

Introduction to Deep Learning

Maxwell Cai, Joris Mollinga
surf.nl

Objectives and Outline

Part I: Fundamentals, ML→DL, Common practices

- Theoretical introduction: 9:00 - 10:00
- Short break: 10:00 - 10:15
- Hands-on: 10:15 - 11:15

Part II: CNN

- Theoretical introduction: 11:15 - 11:45
- Hands-on: 11:45 - 12:30
- Lunch break: 12:30 - 13:30

Part III: RNN, sequential modelling

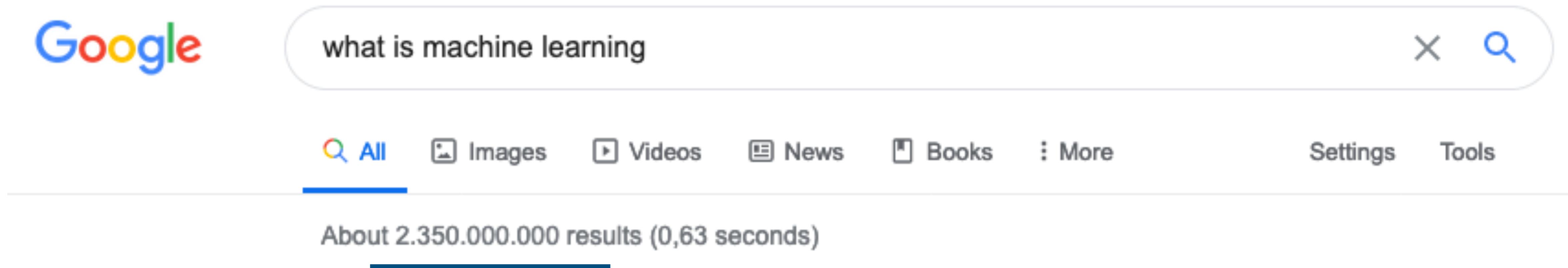
- Theoretical introduction: 13:30-14:15
- Hands-on: 14:15-15:15
- Short break: 15:15 - 15:30

Part IV: High Performance Machine Learning (optional if time permits)

- Theoretical introduction: 15:30 - 16:00
- Hands-on: 16:00 - 17:00

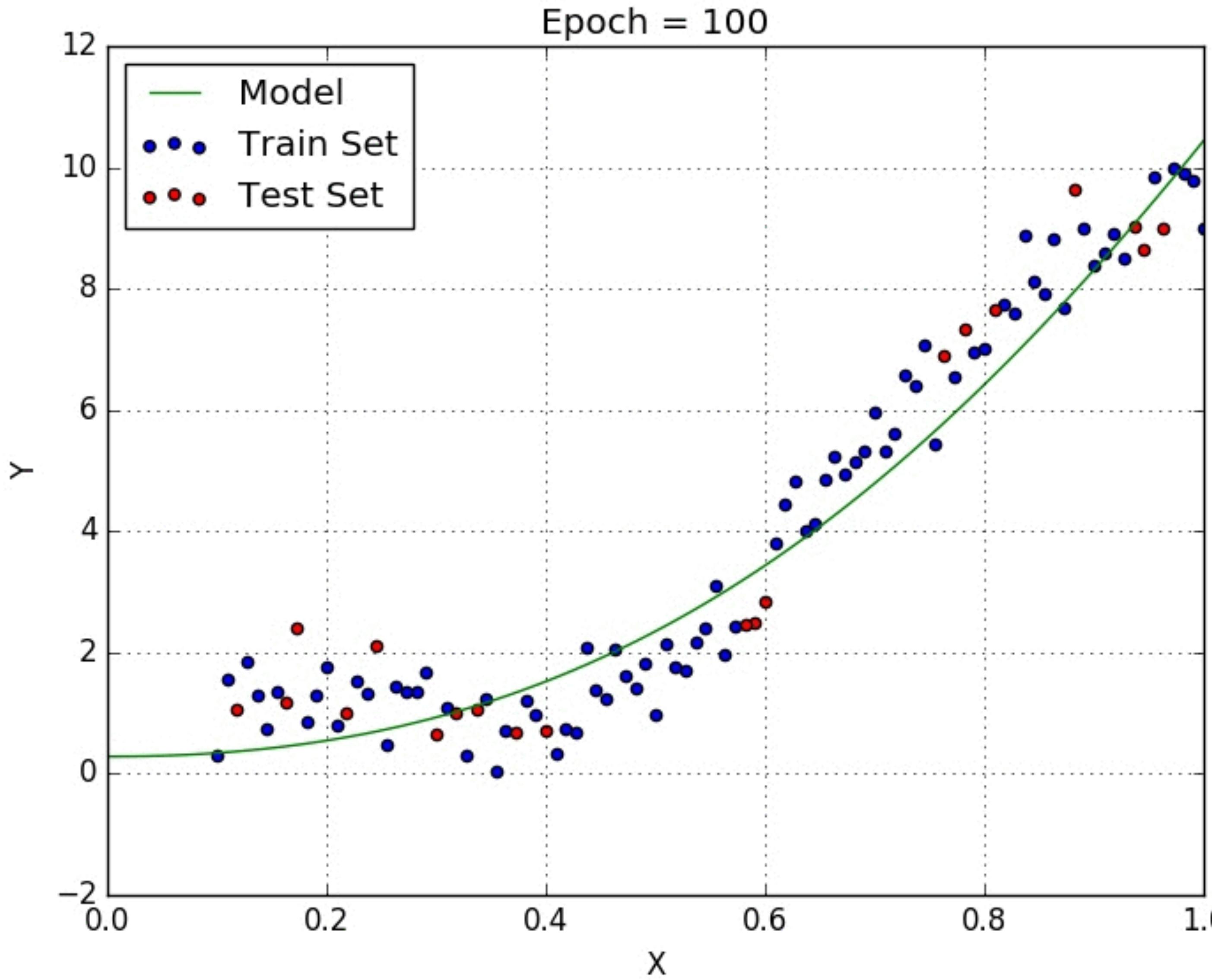
What *is* machine learning?

Search: “what is machine learning”



A screenshot of a Google search results page. The search bar at the top contains the query "what is machine learning". Below the search bar, the "All" tab is selected, followed by "Images", "Videos", "News", "Books", and "More". A horizontal line indicates the start of the search results. The text "About 2.350.000.000 results (0,63 seconds)" is displayed.

What *is* machine learning?



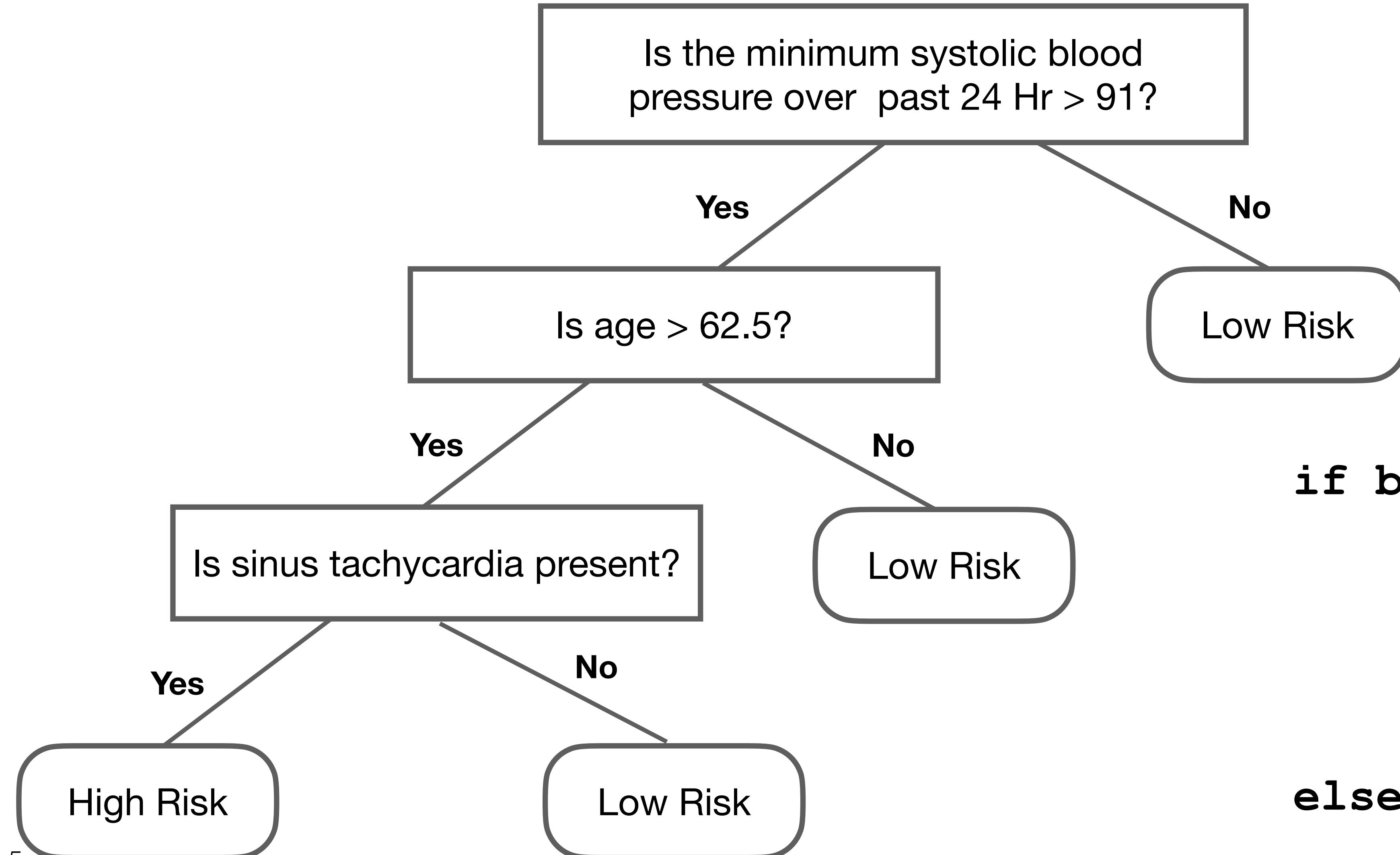
Could be as simple as curve fitting!

Typical applications in general:

- Classification
- Regression
- Dimensionality reduction
- Control

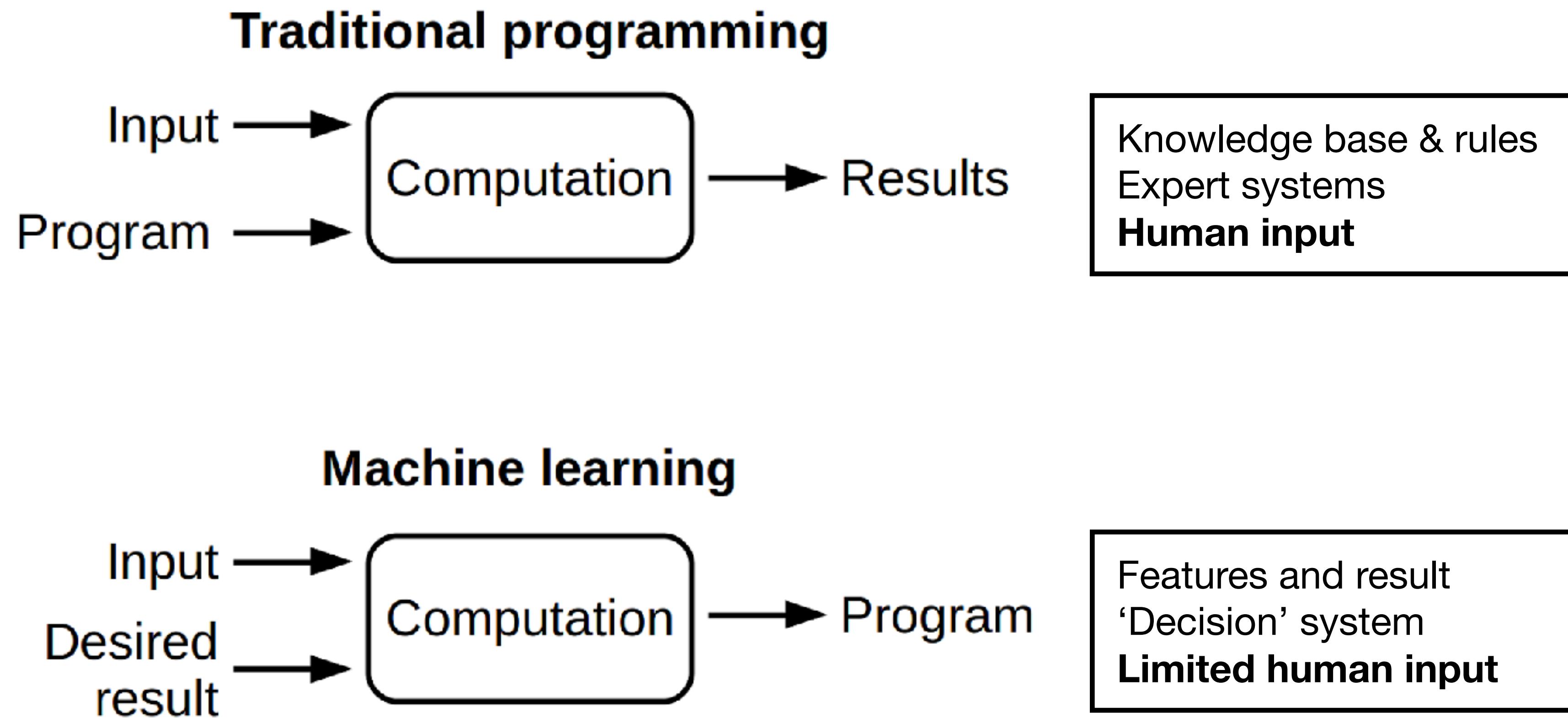
Why Machine Learning?

Think about a simple **decision tree**:



```
if blood_pressure > 91:  
    if age > 62.5:  
        if sinus_tach:  
            ...  
    else:  
        ...  
else:  
    ...
```

What *is* machine learning?



Fla.
70% Dem.

Pa.
89% Dem.

Ohio
54% Rep.

N.C.
66% Dem.

Va.
96% Dem.

Wis.
91% Dem.

Colo.
86% Dem.

Iowa
63% Rep.

Nev.
66% Dem.

N.H.
80% Dem.

Dem Rep

Clinton has **693** ways to win
68% of paths

16 ties
2% of paths

Trump has **315** ways to win
31% of paths

Florida

If Clinton wins Florida...

If Trump wins Florida...

Pennsylvania

Ohio

North Carolina

Virginia

Wisconsin

Colorado

Iowa

Nevada

New Hampshire

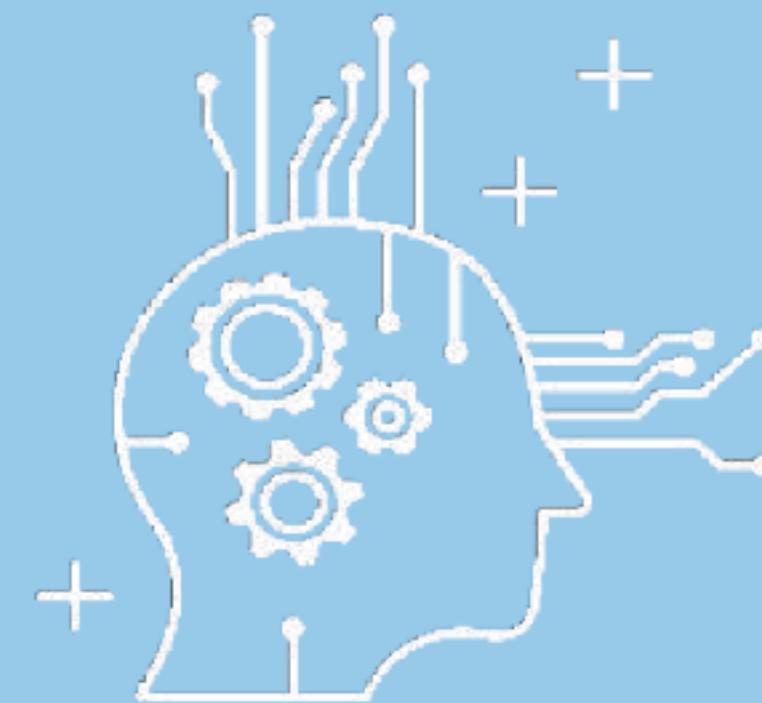
We probably don't want to do code a complex decision tree by hands...

Image source: Rahul/Medium

AI vs ML vs DL

ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



MACHINE LEARNING

Ability to learn without explicitly being programmed



DEEP LEARNING

Extract patterns from data using neural networks

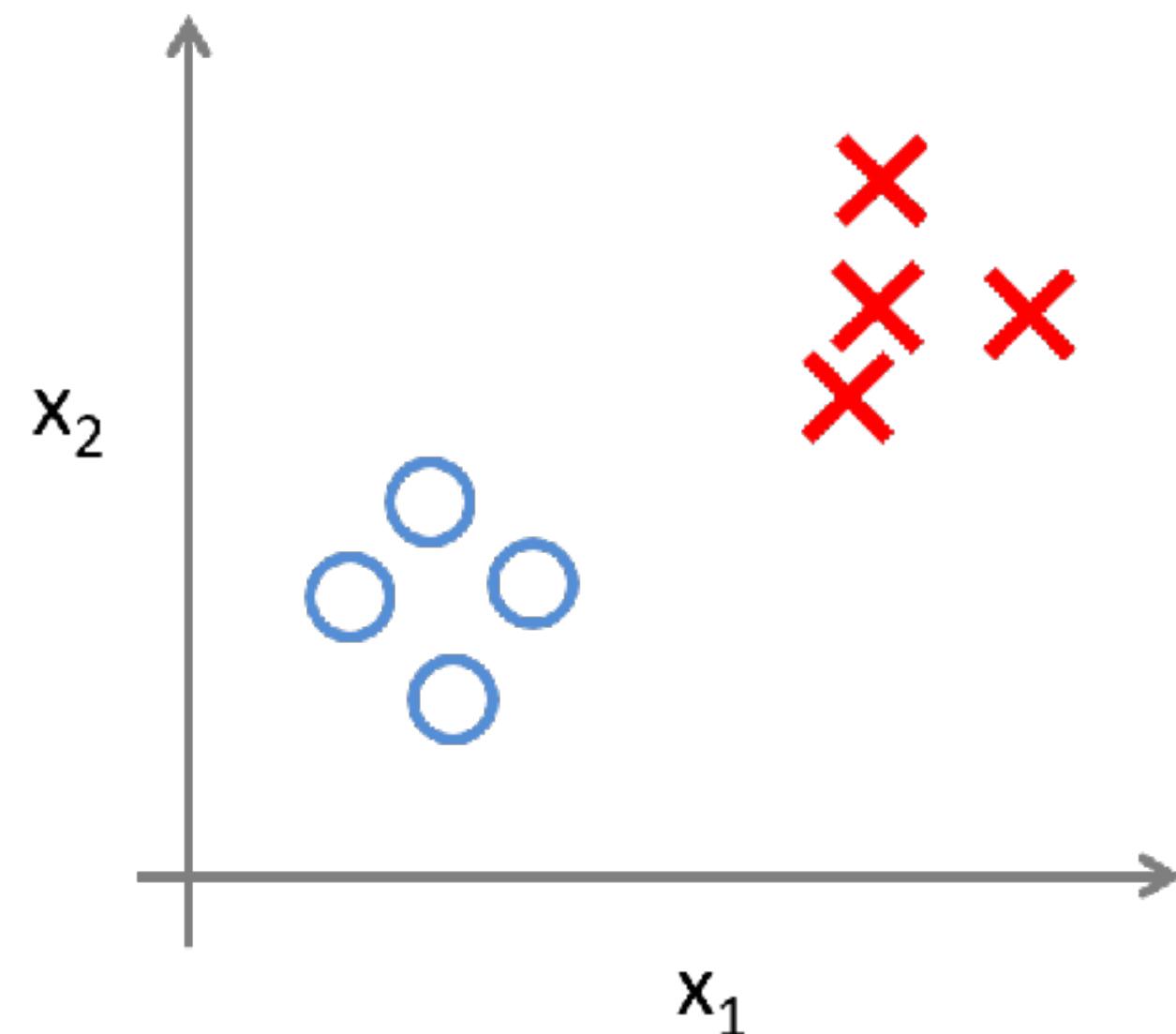


Credit: Alexander Amini/MIT

Categories of machine learning

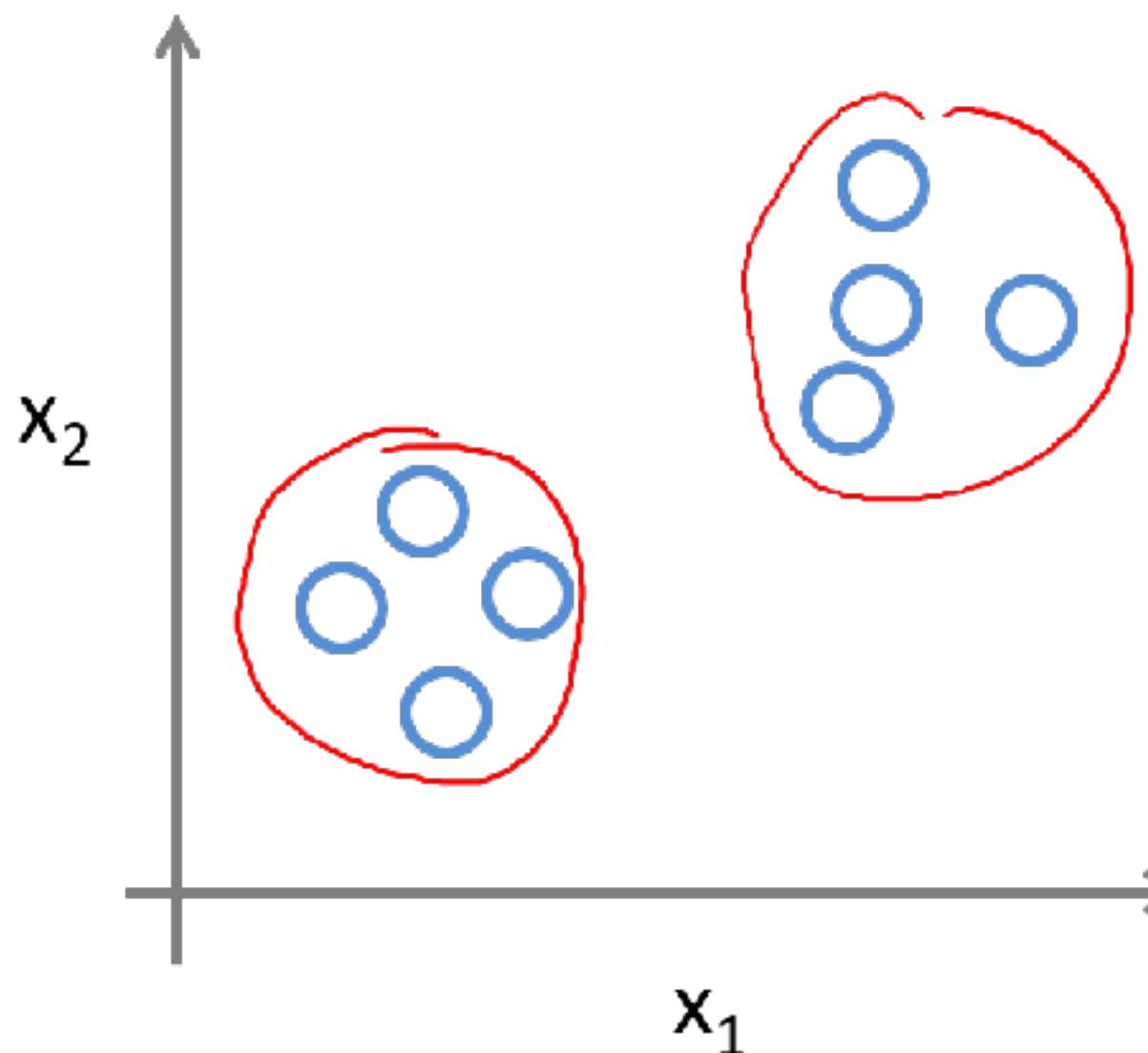
Supervised

Learn from the labels



Unsupervised

Detect patterns in the data



Reinforcement

Learn from mistakes



Categories of machine learning

Supervised

Classification

- Naive bayes
- Support vector machine
- Decision tree
- Random forest
- K-Nearest Neighbor
- Logistic regression

Regression

- Linear
- Generalized

Unsupervised

Clustering

- K-means
- K-medoids

Dimensionality reduction

- PCA
- SVD

Reinforcement

Discrete

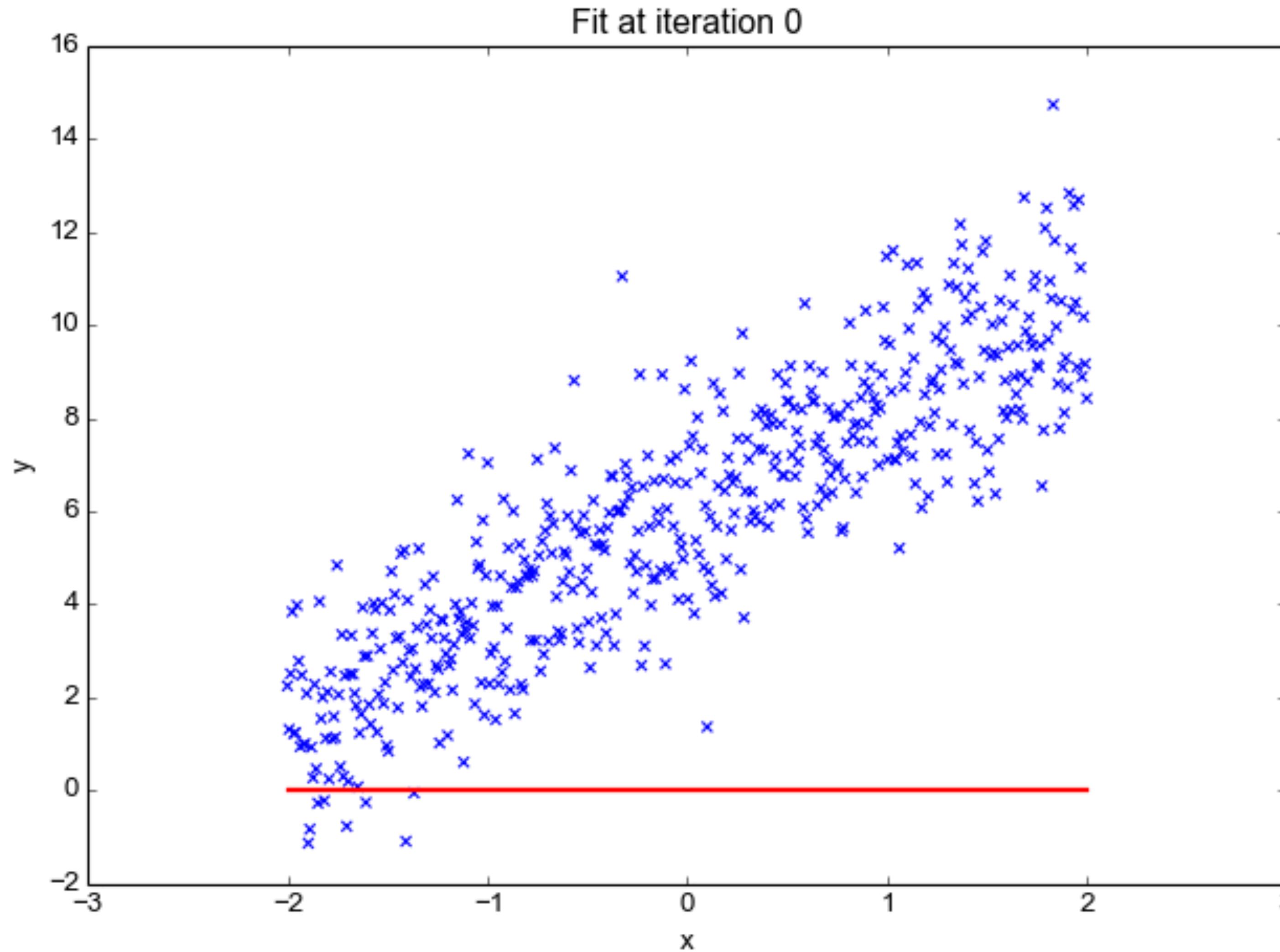
- Markov decision process
- Deep Q Network
- A2C
- A3C

Continuous

- DDPG
- NAF

... and many more... This list is incomplete!

Linear Regression



Iteratively Optimizing the parameters in a linear function

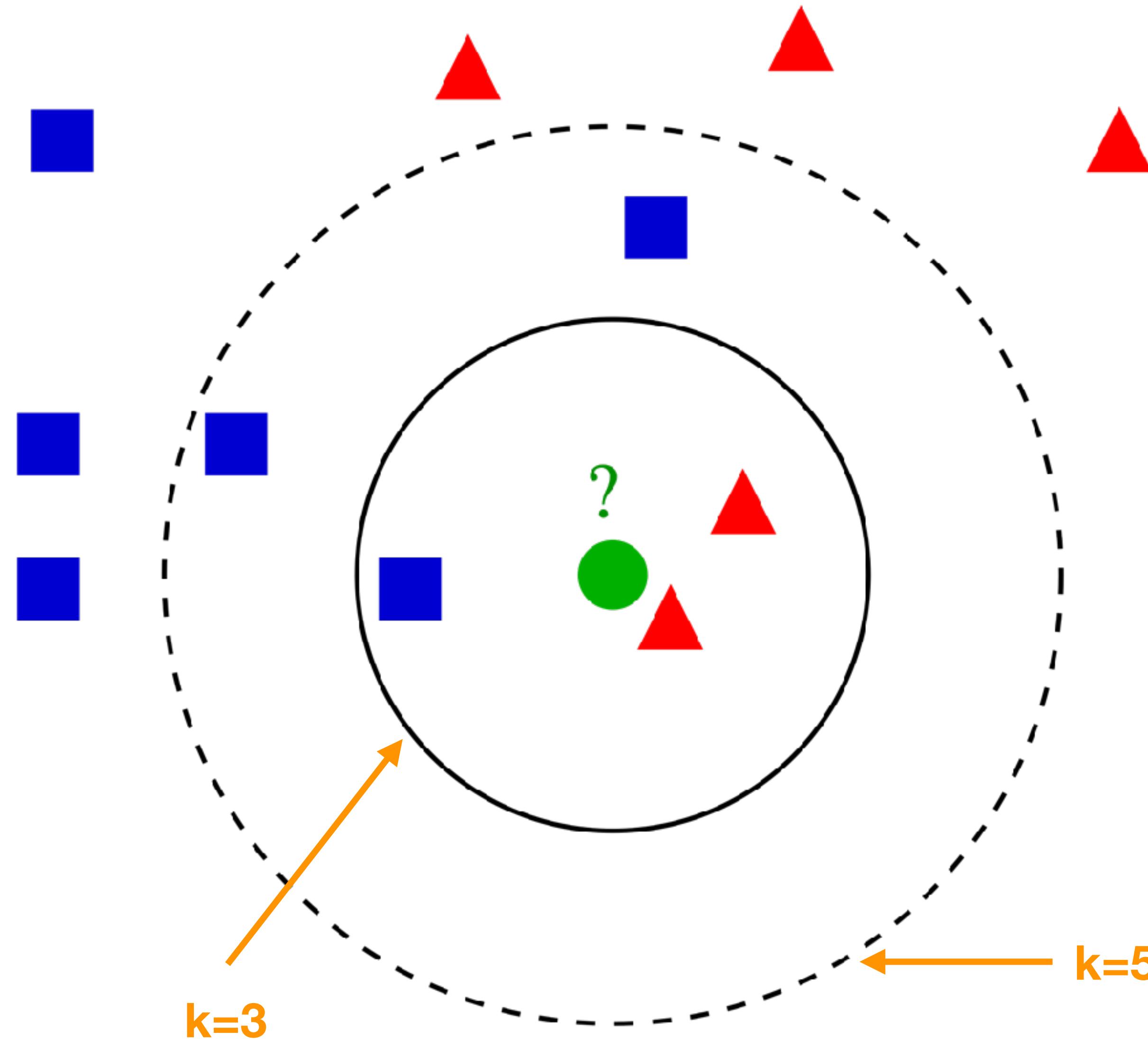
$$\hat{y} \equiv f(x) = wx + b$$

Such that

$$L(\hat{y}, y)$$

Reaches a minimum value.

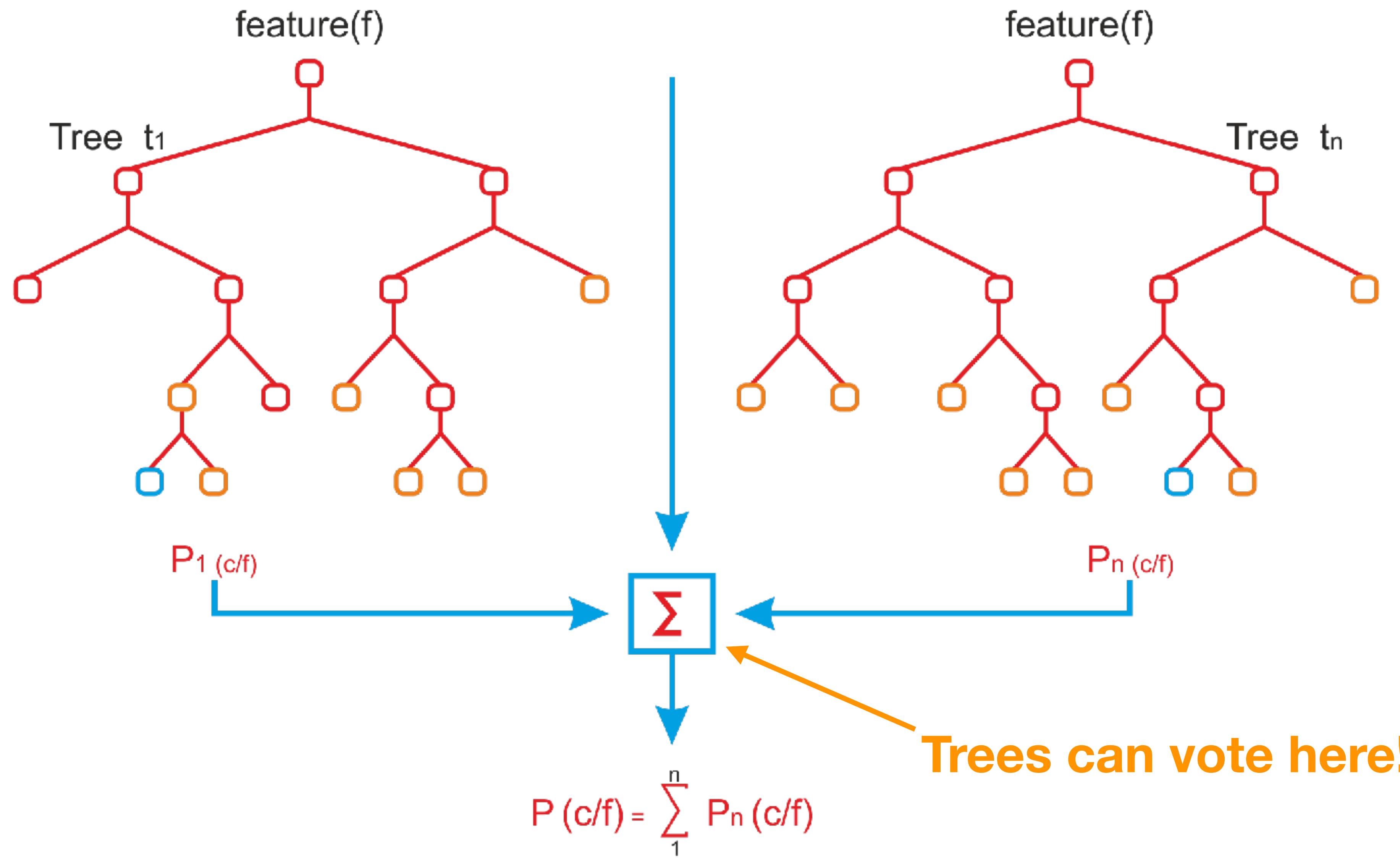
Classification: k -Nearest Neighbor



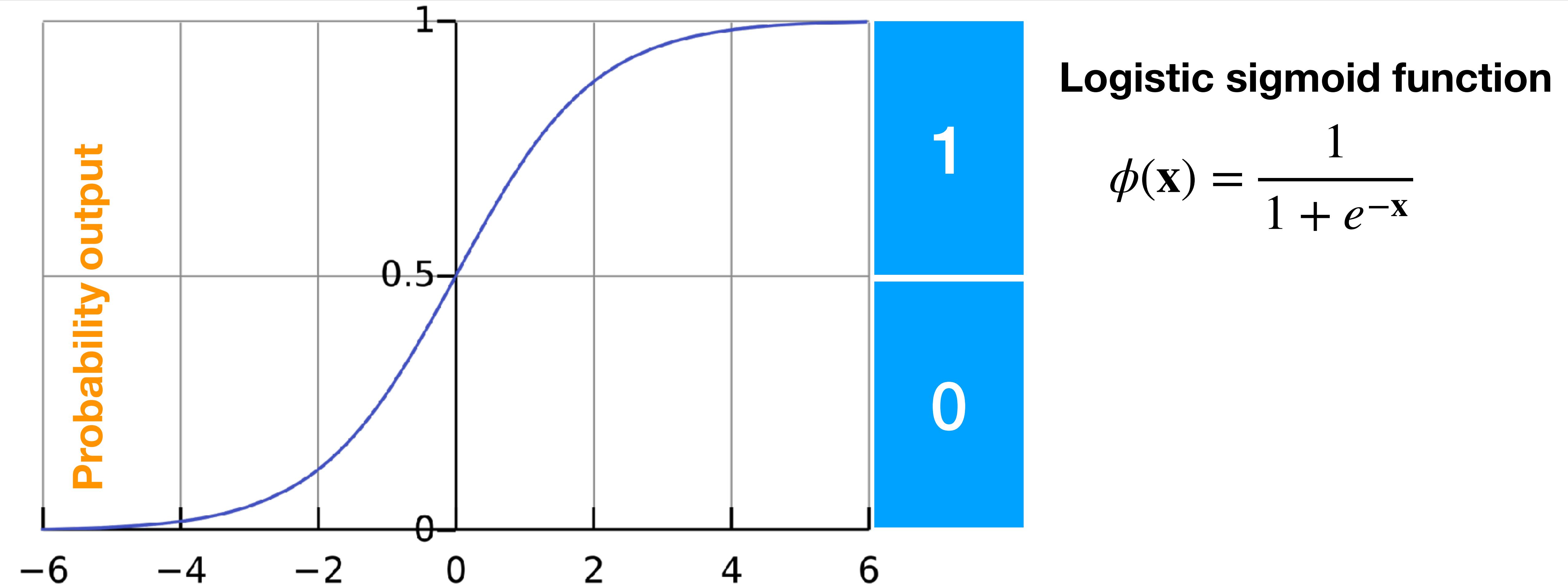
... is a **classification** algorithm

... based on a **voting** process

Classification: Decision Tree → Random Forest

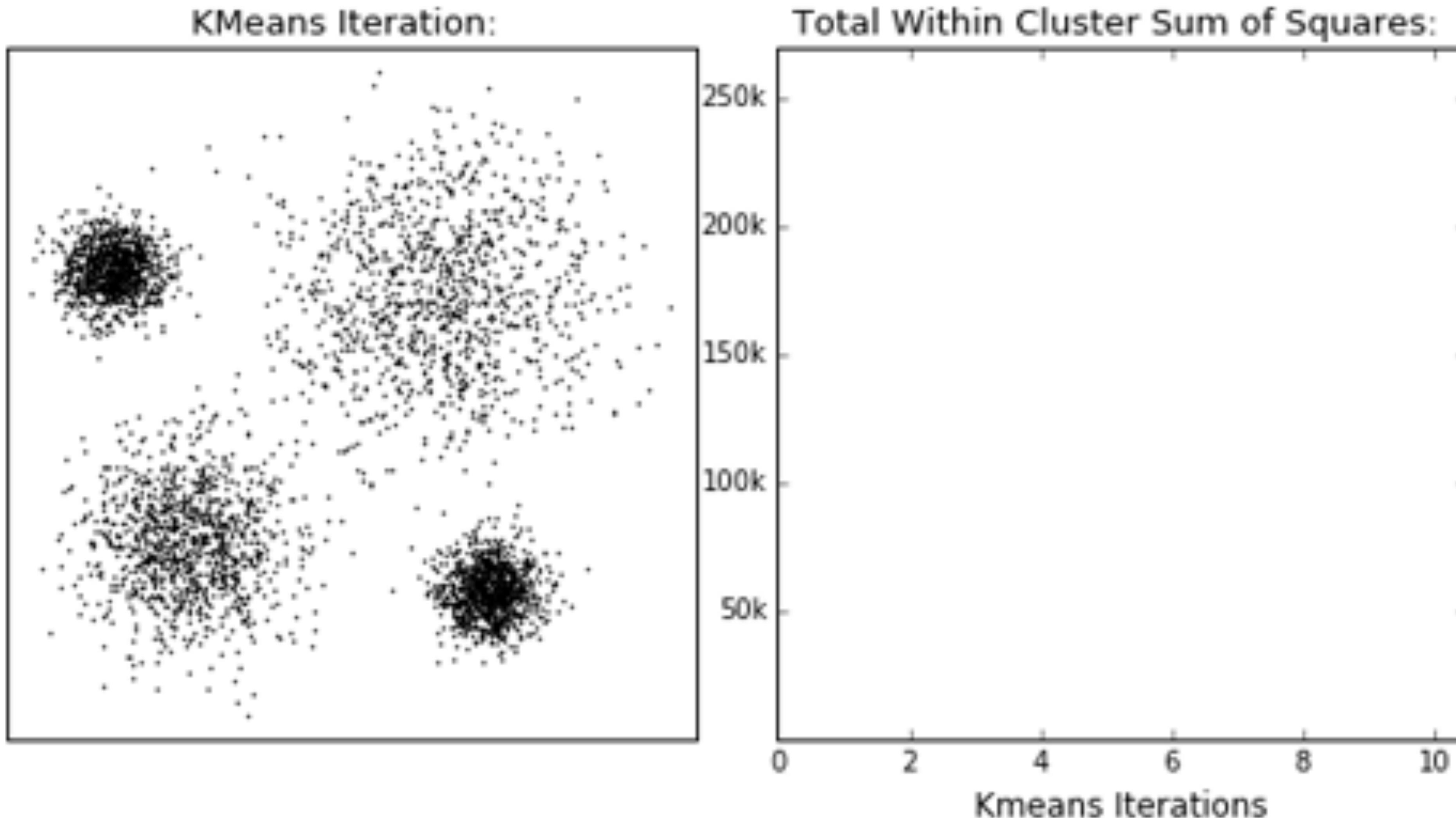


Classification: Logistic Regression

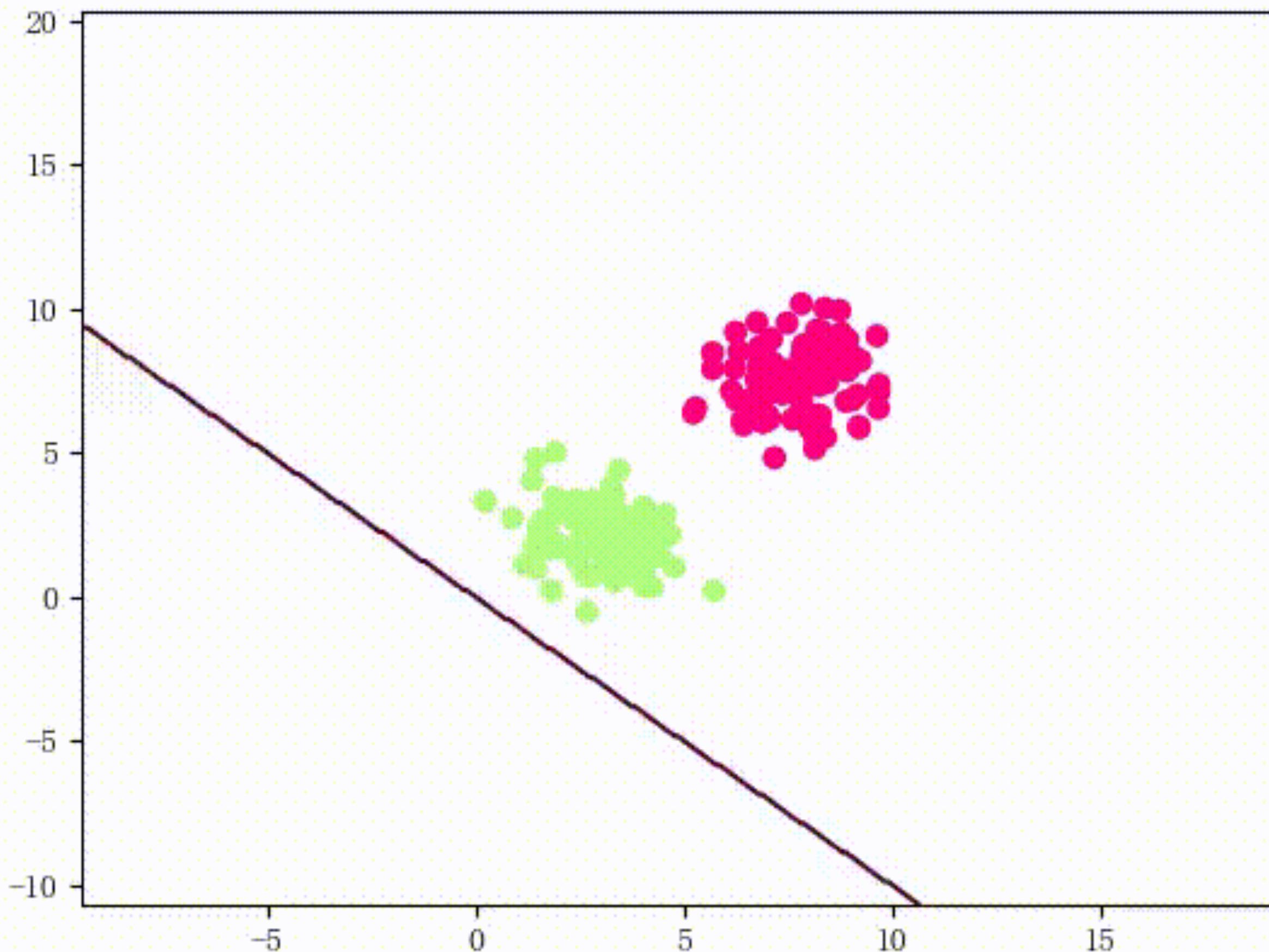


... is a binary **classification** algorithm, not a regression algorithm
... Gives the **probability** of each class

K-means clustering



Support Vector Machine



... is a **classification** algorithm

... does not provide probabilities, but only output class identities.

... aims to find an optimal way (perpendicular to the **support vector**) to separate different classes

What is the machine actually learning?

Machine learning is an optimisation process

Learning: Optimizing the loss function

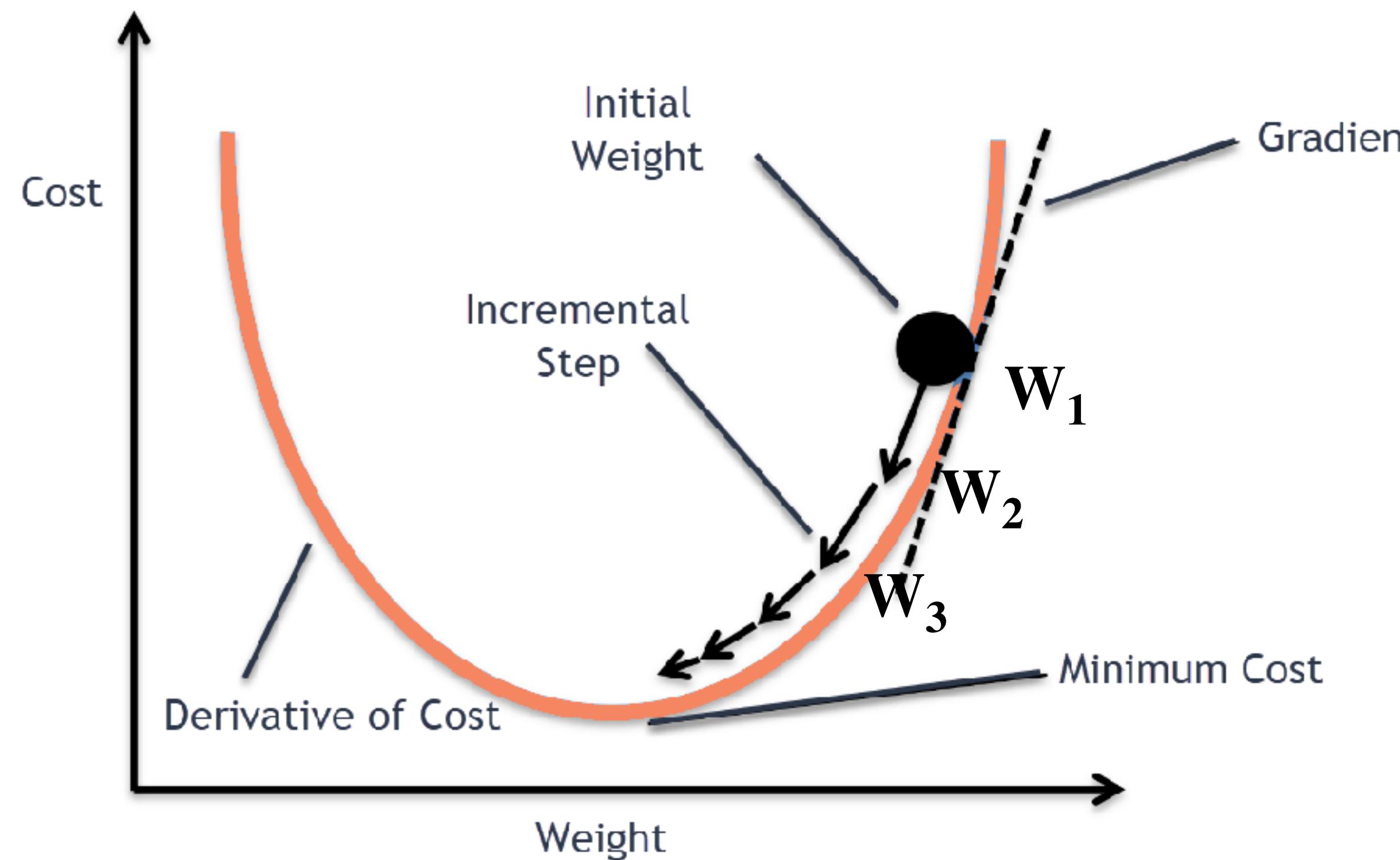
Objective: minimize the **difference** between predicted value the actual value

$$\mathcal{L}(\mathbf{W}, b) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

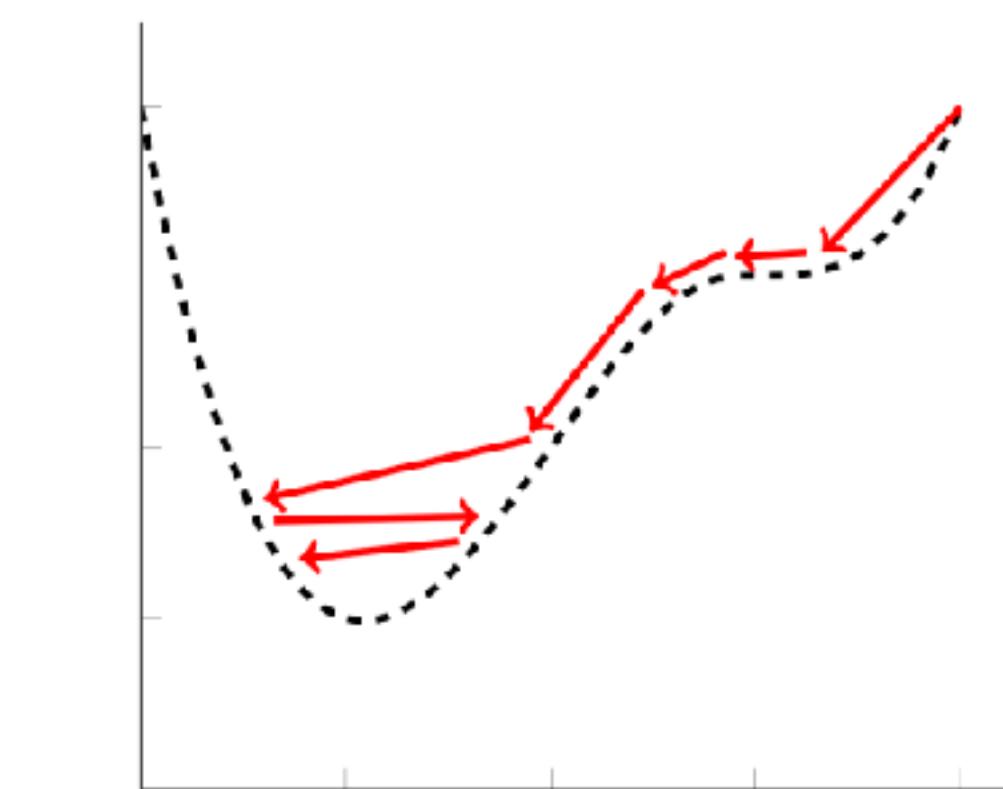
$\nabla \mathcal{L}(\mathbf{W}, b)$

$$\mathbf{W}_{j+1} = \mathbf{W}_j - \alpha \nabla \mathcal{L}(\mathbf{W}_j, b)$$

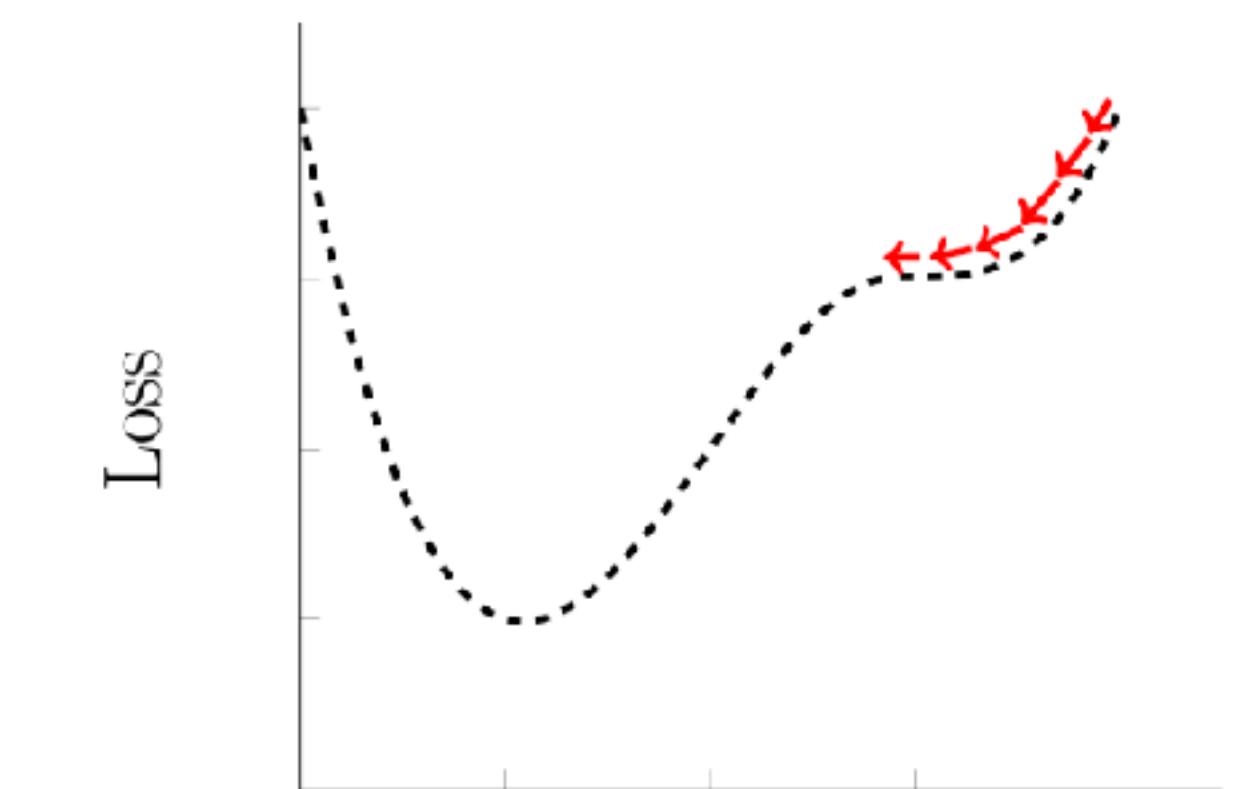
↑
Learning rate



High Learning Rate



Low Learning Rate



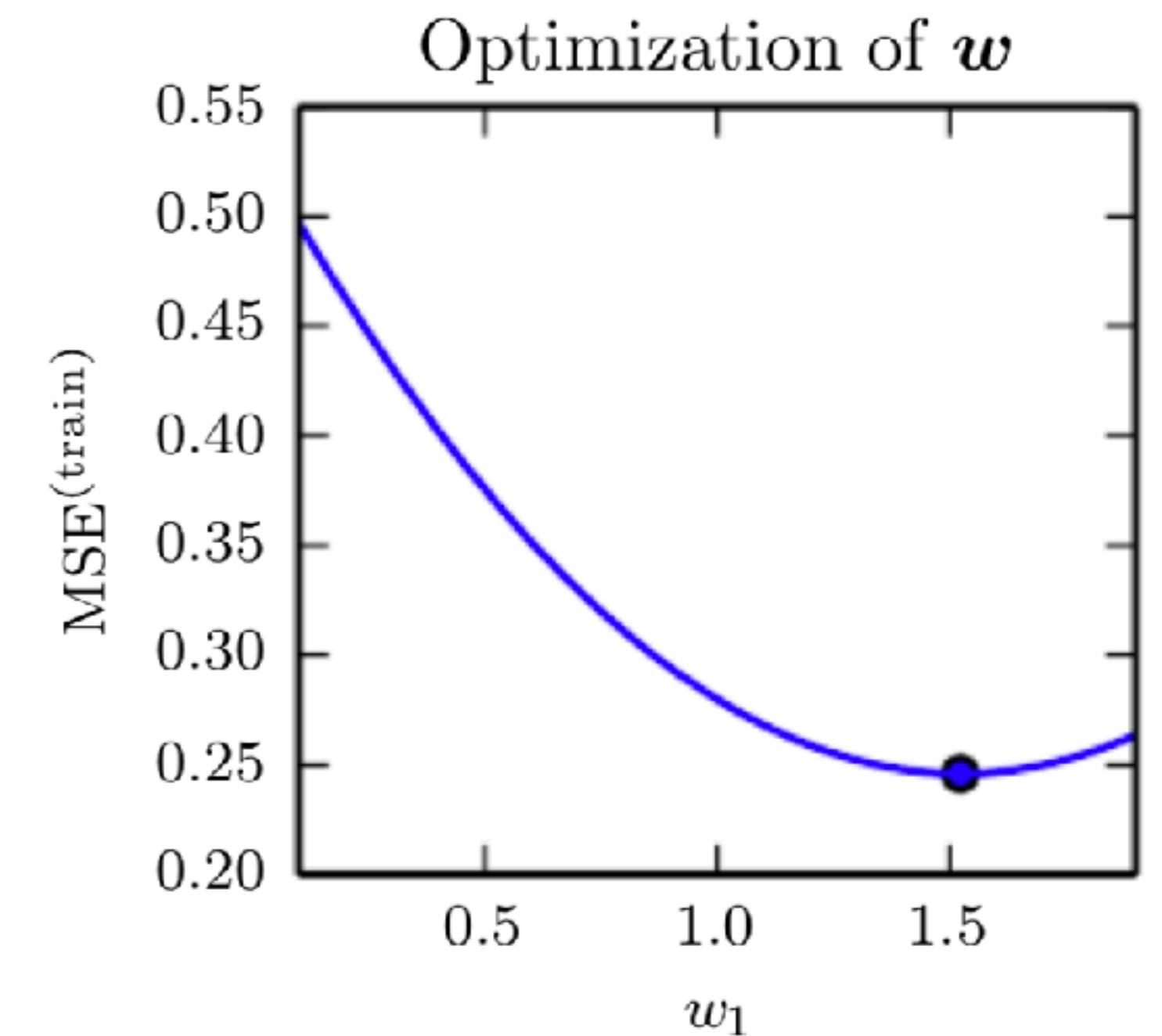
A few more terms...

$$\hat{y} = \mathbf{w}^T \mathbf{x} + b$$

Weights Bias
↓ ↓
Estimate Input features

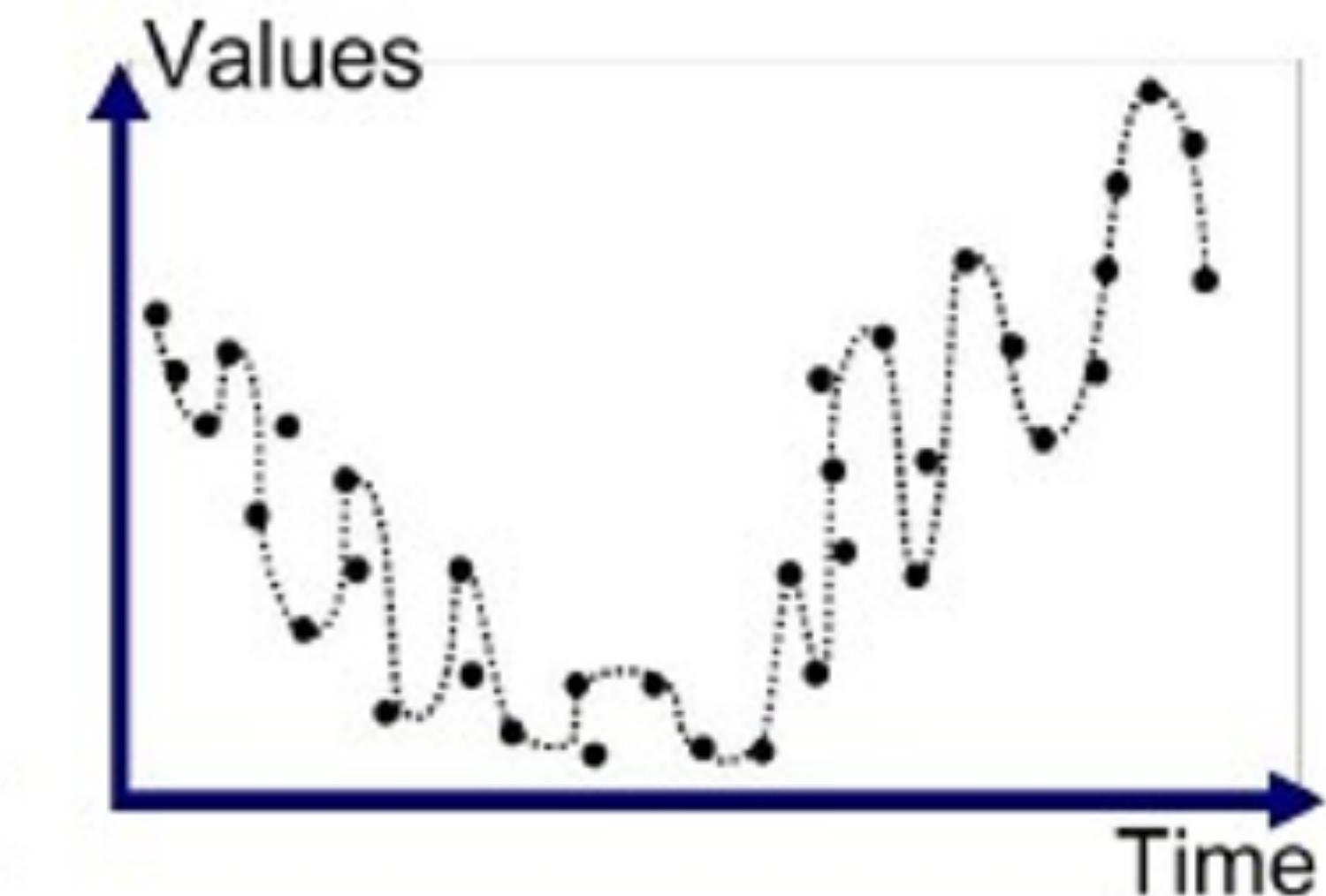
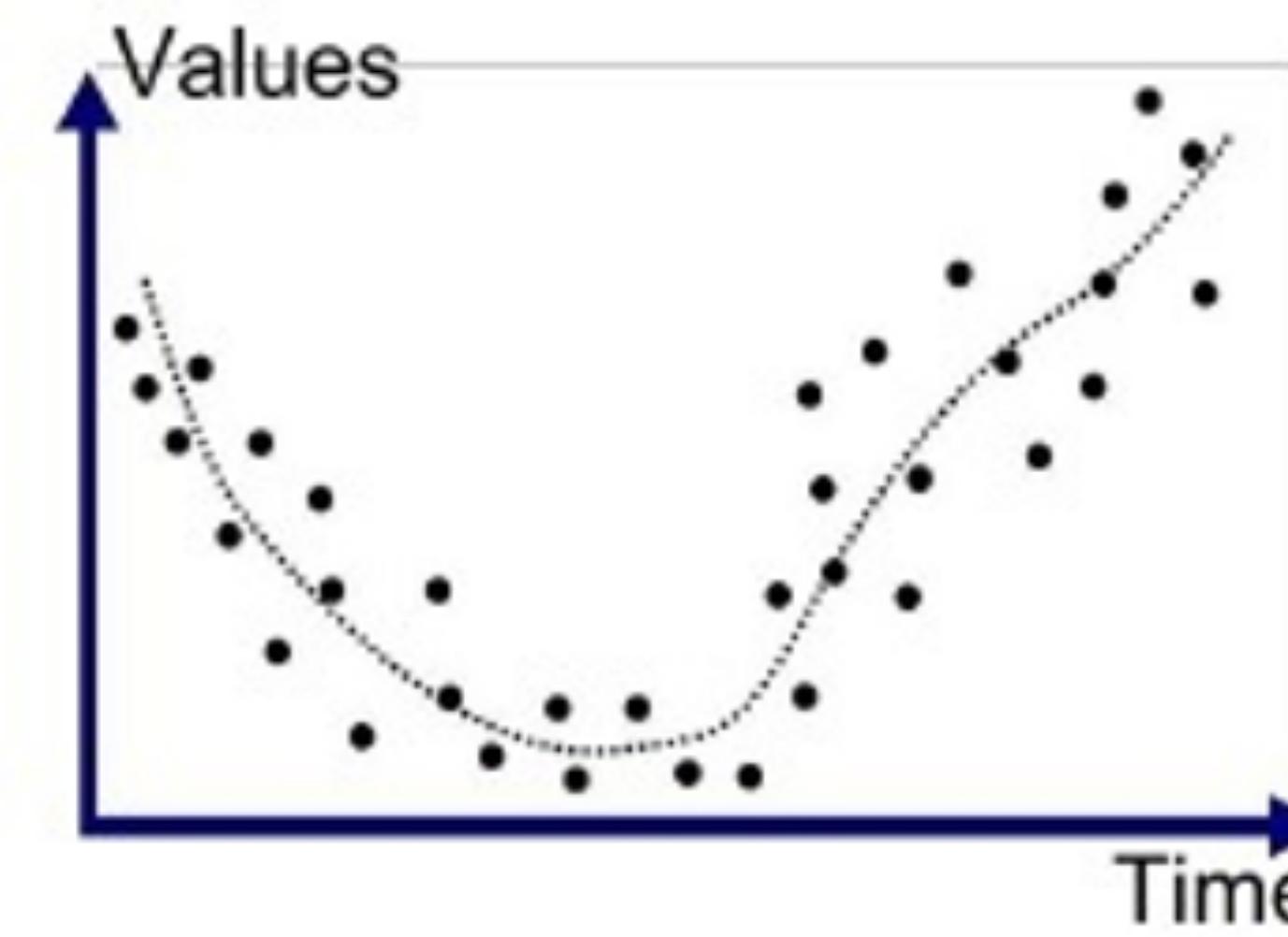
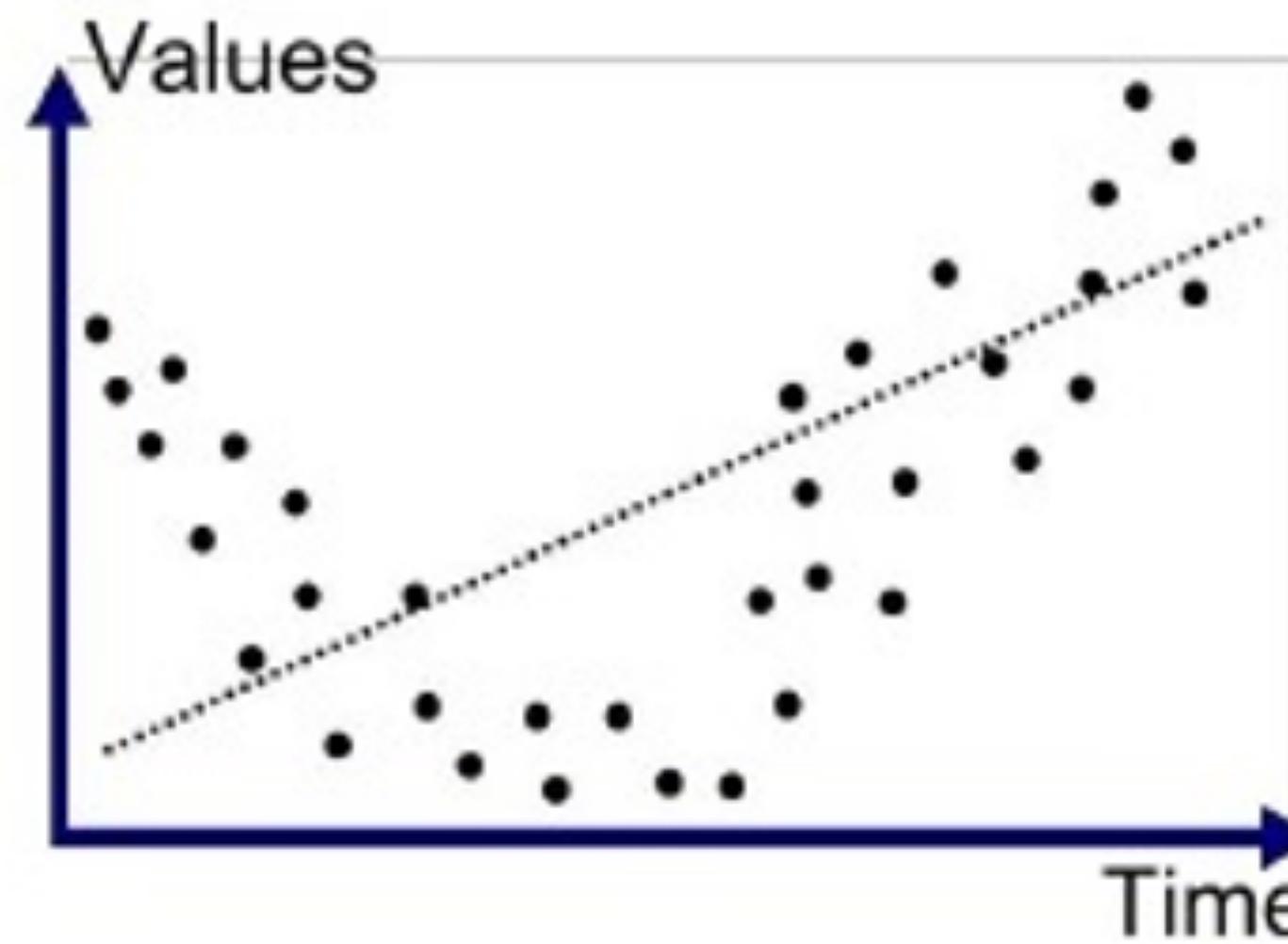
Evaluation with the mean squared error

$$\text{MSE}(\mathbf{W}, b) = \frac{1}{N} \sum_{i=1}^N \left(y_i^{(\text{test})} - \hat{y}_i^{(\text{test})} \right)^2$$



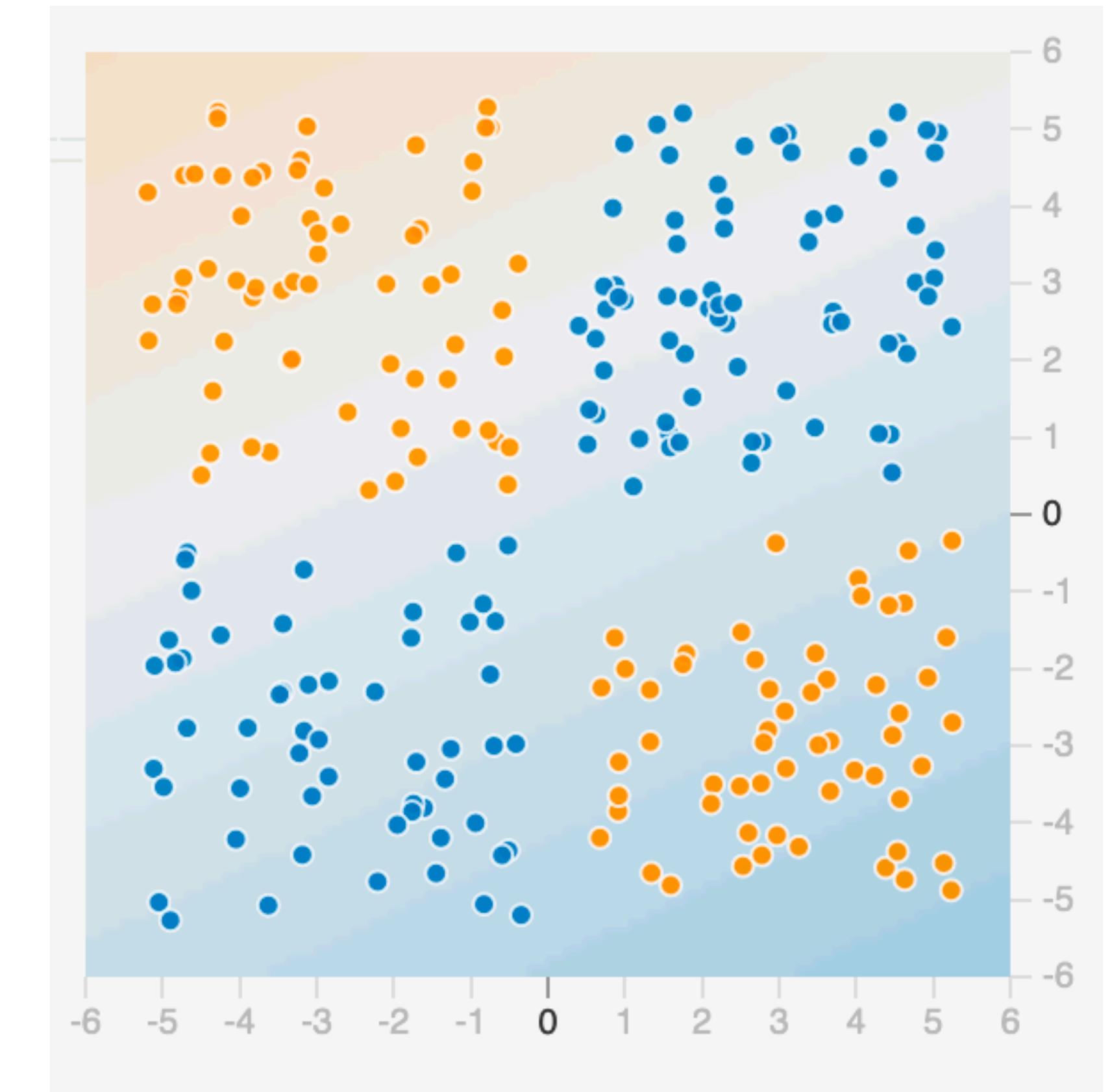
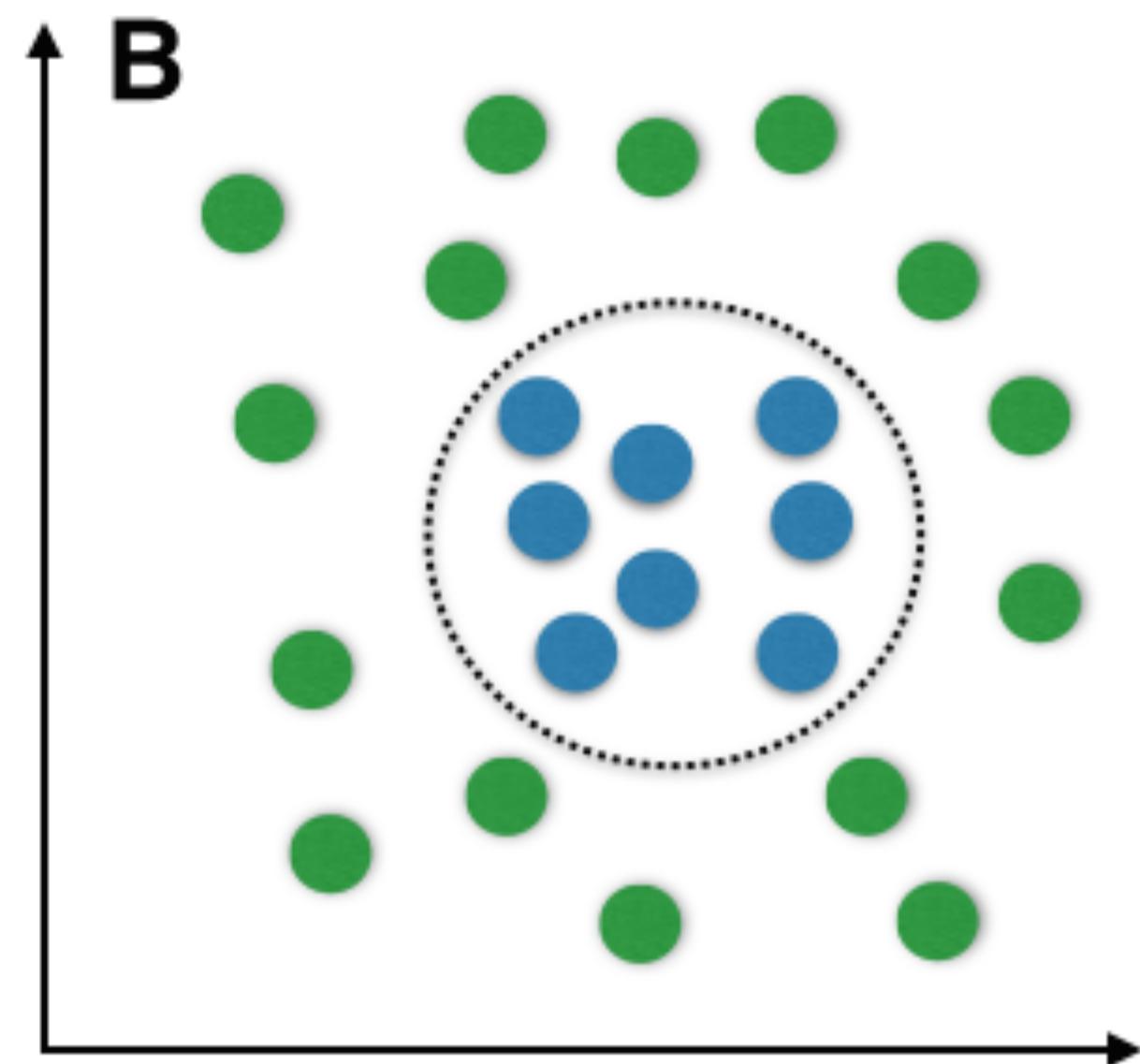
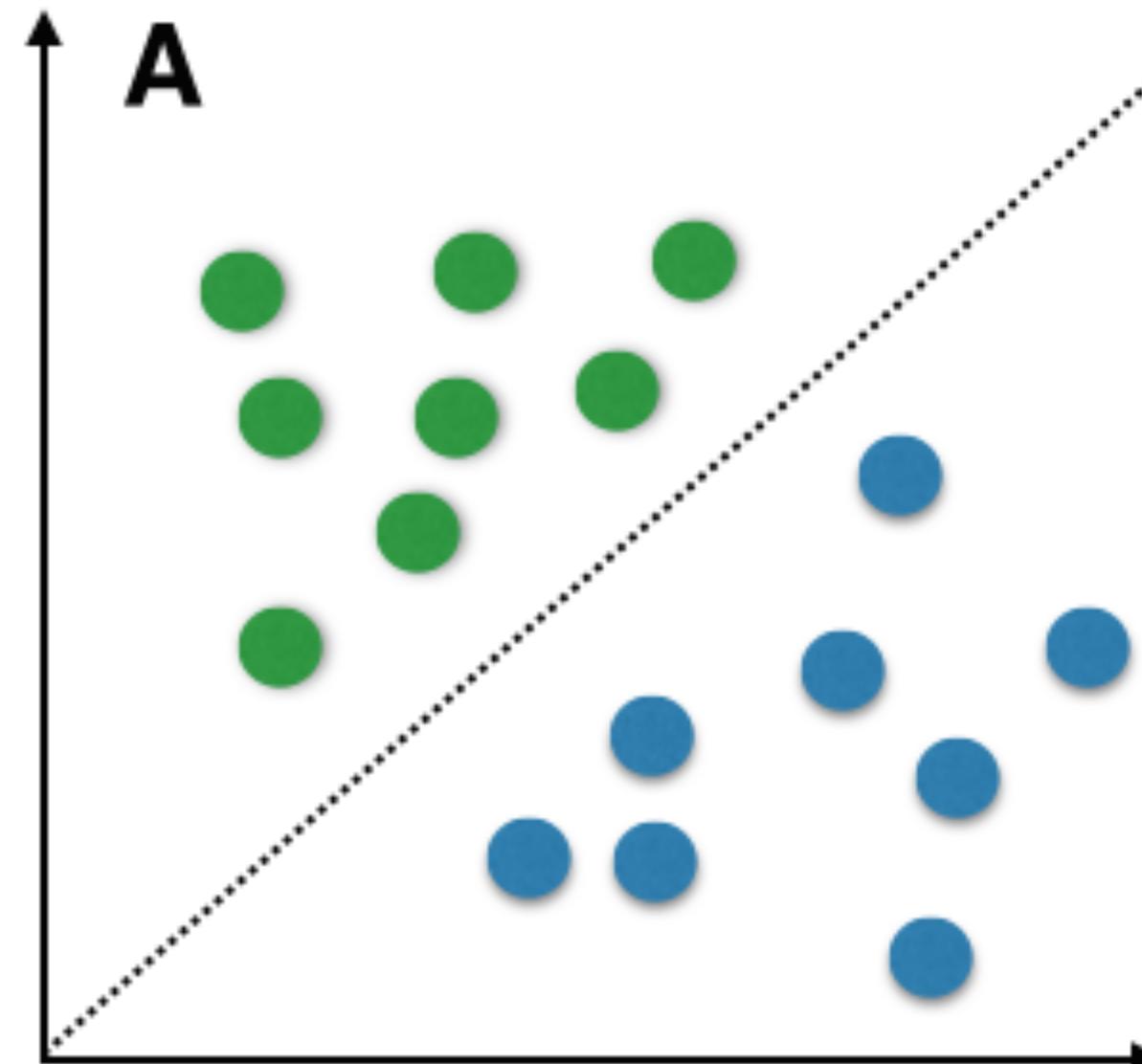
The bias-variance tradeoff

The lower the losses, the better? Not really.



What if the problem is more complicated?

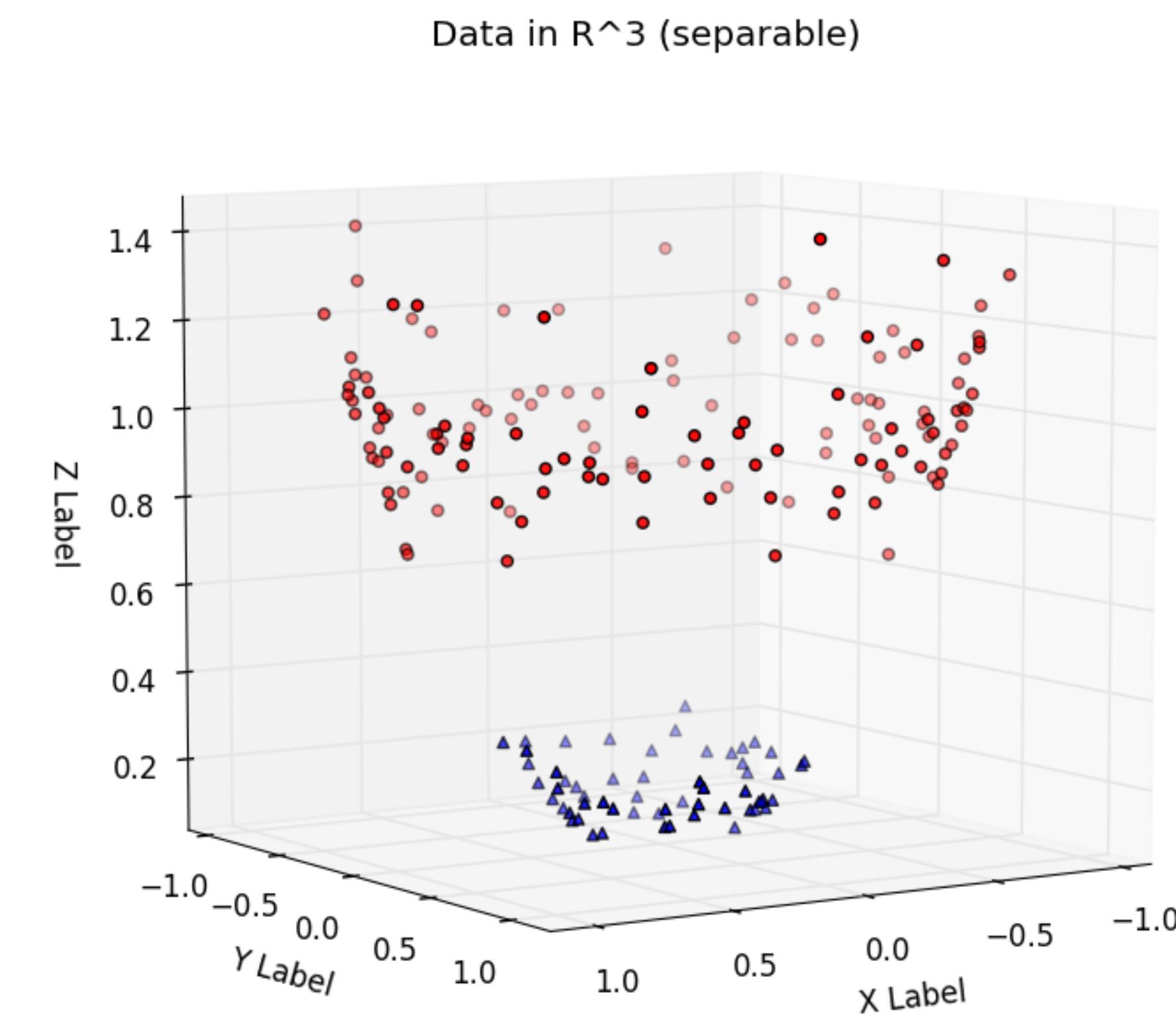
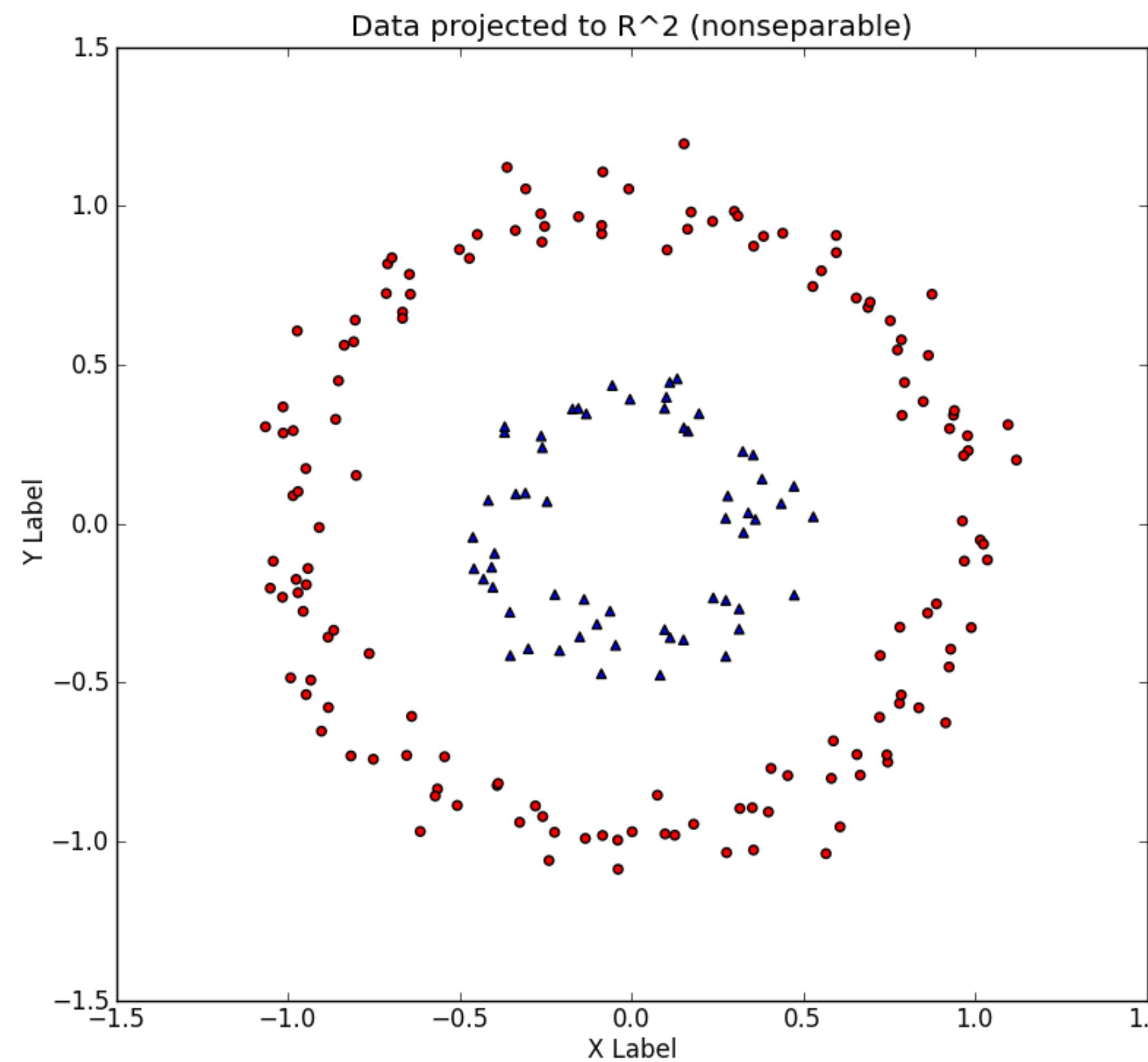
What if the problem is not so simple?



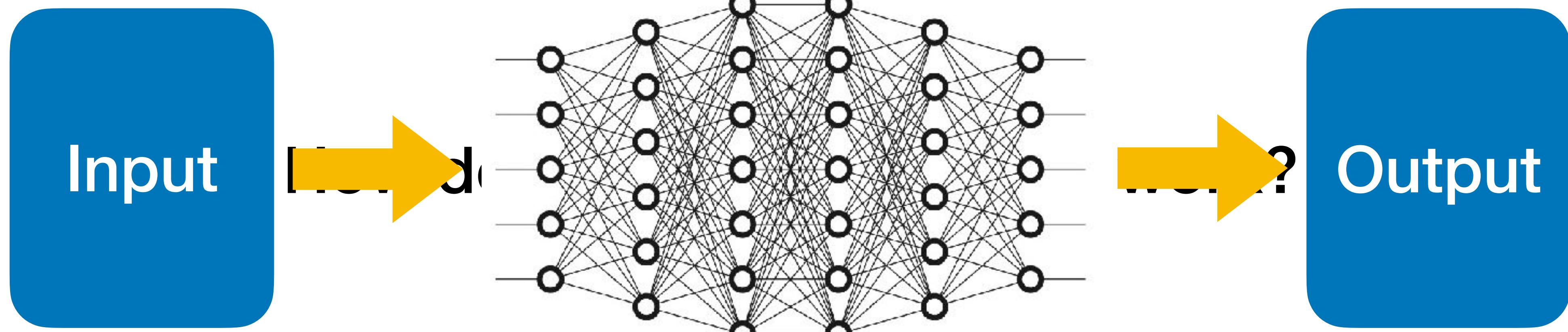
Possible solutions:

- map into another space (kernel trick)
- Add more layers (deep learning)

Kernel Trick

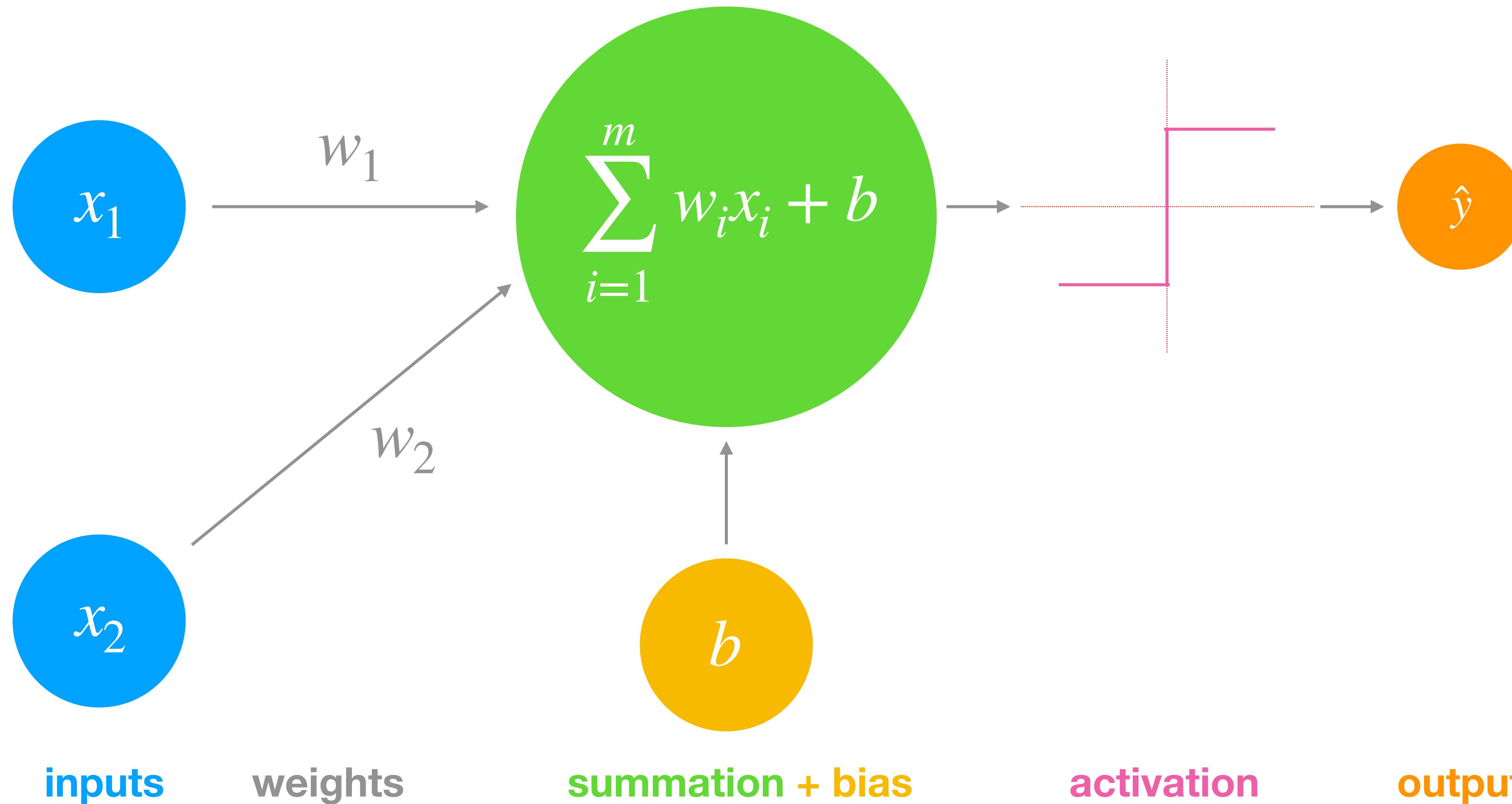


From (traditional) ML to DL

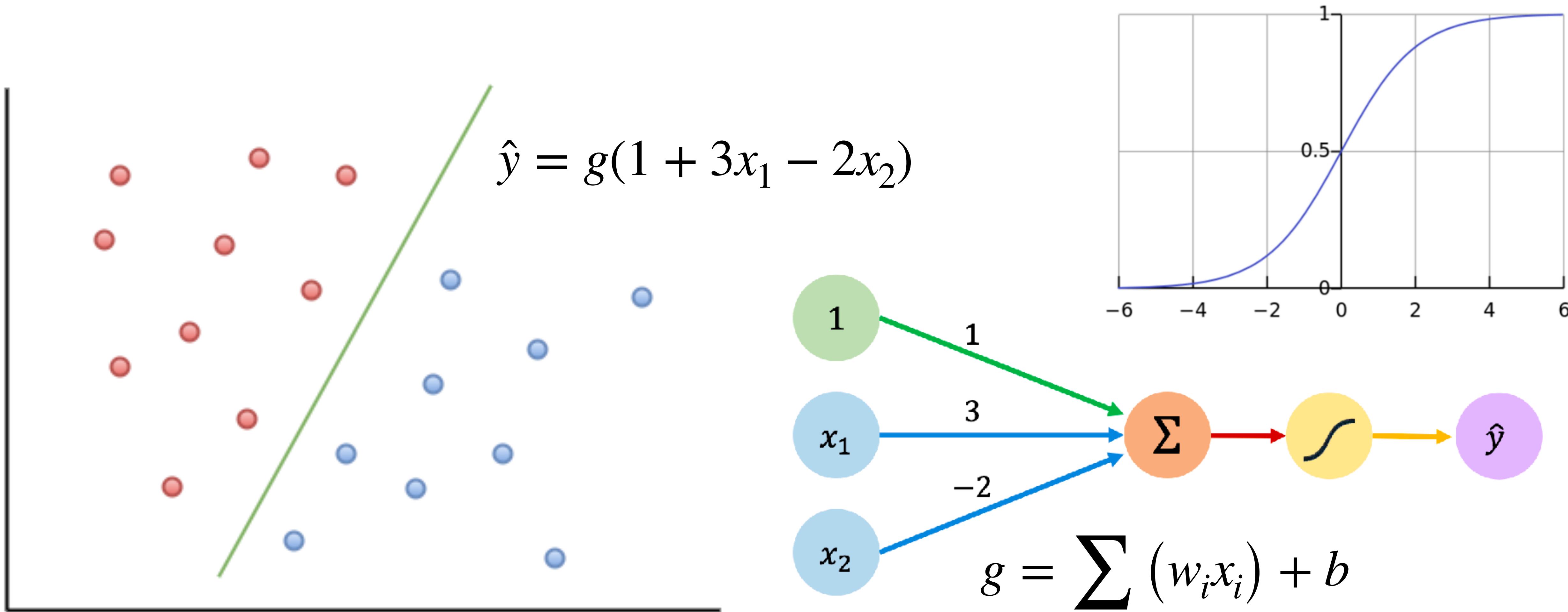


$$\hat{y} = f_{\text{NN}}(x_1, x_2, \dots, x_n)$$

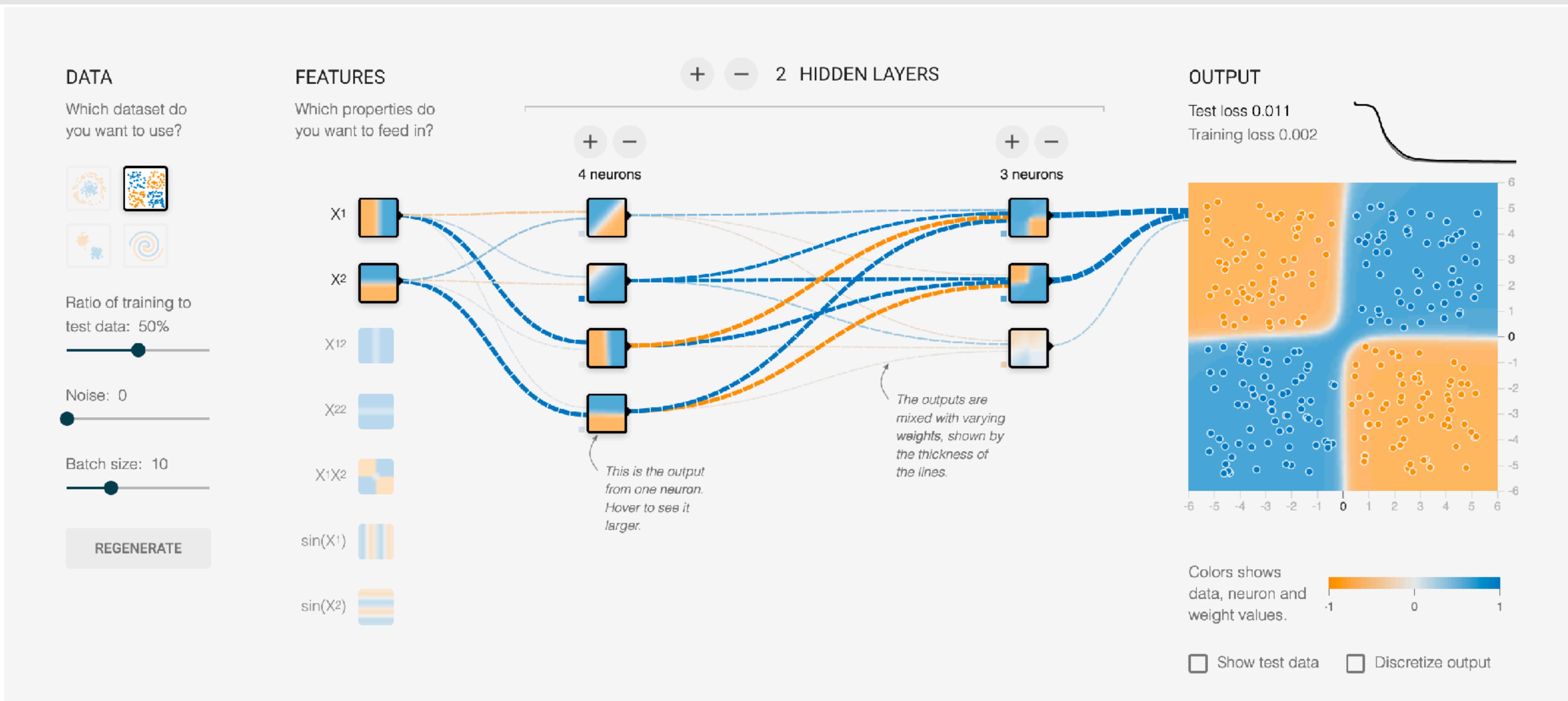
Prediction



Neurons/Perceptrons & Activation Functions



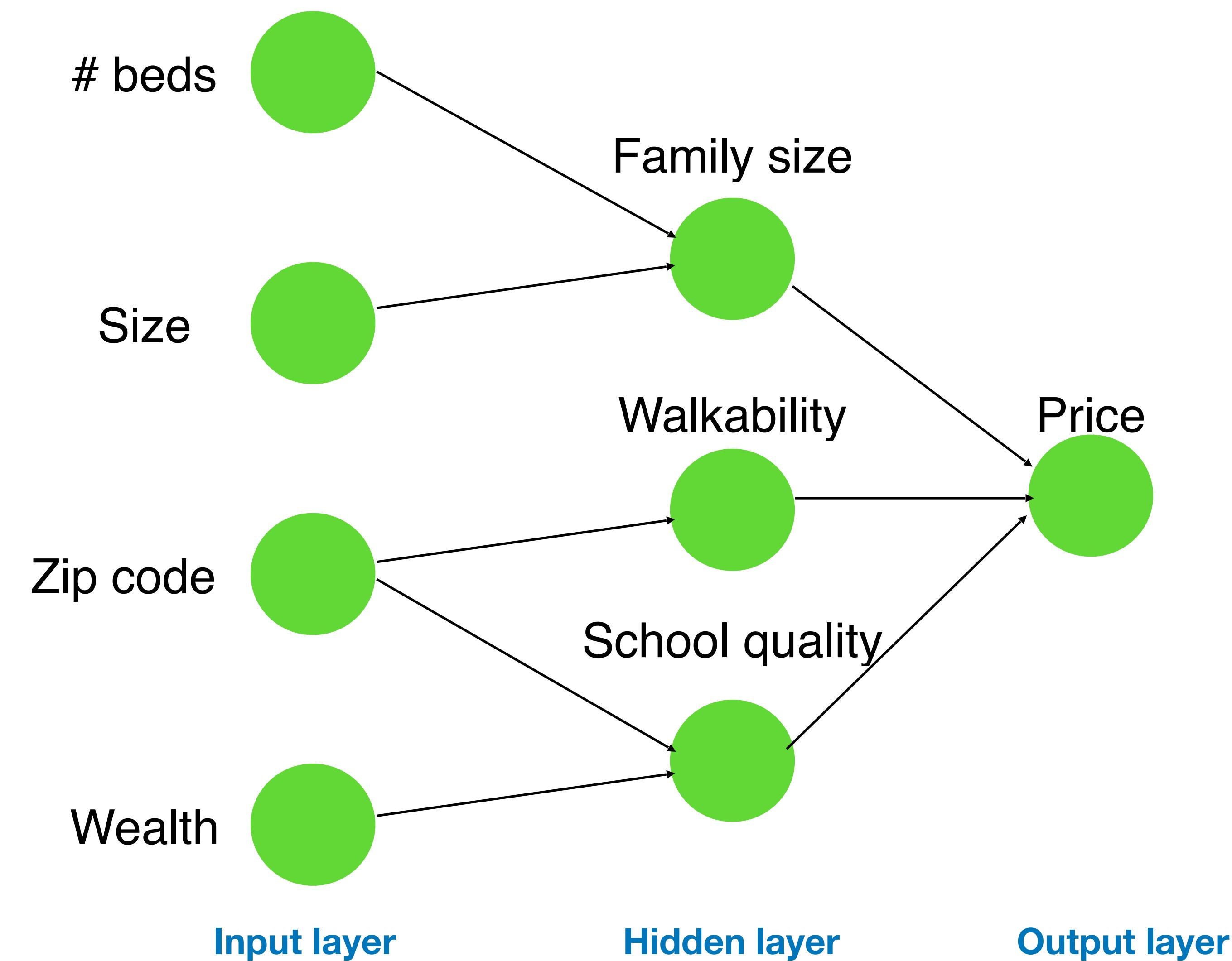
Live demo: Multi-layer Perceptrons



<https://playground.tensorflow.org>

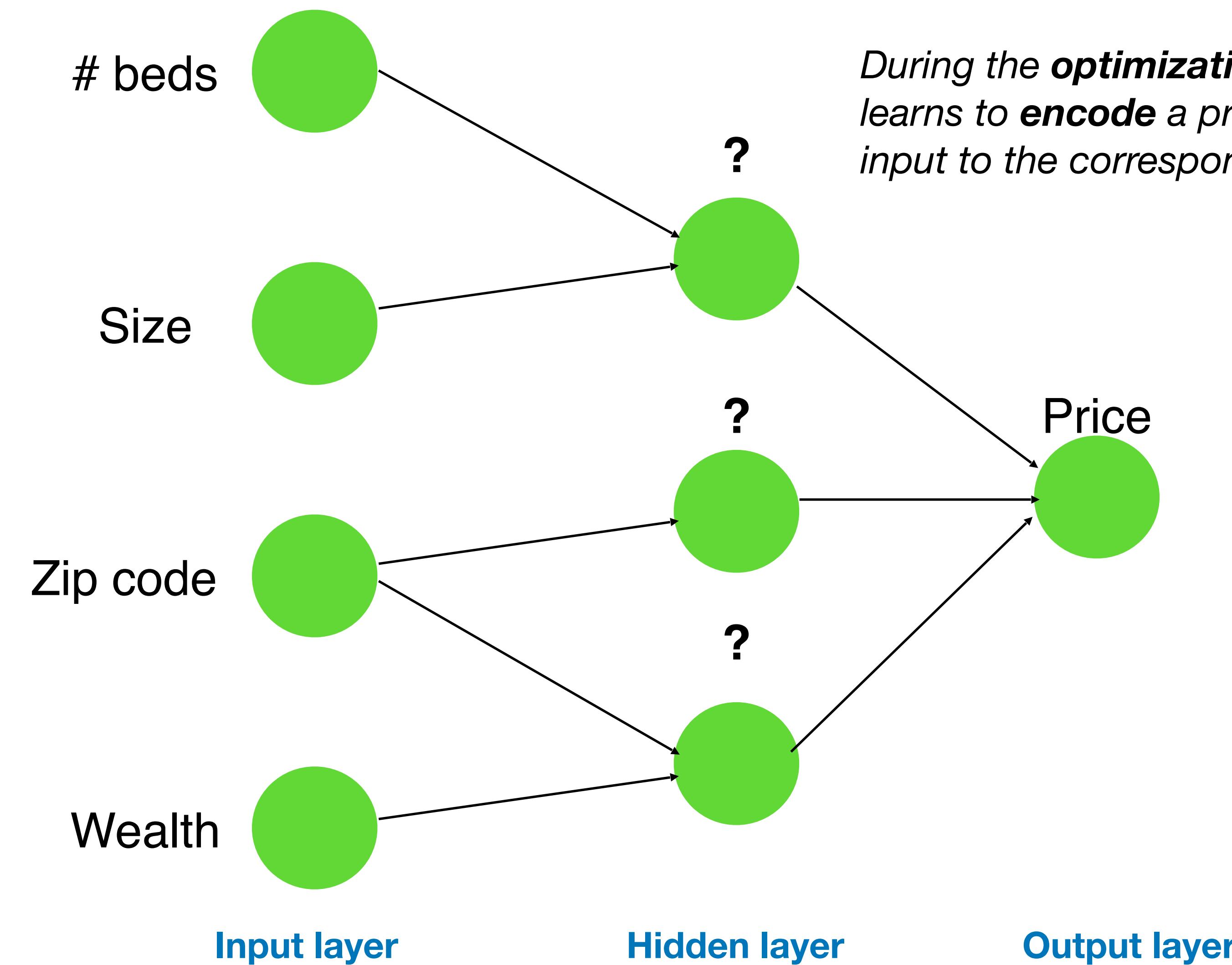
Why multiple layers?

Example: house price prediction model (designed by humans)

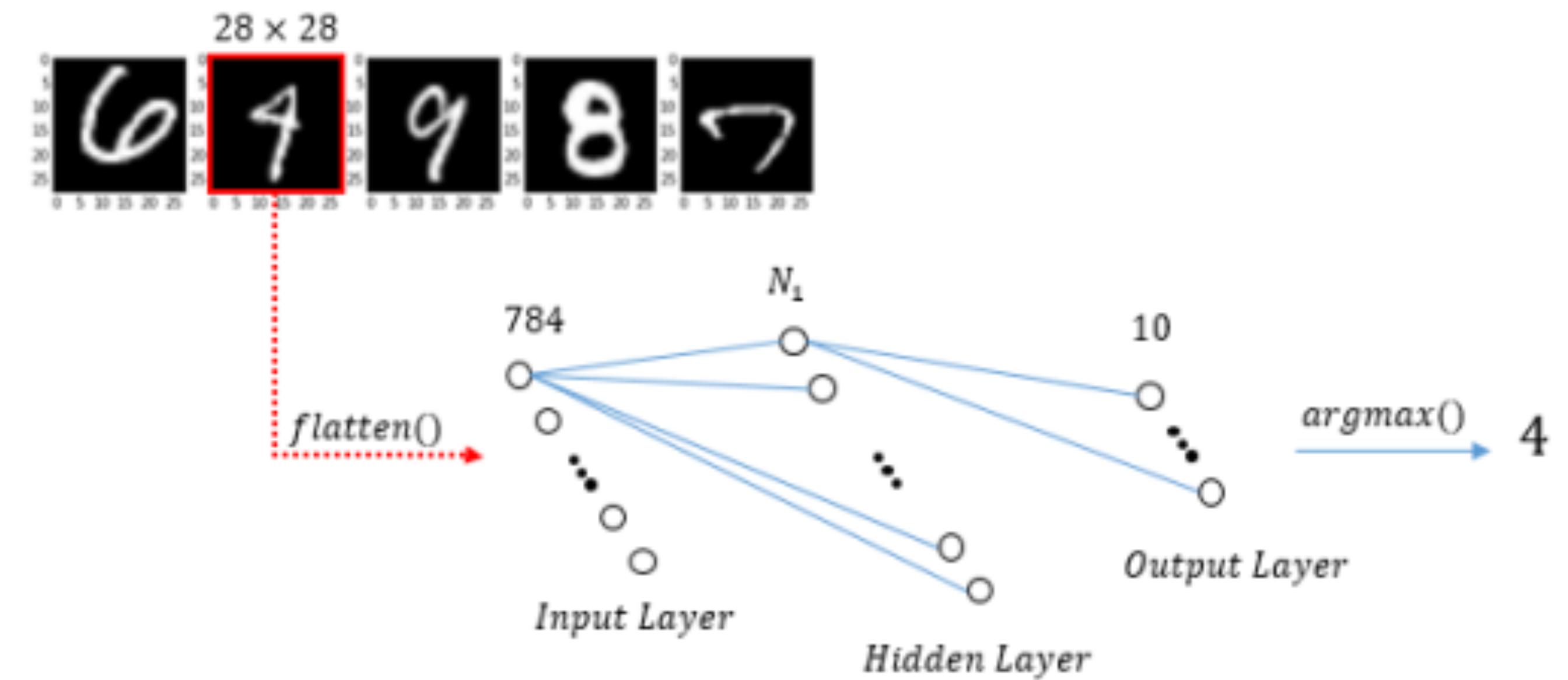


Why multiple layers?

Example: house price prediction model (designed by machines)



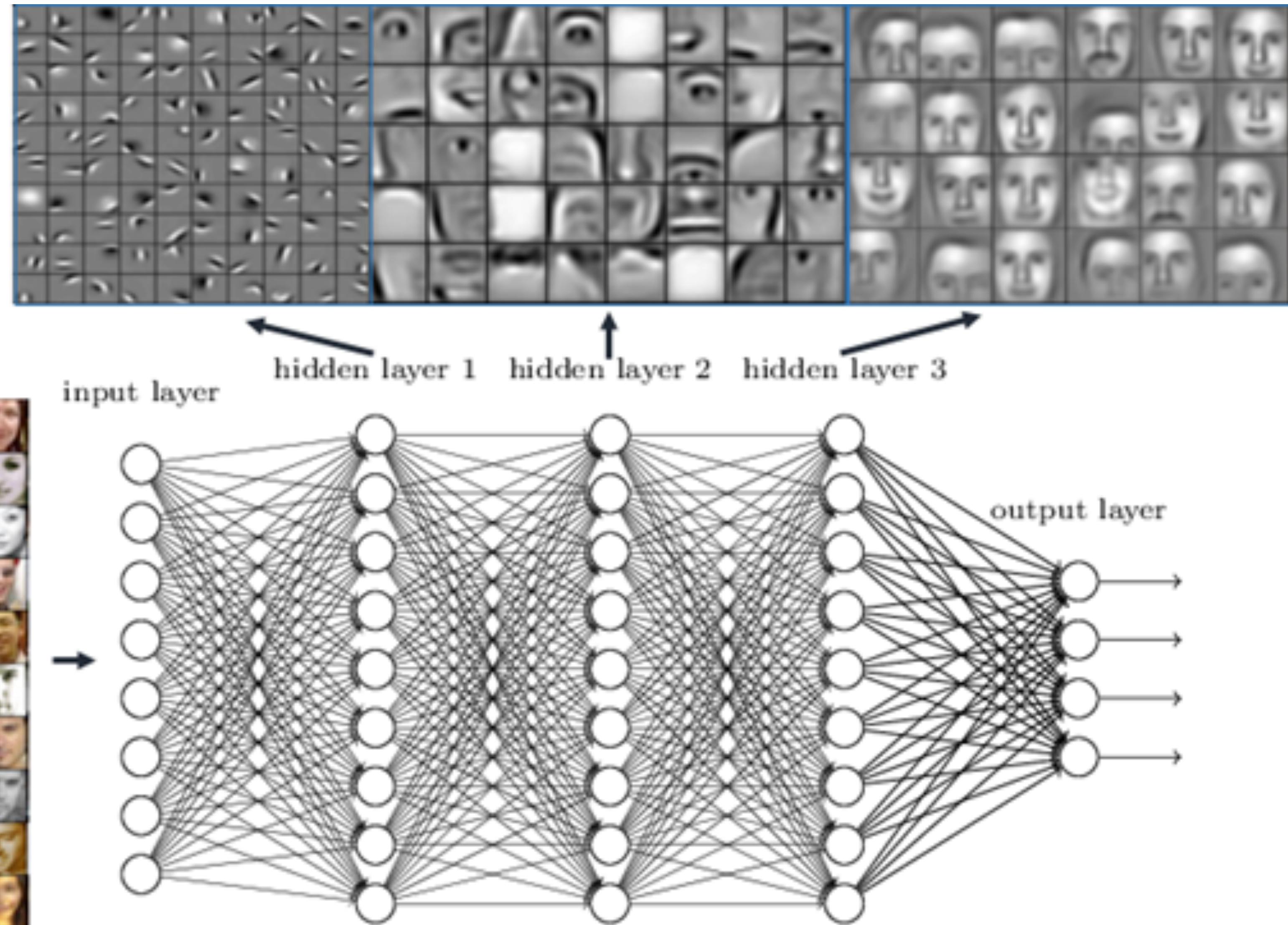
Multi-Layer Perception



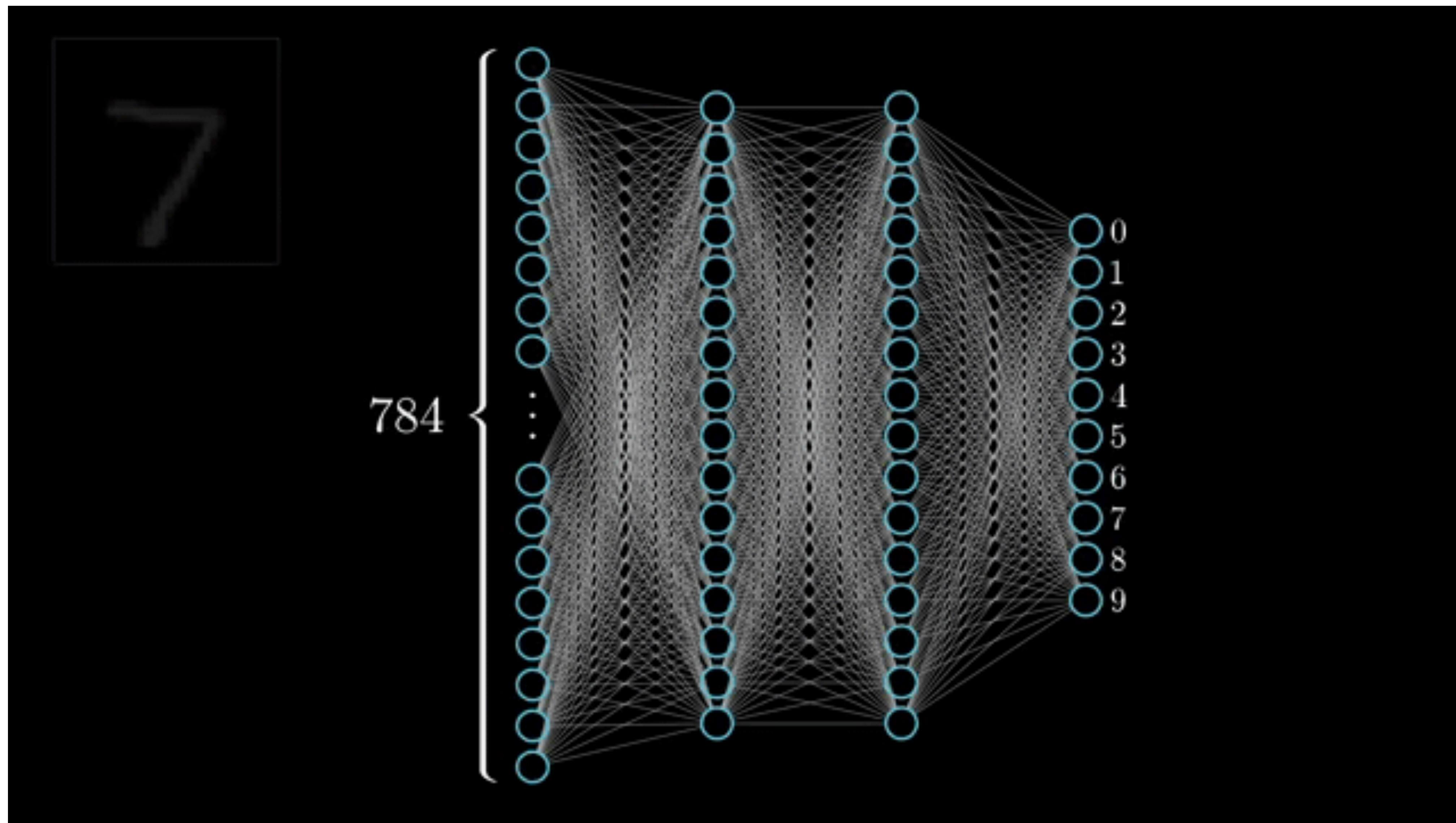
MNIST dataset
(Modified National Institute for Standards and Technology)

A DNN encodes the representation hierarchically

Deep neural networks learn hierarchical feature representations

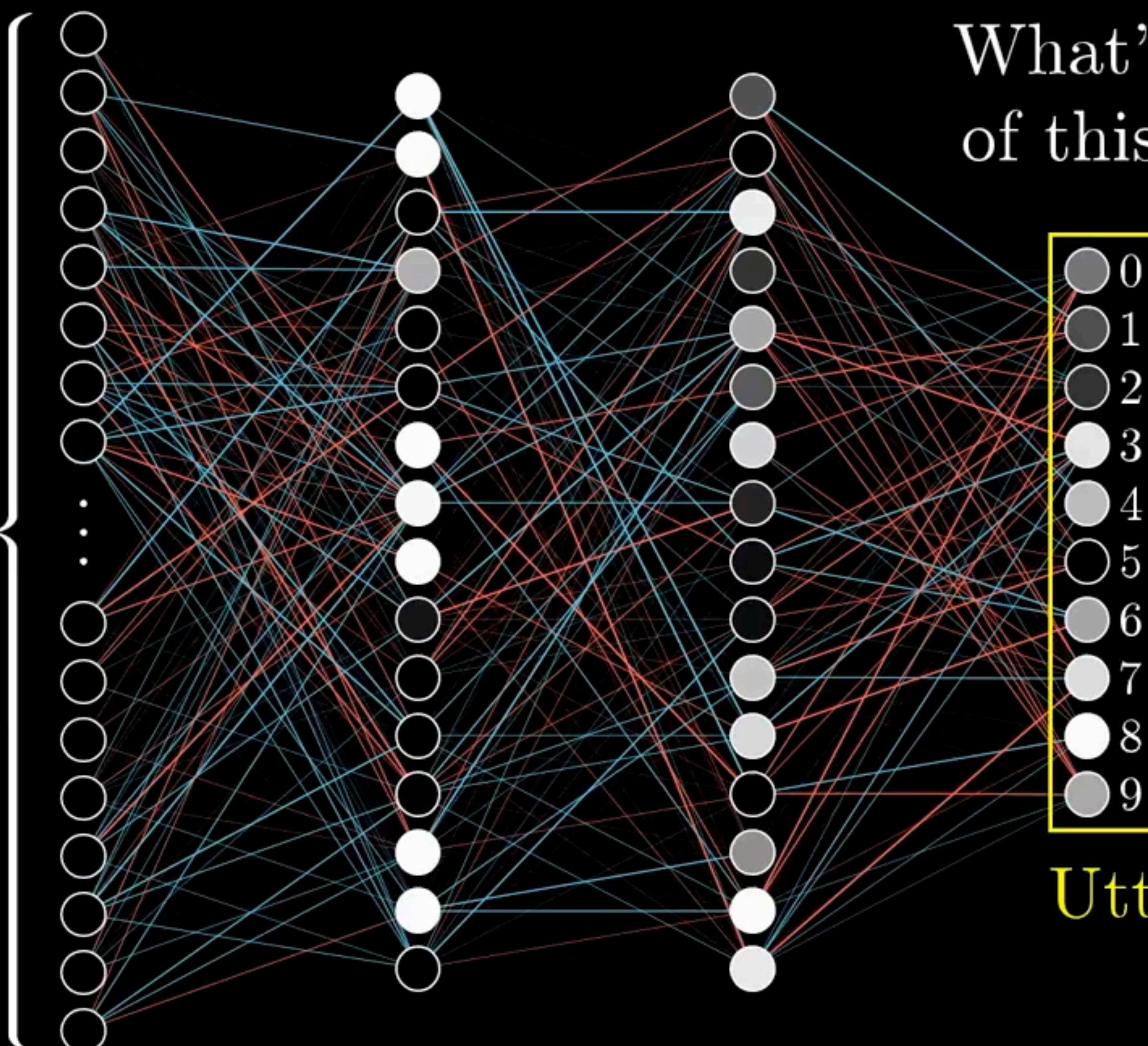


Prediction: forward propagation





784

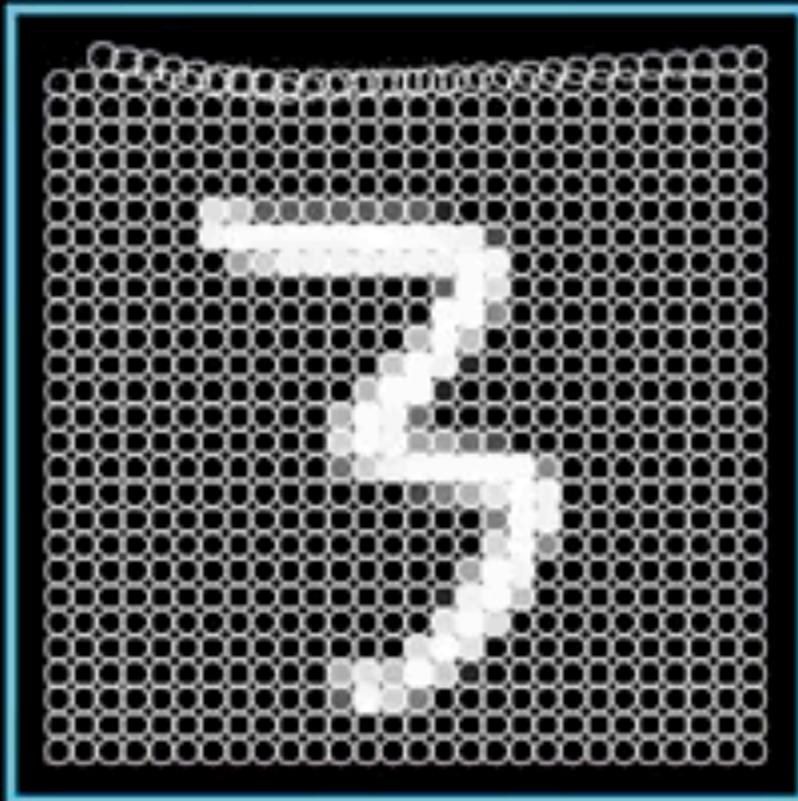


What's the “cost”
of this difference?

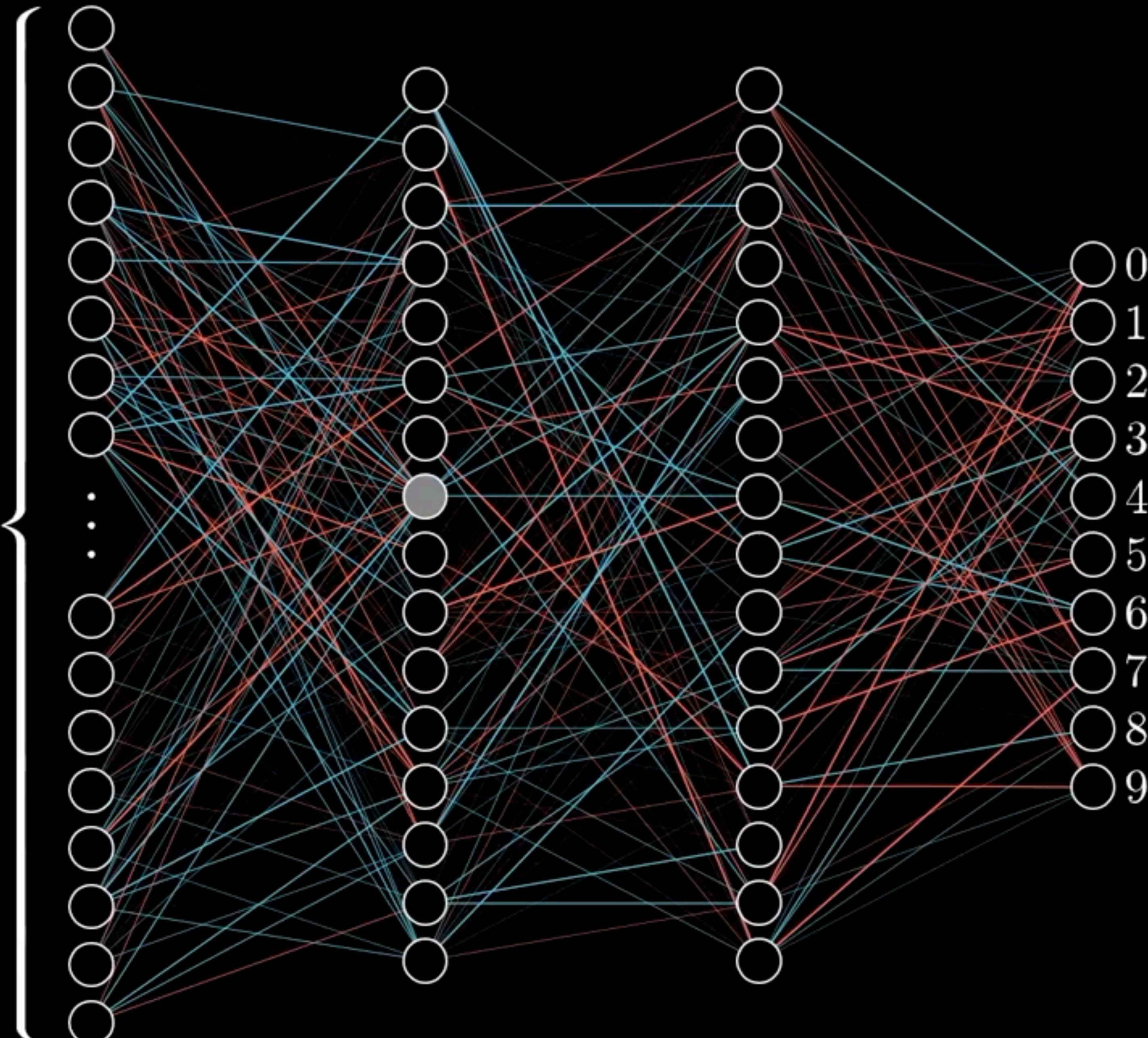
Utter trash

Training: backward propagation

Animation: 3blue1brown



784

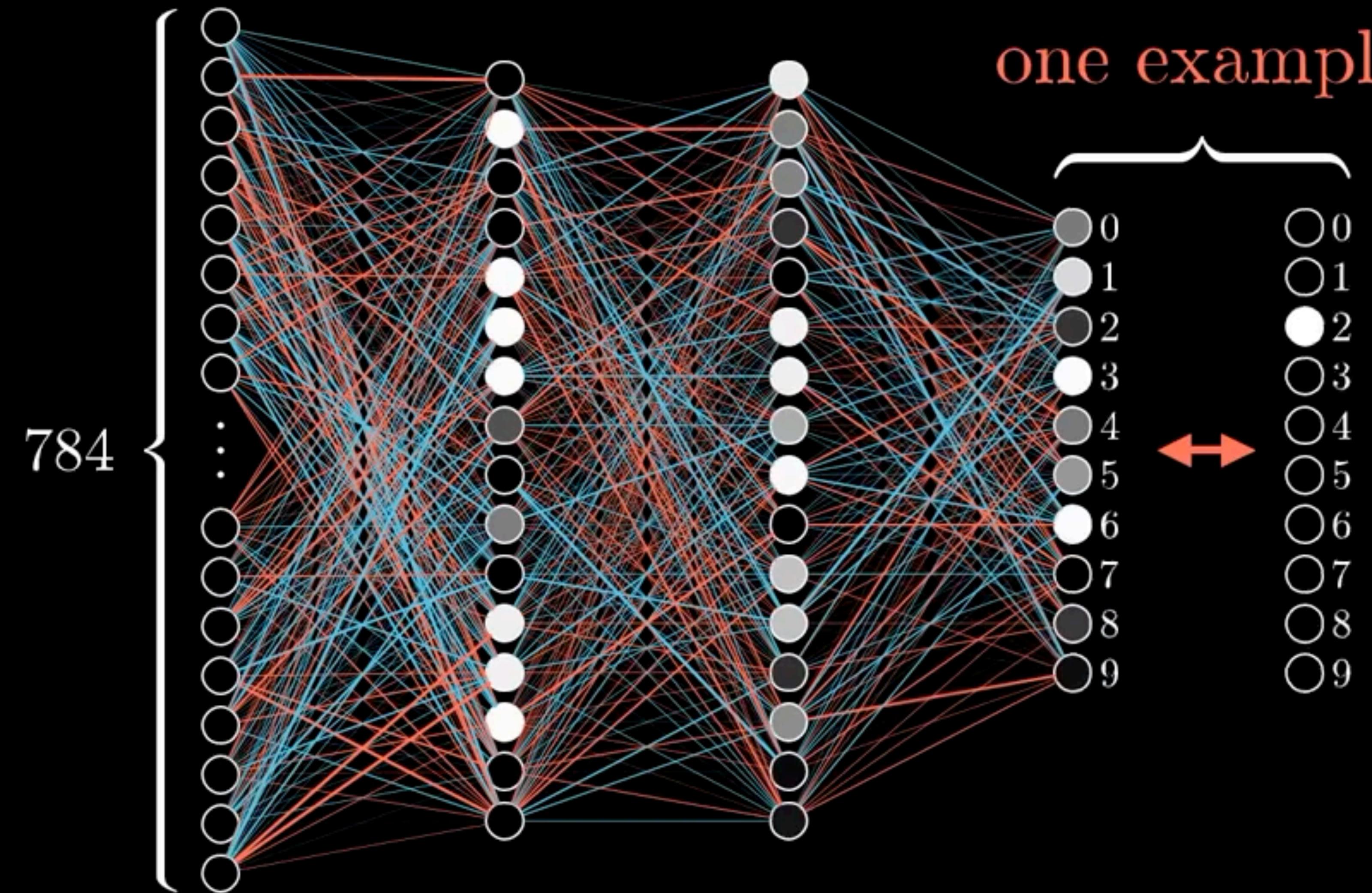


Training: backward propagation

Animation: 3blue1brown

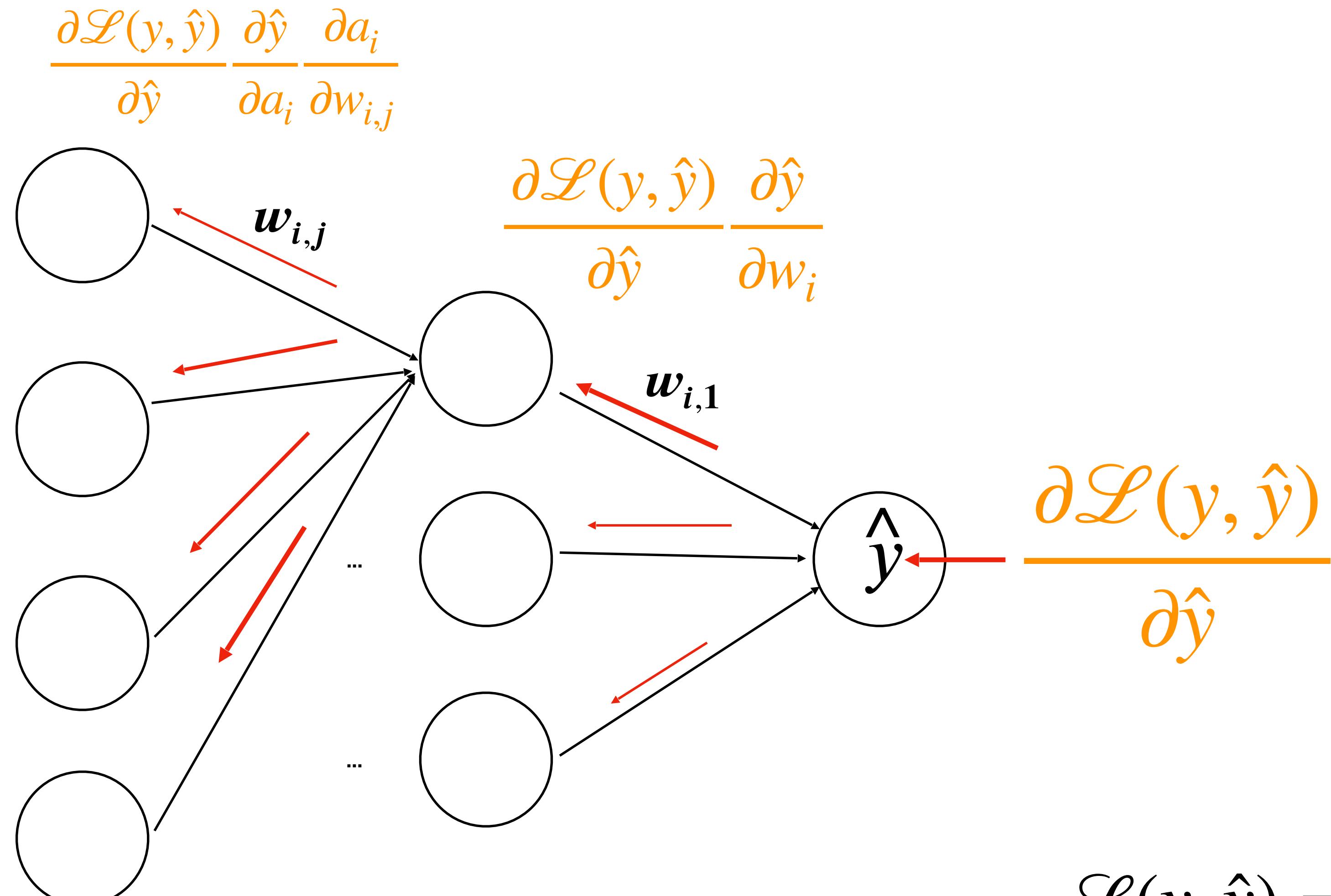


Cost of
one example



Training: backward propagation

Training: backward propagation



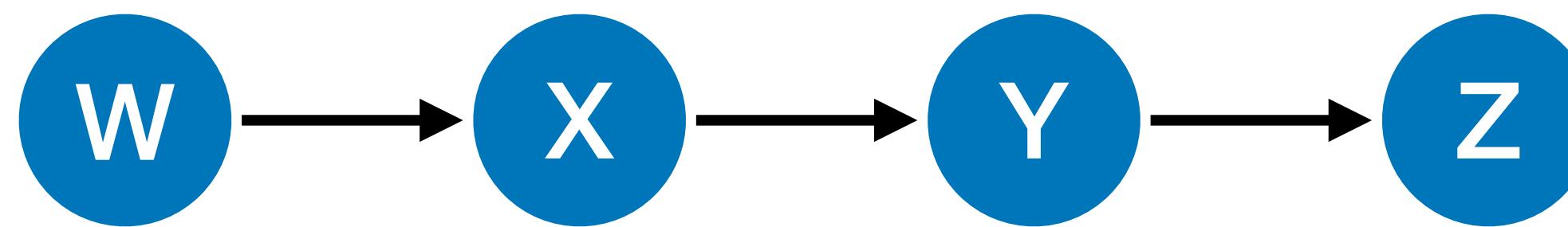
- In the context of DL we need to compute the gradient for each layer.
- We do this by applying the **chain rule** of derivatives.
- This algorithm is known as **backpropagation**.

$$\mathcal{L}(y, \hat{y}) = L(\mathbf{W}, b) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Neural Network: the deeper, the better?

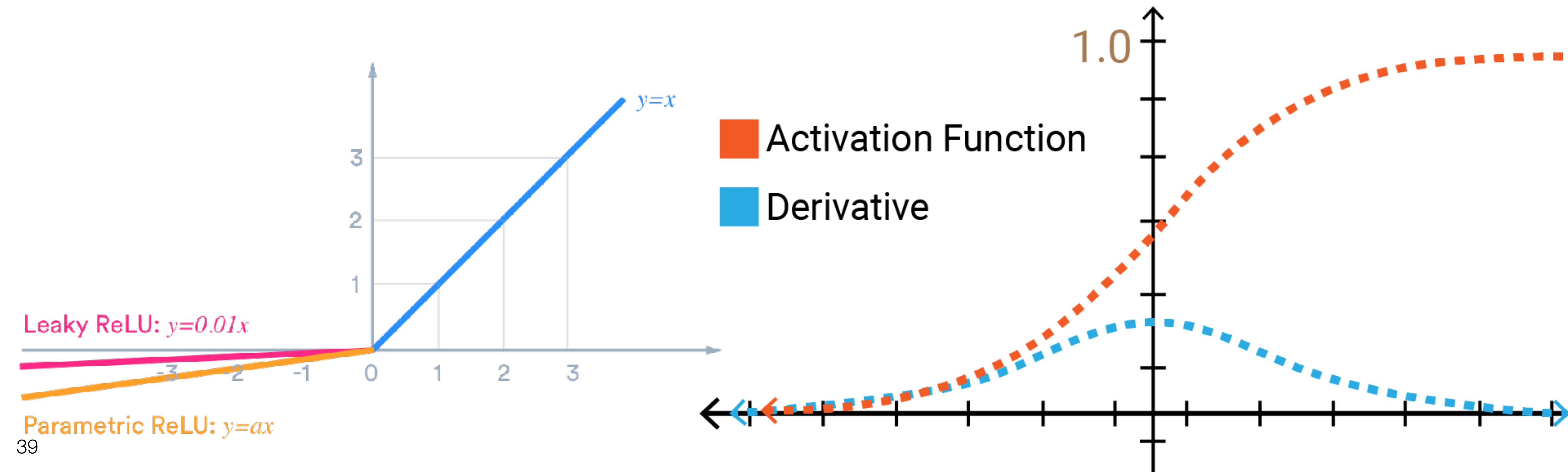
Not really.

The vanishing gradient problem



Chain rule

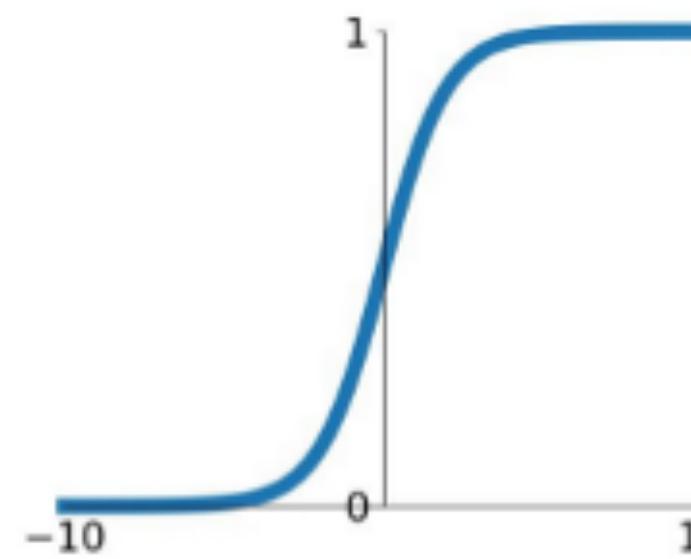
$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w}$$



Activation Functions

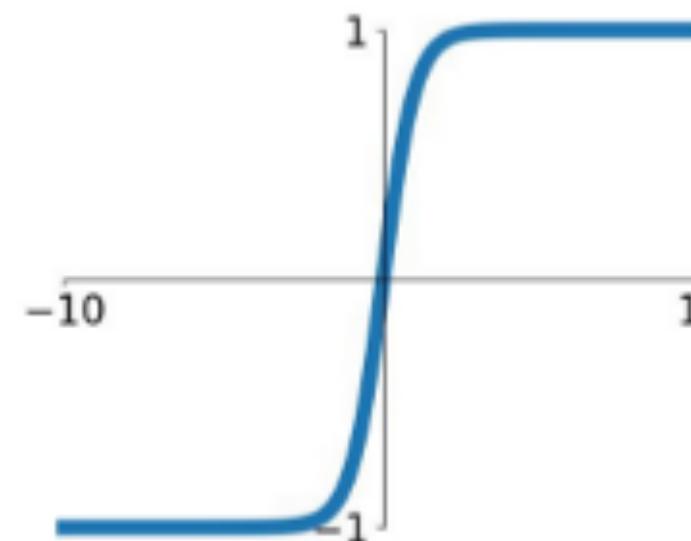
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



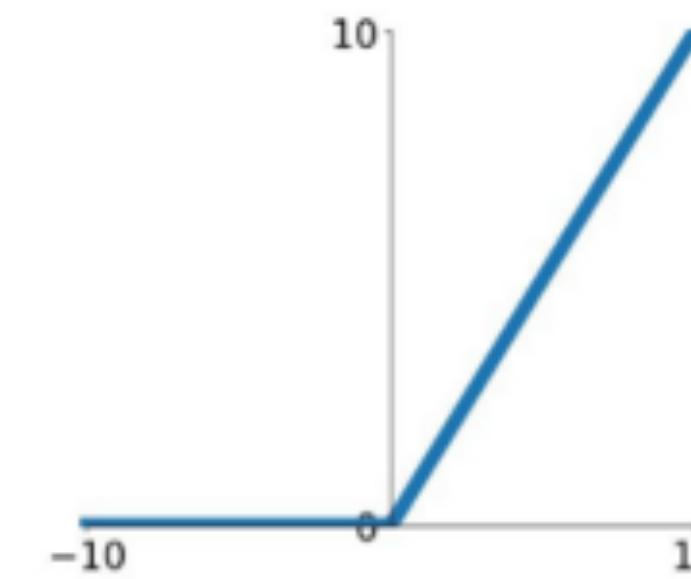
tanh

$$\tanh(x)$$



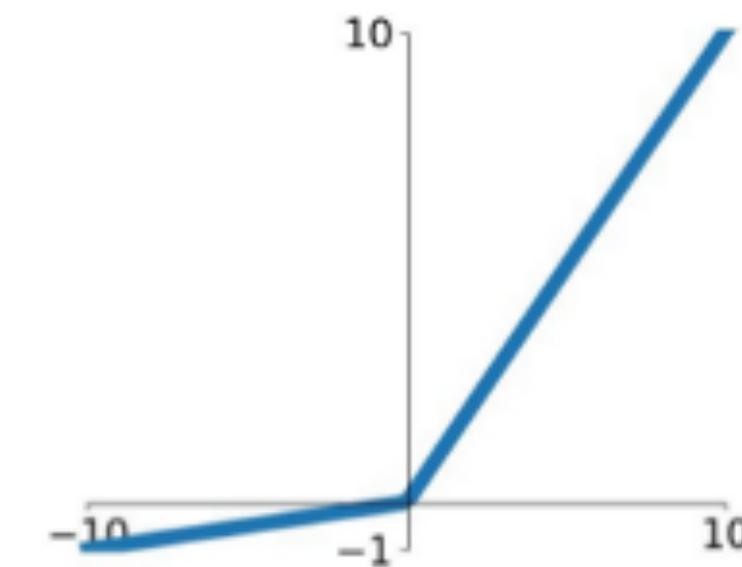
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

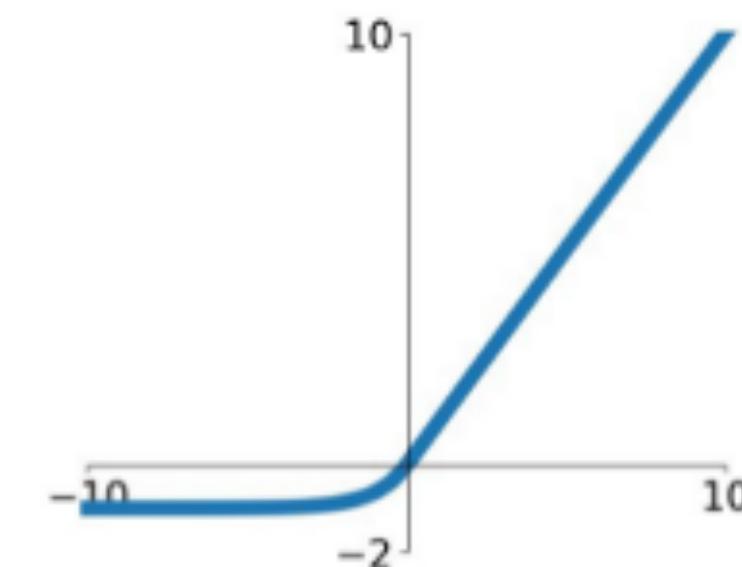


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

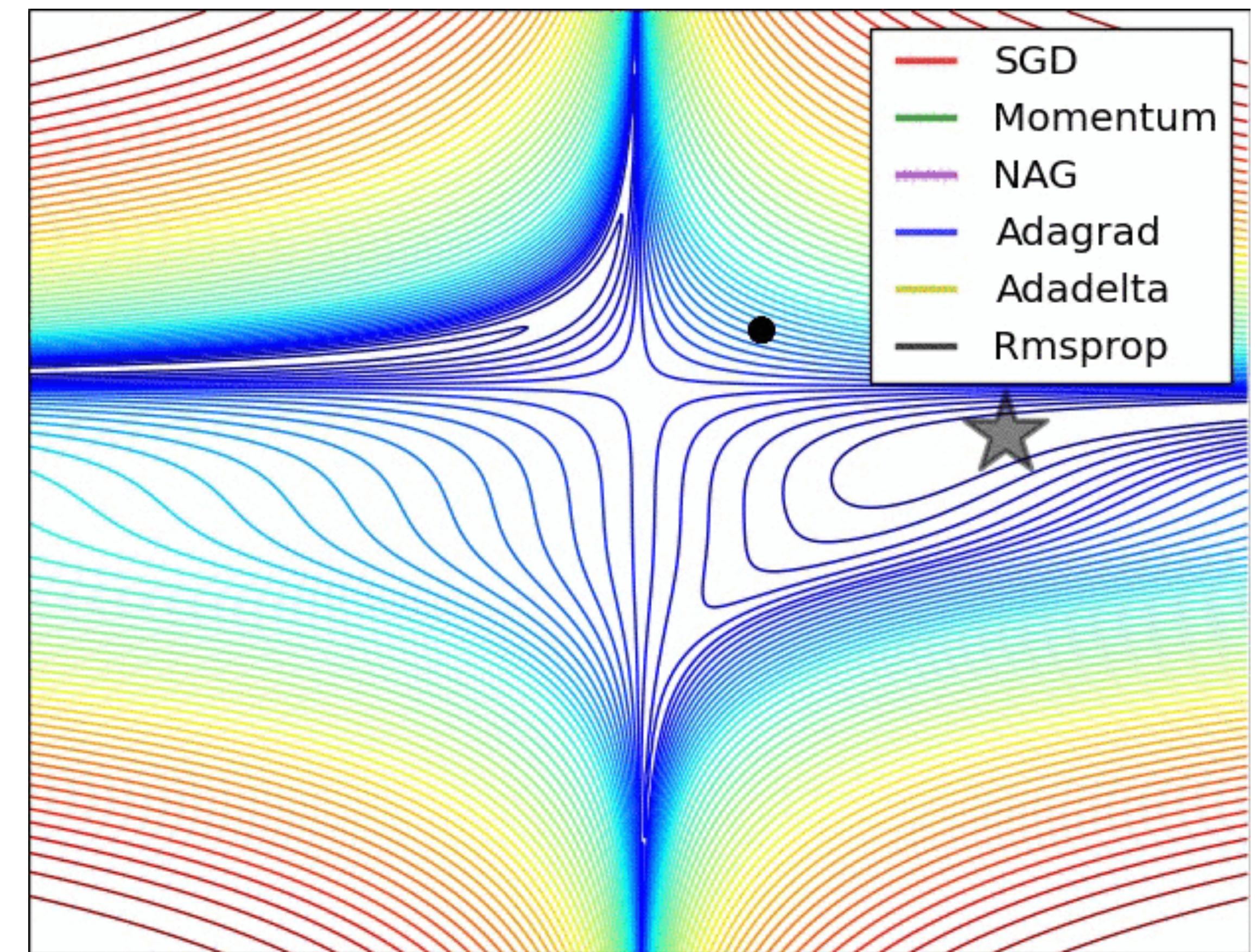
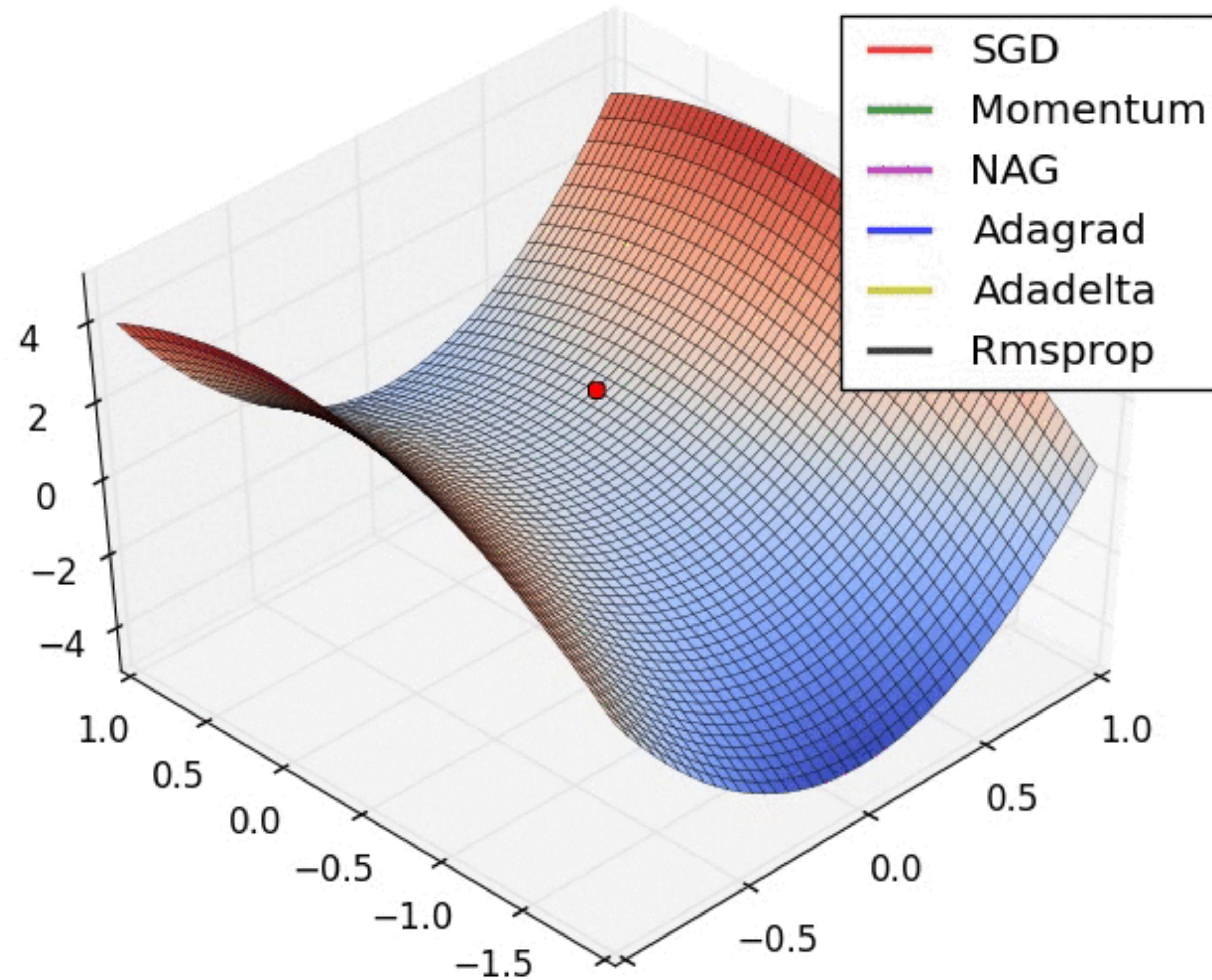
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

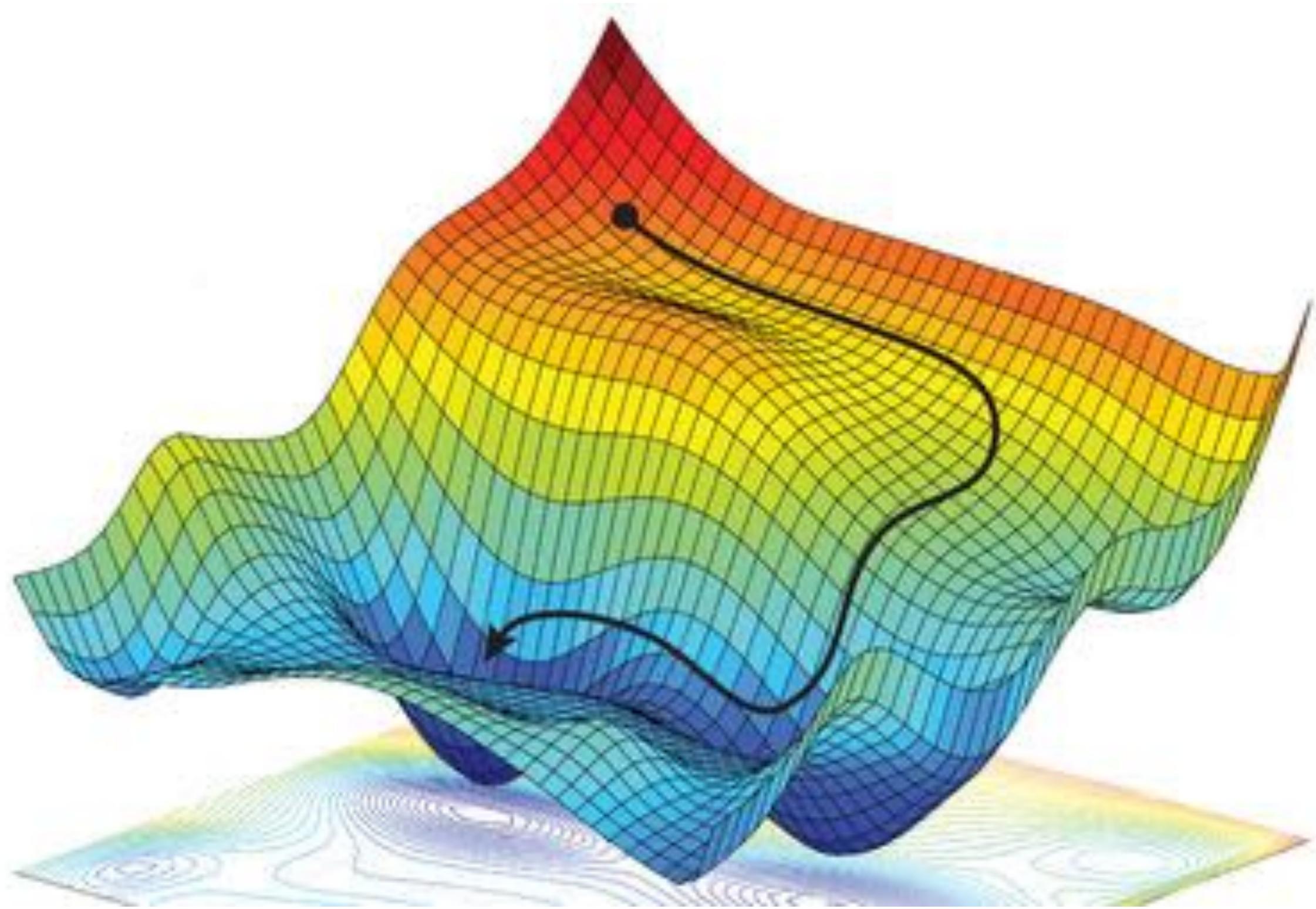


But one can design their own activation functions!

Optimizers



Source: medium.com/analytics-vidhya



Common practice for loss functions

Regression

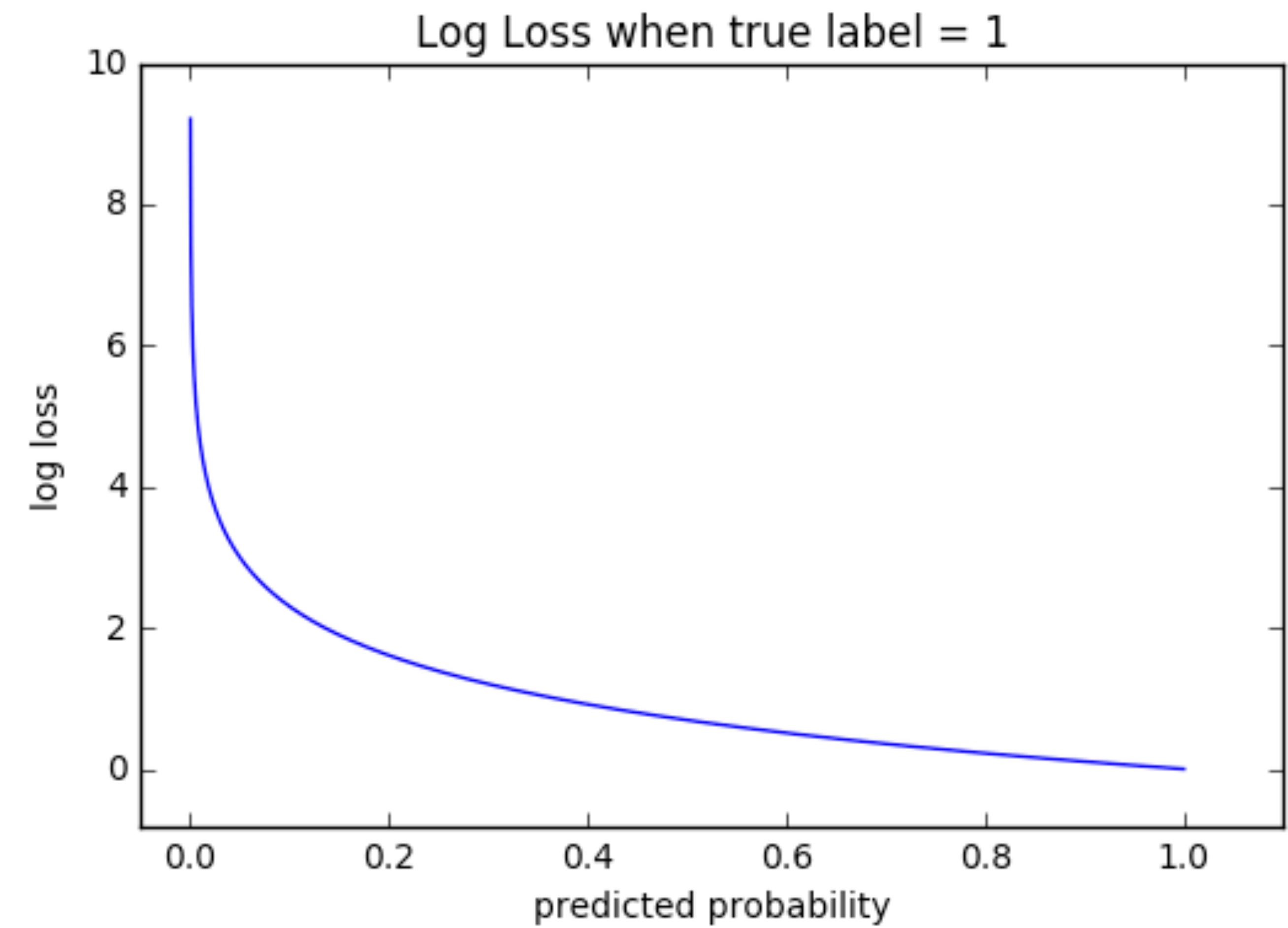
- Mean squared error
- Mean squared logarithmic error
- Mean absolute error

Binary Classification

- Binary cross-entropy
- Hinge loss
- Squared hinge loss

Multi-Class Classification

- Multi-class cross-entropy
- Sparse multi-class cross-entropy
- Kullback-Leibler divergence



DL frameworks

In DL, you need to

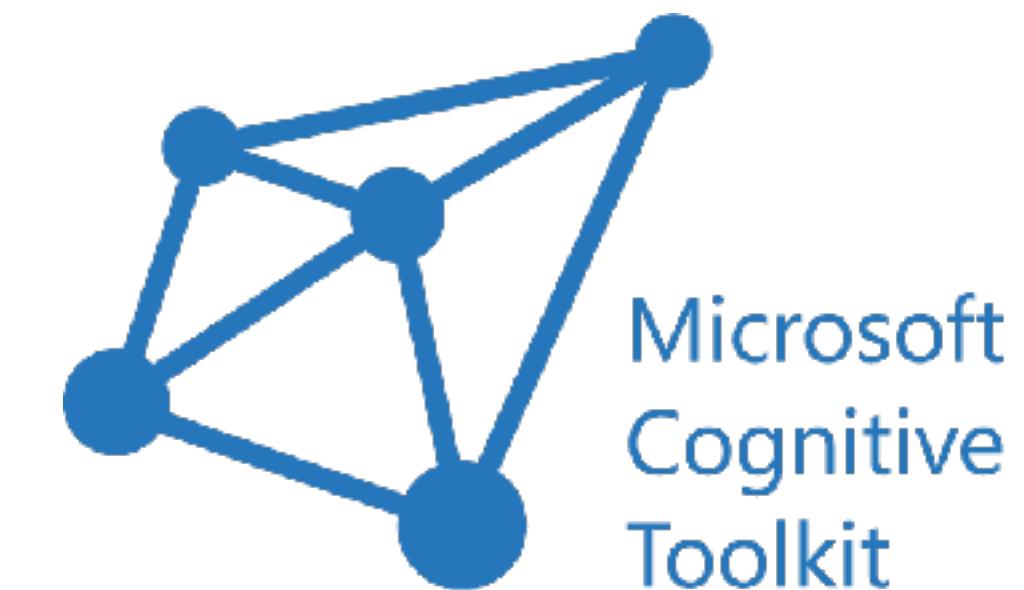
- Define neurons and layers
 - Define loss function
 - Calculate losses
 - Calculate gradient (multivariate calculus)
 - Backward propagation
 - Update weights
-
- Many frameworks exist; **TensorFlow**, **CNTK**, **Torch**, **Keras**, **Theano**, **Caffe**, ...
 - We will use **TensorFlow/Keras**
 - Keras used to call TensorFlow as a *backend*, but is now fully integrated in TensorFlow.



theano



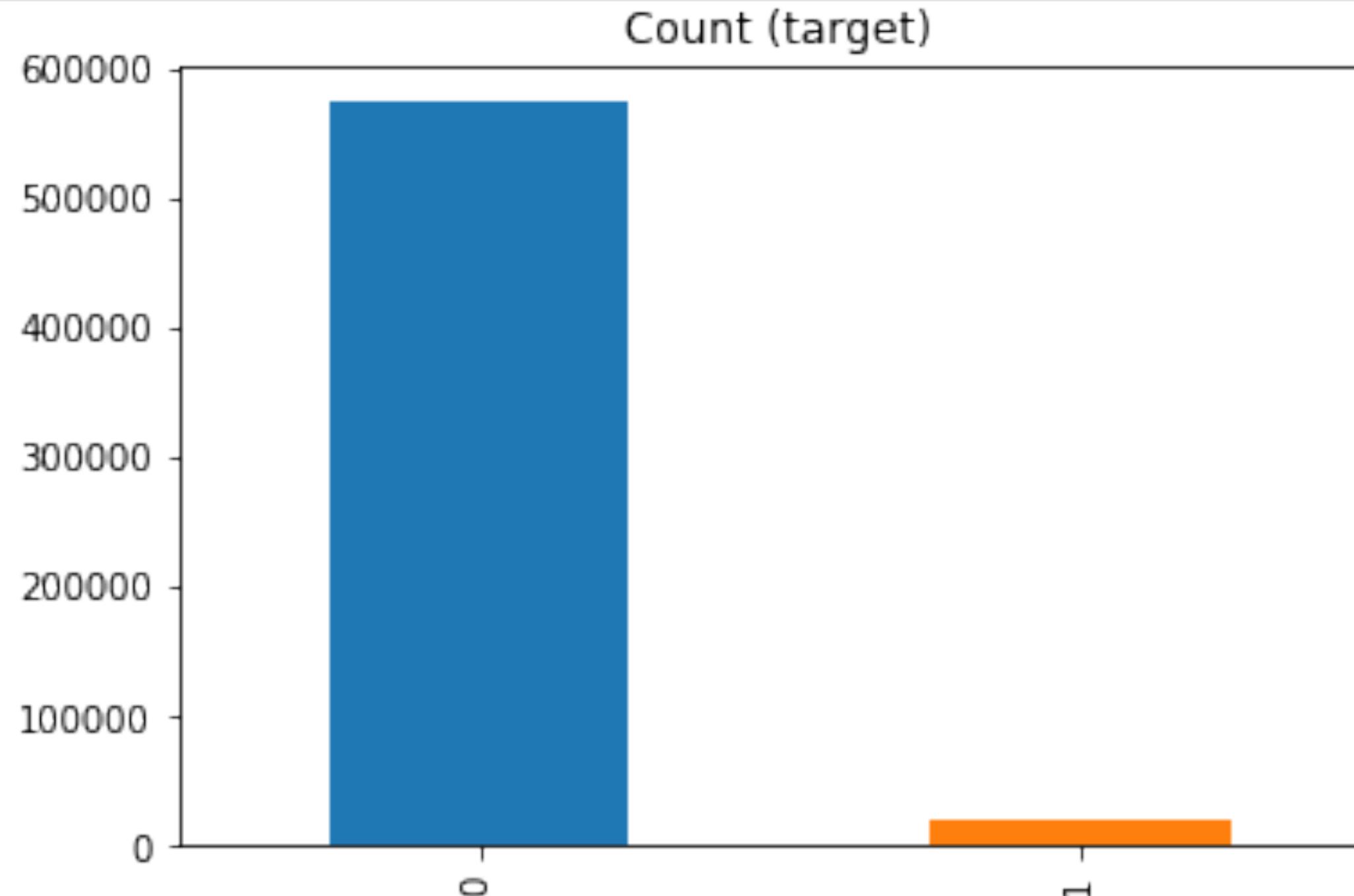
TensorFlow



PyTorch

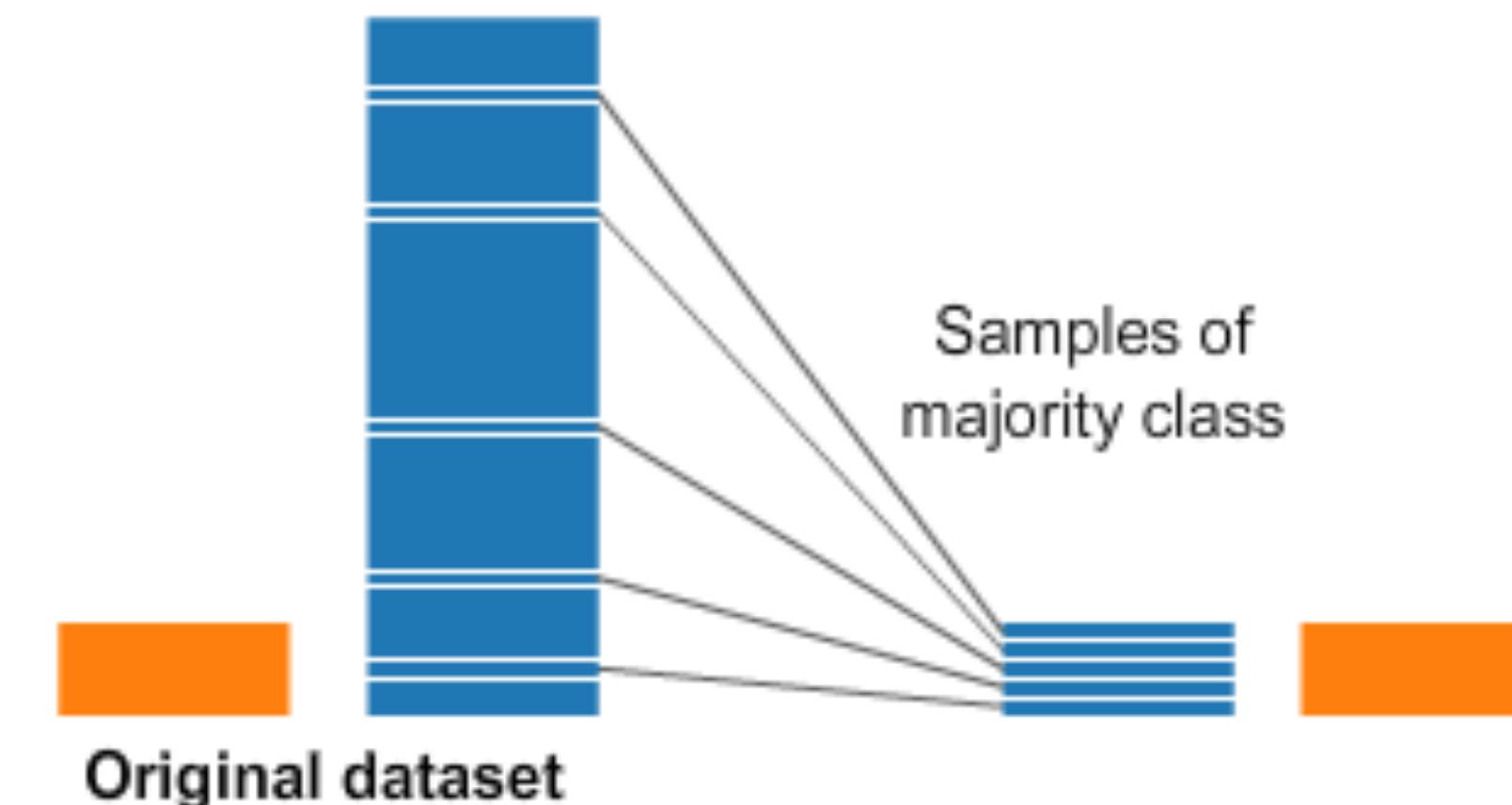
Model Evaluation

Balanced/Imbalanced training set



Pay attention to your data: they may fool your model.

Undersampling

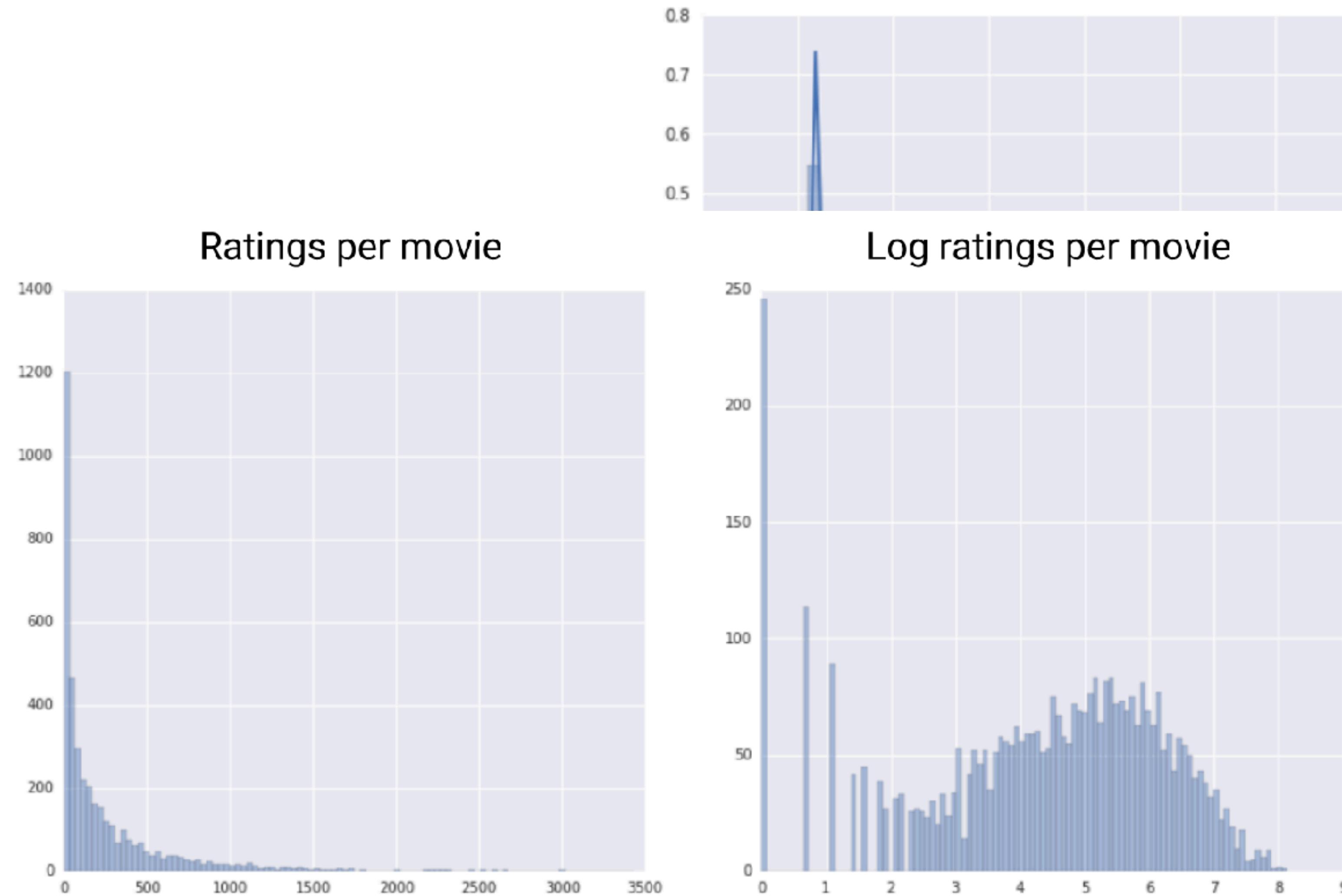


Oversampling



Data Normalization

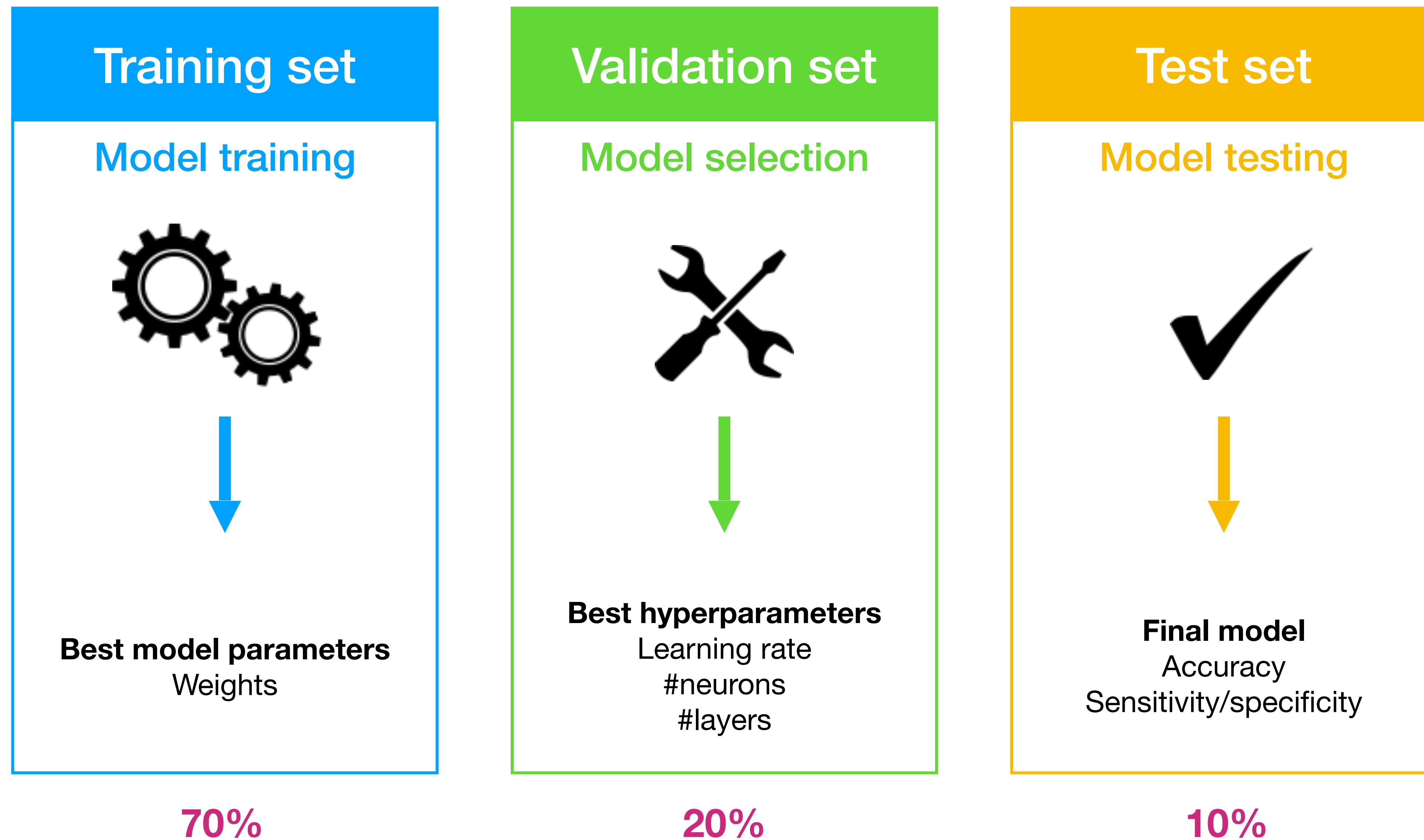
A process to transform the input **data** in a **well-behaved** form.



Further reading: [sklearn scalers](#)

Source: developers.google.com

Dataset splitting



Confusion Matrix

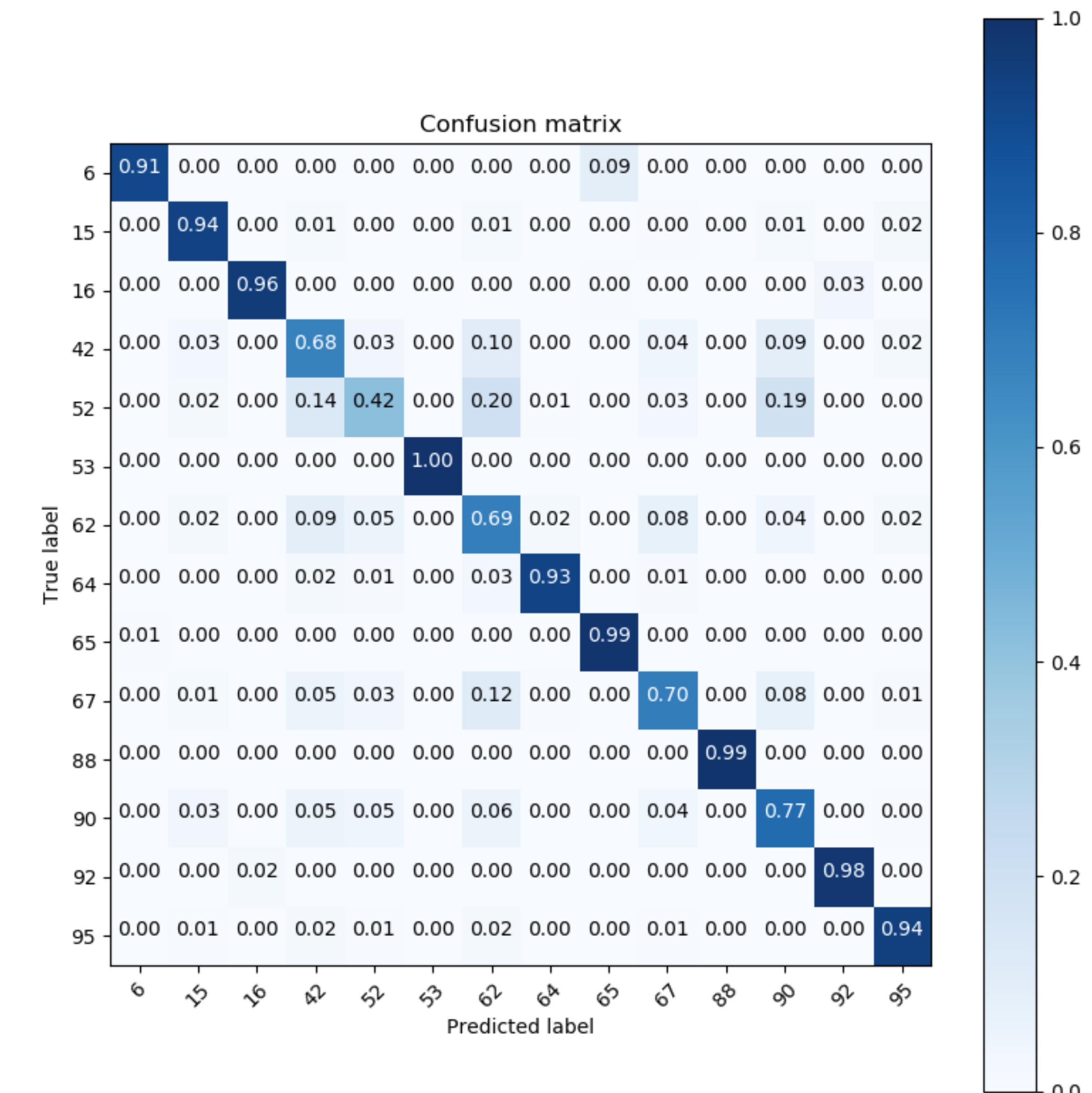
		Prediction outcome		actual value	total
		p	n		
p'	True Positive		False Negative	P'	
	False Positive		True Negative	N'	
total		P	N		

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{FN} + \text{FP} + \text{TN})$$

$$\text{Precision (p)} = \text{TP} / (\text{TP} + \text{FP})$$

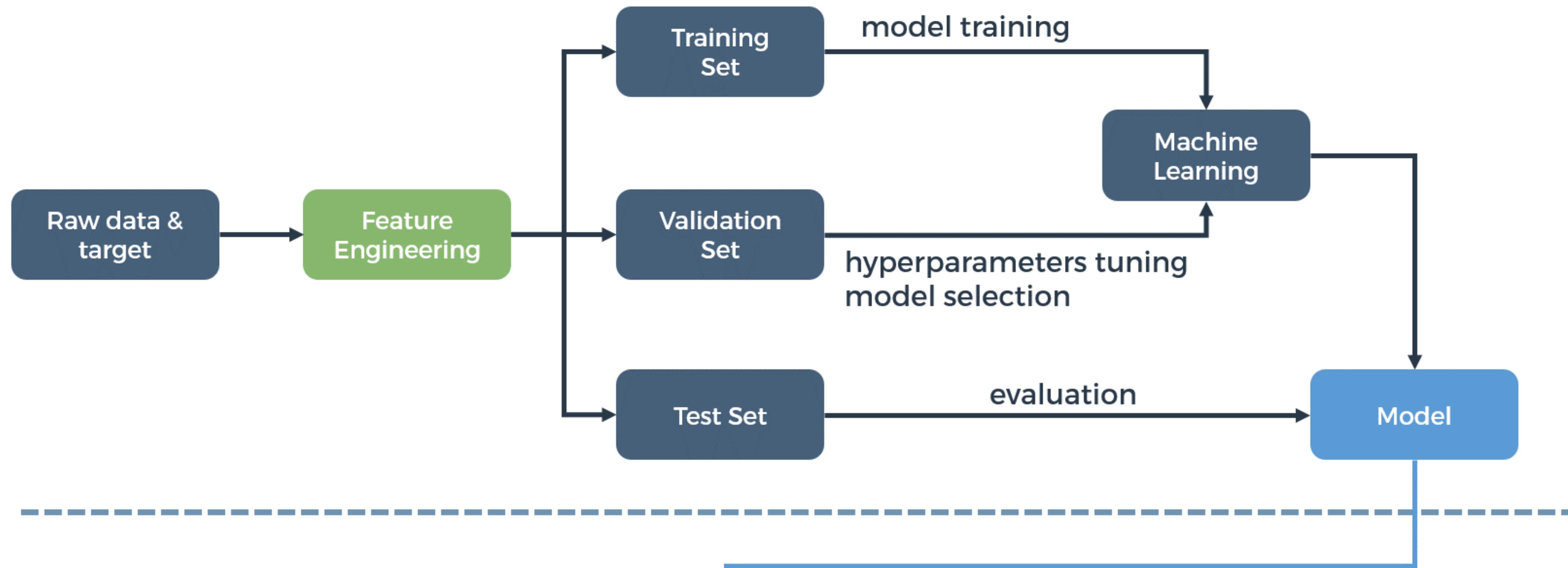
$$\text{Recall (r)} = \text{TP} / (\text{TP} + \text{FN})$$

$$F_1 = \frac{2}{r^{-1} + p^{-1}}$$



General Workflow of ML/DL

TRAINING



PREDICTING



Open Datasets

Datasets

Find and use datasets or complete tasks. [Learn more.](#)

Processed, balanced, well-behaved, labeled datasets to benchmark your networks!

Help the community by creating and solving Tasks on datasets!



Search 29,853 datasets

Feedback Filter

PUBLIC

Sort by: Hottest

 Hotel booking demand
Jesse Mostipak
19 days 1 MB 10.0 1 File (CSV) 1 Task

 Big Five Personality Test
Bojan Tunguz
14 days 159 MB 9.7 3 Files (CSV, other)

 StartUp Investments (Crunchbase)
Andy_M
14 days 3 MB 8.8 1 File (CSV)

Can we predict the possibility of a bo...
0 Submissions · In Hotel booking demand

Visualize US Accidents Dataset
12 Submissions · In US Accidents (3.0 million...)

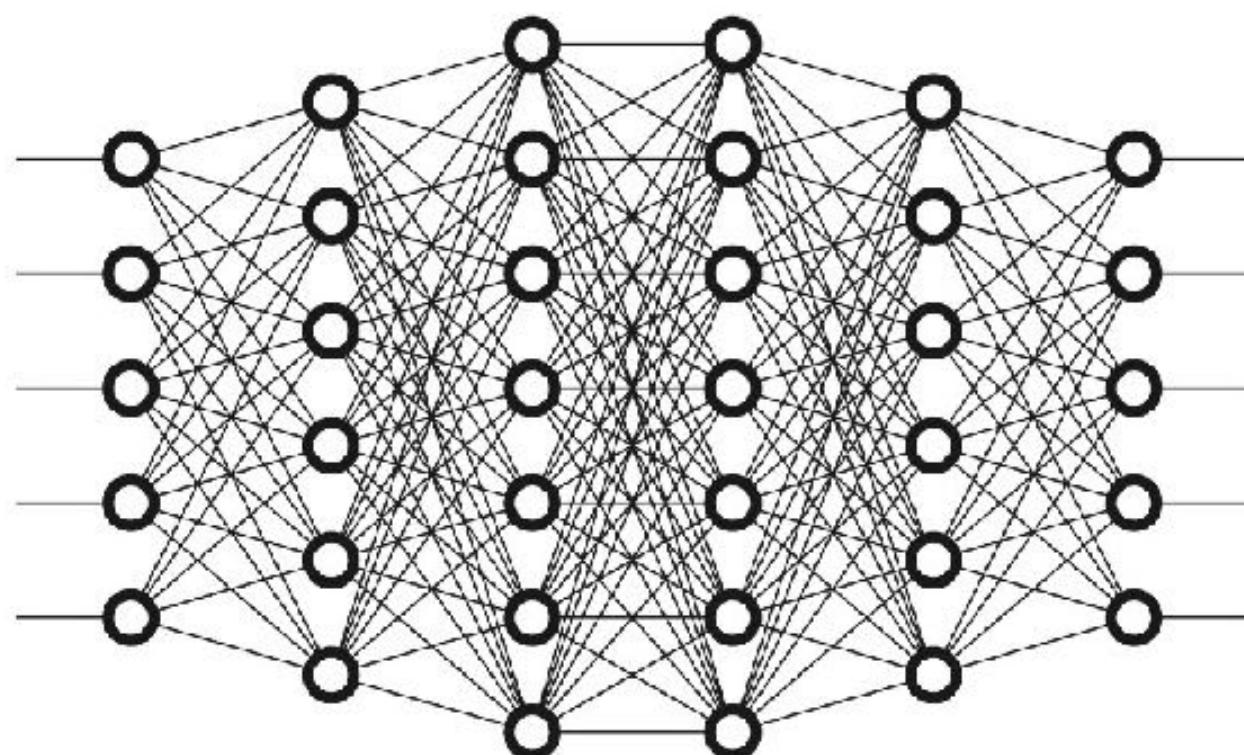
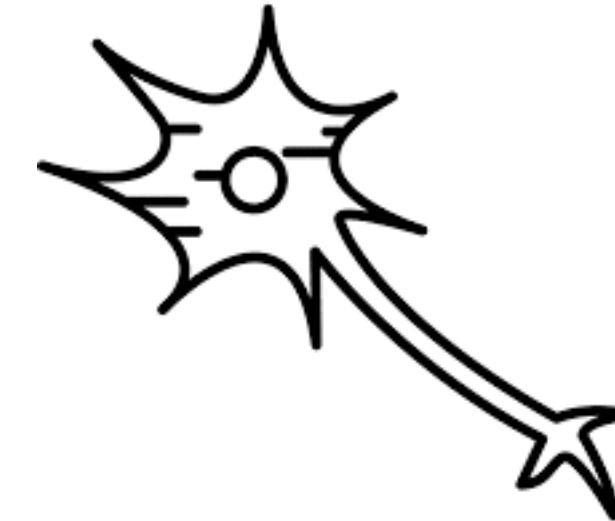
What to watch on Netflix ?
4 Submissions · In Netflix Movies and TV Sh...

<https://www.tensorflow.org/datasets>
<https://www.kaggle.com/datasets>
<http://topepo.github.io/caret/data-sets.html>
<https://github.com/awesomedata/awesome-public-datasets>

Take-home messages

In a neuron:

- ... the main job is to calculate a **weighted average**
- ... the **decision** is made through the **activation** function

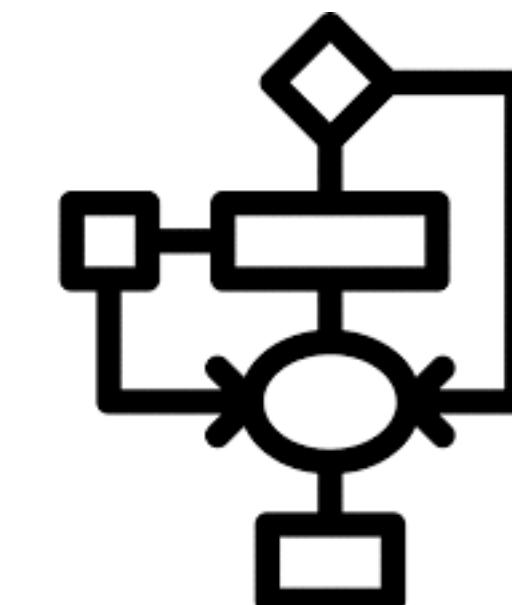


In a neural network:

- ... losses are calculated using the loss function
- ... losses are calculated by **comparing** the truths and the prediction
- ... **predictions** are made through **forward** propagation
- ... weights are **updated** through the **backward** propagation process
- ... **optimizers** are used to decide the weights updating **strategies**

In a deep learning workflow:

- ... the heavy lifting is mostly done by **DL frameworks**
- ... open datasets are crucial for benchmarking and bootstrapping DNNs





Go to <https://jupyter.lisa.surfsara.nl/jhsrf003>

Select “SURF jupyterhub - course hours” profile

Notebook: `notebook_1_mnist.ipynb`

Hacking until 11:15