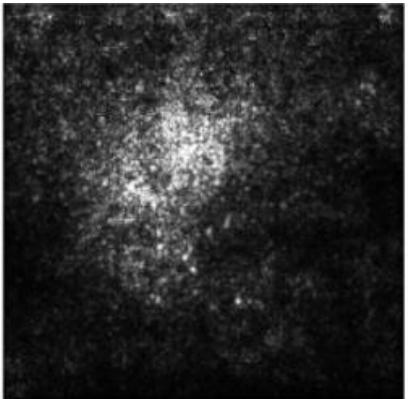
Attribution Methods

Gradient Based

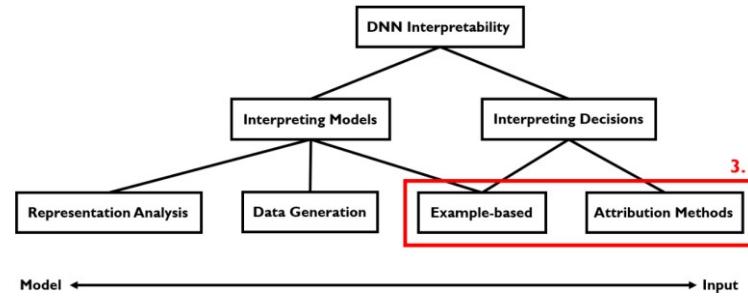
Backprop. Based

Summary

- Attribution method assigns “attribution score” to each input pixel.
- Baseline attribution method Saliency Map is noisy.
- Hypothesis 1:Saliency maps are truthful.
- Hypothesis 2:Gradients are discontinuous.
- Hypothesis 3: $f(x)$ saturates.
- Two solution approaches: Gradient-based method and Backprop-based method.



Types of DNN Interpretability



Example-based

Attribution Methods

Gradient Based

Backprop. Based

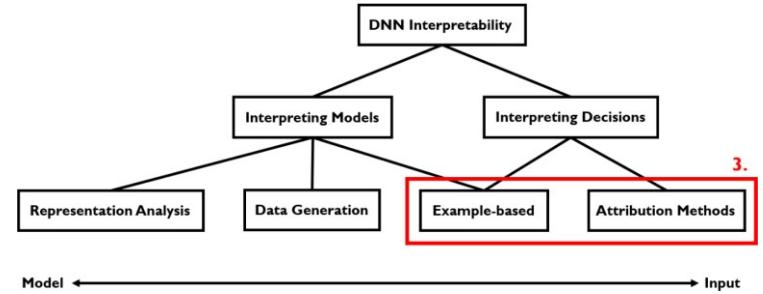
Other Attribution Methods

- Gradient * Input – <https://arxiv.org/pdf/1704.02685.pdf>
- Integrated Gradient – <https://arxiv.org/pdf/1703.01365.pdf>
- Layer-wise Relevance Propagation – <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0130140>
- Deep Taylor Decomposition – <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0130140>
- DeepLIFT – <https://arxiv.org/pdf/1704.02685.pdf>
- PatternNet and PatternAttribution – <https://arxiv.org/pdf/1705.05598.pdf>

2c – Interpreting decisions:

- (i) Example based
- (ii) Attribution Methods: why are gradients noisy?
- (iii) Gradient-based Attribution: SmoothGrad, Interior Gradient
- (iv) Backprop-based Attribution: Deconvolution, Guided Backprop

Types of DNN Interpretability



Example-based

Attribution Methods

Gradient Based

Backprop. Based

I. SmoothGrad Hypothesis 2 – *Gradients are discontinuous: Just smooth the gradient!*

$$\text{SmoothGrad}(\mathbf{x}) := \frac{1}{n} \int_1^n \frac{\partial f(\mathbf{x}^*)}{\partial \mathbf{x}^*}, \quad \mathbf{x}^* = \mathbf{x} + \mathcal{N}(0, \sigma^2)$$

Gaussian smoothing

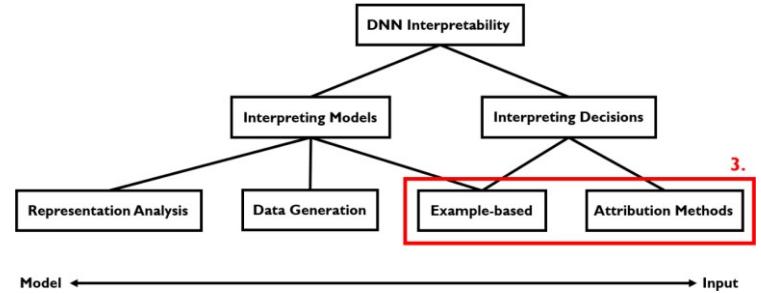
Function's view (cartoon)



shallow

deep

Types of DNN Interpretability



Example-based

Attribution Methods

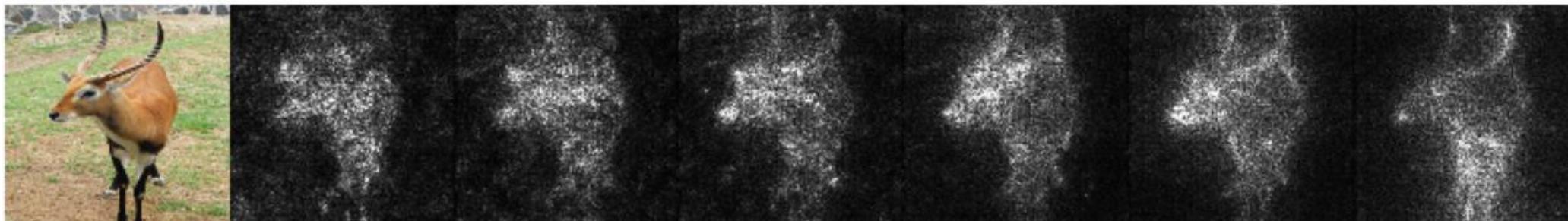
Gradient Based

Backprop. Based

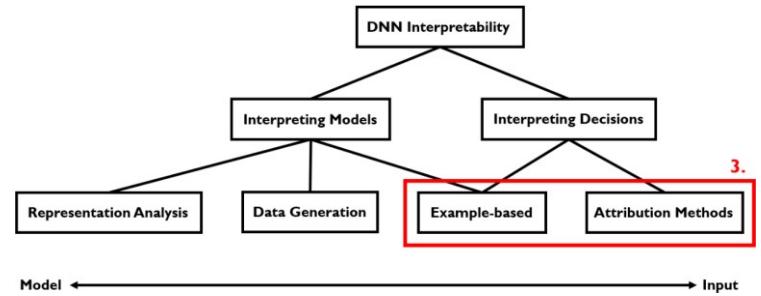
I. SmoothGrad Hypothesis 2 – Gradients are discontinuous: Just smooth the gradient!

$$\text{SmoothGrad}(\mathbf{x}) := \frac{1}{n} \int_1^n \frac{\partial f(\mathbf{x}^*)}{\partial \mathbf{x}^*}, \quad \mathbf{x}^* = \mathbf{x} + \mathcal{N}(0, \sigma^2)$$

Noise level:



Types of DNN Interpretability



Example-based

Attribution Methods

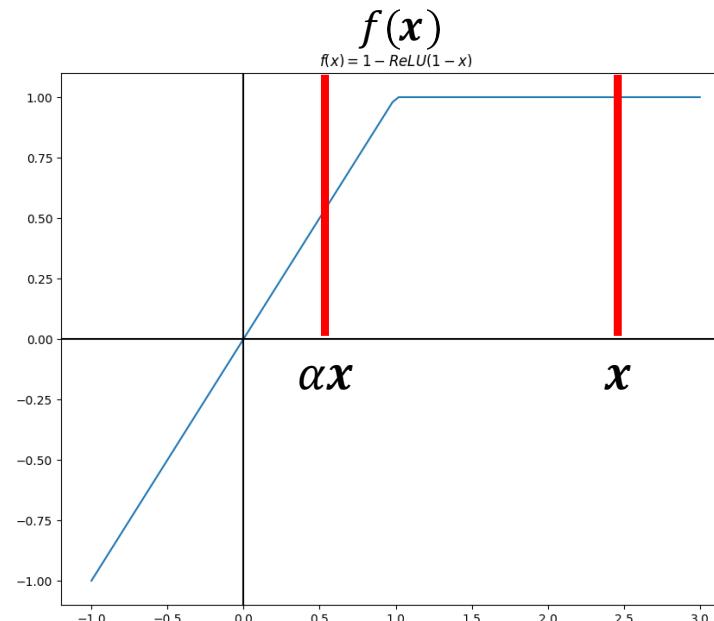
Gradient Based

Backprop. Based

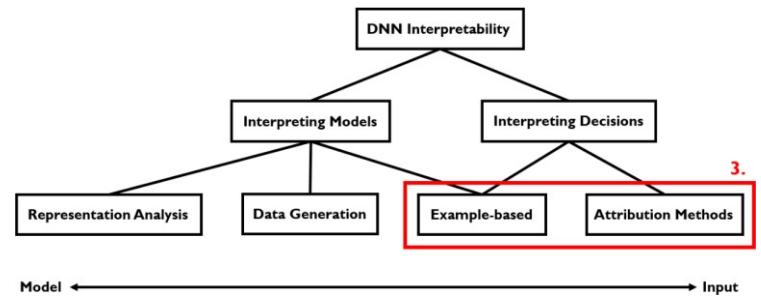
2. Interior Gradient Hypothesis 3 – $f(x)$ saturates

$$\text{IntGrad}(x) := \frac{\partial f(x^*)}{\partial x^*}, \quad x^* = \alpha x, \quad 0 < \alpha \leq 1$$

-Appropriate α will trigger informative activation functions



Types of DNN Interpretability



Example-based

Attribution Methods

Gradient Based

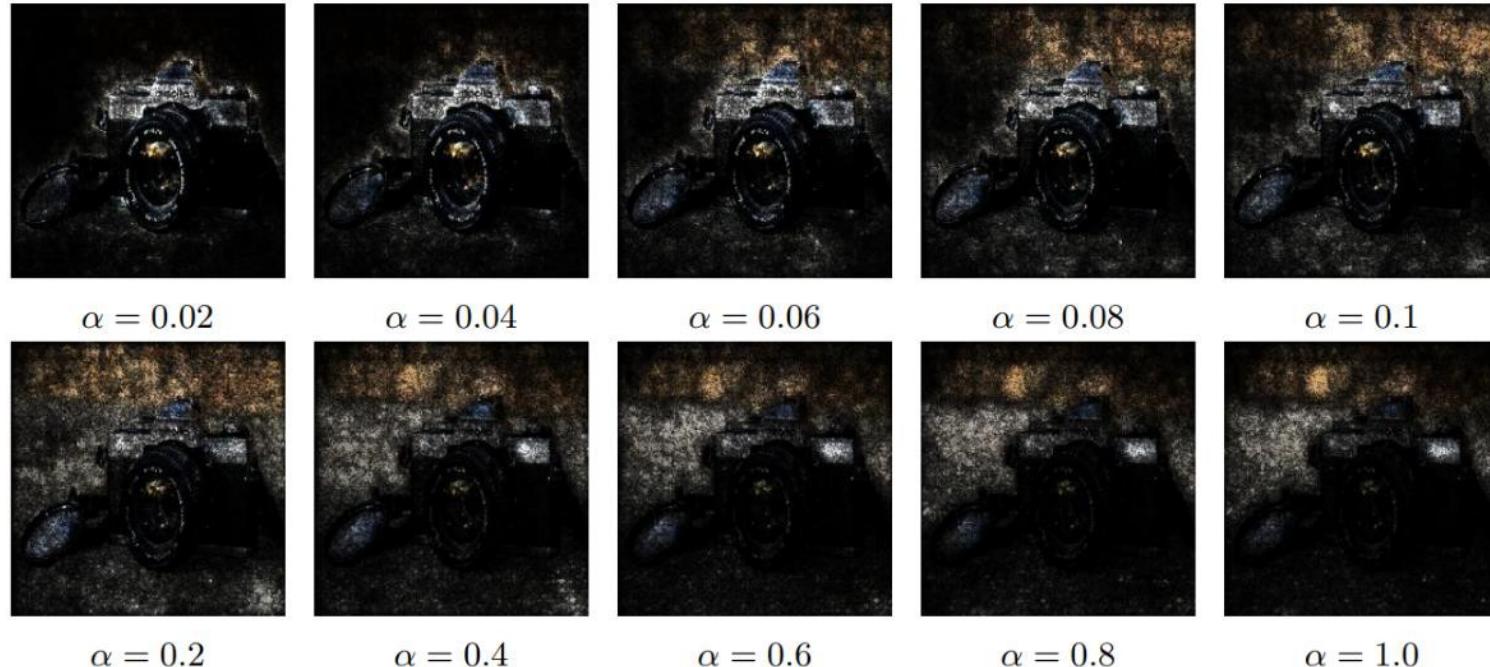
Backprop. Based

2. Interior Gradient

$$\text{IntGrad}(x) := \frac{\partial f(x^*)}{\partial x^*},$$

$$x^* = \alpha x,$$

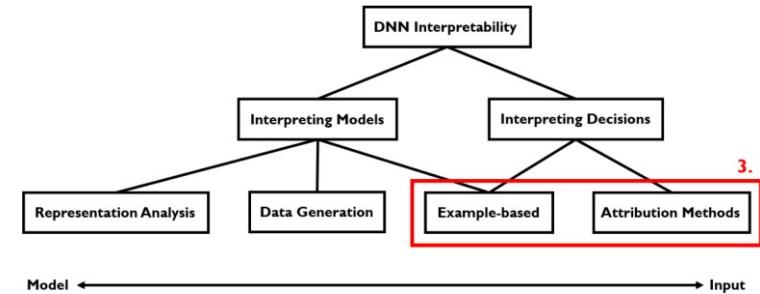
$$0 < \alpha \leq 1$$



2c – Interpreting decisions:

- (i) Example based
- (ii) Attribution Methods: why are gradients noisy?
- (iii) Gradient-based Attribution: SmoothGrad, Interior Gradient
- (iv) Backprop-based Attribution: Deconvolution, Guided Backprop

Types of DNN Interpretability



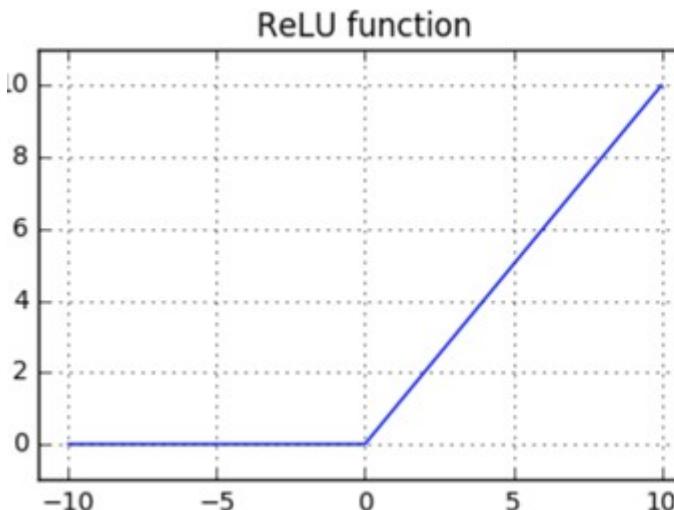
Example-based

Attribution Methods

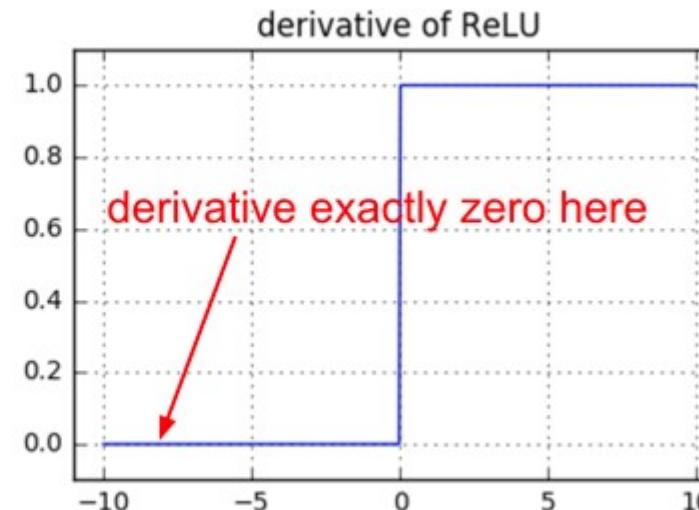
Gradient Based

Backprop. Based

Review: Backpropagation at ReLU



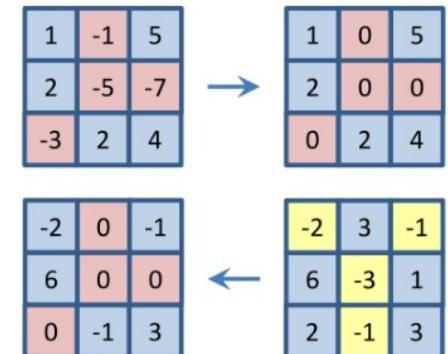
$$ReLU(z) = \max(0, z)$$



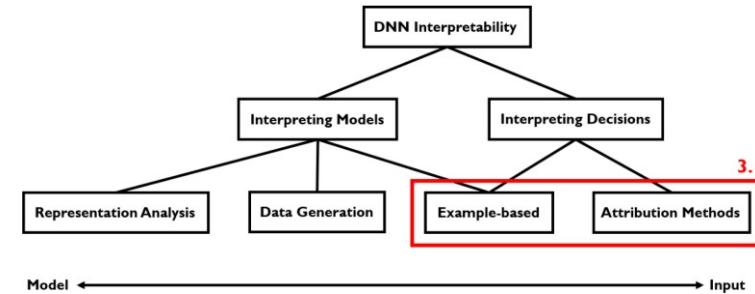
$$ReLU'(z) = (z > 0)$$

Forward pass

Backward pass:
backpropagation



Types of DNN Interpretability



Example-based

Attribution Methods

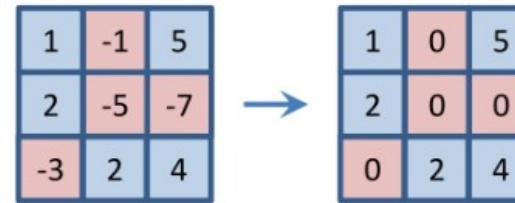
Gradient Based

Backprop. Based

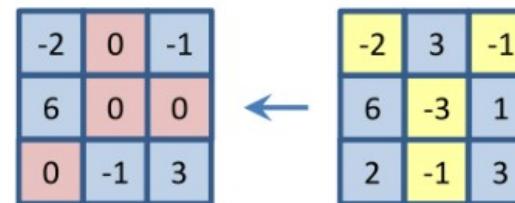
I. Deconvnet

- Maps feature pattern to input space (image reconstruction)
 - To obtain valid feature reconstruction, pass the reconstructed signal through ReLUs
 - Removing noise by removing negative gradient

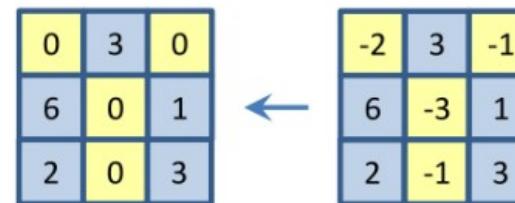
Forward pass

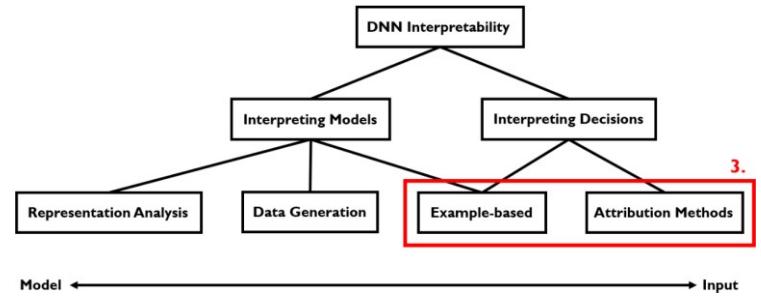


Backward pass:
backpropagation



Backward pass “deconvnet”





Types of DNN Interpretability

Example-based

Attribution Methods

Gradient Based

Backprop. Based

2. Guided Backpropagation

- Combine Deconvnet with Backpropagation
- Removing negative gradient + consider forward activations

Forward pass

$$\begin{matrix} 1 & -1 & 5 \\ 2 & -5 & -7 \\ -3 & 2 & 4 \end{matrix} \rightarrow \begin{matrix} 1 & 0 & 5 \\ 2 & 0 & 0 \\ 0 & 2 & 4 \end{matrix}$$

Backward pass:
backpropagation

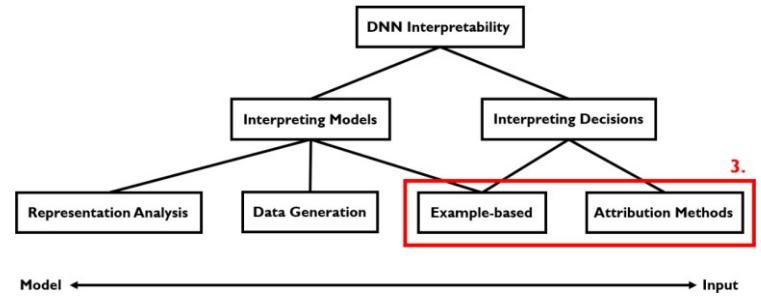
$$\begin{matrix} -2 & 0 & -1 \\ 6 & 0 & 0 \\ 0 & -1 & 3 \end{matrix} \leftarrow \begin{matrix} -2 & 3 & -1 \\ 6 & -3 & 1 \\ 2 & -1 & 3 \end{matrix}$$

Backward pass:
“deconvnet”

$$\begin{matrix} 0 & 3 & 0 \\ 6 & 0 & 1 \\ 2 & 0 & 3 \end{matrix} \leftarrow \begin{matrix} -2 & 3 & -1 \\ 6 & -3 & 1 \\ 2 & -1 & 3 \end{matrix}$$

Backward pass:
*guided
backpropagation*

$$\begin{matrix} 0 & 0 & 0 \\ 6 & 0 & 0 \\ 0 & 0 & 3 \end{matrix} \leftarrow \begin{matrix} -2 & 3 & -1 \\ 6 & -3 & 1 \\ 2 & -1 & 3 \end{matrix}$$



Types of DNN Interpretability

Example-based

Attribution Methods

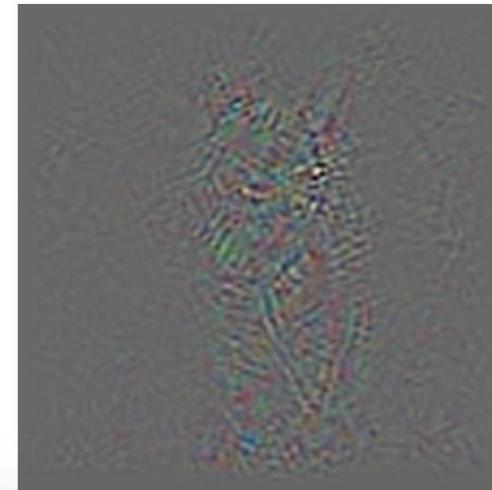
Gradient Based

Backprop. Based

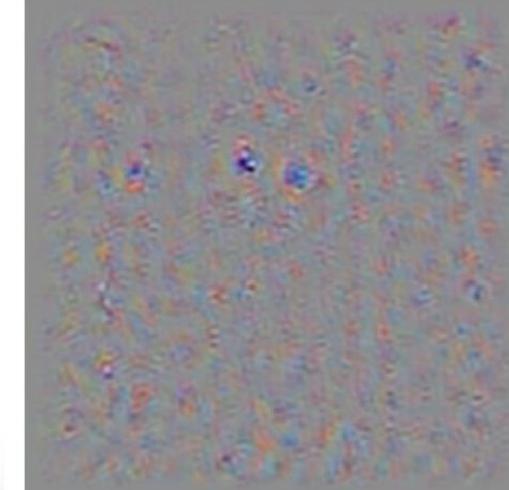
Input image



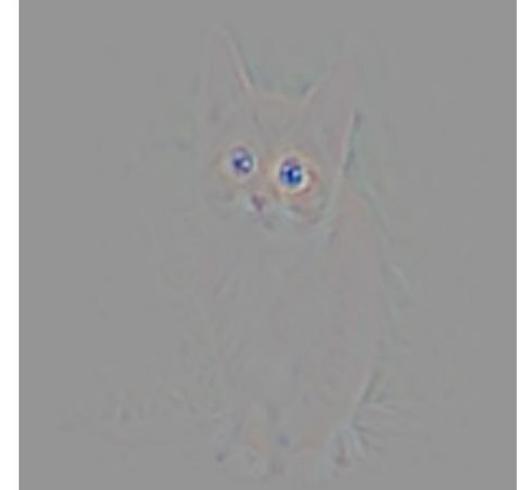
Backpropagation



Deconvolution



Guided Backprop



Observation: Removing **more** gradient leads to **sharper** visualizations

Interpretable Deep Learning

1. Intro to Interpretability

- 1a. **Interpretability definition:** Convert implicit NN information to human-interpretable information
- 1b. **Motivation:** Verify model works as intended; debug classifier; make discoveries; Right to explanation
- 1c. **Ante-hoc** (train interpretable model) vs. **Post-hoc** (interpret complex model; degree of “locality”)

2. Interpreting Deep Neural Networks

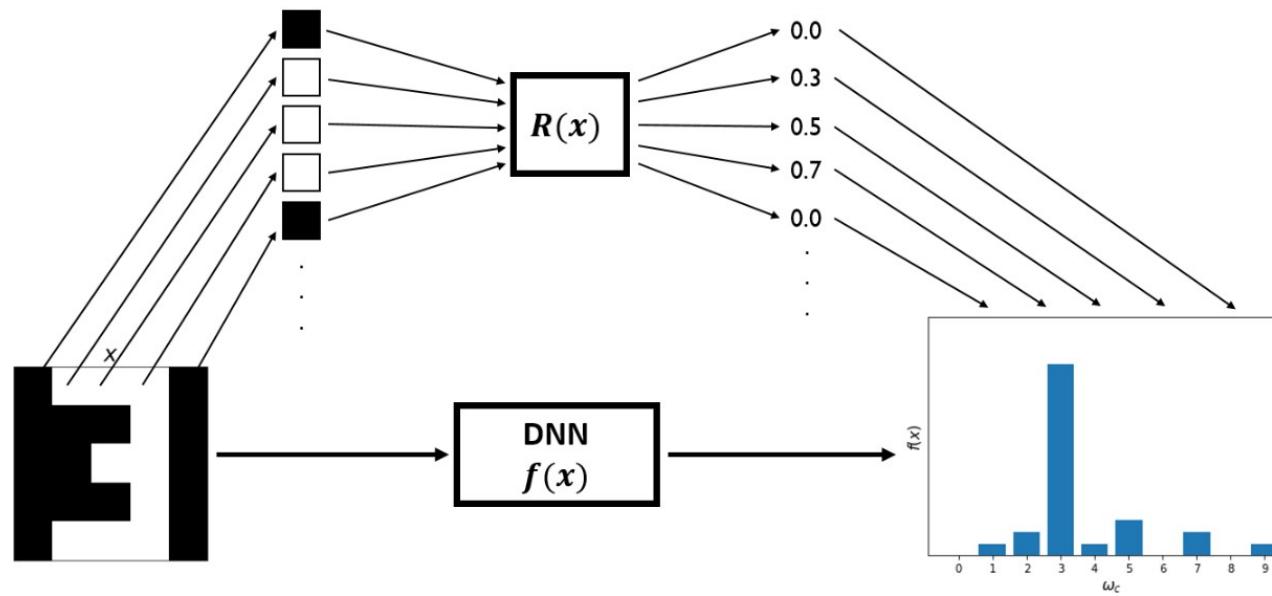
- 2a. **Interpreting Models** (macroscopic, understand internals) vs. **decisions** (microscopic, practical applications)
- 2b. **Interpreting Models:** Weight visualization, Surrogate model, Activation maximization, Example-based
- 2c. **Interpreting Decisions:**
 - Example-based
 - Attribution Methods: why are gradients noisy?
 - Gradient-based Attribution: SmoothGrad, Interior Gradient
 - Backprop-based Attribution: Deconvolution, Guided Backpropagation

3. Evaluating Attribution Methods

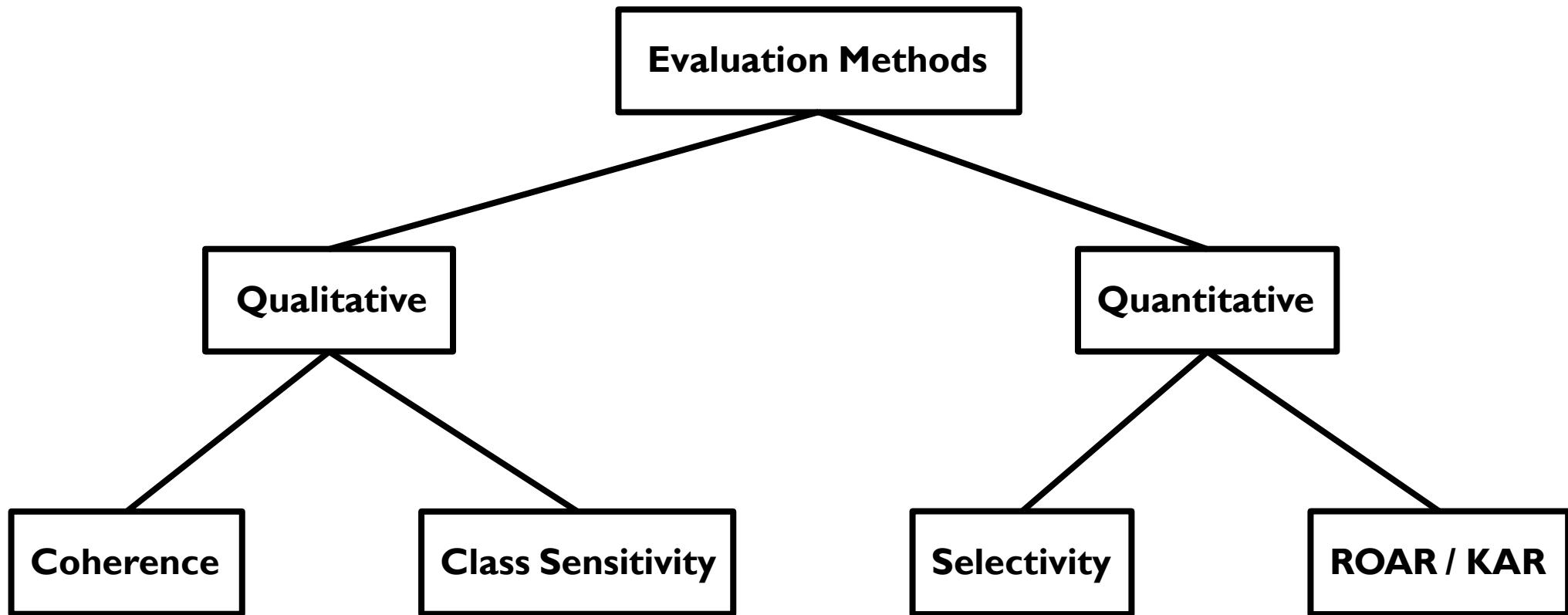
- 3a. **Qualitative: Coherence:** Attributions should highlight discriminative features / objects of interest
- 3b. **Qualitative: Class Sensitivity:** Attributions should be sensitive to class labels
- 3c. **Quantitative: Sensitivity:** Removing feature with high attribution → large decrease in class probability
- 3d. **Quantitative: ROAR & KAR.** Low class prob cuz image unseen → remove pixels, retrain, measure acc. drop

Attribution Method Review

Given an image $x \in \mathbb{R}^n$ and a decision $f(x)$,
assign to each pixel x_1, x_2, \dots, x_n **attribution values** $R_1(x), R_2(x), \dots, R_n(x)$.



Evaluating Attribution Methods



3a – Qualitative: coherence

Evaluating Attribution Methods

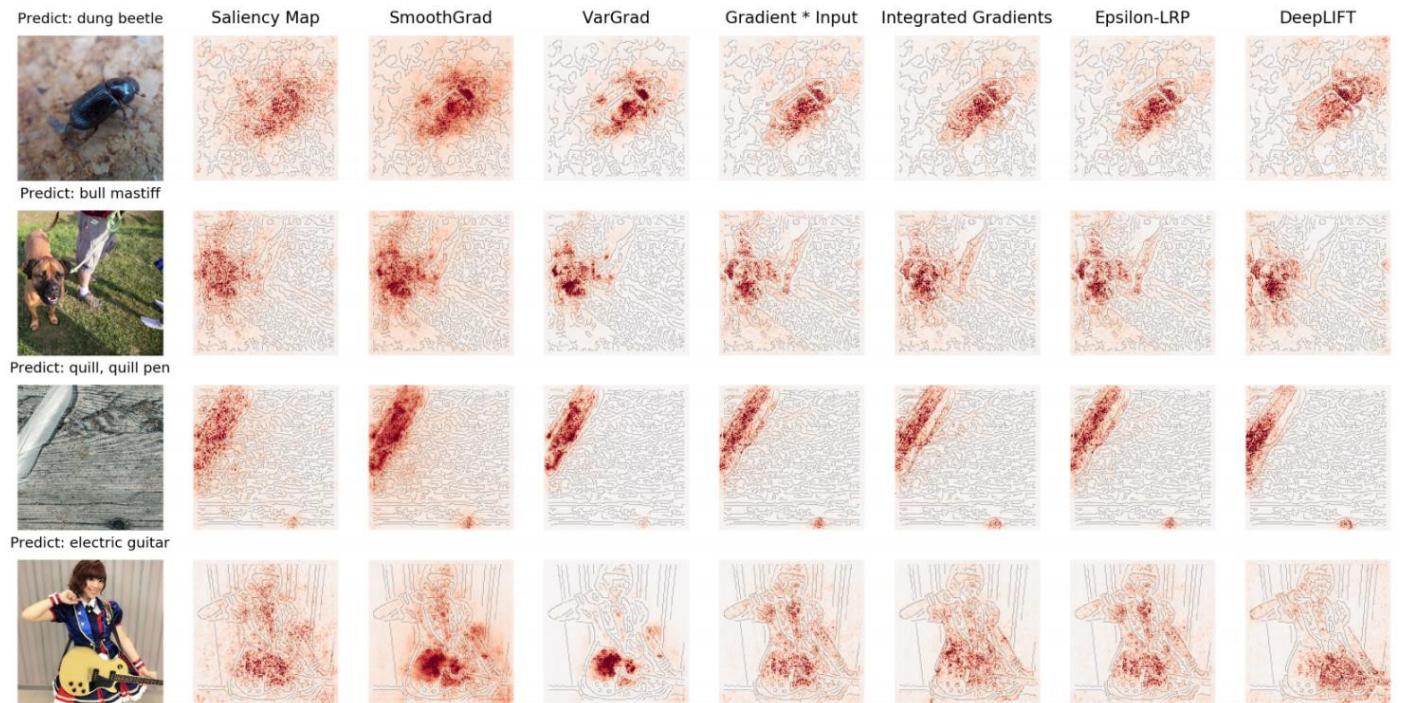
Coherence

Class Sensitivity

Selectivity

ROAR / KAR

-Attributions should fall on discriminative features (e.g. the object of interest)



Evaluating Attribution Methods

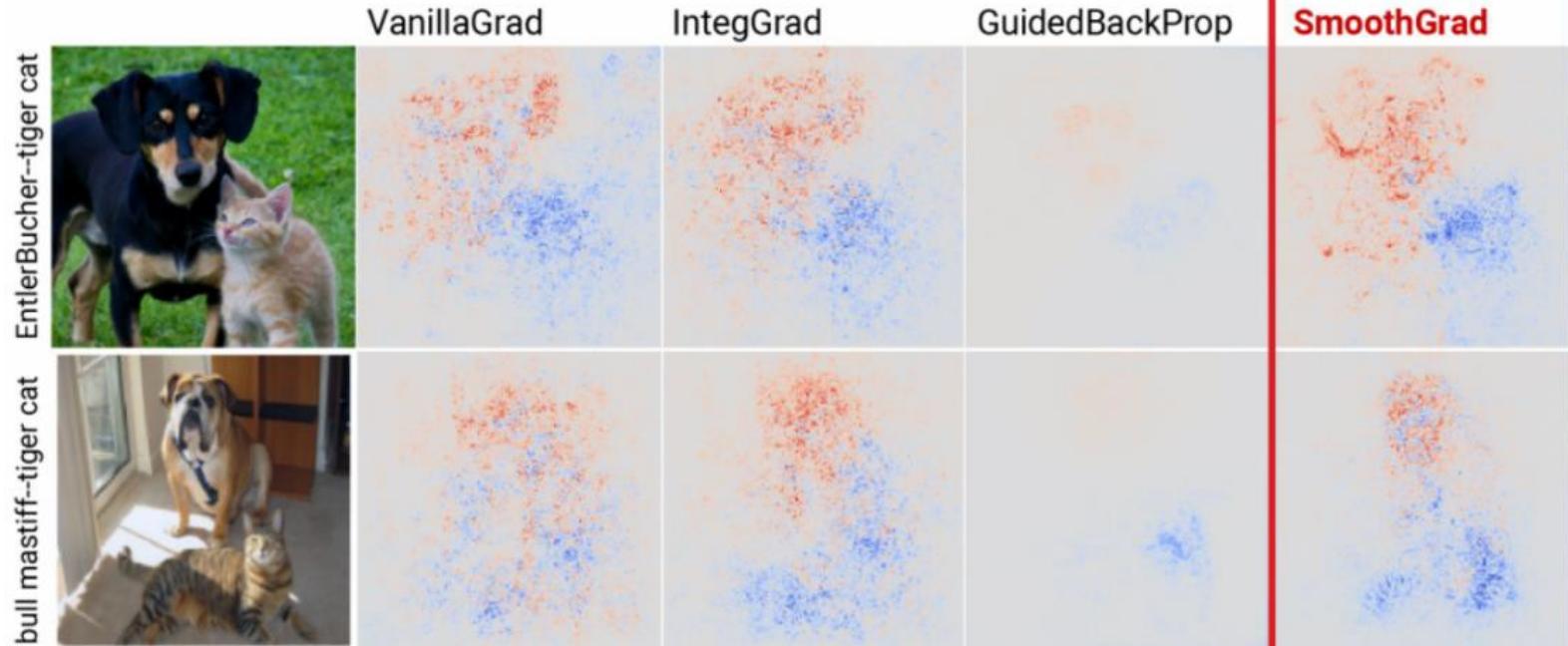
Coherence

Class Sensitivity

Selectivity

ROAR / KAR

-Attributions should be sensitive to class labels



Evaluating Attribution Methods

Coherence

Class Sensitivity

Selectivity

ROAR / KAR

- Removing feature with high attribution should cause large decrease in class probability

Algorithm

Sort pixel attribution values $R_i(x)$

Iterate:

 Remove pixels

 Evaluate $f(x)$

Measure decrease of $f(x)$

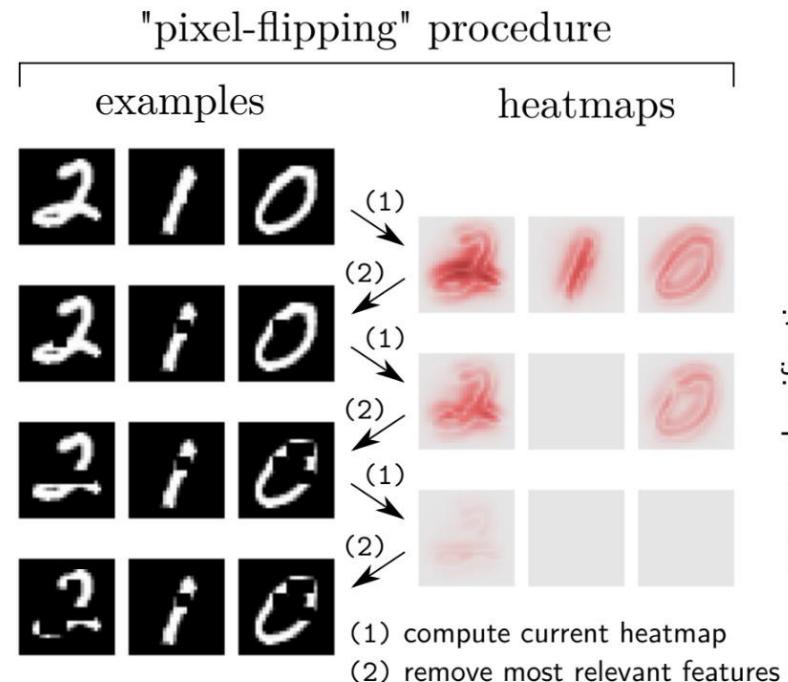
Evaluating Attribution Methods

Coherence

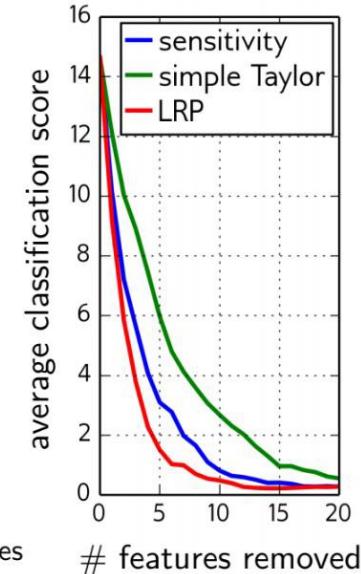
Class Sensitivity

Selectivity

ROAR / KAR

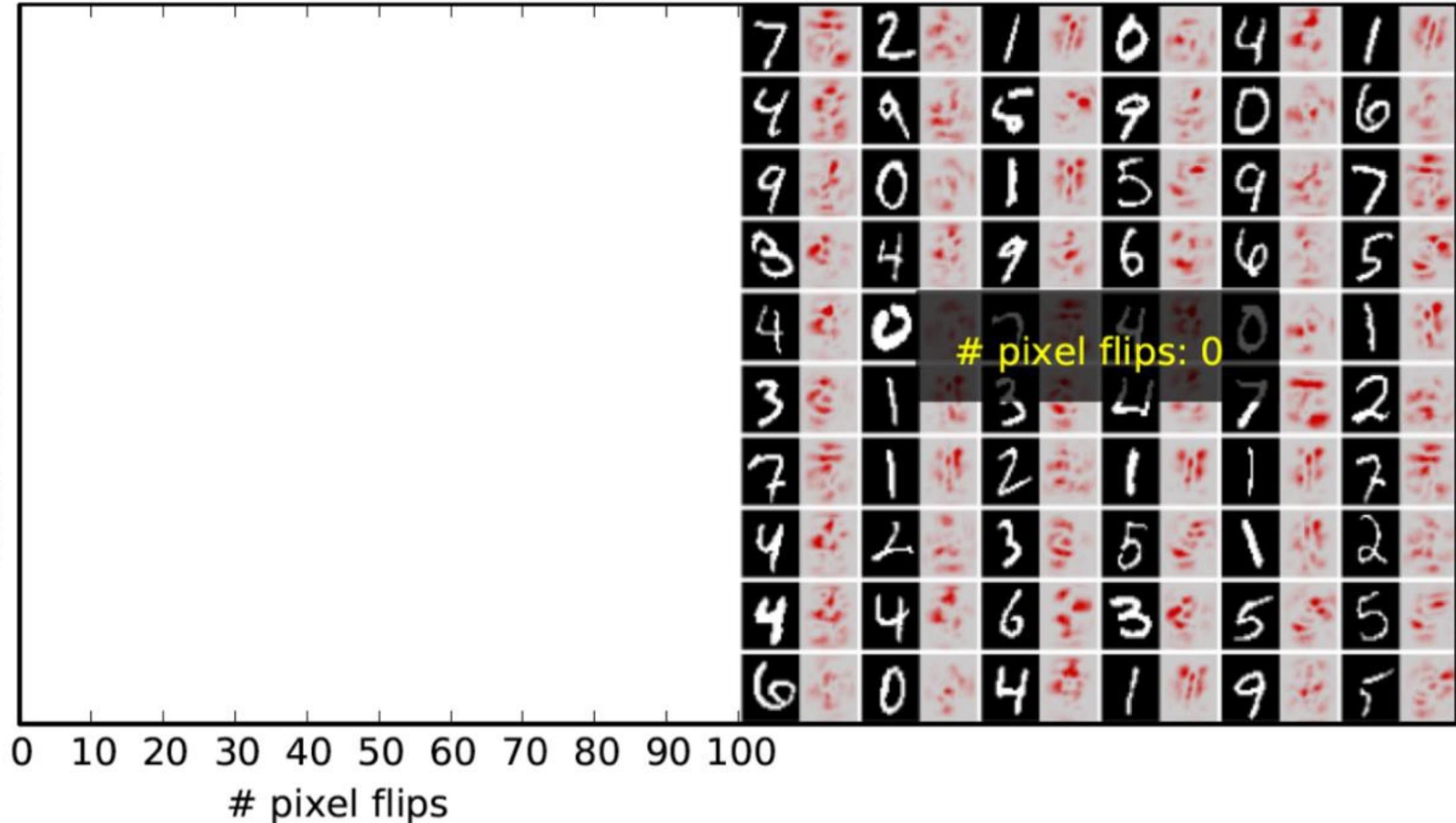


comparing explanation techniques

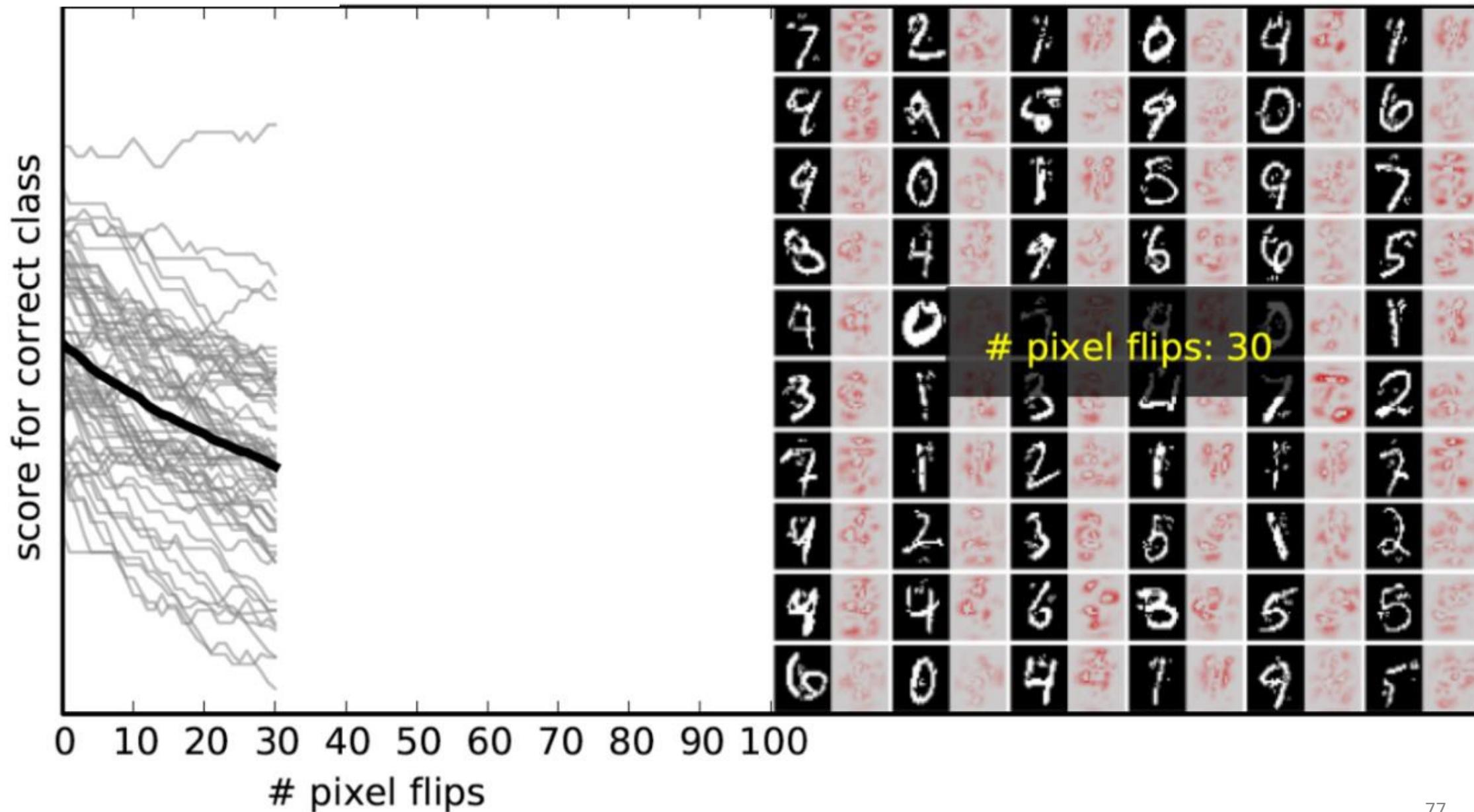


Selectivity on Saliency Map

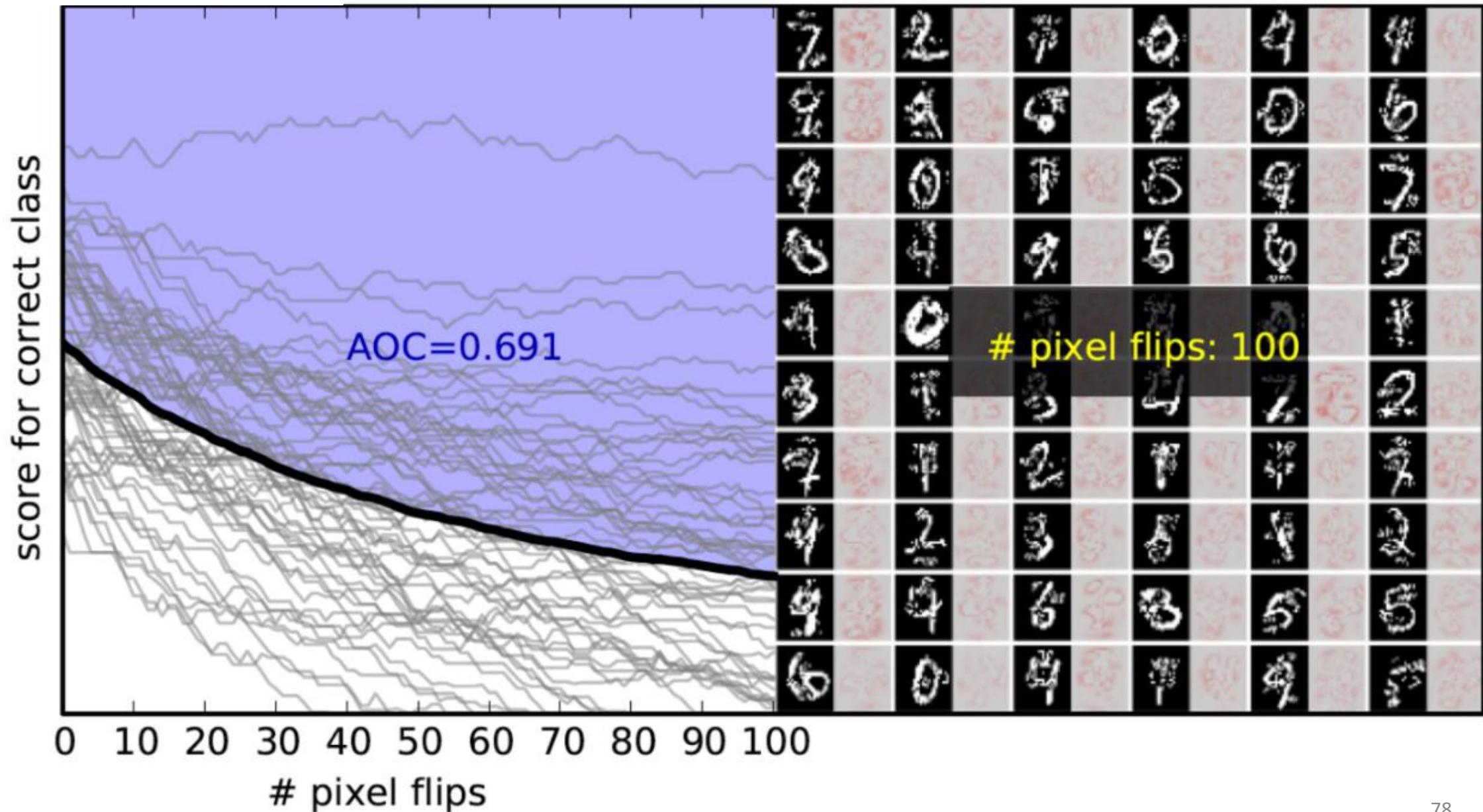
score for correct class



Selectivity on Saliency Map



Selectivity on Saliency Map



Evaluating Attribution Methods

Coherence

Class Sensitivity

Selectivity

ROAR / KAR

- Sensitivity may not be accurate
- Class probability may decrease because the DNN has never seen such image

Remove and Retrain (ROAR) / Keep and Retrain (KAR)

Measure how the performance of the classifier changes as features are removed based on the attribution method

- ROAR: replace $N\%$ of pixels estimated to be *most* important
- KAR: replace $N\%$ of pixels estimated to be *least* important
- Retrain DNN and measure change in test accuracy

Evaluating Attribution Methods

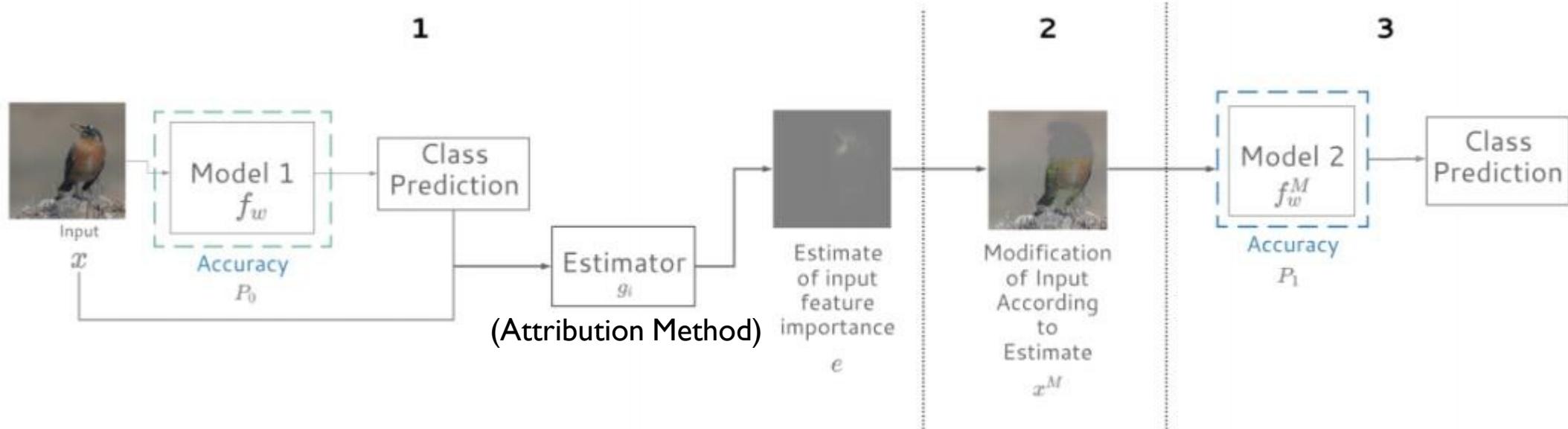
Coherence

Class Sensitivity

Selectivity

ROAR / KAR

ROAR – RemOve And Retrain



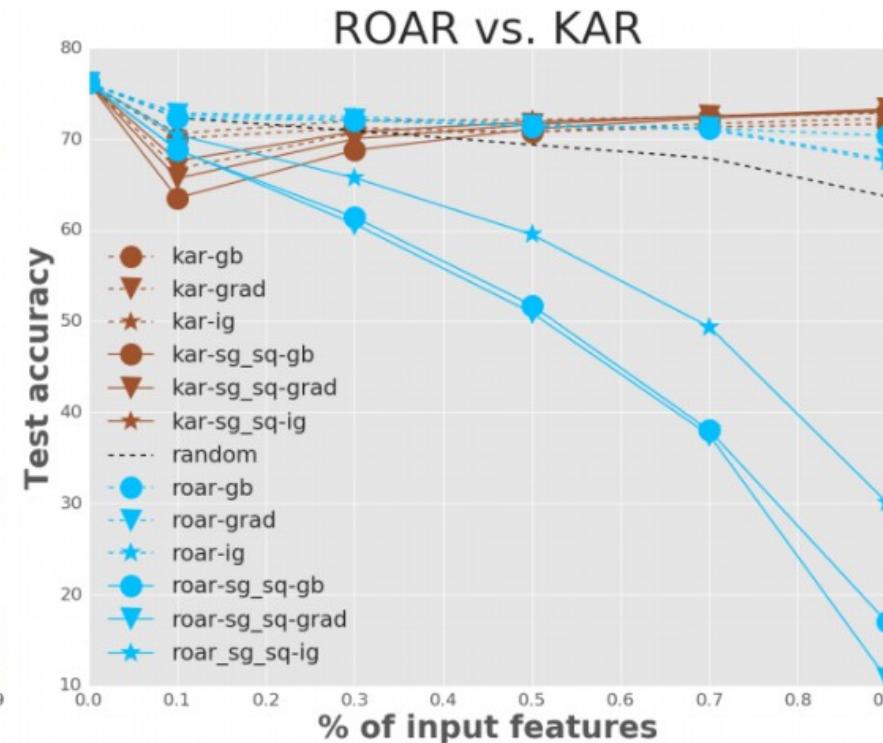
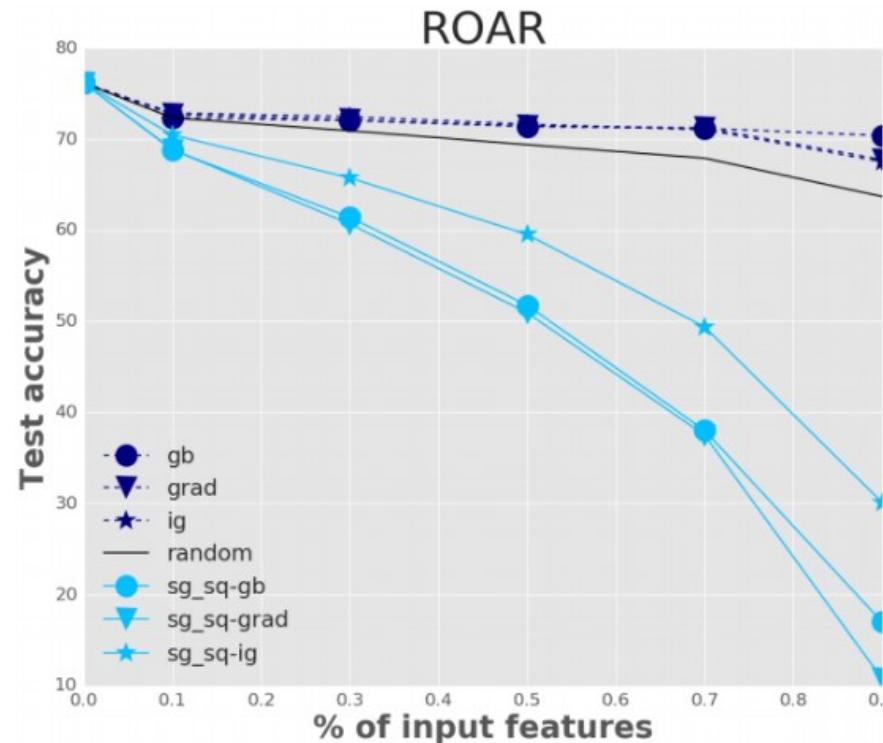
Evaluating Attribution Methods

Coherence

Class Sensitivity

Selectivity

ROAR / KAR



Interpretable Deep Learning

1. Intro to Interpretability

- 1a. **Interpretability definition:** Convert implicit NN information to human-interpretable information
- 1b. **Motivation:** Verify model works as intended; debug classifier; make discoveries; Right to explanation
- 1c. **Ante-hoc** (train interpretable model) vs. **Post-hoc** (interpret complex model; degree of “locality”)

2. Interpreting Deep Neural Networks

- 2a. **Interpreting Models** (macroscopic, understand internals) vs. **decisions** (microscopic, practical applications)
- 2b. **Interpreting Models:** Weight visualization, Surrogate model, Activation maximization, Example-based
- 2c. **Interpreting Decisions:**
 - Example-based
 - Attribution Methods: why are gradients noisy?
 - Gradient-based Attribution: SmoothGrad, Interior Gradient
 - Backprop-based Attribution: Deconvolution, Guided Backpropagation

3. Evaluating Attribution Methods

- 3a. **Qualitative: Coherence:** Attributions should highlight discriminative features / objects of interest
- 3b. **Qualitative: Class Sensitivity:** Attributions should be sensitive to class labels
- 3c. **Quantitative: Sensitivity:** Removing feature with high attribution → large decrease in class probability
- 3d. **Quantitative: ROAR & KAR.** Low class prob cuz image unseen → remove pixels, retrain, measure acc. drop

Summary

1. Introduction to Interpretability

- Interpretability is converting implicit information in DNN to (human) interpretable information
- Ante-hoc Interpretability vs. Post-hoc Interpretability
- Post-hoc interpretability techniques can be classified by degree of “locality”

2. Interpreting Deep Neural Networks

- Interpreting Models vs. Interpreting Decisions
- Interpreting Models: weight visualization, surrogate model, activation maximization, example-based
- Interpreting Decisions: example-based, attribution methods

3. Evaluating Attribution Methods

- Qualitative Evaluation Methods: coherence, class sensitivity
- Quantitative Evaluation Methods: Sensitivity, ROAR & KAR