

# Explainable AI (XAI)

# 11A

**CS 3244**  
**Machine Learning**



**NUS** | Computing

# Week 11A: Learning Outcomes

1. Describe multiple methods to interpret **feature importance**
2. Appropriately **interpret** feature attributions from each type of explanation
3. Describe how LIME explanations are generated
4. Describe how Grad-CAM explanations are generated

# Week 11A: Lecture Outline

1. Introduction
  1. Motivation for Explainable AI (XAI)
  2. Explaining Why: Feature Importance
2. Explanation techniques
  1. Glassbox Models (Linear Regression, Logistic Regression)
  2. Model-Agnostic Explanations (LIME)
  3. Model-Specific Explanations (Grad-CAM)
3. Human-centered XAI (not in exam)



# EXPLAINABLE AI EXPLAINED!

by [illegible]

# Case 1:

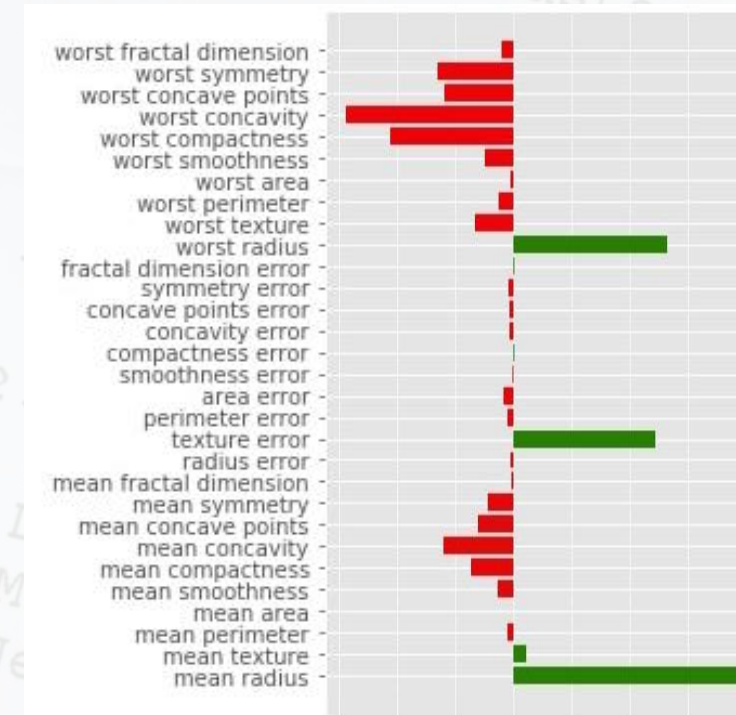
## Does patient have cancer?

Feature	Value
worst area	1315.00
mean radius	16.13
worst radius	20.96
area error	54.18
worst perimeter	136.80
worst texture	31.48
mean perimeter	108.10
smoothness error	0.01
mean area	798.80
mean concave points	0.10

**Prediction:** Cancer

Further reading: <https://coderzcolumn.com/tutorials/machine-learning/how-to-use-lime-to-understand-sklearn-models-predictions>

## Why??



Evidence for  
Cancer

Evidence for  
No Cancer

**Explanation:** Feature Attributions



## Case 2:

Is this skin cancer?

Why??



**Prediction:** Skin Cancer



**Explanation:** Highlighted Salient Region

Further reading: <https://towardsdatascience.com/medical-image-analysis-using-probabilistic-layers-and-grad-cam-42cc0118711f>

Image credit: <https://news.yale.edu/2019/11/13/yale-study-reveals-hyperhotspots-identifying-skin-cancer-risk>

# Why? Feature Importance

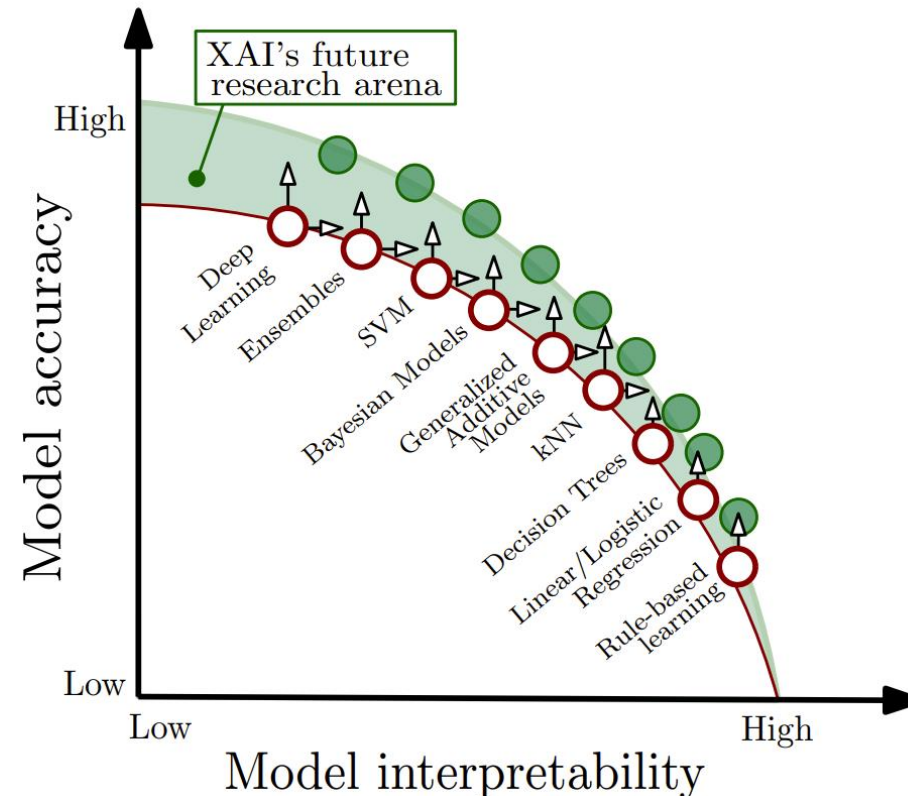
- Question
  - **Why?**
  - What caused this prediction?
- Explains
  - Which features are **important** for the prediction
  - In what way the features **influenced** the prediction
- Implementation
  - Weights in Linear / Logistic Regression
  - Surrogate Weights from LIME
  - Saliency Map of CNN

# Performance-Interpretability Trade-off

## High performance

### (Low interpretability)

- **Non-Linear** relationships
  - $y = w_1x_1^2 + \log(x_2)$
- **Interacting** features
  - $y = x_1 + x_2 + x_1x_2$
- **Many** features or parameters
- **Unrelatable** features



## High interpretability

### (Lower performance)

- **Linear** relationships
  - $y = w_1x_1 + w_2x_2$
- **Independent** features
  - $y = x_1 + x_2$
  - $y = f_1(x_1) + f_2(x_2)$
- **Few** features ( $x_r$ ), and parameters ( $w_r, \theta_r$ )
- **Semantically** meaningful features

Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., ... & Herrera, F. (2020). [Explainable Artificial Intelligence \(XAI\): Concepts, taxonomies, opportunities and challenges toward responsible AI](#). *Information Fusion*, 58, 82-115.



# Explainable AI (XAI) Categorization

		Explanability			Scope	
		Glassbox	Model-Agnostic	Model-Specific		Blackbox
Global	<ul style="list-style-type: none"><li>Linear Regression</li><li>Logistic Regression</li><li>Decision Tree</li></ul>	<ul style="list-style-type: none"><li>Collection of local explanations</li></ul>		<ul style="list-style-type: none"><li>Deep Neural Networks</li><li>Highly non-linear models</li></ul>		
Local	<ul style="list-style-type: none"><li>Examples (e.g., kNN)</li></ul>	<ul style="list-style-type: none"><li>LIME</li></ul>	<ul style="list-style-type: none"><li>Grad-CAM</li></ul>			

## Definitions

### **Glassbox** model

- Prediction model is implicitly interpretable

### **Model-Agnostic** explanation

- Uses surrogate model to provide simple explanations

### **Model-Specific** explanation

- Derived from calculations in specific prediction model

### **Blackbox** model

- Very uninterpretable, not transparent

### **Global** explanation

- Explains general behavior for all instances
- Does not change for different instances

### **Local** explanation

- Explains behavior for specific instance



# Interpreting Linear Regression

How would you interpret?  
Linear Regression

$$\begin{aligned}\hat{y} &= w_0 + w_1x_1 + w_2x_2 + \cdots + w_nx_n = \sum_{r=0}^n w_rx_r \\ &= \mathbf{w} \cdot \mathbf{x} = \mathbf{w}^\top \mathbf{x}\end{aligned}$$

# How would you interpret? Linear Regression

$$\begin{aligned}\hat{y} &= w_0 + w_1x_1 + w_2x_2 + \cdots + w_nx_n = \sum_{r=0}^n w_rx_r \\ &= \mathbf{w} \cdot \mathbf{x} = \mathbf{w}^\top \mathbf{x}\end{aligned}$$

$$\hat{y} \propto w_rx_r, \forall r$$

If  $w_1 = kw_2$ ,  
then  $w_1$  is  $k$  times more important than  $w_2$

## Weighted Sum Interpretation

**Bigger**  $w_r$  means

- **Larger** weight
- More **importance** for  $x_r$
- Direction? **Supportive (positive)** or **opposing (negative) influence**

## Gradient Interpretation

**Bigger**  $w_r$  means

- **Steeper** slope for  $x_r$  axis
  - Changes in  $x_r$  lead to bigger in  $\hat{y}$  changes
- More **importance** for  $x_r$
- Direction indicates **increasing** or **decreasing influence**

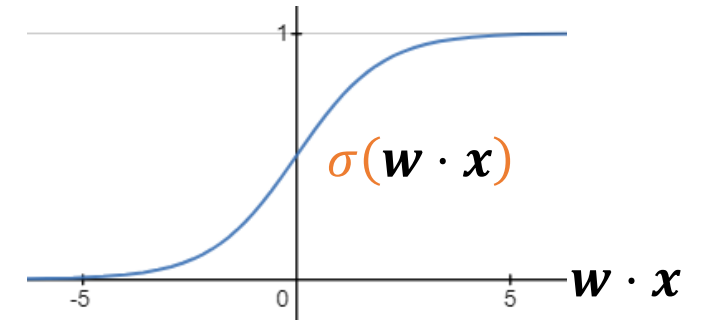




# Interpreting Logistic Regression

How would you interpret?  
Logistic Regression

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$
$$z = \mathbf{w} \cdot \mathbf{x} = \sum_{r=0}^n w_r x_r$$



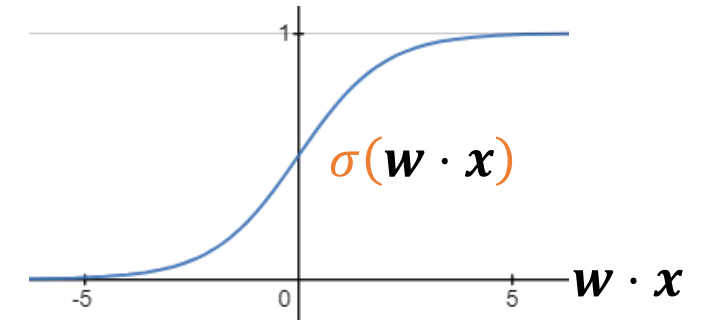
# How would you interpret? Logistic Regression

Not proportional,  
since non-linear  
If  $w_1 = kw_2$ ,  
then what is the  
relationship  
between  $w_1$  and  $w_2$ ?

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$z = \mathbf{w} \cdot \mathbf{x} = \sum_{r=0}^n w_r x_r$$

$$\hat{y} \not\propto w_r x_r, \forall r$$



## Weighted Sum Interpretation

**Bigger**  $w_r$  means

- **Larger** importance
- Direction indicates **influence**

$\hat{y} = \sigma(f)$  is  
positively monotonic, i.e.,  
 $f_1 > f_2 \Rightarrow \sigma(f_1) > \sigma(f_2)$   
Bigger  $f \Rightarrow$  bigger  $\hat{y}$

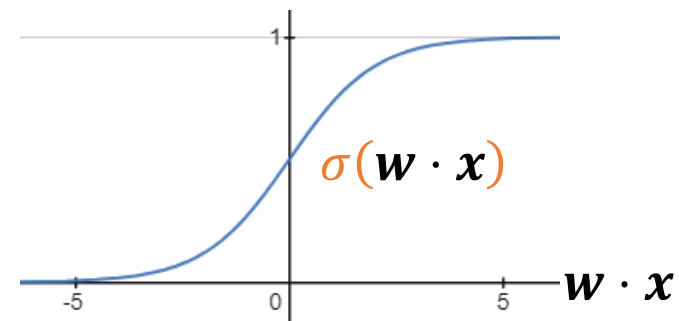
# Interpreting weights in Logistic Regression

Probability

$$P(\hat{y}) = p = \hat{y} = \frac{1}{1 + e^{-w \cdot x}} = \sigma(p)$$

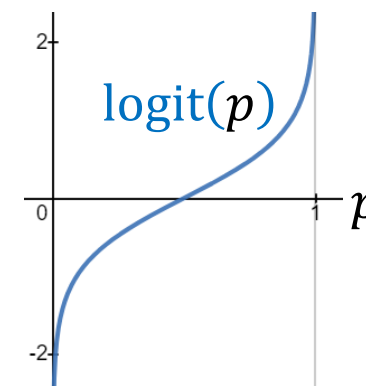
$$\frac{1}{p} = 1 + e^{-w \cdot x}$$

$$\frac{1-p}{p} = \frac{1}{p} - 1 = \frac{1}{1 + e^{-w \cdot x}} - 1 = \frac{1}{e^{w \cdot x}}$$



Odds Ratio

$$\frac{p}{1-p} = e^{w \cdot x}$$



Log Odds Ratio  $\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = w \cdot x$

$$\text{logit}(P(\hat{y})) \propto w_r x_r, \forall r$$

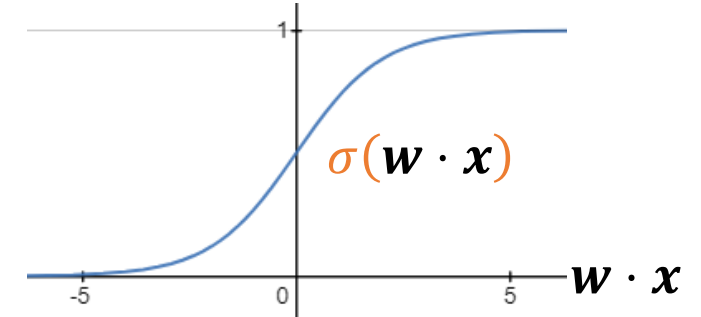
If  $w_1 = k w_2$ ,  
then log odds ratio of  $w_1$  is  
 $k$  times bigger than of  $w_2$



# How would you interpret? Logistic Regression

$$\hat{y} = \sigma(f) = \frac{1}{1 + e^{-f}}$$

$$f = \mathbf{w} \cdot \mathbf{x} = \sum_{r=0}^n w_r x_r$$



$$\text{logit}(P(\hat{y})) \propto w_r x_r, \forall r$$

If  $w_1 = k w_2$ ,  
then log odds ratio of  $w_1$  is  
 $k$  times bigger than of  $w_2$

## Weighted Sum Interpretation

**Bigger**  $w_r$  means

- **Larger** importance
- Direction indicates **influence**

## Gradient Interpretation

**Bigger**  $w_r$  means?

- Steepness? Sigmoid bounded between 0 and 1
- Direction in 2D (or higher)?



1  $f(x) = \frac{1}{1 + e^{-(w_0 + w_1 x)}}$



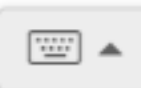
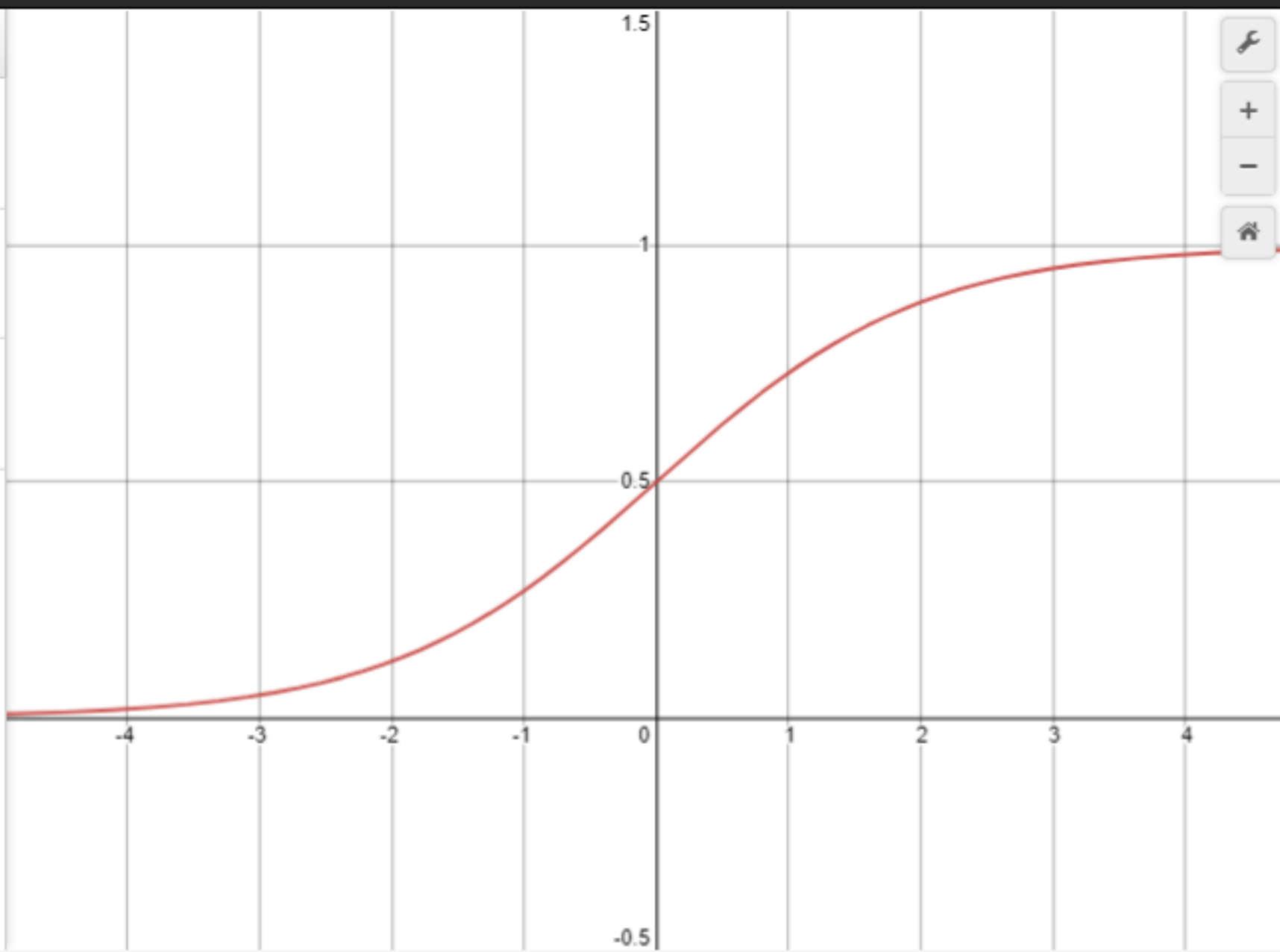
2  $w_1 = 1$

-10 10



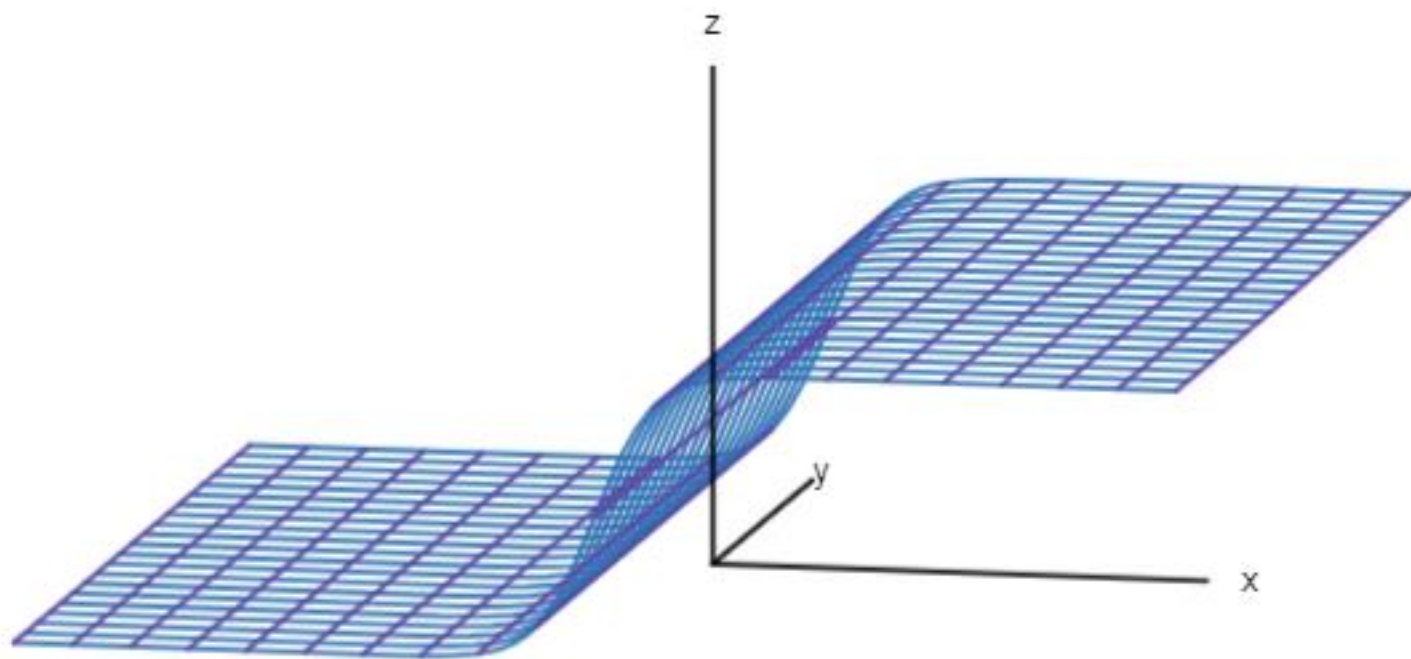
3  $w_0 = 0$

-10 10





- 1 Display controls
- 8
- 9 This graph has lots of different ways of plotting 3D objects.
- 10 To turn on a plotter, click on its folder icon on the left hand side. Settings for the plotters can be configured inside the respective folders.
- 11 To plot more than one function at once, simply replace the single function with a list of functions. E.g. below, change  $g(x, y) = \cos(xy+T)$  to  $g(x, y) = [\cos(xy+T), x^2]$ .
- 12 Some things work very well and run very smoothly, others not so much. If something is running particularly slowly, then turn off all unnecessary features i.e. arrowheads on the axes or on vector fields. If this doesn't help, then try deleting unnecessary parts of the program e.g. if you are only using the function plotter then delete the vector field plotter.



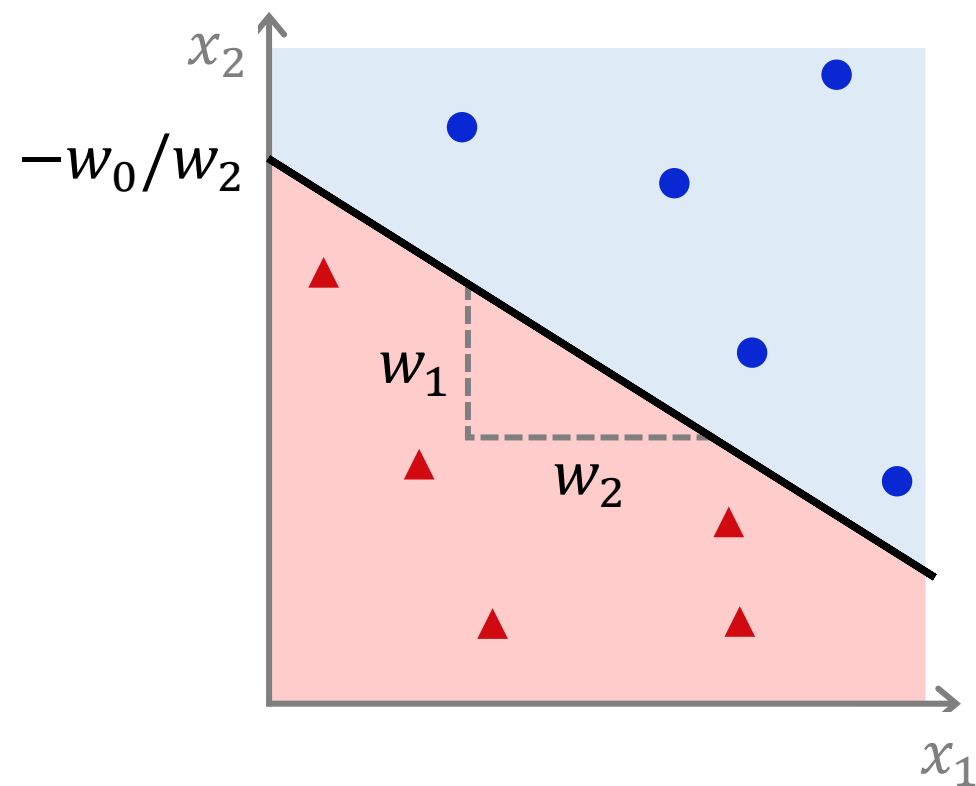
# Linear Classification

$$w_2x_2 + w_1x_1 + w_0 = 0$$

$$\sum_{r=0}^n w_r x_r = 0 \quad x_0 = 1$$

$$\sum_{r=0}^n w_r x_r > 0 \quad \sum_{r=0}^n w_r x_r \leq 0$$

$$\hat{y} = \sigma \left( \sum_{r=0}^n w_r x_r \right)$$







$$w_0 + w_1x + w_2y = 0$$



$$w_0 = -1.$$

-10  10



$$w_1 = 1$$

-10  10



$$w_2 = 1$$

-10  10

5



powered by  
desmos

-3 -2 -1 0 1 2 3 4 5 6 7

-1

-2

4

3

2

1

# Linear Classification

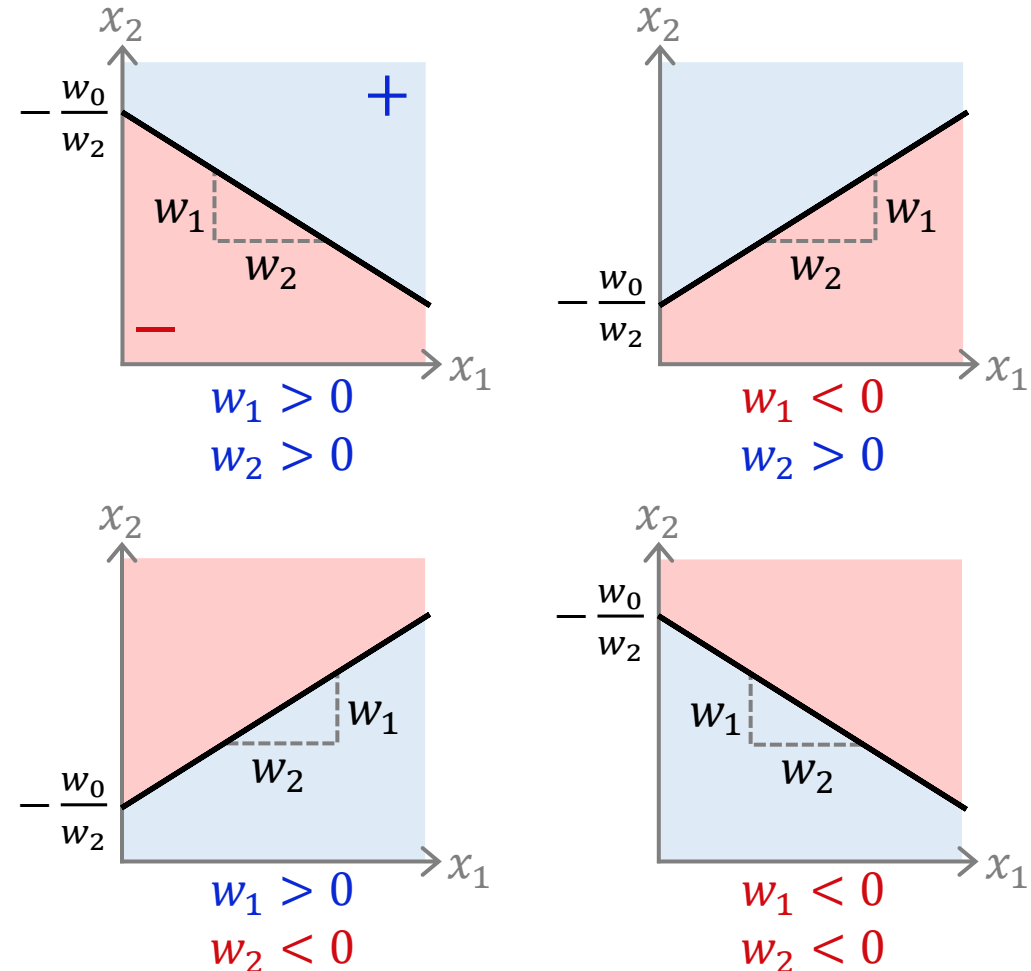
$$w_2x_2 + w_1x_1 + w_0 = 0$$

$$\sum_{r=0}^n w_r x_r = 0 \quad x_0 = 1$$

$$\sum_{r=0}^n w_r x_r > 0 \quad \sum_{r=0}^n w_r x_r \leq 0$$

$$\hat{y} = \sigma \left( \sum_{r=0}^n w_r x_r \right)$$

**Weight sign** indicates direction towards **pos/neg** prediction

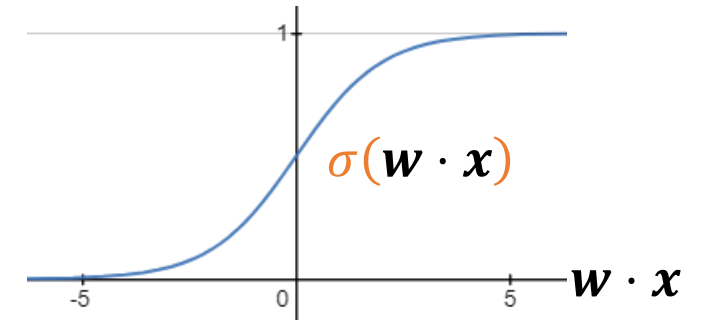


# How would you interpret?

## Logistic Regression

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$f = \mathbf{w} \cdot \mathbf{x} = \sum_{r=0}^n w_r x_r$$



$$\text{logit}(P(\hat{y})) \propto w_r x_r, \forall r$$

If  $w_1 = k w_2$ , then **log odds ratio** of  $w_1$  is  $k$  times bigger than of  $w_2$

### Weighted Sum Interpretation

**Bigger**  $w_r$  means

- **Larger** importance
- Direction indicates **influence**

### Gradient Interpretation

**Bigger**  $w_r$  means?

- **Steeper** slope for  $x_r$  near decision boundary
- **Decision boundary** more *perpendicular* to  $x_r$
- Weight sign indicates direction of **pos/neg** prediction



*Questions!*



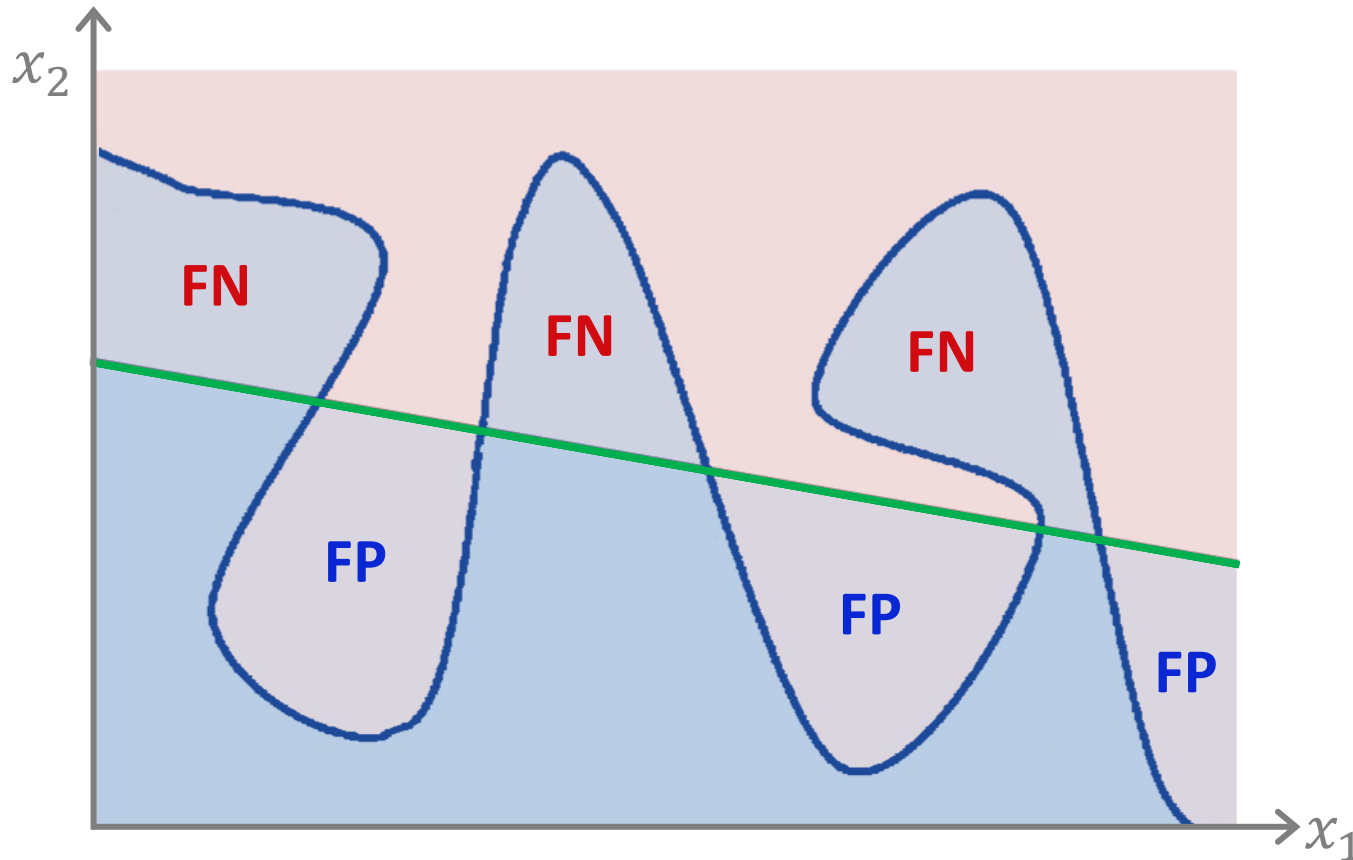




# Local Interpretable Model-agnostic Explanations LIME

How to describe with just  $x_1$  and  $x_2$ ?

# Non-Linear Decision Boundary $f(\mathbf{x})$



## Prediction Model

$f(\mathbf{x})$

- Non-linear model of  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \end{pmatrix}$
- Shown as curvy decision boundary

## Explanation: Linear Model

$$g(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} = \sum_{r=0}^n w_r x_r$$

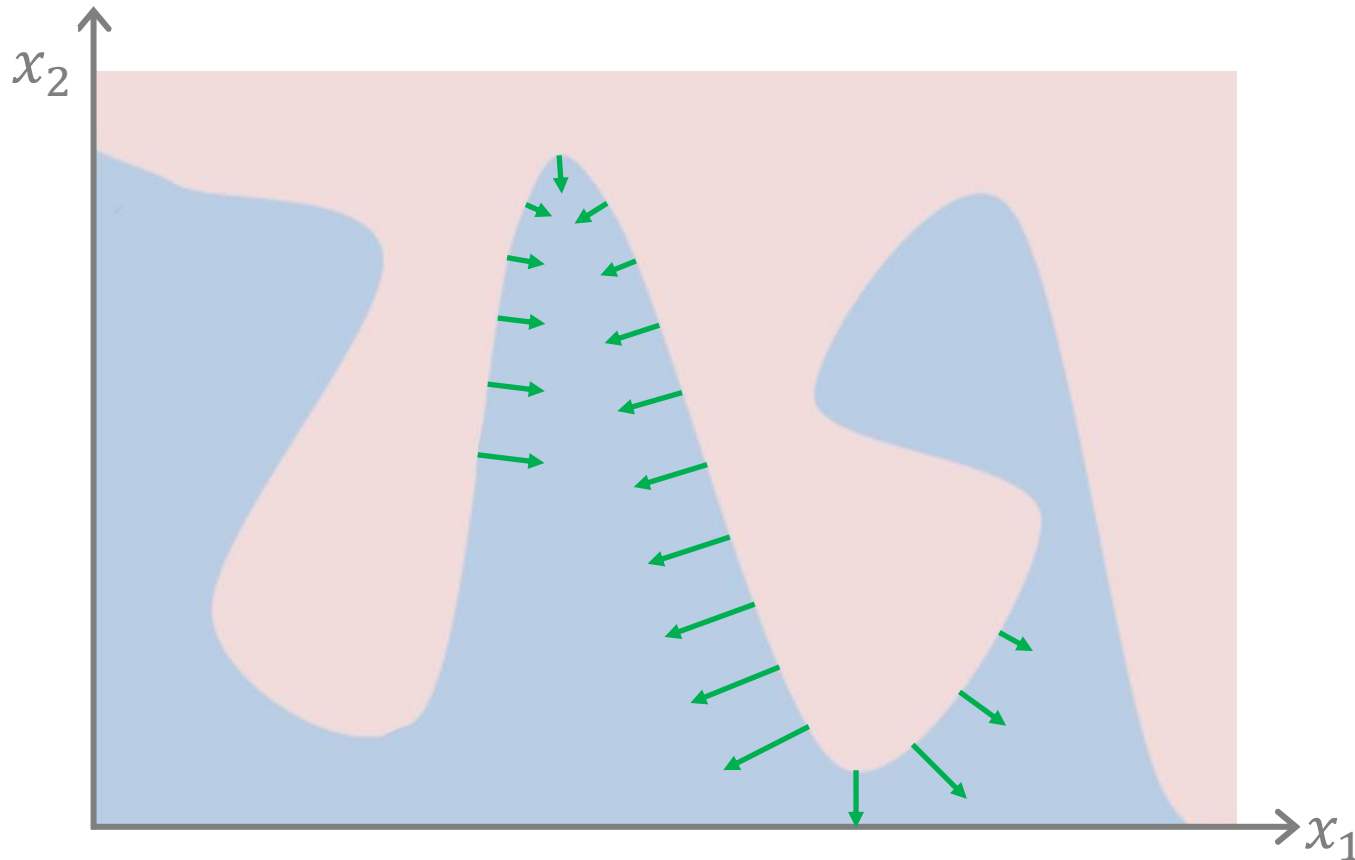
- Simple to interpret
- But **too many errors** between  $g$  and  $f$

$$L = f(\mathbf{x}) - g(\mathbf{x})$$

# How to describe with just $x_1$ and $x_2$ ? Non-Linear Decision Boundary

$$\nabla f(\mathbf{x}) \neq \begin{pmatrix} w_1 \\ w_2 \\ \vdots \end{pmatrix} = \mathbf{w}$$

- Since  $f$  is not linear



## Prediction Model

$$f(\mathbf{x})$$

- Non-linear model of  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \end{pmatrix}$
- Shown as curvy decision boundary

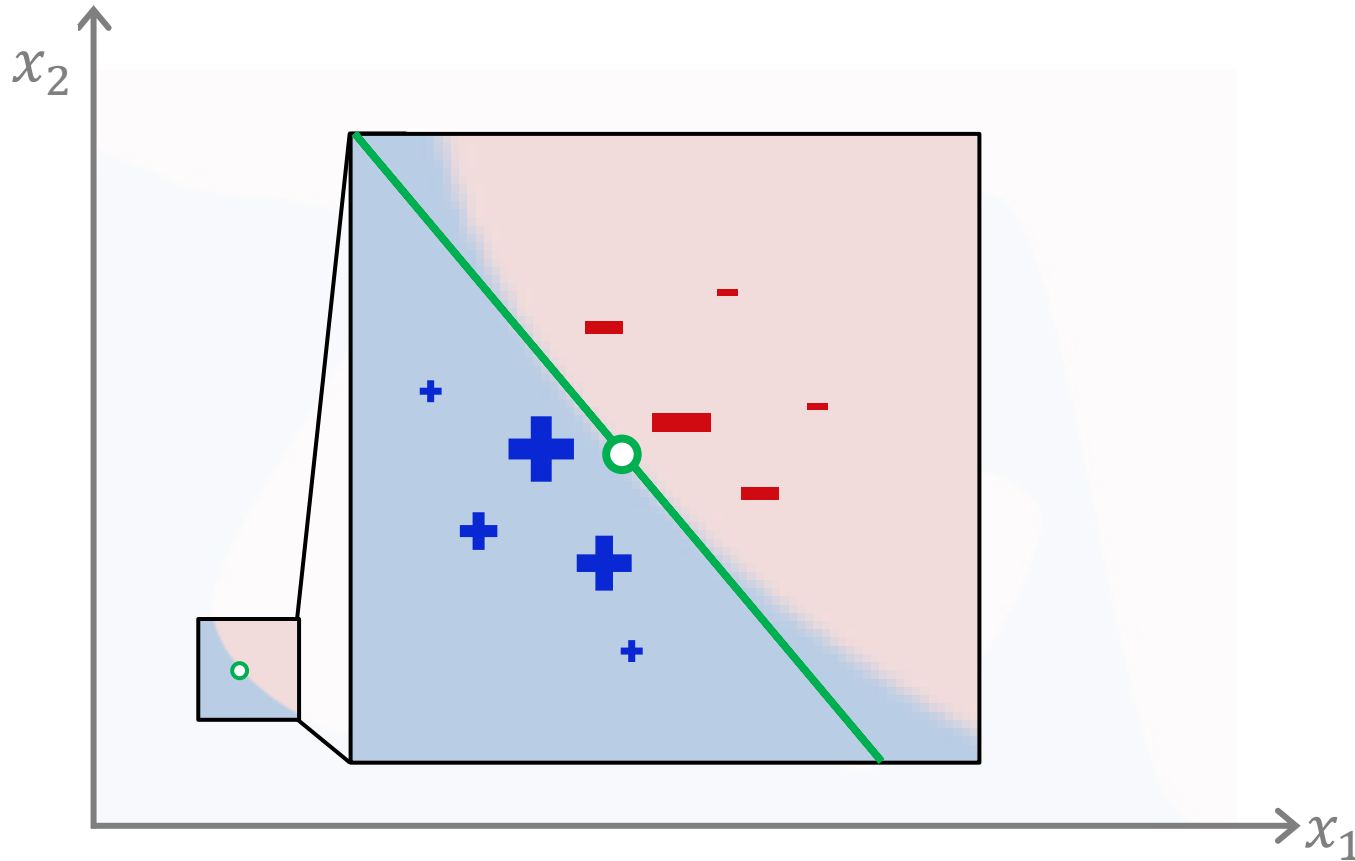
## Explanation: Gradients

$$g(\mathbf{x}) = \nabla f(\mathbf{x}) = \frac{\partial f}{\partial \mathbf{x}} = \begin{pmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \vdots \end{pmatrix}$$

- **Steepness** for each feature  $x_r$
- **Difficult to remember**, since gradients are different for each instance (point)

# LIME

## Local Interpretable Model-agnostic Explanations



### Prediction Model

$$f(\mathbf{x})$$

- Non-linear model of  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \end{pmatrix}$
- Shown as curvy decision boundary

### Explanation: LIME

1. Starting with **instance**  $\mathbf{x}$  to explain
2. Focus on **Local** region
3. Training set as **neighbors**  $\mathbf{x}^{(\eta)} \in X^{(\eta)}$
4. Train **surrogate model**, e.g., linear:

$$g(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} = \sum_{r=0}^n w_r x_r$$

# LIME

Python API:

<https://github.com/marcotcr/lime>

Find “best” explainer  $g$  that minimizes  $\xi(\mathbf{x})$

$$\xi(\mathbf{x}) = \underset{g \in G}{\operatorname{argmin}} \left( \overbrace{L(f, g, \pi_x)}^{\text{“Faithful”}} + \overbrace{\Omega(g)}^{\text{“Simple”}} \right)$$

$f$  is the **predictor** function (model)

$g$  is the **explainer** function (model)

$\pi_x(\mathbf{x}^{\langle \eta \rangle})$  is the neighbor **proximity** function

- E.g., exponential decay  $\exp\left(-\left(d(\mathbf{x}, \mathbf{x}^{\langle \eta \rangle})\right)^2\right)$

$L(f, g, \pi_x)$  is the **locally-weighted error loss** function between predictor  $f$  and explainer  $g$

$$L(f, g, \pi_x) = \sum_{\mathbf{x}^{\langle \eta \rangle} \in X^{\langle \eta \rangle}} \pi_x(\mathbf{x}^{\langle \eta \rangle}) \left( f(\mathbf{x}^{\langle \eta \rangle}) - g(\mathbf{x}^{\langle \eta \rangle}) \right)^2$$

$\Omega(g)$  is the **sparsity regularization**

- Want simpler explanation
  - $\Rightarrow$  *fewer weights*
  - $\Rightarrow$  Lasso (L1 norm)
- Penalizes if total weights is too large
- $\lambda$  is hyperparameter on how much to penalize

$$\Omega(g) = \lambda \|w_r\|_1 = \lambda \sum_{r=1}^n w_r$$

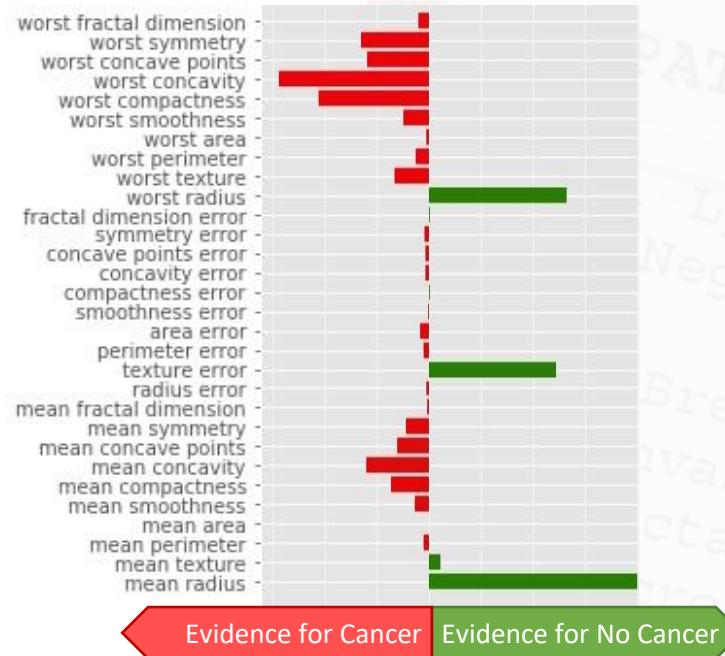
# Case 1: Does patient have cancer?

Why do the two set of weights **differ**?

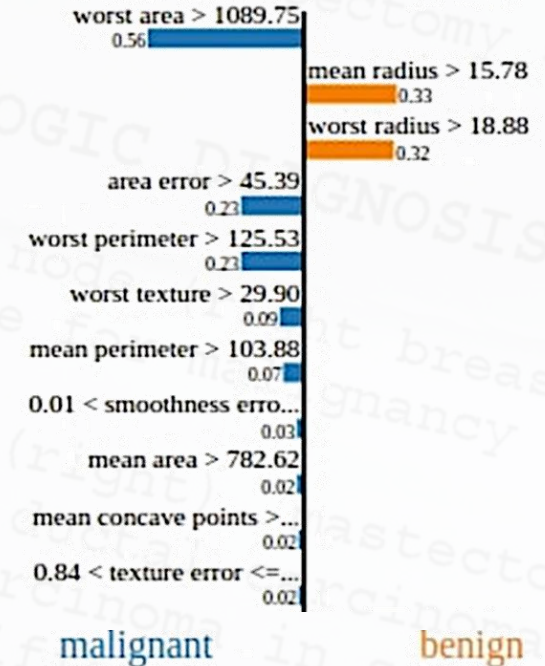
Instance  $x$

Feature	Value
worst area	1315.00
mean radius	16.13
worst radius	20.96
area error	54.18
worst perimeter	136.80
worst texture	31.48
mean perimeter	108.10
smoothness error	0.01
mean area	798.80
mean concave points	0.10

Logistic Regression



LIME



Prediction  $\hat{y} = \underline{\text{Cancer}}$

Weights  $w$  of  
surrogate explanation  $f$

Weights  $w$  of  
surrogate explanation  $g$

Further reading: <https://coderzcolumn.com/tutorials/machine-learning/how-to-use-lime-to-understand-sklearn-models-predictions>





*Questions!*



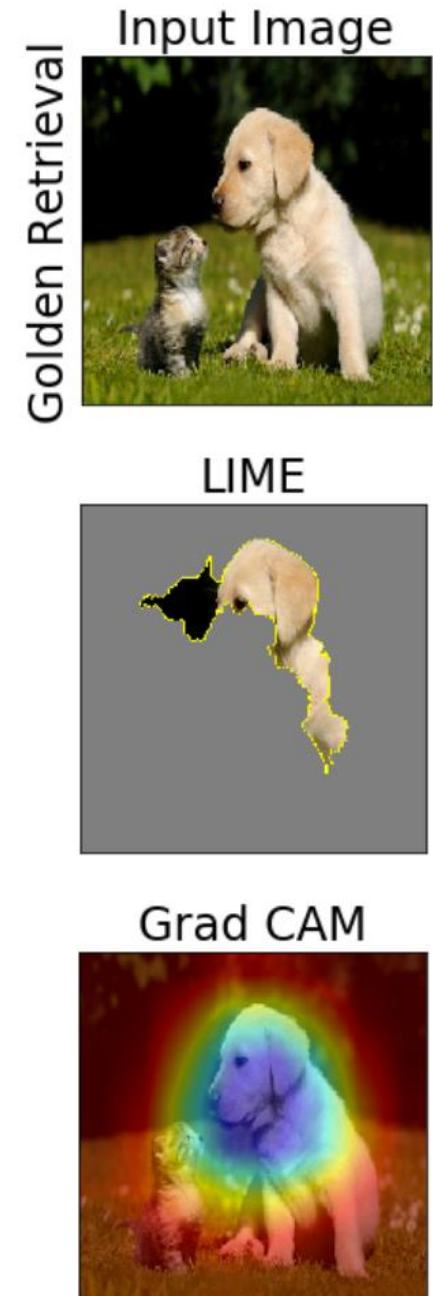


# Gradient-weighted Class Activation Mapping Grad-CAM

# Explaining Image Predictions

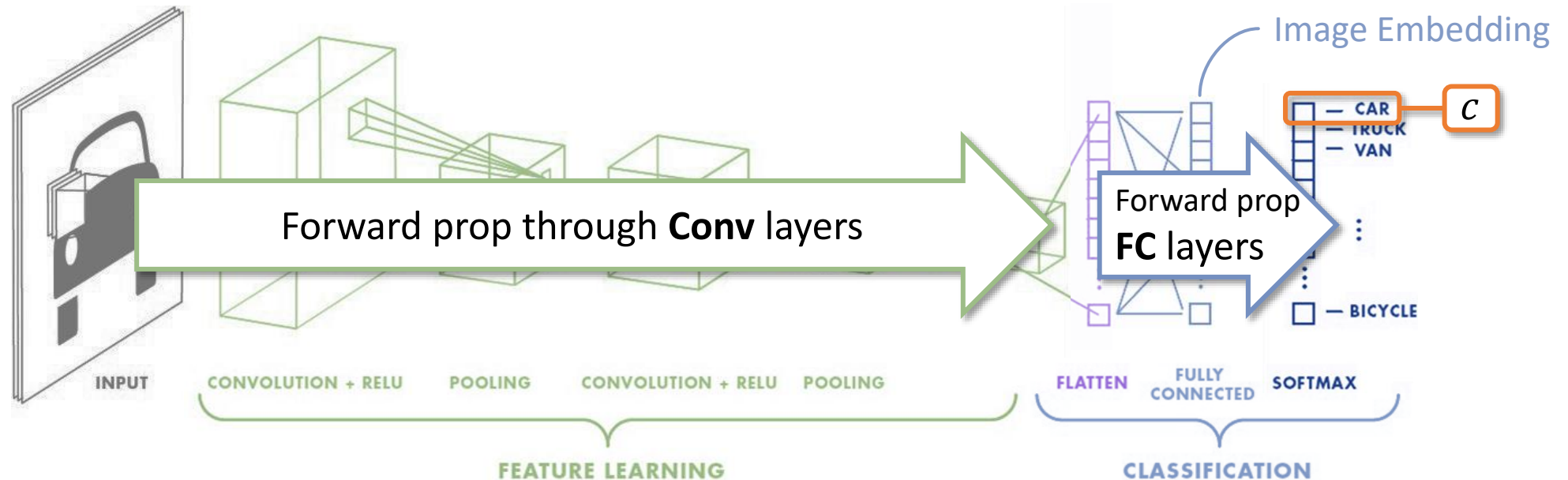
- LIME to explain image prediction?
- What are the input features?
  - Feature = Pixels?
    - **Too many** features
  - Need “super pixels”
- Another way: Attribution → Saliency Map
  - Feature = Activation Map
  - Grad-CAM

Image credit: <https://arxiv.org/pdf/1908.04389.pdf>





# Convolutional Neural Network



## Key concepts

### ① Learn Spatial Feature

- Series of multiple convolution + pooling layers
- Progressively learn more diverse and higher-level features

### ② Flattening

- Convert to fixed-length 1D vector

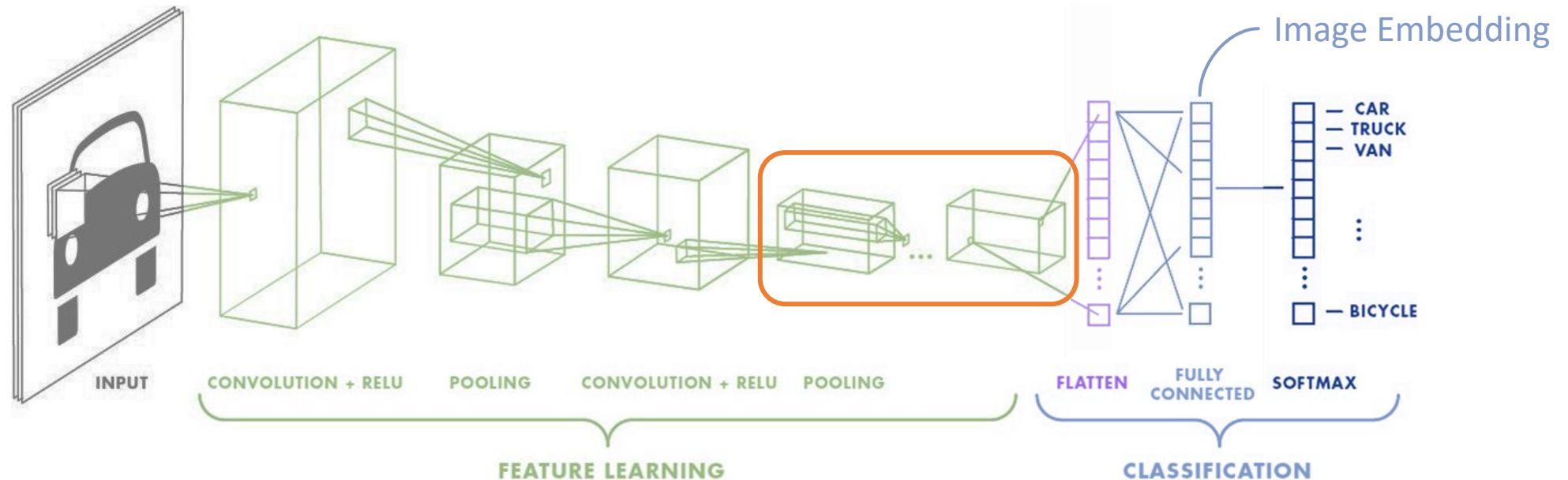
### ③ Learn Nonlinear Features

- With fully connected layers (regular neurons)
- Learns nonlinear relations with multiple layers

### ④ Classification

- Softmax := Multiclass Logistic Regression
- Feature input = image embedding vector (typically large vector)

# Convolutional Neural Network



## Key concepts

### ① Learn Spatial Feature

- Series of multiple convolution + pooling layers
- Progressively learn more diverse and higher-level features

### ② Flattening

- Convert to fixed-length 1D vector

### ③ Learn Nonlinear Features

- With fully connected layers (regular neurons)
- Learns nonlinear relations with multiple layers

### ④ Classification

- Softmax := Multiclass Logistic Regression
- Feature input = image embedding vector (typically large vector)

# Convolutional Layer: Feature Kernels & Feature Maps

$$\mathbf{X}^{[0]} \rightarrow g^{[1]}(\mathbf{W}^{[1]} * \mathbf{X}^{[0]}) = \mathbf{A}^{[1]} \xrightarrow{\text{Pooling}} g^{[2]}(\mathbf{W}^{[2]} * \mathbf{X}^{[1]}) = \mathbf{A}^{[2]}$$

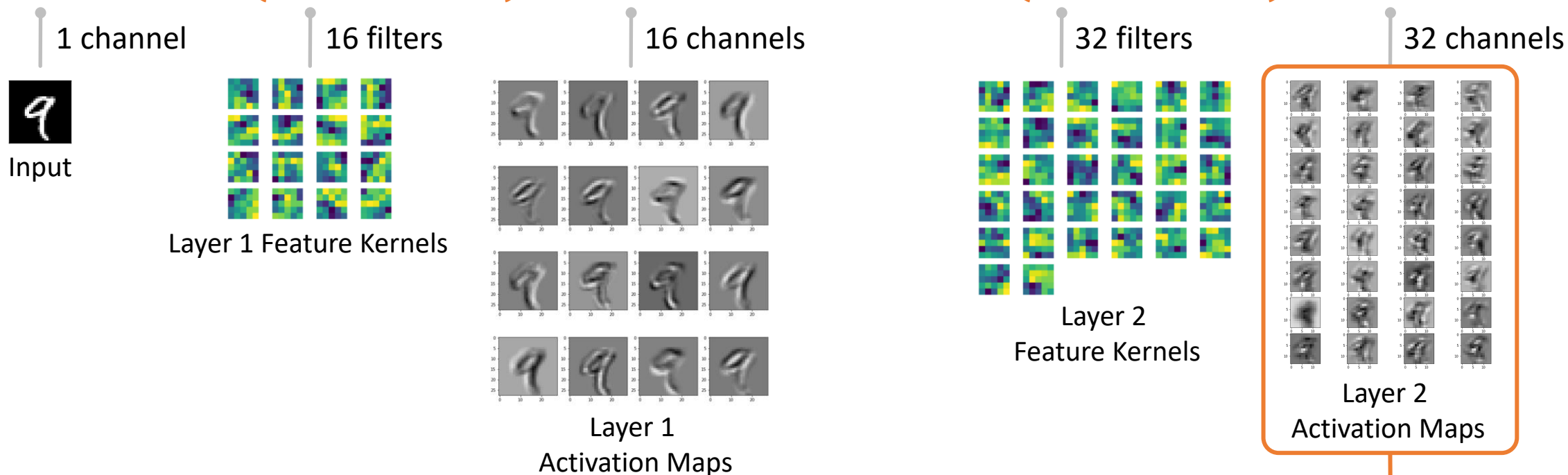
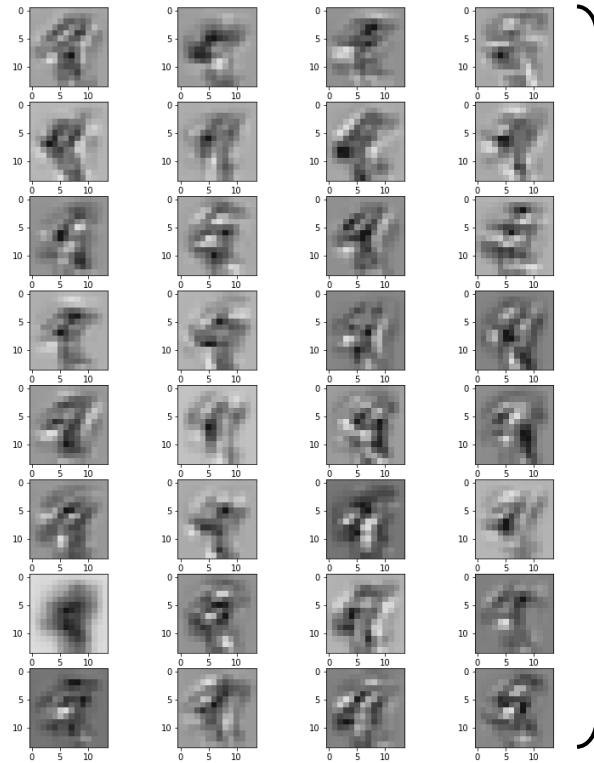


Image credit: <https://medium.com/dataseries/visualizing-the-feature-maps-and-filters-by-convolutional-neural-networks-e1462340518e>



# Multi-Channel Activation Maps (layers diagram)

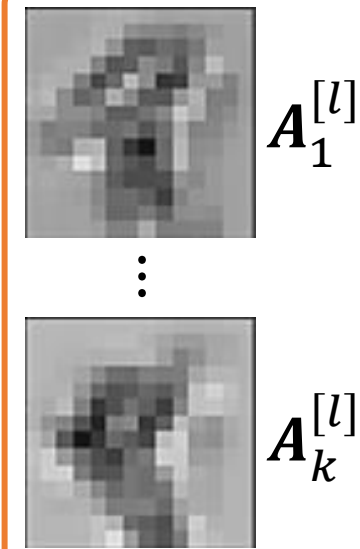


$$\mathbf{A}^{[l]} = (\mathbf{A}_1^{[l]} \quad \dots \quad \mathbf{A}_k^{[l]})$$

$$\mathbf{A}^{[l]} = \begin{pmatrix} a_{111}^{[l]} & \dots & a_{1w1}^{[l]} \\ \vdots & \ddots & \vdots \\ a_{h11}^{[l]} & \dots & a_{hw1}^{[l]} \\ \vdots & \ddots & \vdots \\ a_{11k}^{[l]} & \dots & a_{1wk}^{[l]} \\ \vdots & \ddots & \vdots \\ a_{h1k}^{[l]} & \dots & a_{hwk}^{[l]} \end{pmatrix}$$

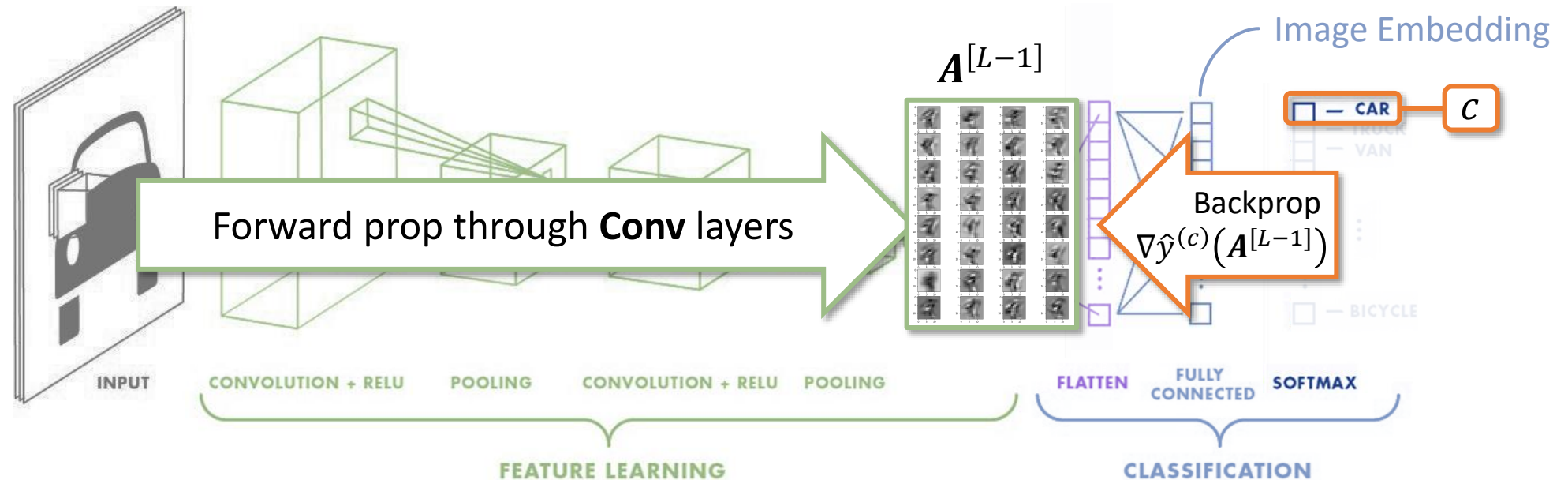
We know the **Activation Maps**,  
But which is more **important**?

Calculate with **Gradients**:  
Steeper  $\Rightarrow$  More **important**



**Activation** of  $k$ th channel at pixel  
position  $(h, w)$  in layer  $l$

# Convolutional Neural Network



## Key concepts

### ① Learn Spatial Feature

- Series of multiple convolution + pooling layers
- Progressively learn more diverse and higher-level features

### ② Flattening

- Convert to fixed-length 1D vector

### ③ Learn Nonlinear Features

- With fully connected layers (regular neurons)
- Learns nonlinear relations with multiple layers

### ④ Classification

- Softmax := Multiclass Logistic Regression
- Feature input = image embedding vector (typically large vector)

# Grad-CAM

## Gradient-Weighted Class Activation Maps

$$\frac{\partial \hat{y}^{(c)}}{\partial a_{hwk}^{[L]}} = \frac{\partial f^{[L]}}{\partial a_{hwk}^{[L]}} \frac{\partial g^{[L]}}{\partial f^{[L]}} \frac{\partial \hat{y}^{(c)}}{\partial g^{[L]}}$$

$$\frac{\partial \hat{y}^{(c)}}{\partial \mathbf{A}_k^{[L]}} = \begin{pmatrix} \partial \hat{y}^{(c)} / \partial a_{11k}^{[L]} & \dots & \partial \hat{y}^{(c)} / \partial a_{1wk}^{[L]} \\ \vdots & \ddots & \vdots \\ \partial \hat{y}^{(c)} / \partial a_{h1k}^{[L]} & \dots & \partial \hat{y}^{(c)} / \partial a_{hwk}^{[L]} \end{pmatrix}$$

$$\left\| \frac{\partial \hat{y}^{(c)}}{\partial \mathbf{A}_k^{[L]}} \right\|_{1,1} = \sum_{ij} \frac{\partial \hat{y}^{(c)}}{\partial a_{ijk}^{[L]}} = \frac{\partial \hat{y}^{(c)}}{\partial a_{11k}^{[L]}} + \dots + \frac{\partial \hat{y}^{(c)}}{\partial a_{hwk}^{[L]}}$$

**CAM**  
2D Matrix  
 $\{h^{[L]} \times w^{[L]}\}$

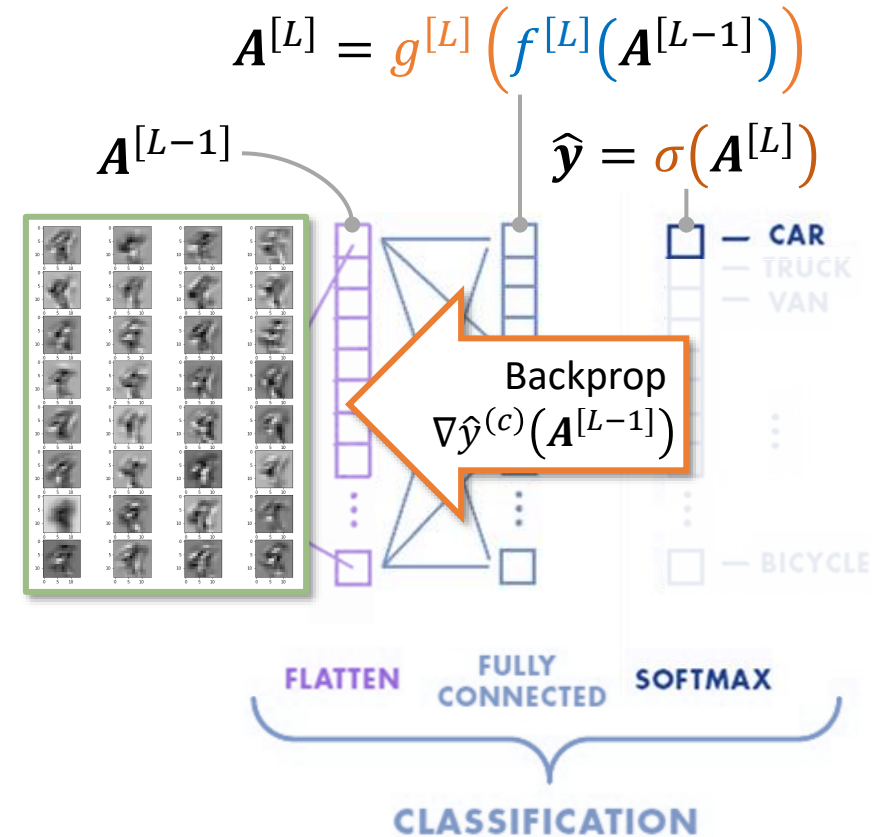
$$g(\hat{y}^{(c)}) = \text{ReLU} \left( \sum_k \alpha_k^{(c)} \mathbf{A}_k^{[L-1]} \right),$$

**Weighted Sum**  
Keep positive activations only

$$\alpha_k^{(c)} = \left\| \frac{\partial \hat{y}^{(c)}}{\partial \mathbf{A}_k^{[L-1]}} \right\|_{1,1}, \mathbf{A}_k^{[L]} = \begin{pmatrix} a_{11k}^{[L-1]} & \dots & \dots \\ \vdots & \ddots & \vdots \\ \dots & \dots & a_{hwk}^{[L-1]} \end{pmatrix}$$

**Importance Weight**  
of  $k$ th filter for  $c$ th class

**Activation Map**  
of  $k$ th filter in the  $L$ th layer



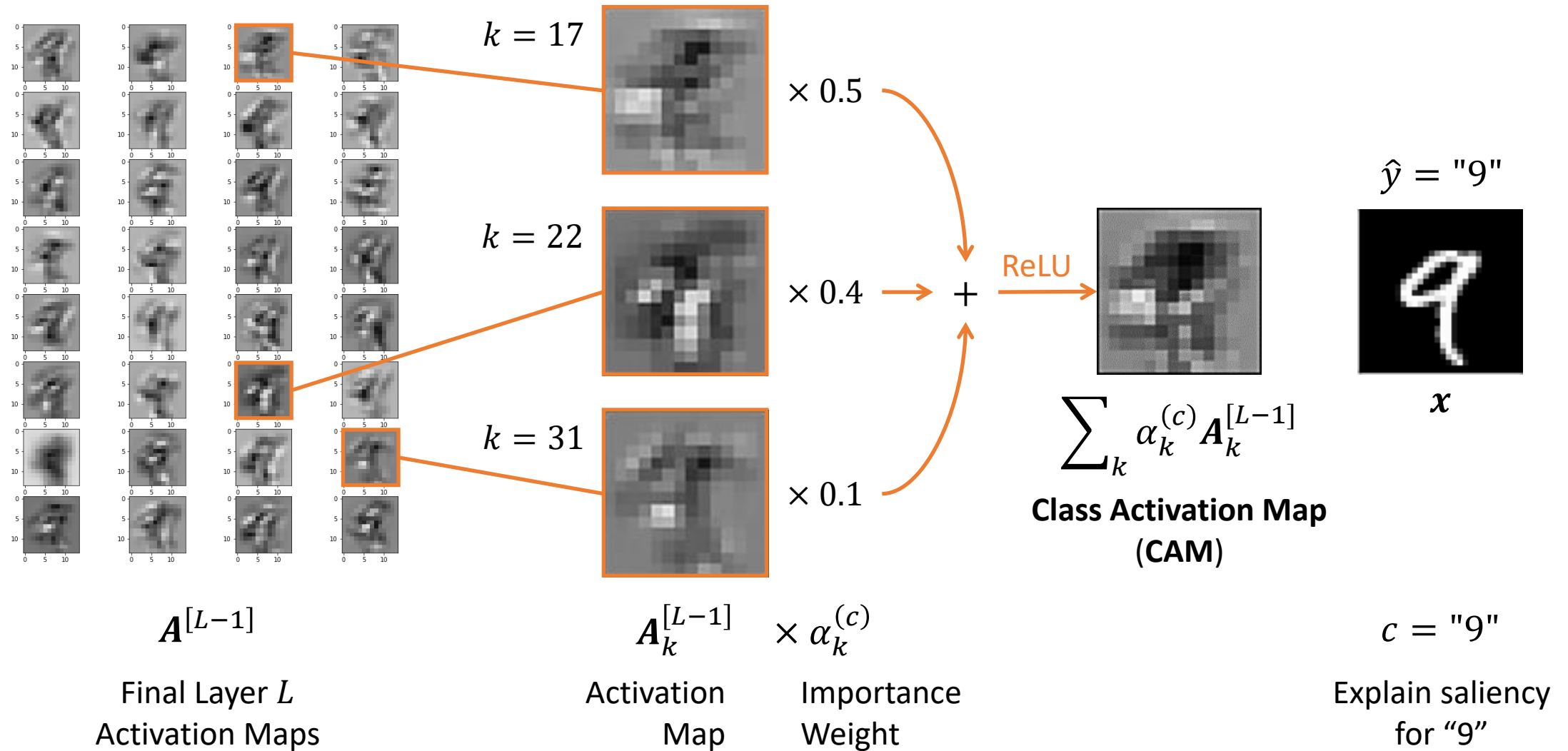
# Grad-CAM Steps

1. Compute Activation Maps  $\mathbf{A}^{[L]}$  of last conv layer  $L$ 
  1. via Forward Propagation
2. Choose class label  $c$  to explain about (e.g., predict “9”, “car”)
3. Filter prediction  $\hat{\mathbf{y}}$  to be about class  $c$

1. Given:  $\hat{\mathbf{y}} = \begin{pmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \hat{y}^{(c)} \\ \hat{y}^{(n)} \end{pmatrix}$ ,  $\mathbf{e}^{(c)} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$ , then  $\hat{\mathbf{y}}^{(c)} = \hat{\mathbf{y}} \circ \mathbf{e}^{(c)} = \begin{pmatrix} 0 \\ 0 \\ y^c \\ 0 \end{pmatrix}$

2. To generate explanation only for that class  $c$
4. Compute importance weight  $\alpha_k^{(c)}$  for each Activation Map  $\mathbf{A}_k^{[L]}$ 
  1. Backprop from  $\hat{\mathbf{y}}^{(c)}$  to get gradients at last conv layer
    - Note that gradient is **relative to activations**, not weights
5. Compute weighted sum with ReLU to get **Class Activation Map**

# Grad-CAM example: Why did the CNN predict “9”?



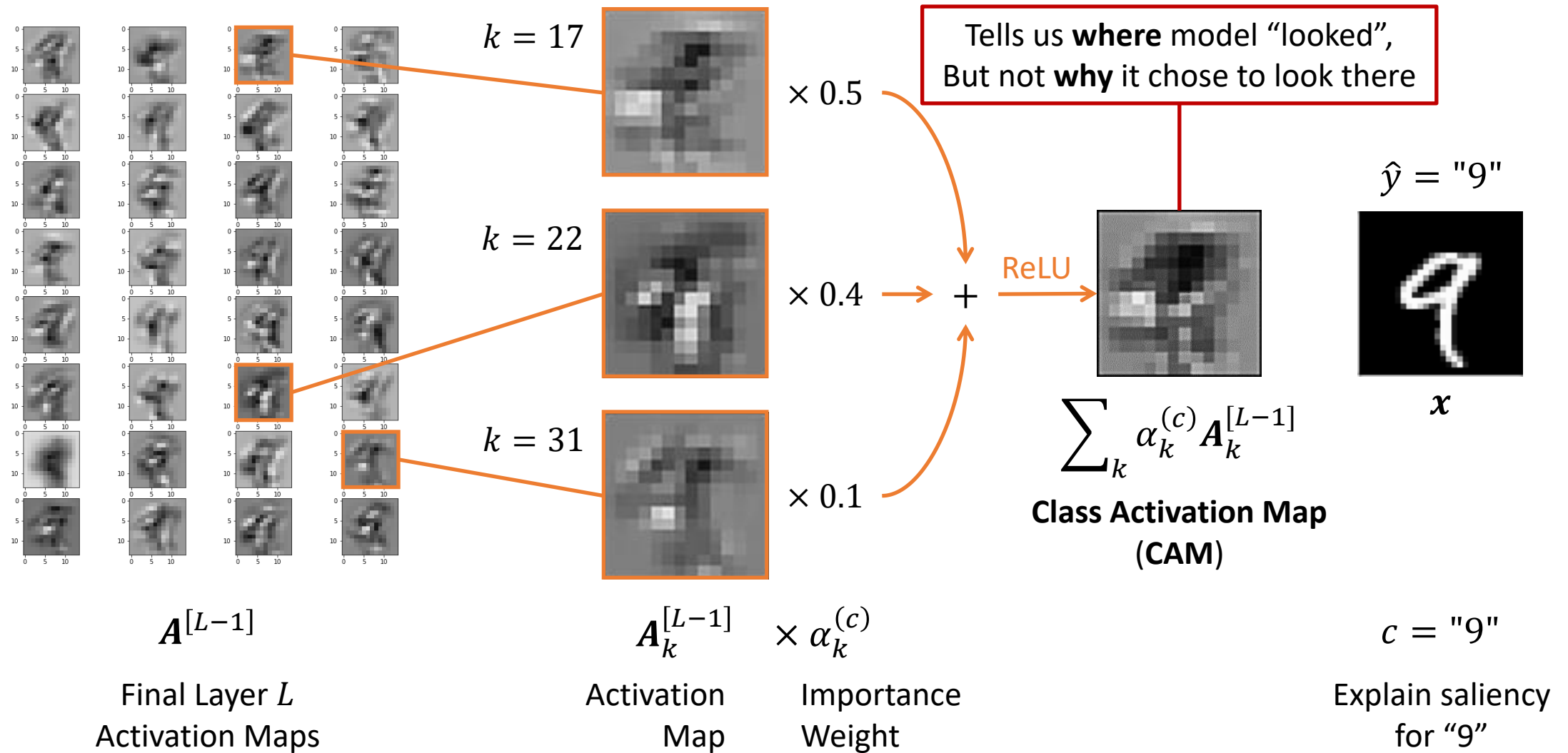


*Questions!*





# Grad-CAM example: Why did the CNN predict “9”?



# W10 Pre-Lecture Task (due before next Mon)

## Watch

- [Who Invented A.I.? - The Pioneers of Our Future](#) by [ColdFusion](#)

## Play

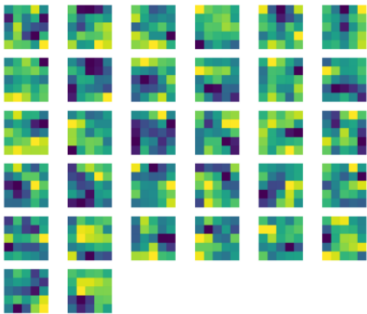
- <https://distill.pub/2018/building-blocks/>
  - Don't worry about reading the whole article


## Discuss

1. Identify what is strange, funny, or erroneous in the deep learning model in Building-Blocks
2. Take a screenshot of the issue and share with your tutorial mates
3. Try to explain why the model was behaving as identified
3. Post a 2–3 sentence description to the topic in your tutorial group: [#tg-xx](#)

# Understanding Filters

Hard to interpret kernels.  
They are just **matrices** used  
for **convolution**.





MIXED3B

MIXED4A

MIXED4B

OUTPUT CLASSES

Otterhound

Bucket

Soccer Ball

Beagle


Golden Retriever

Showing 3 of 480

Showing 3 of 508

Showing 3 of 512

Types of Dogs



MIXED3B

OUTPUT CLASSES

Bow Tie

Trench Coat

Windsor Tie


Suit

Military Uniform

Showing 3 of 480

Bowtie?

For instance, by combining feature visualization (what is a neuron looking for?) with attribution (how does it affect the output?), we can explore how the network decides between labels like **padlock** and **chain**.



CHANNELS THAT MOST SUPPORT ...

PADLOCK

CHAIN

feature visualization of channel

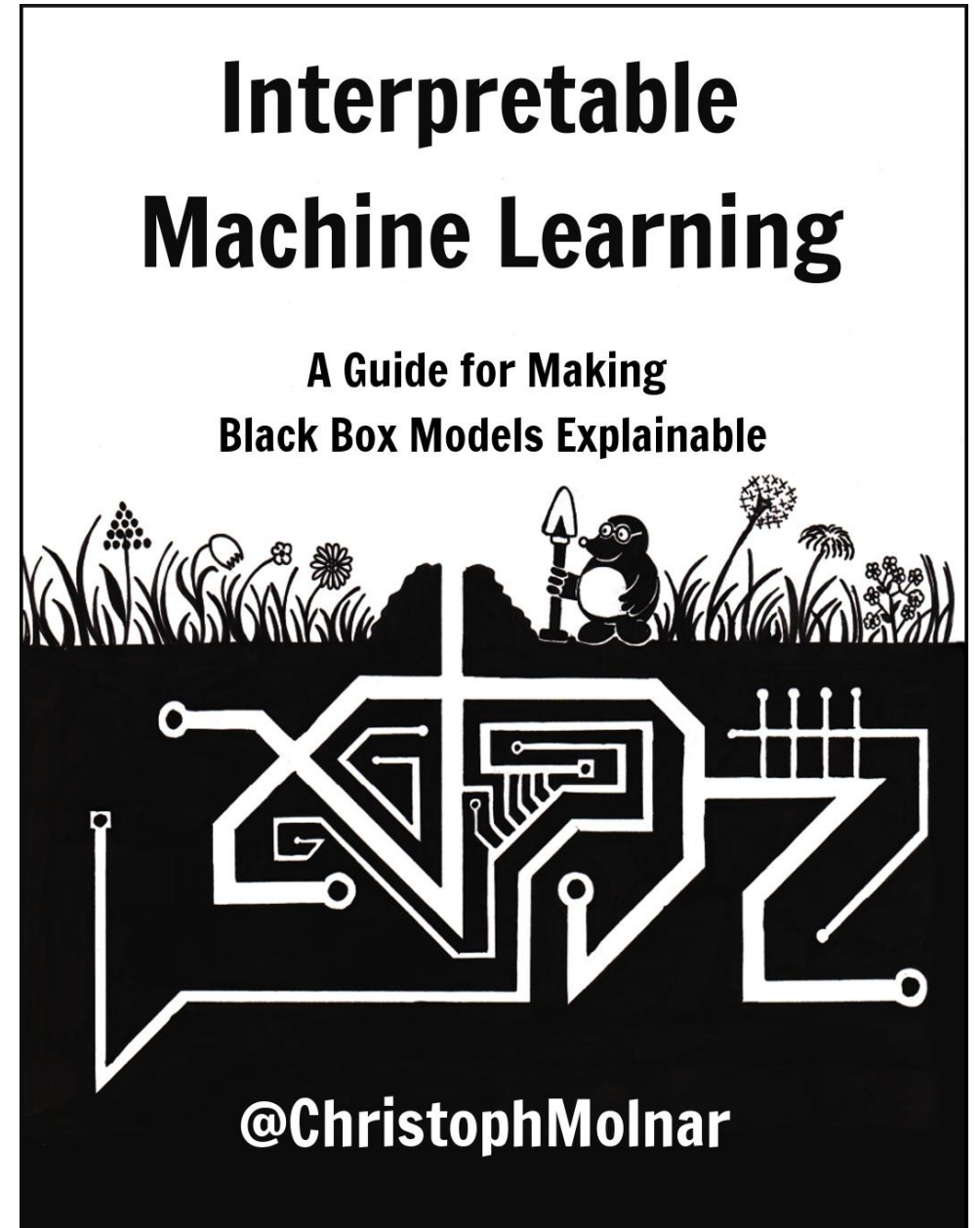
hover for attribution maps →

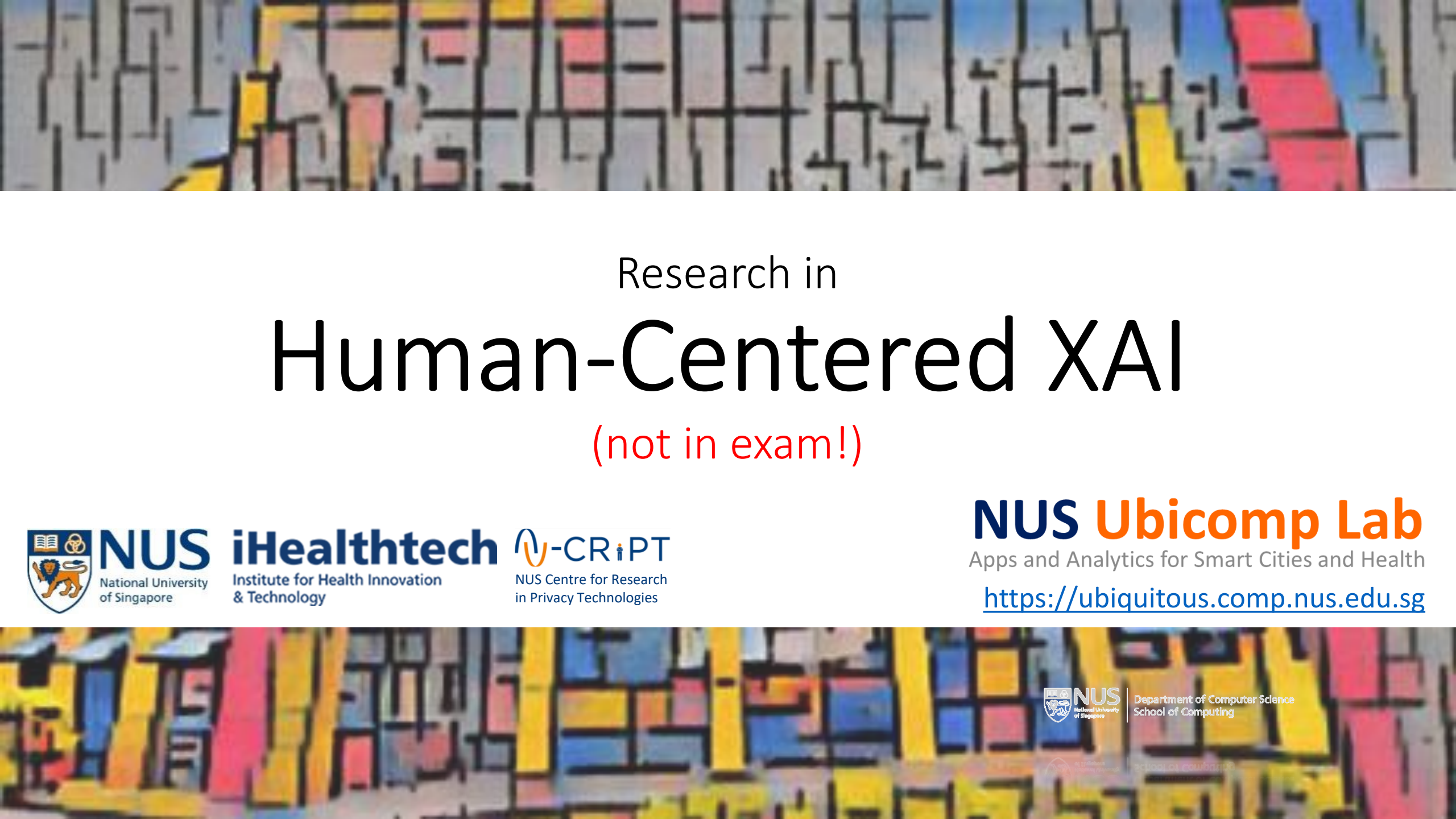
net evidence	0.75	0.58	0.52	0.64	0.76	0.81
for "padlock"	-0.96	0.47	0.14	-0.07	-0.17	-0.66
for "chain"	-1.71	-0.11	-0.38	0.57	0.59	0.15

Chain?

# Further Reading

- <https://christophm.github.io/interpretable-ml-book>





Research in

# Human-Centered XAI

(not in exam!)



**NUS**  
National University  
of Singapore

**iHealthtech**  
Institute for Health Innovation  
& Technology

**U-CRiPT**  
NUS Centre for Research  
in Privacy Technologies

**NUS Ubicomp Lab**

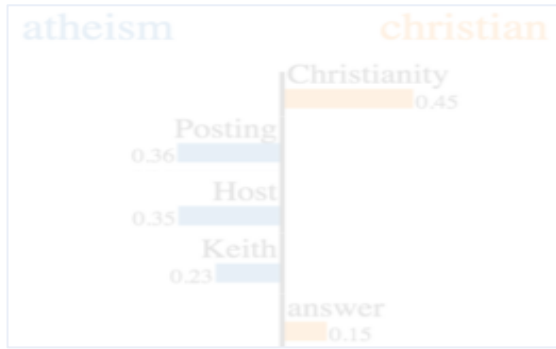
Apps and Analytics for Smart Cities and Health

<https://ubiquitous.comp.nus.edu.sg>

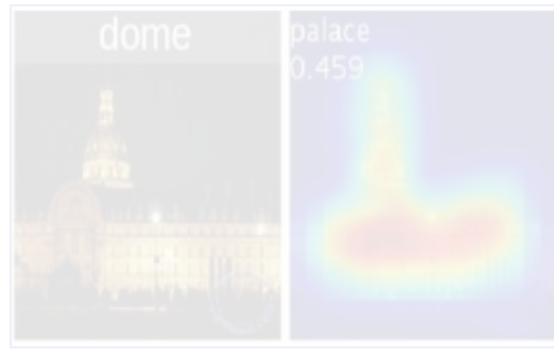


Department of Computer Science  
School of Computing





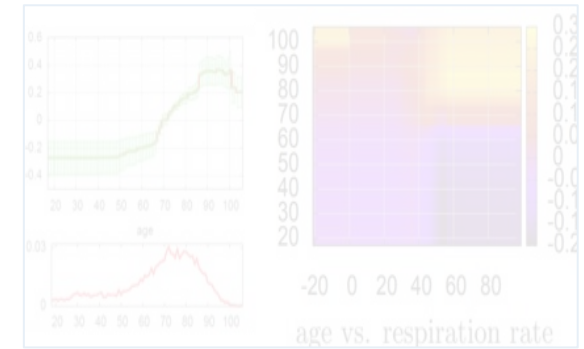
[Ribeiro 2016]



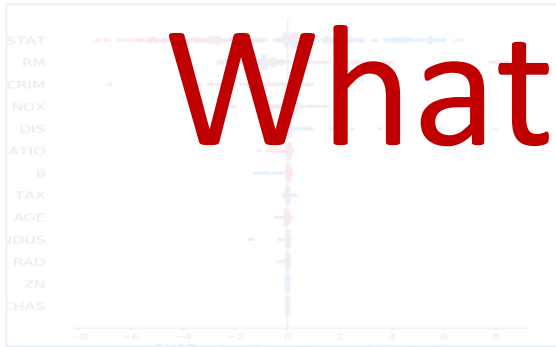
[Zhou 2016]

If Respiratory-Illness=Yes and Smoker=Yes and Age  $\geq$  50 then Lung Cancer  
 If Risk-LungCancer=Yes and Blood-Pressure  $\geq$  0.3 then Lung Cancer  
 If Risk-Depression=Yes and Past-Depression=Yes then Depression  
 If BMI  $\geq$  0.3 and Insurance=None and Blood-Pressure  $\geq$  0.2 then Depression  
 If Smoker=Yes and BMI  $\geq$  0.2 and Age  $\geq$  60 then Diabetes  
 If Risk-Diabetes=Yes and BMI  $\geq$  0.4 and Prob-Infections  $\geq$  0.2 then Diabetes  
 If Doctor-Visits  $\geq$  0.4 and Childhood-Obesity=Yes then Diabetes

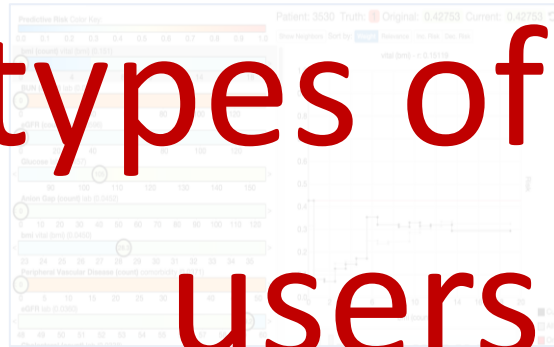
[Lakkaraju 2016]



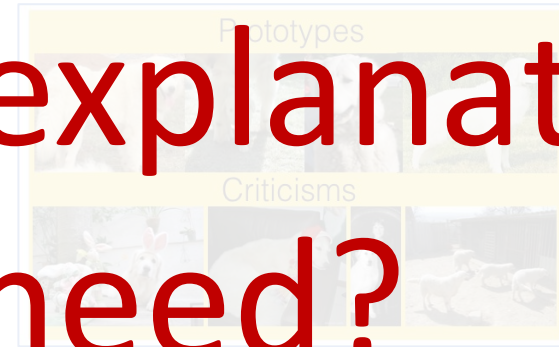
[Caruana 2016]



[Lundberg 2017]



[Krause 2016]



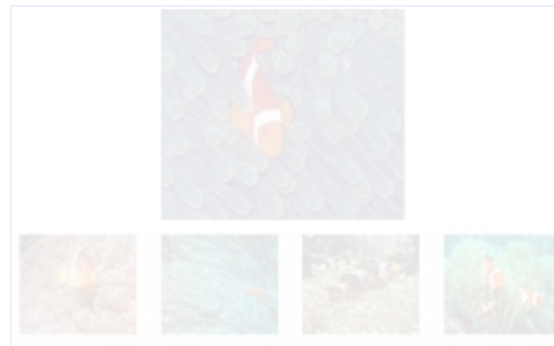
[Kim 2016]

**Saliency Map**  
 Long reaches for the propeller of a vintage  
 aircraft by leaning in a way  
 a brown a dog and a black dog in the edge of the ocean with a  
 wave under them boats are on the water in the background.  
 the pets are sleeping on the grass..  
 man in a blue shirt standing in front of a structure painted with  
 geometric designs.  
 a man is wearing a blue shirt.  
 a man is wearing a black shirt.

[Wallace 2018]



[Cai 2019]



[Koh 2017]



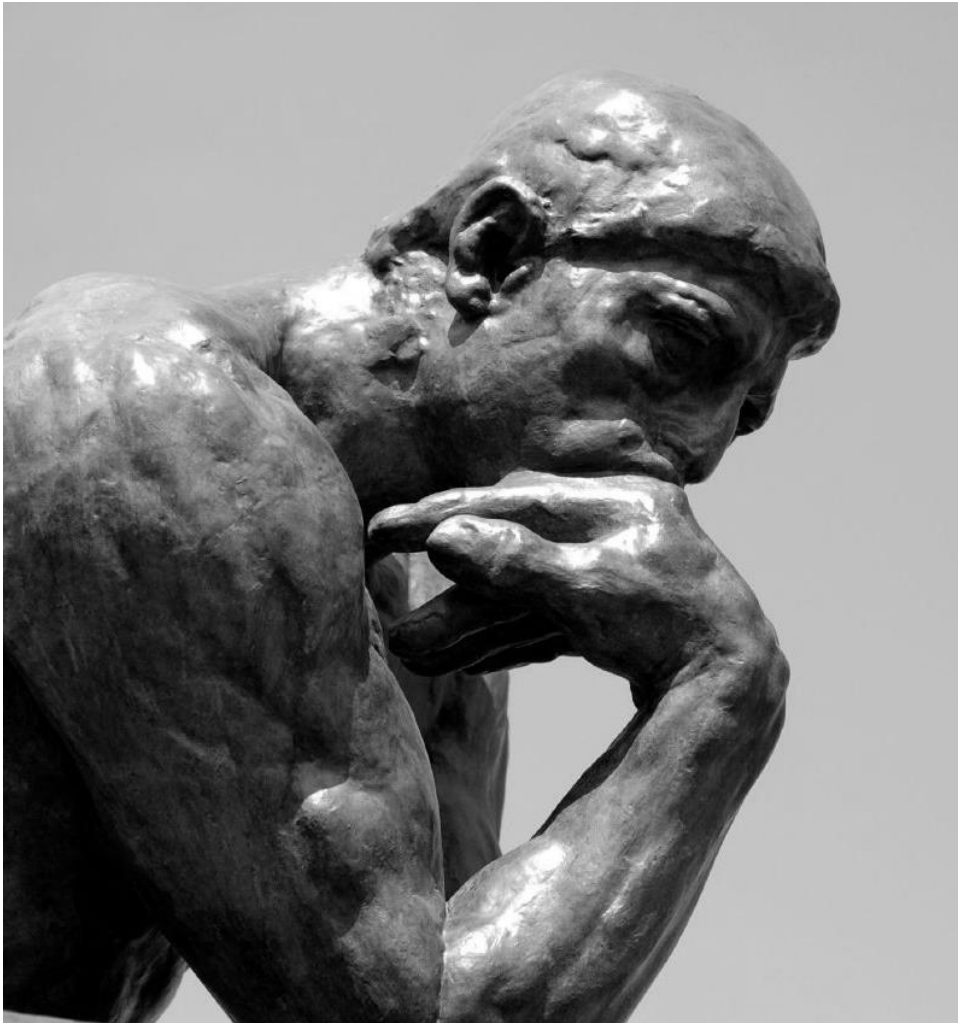
[Kim 2019]

...

What types of explanation do users need?



# How do people think and explain?

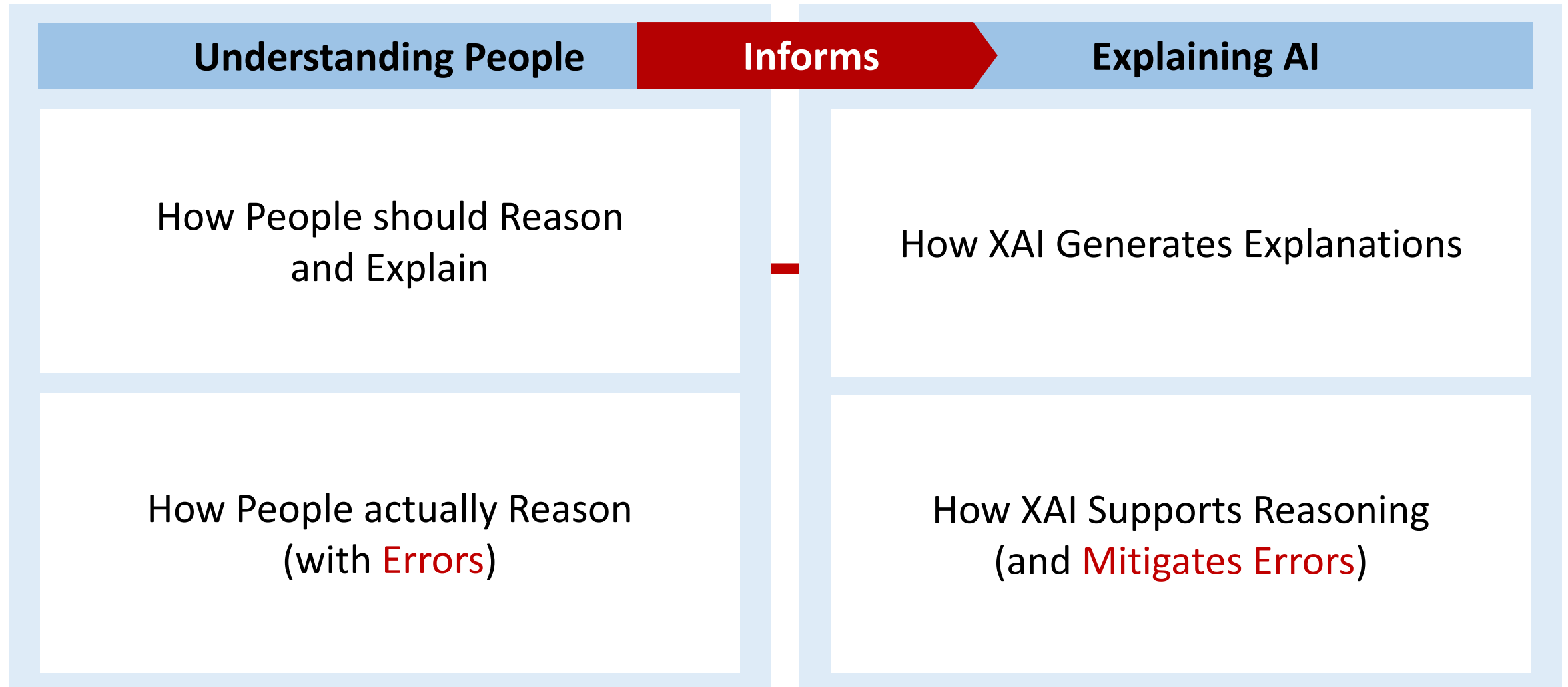


Philosophy

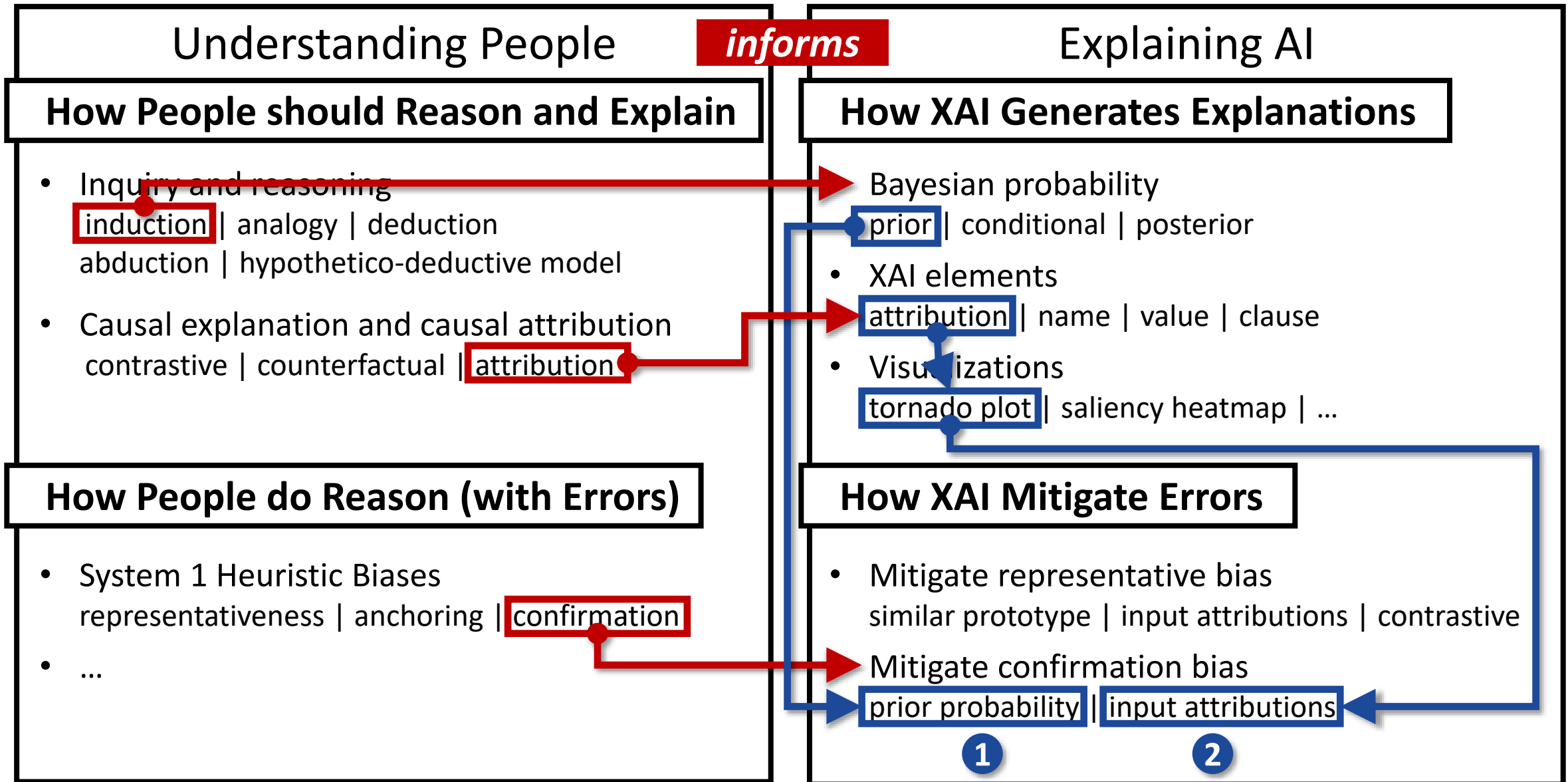


Psychology

# Human reasoning theories to inform XAI applications



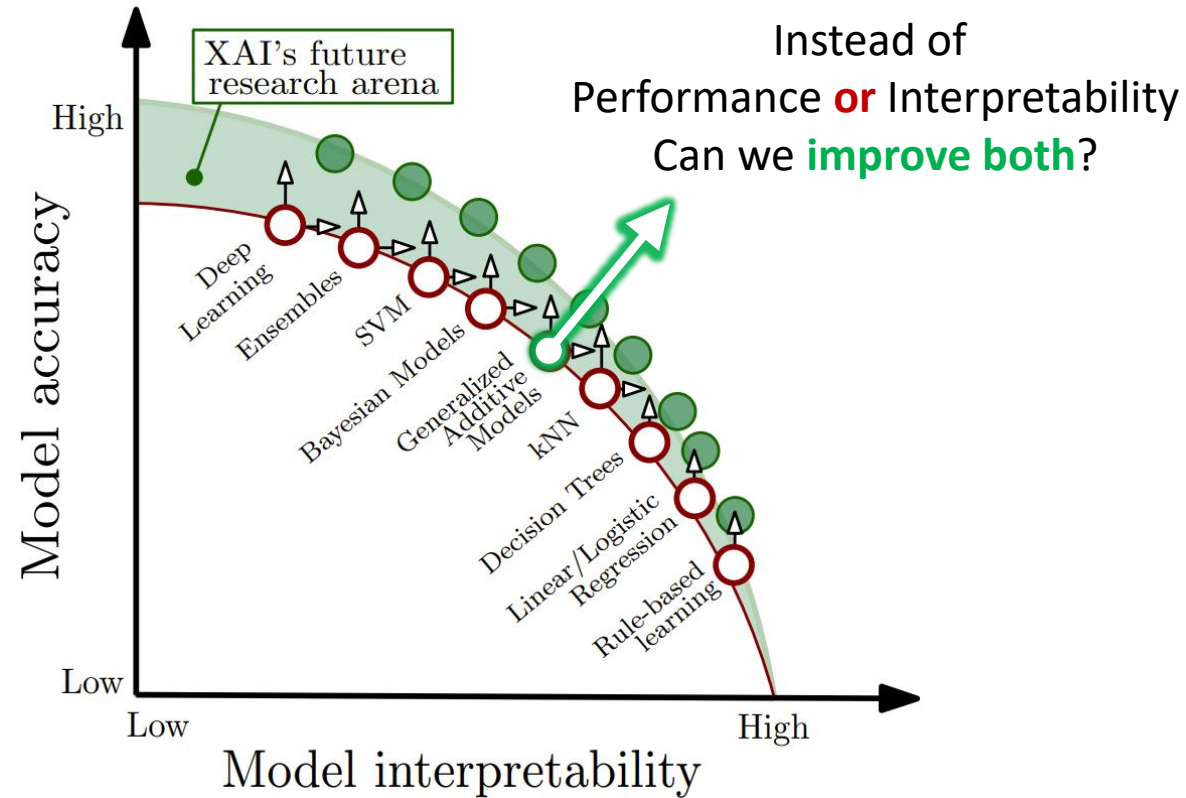
# Human reasoning theories to inform XAI applications



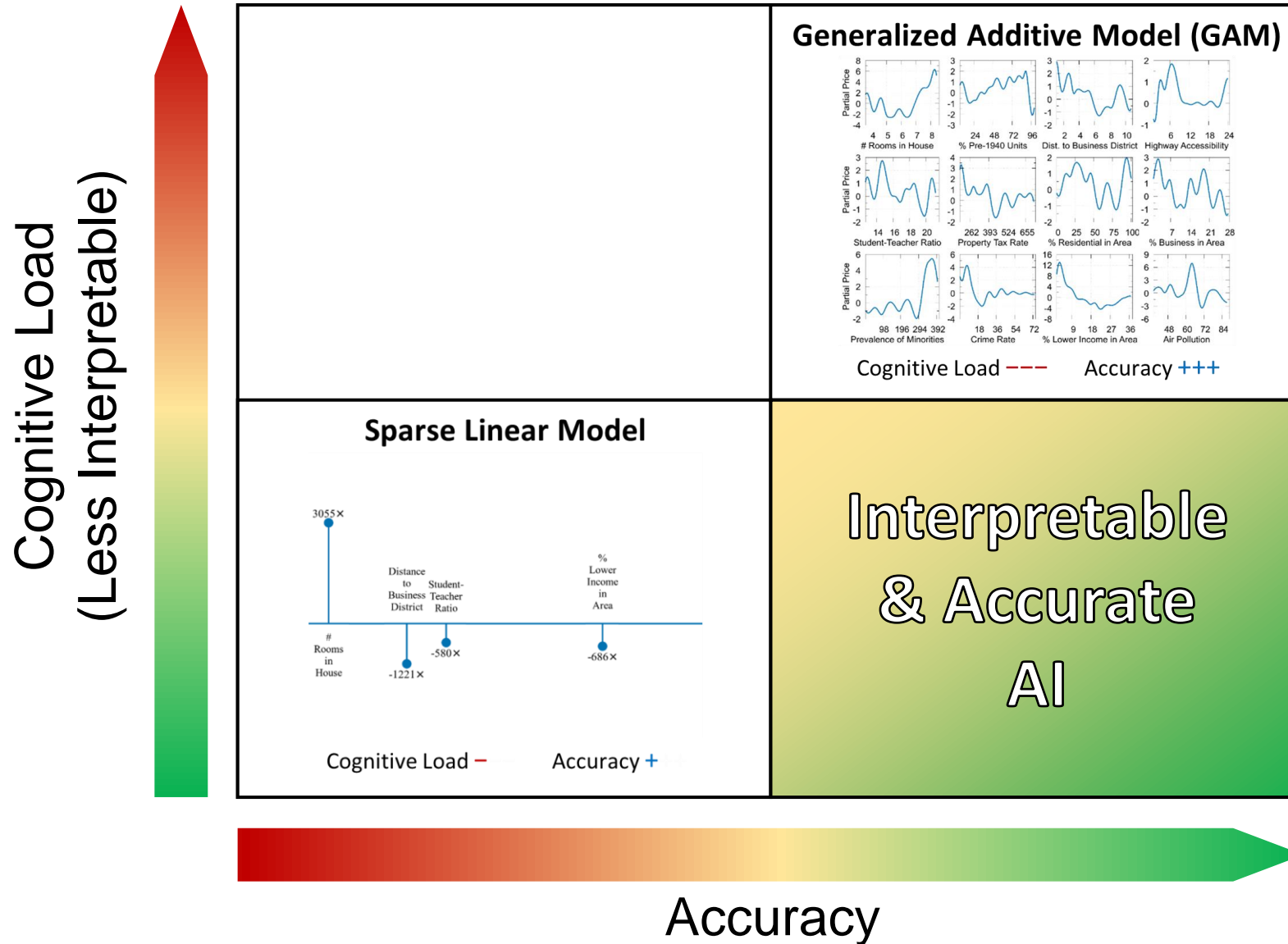
# User-Centered Explanations

- **Why?** Attribution Explanations
- **Why Not?** Contrastive Explanations
- **How To?** Counterfactual Explanations

# Performance-Interpretability Trade-off

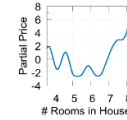


# Moderating **cognitive load** in Model Explanations

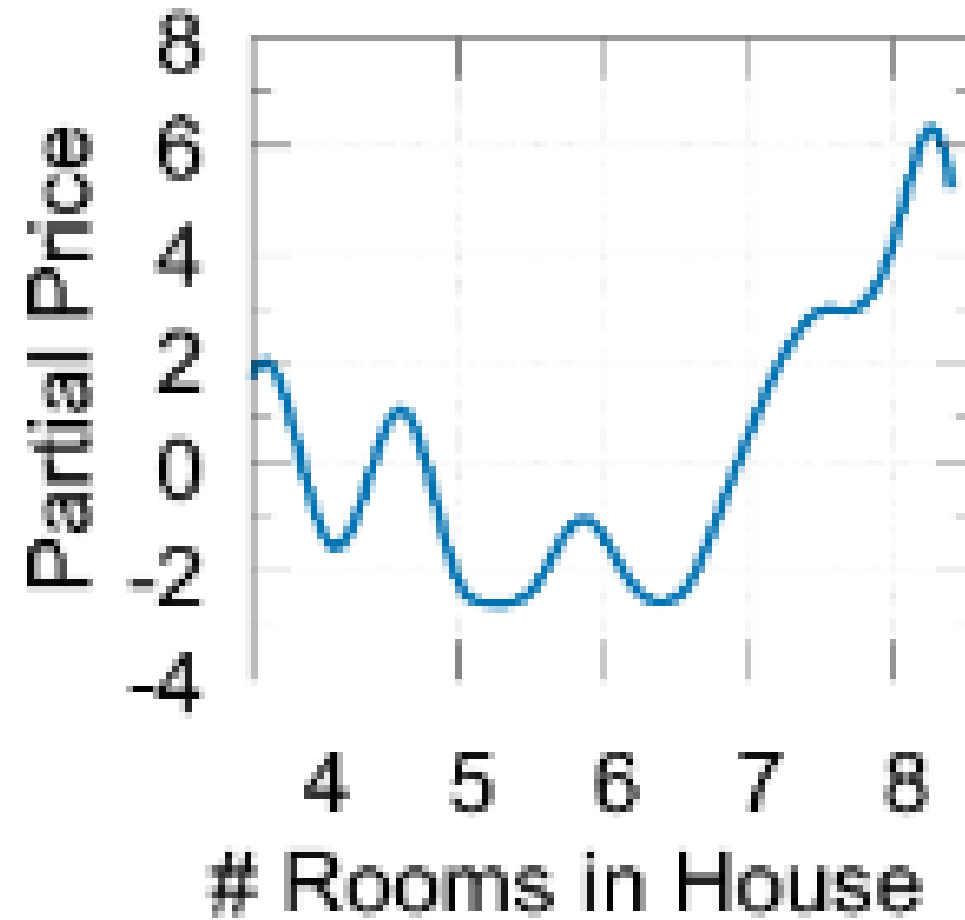




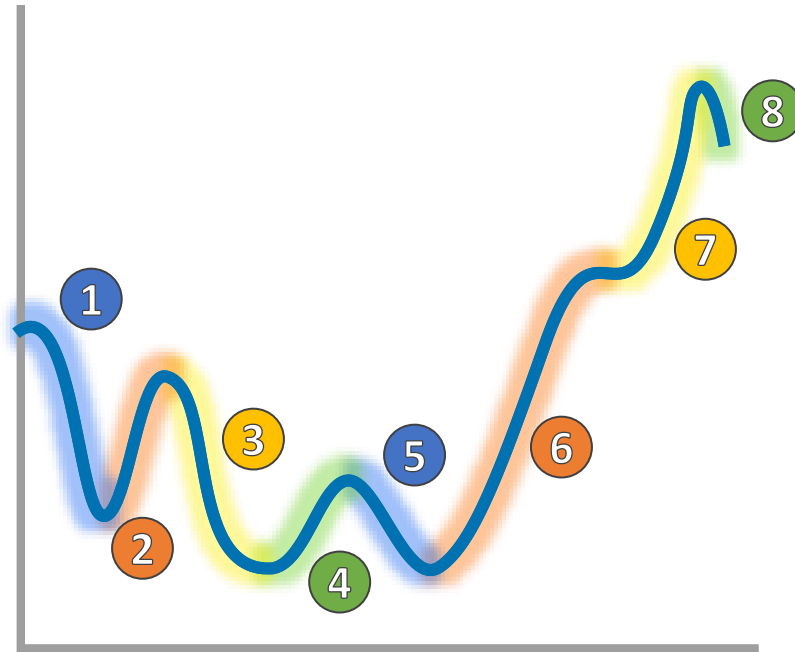
# Moderating **cognitive load** in Model Explanations



# Moderating **cognitive load** in Model Explanations

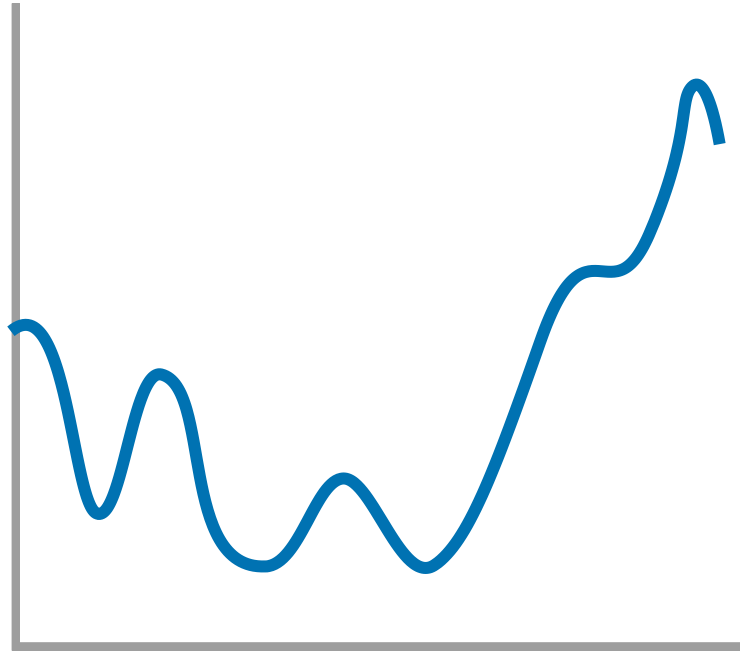


# Moderating **cognitive load** in Model Explanations



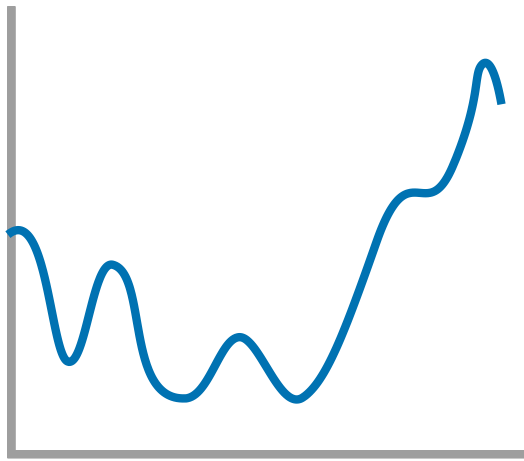
Cognitive Load = Number of **Visual “Chunks”**

# Moderating **cognitive load** in Model Explanations

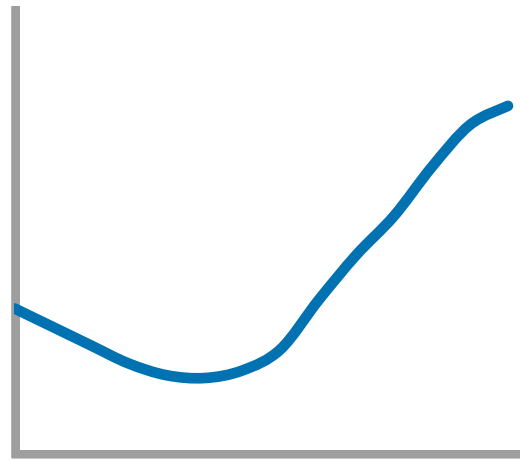


Cognitive Load = Number of **Visual “Chunks”**

# Moderating **cognitive load** in Model Explanations



8 Chunks



2 Chunks



1 Chunk

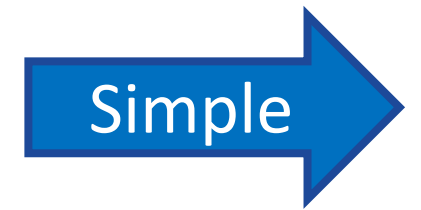
0 Chunks



Minimize

$$\sum_j^n \left( y_i - \sum_i^d f_i \left( x_j^{(i)} \right) \right)^2 + \lambda \sum_i^d \int f_i''(t)^2 dt$$

Less Curviness



Cognitive Load  
(Less Interpretable)



**COGAM: Measuring and Moderating Cognitive Load in Machine Learning Model Explanations.** *CHI '20.*





# Wrapping Up

# What did we learn?

## Feature Importance Explanations

### Feature Attribution



Weights  $w$  of  $f$

**Glassbox** model  $f$   
e.g. Linear Regression,  
Logistic Regression

Weights  $w$  of  $g$

**Model-agnostic**  
explainer model  $g$   
e.g. LIME

**Blackbox**  
nonlinear model  $f$

### Saliency Map



**Grad-CAM** explanation  $g$

**Blackbox**  
CNN model  $f$





# Next week: Unsupervised Learning

Image credit: <https://hip2save.com/2019/11/27/lego-classic-creative-fun-900-piece-set-only-20-at-walmart-regularly-40/>

# W12 Pre-Lecture Task (due before next Mon)

## Read

1. [Clustering With More Than Two Features? Try This To Explain Your Findings](#) by [Mauricio Letelier](#)

## Task

1. Describe other use cases where you need to **apply domain knowledge** with data-driven **unsupervised learning** to better understand your business or engineering problem

**Tip:** you can your own projects too; you don't have to be correct

2. Post a 1–2 sentence answer to the topic in your tutorial group: [#tg-xx](#)