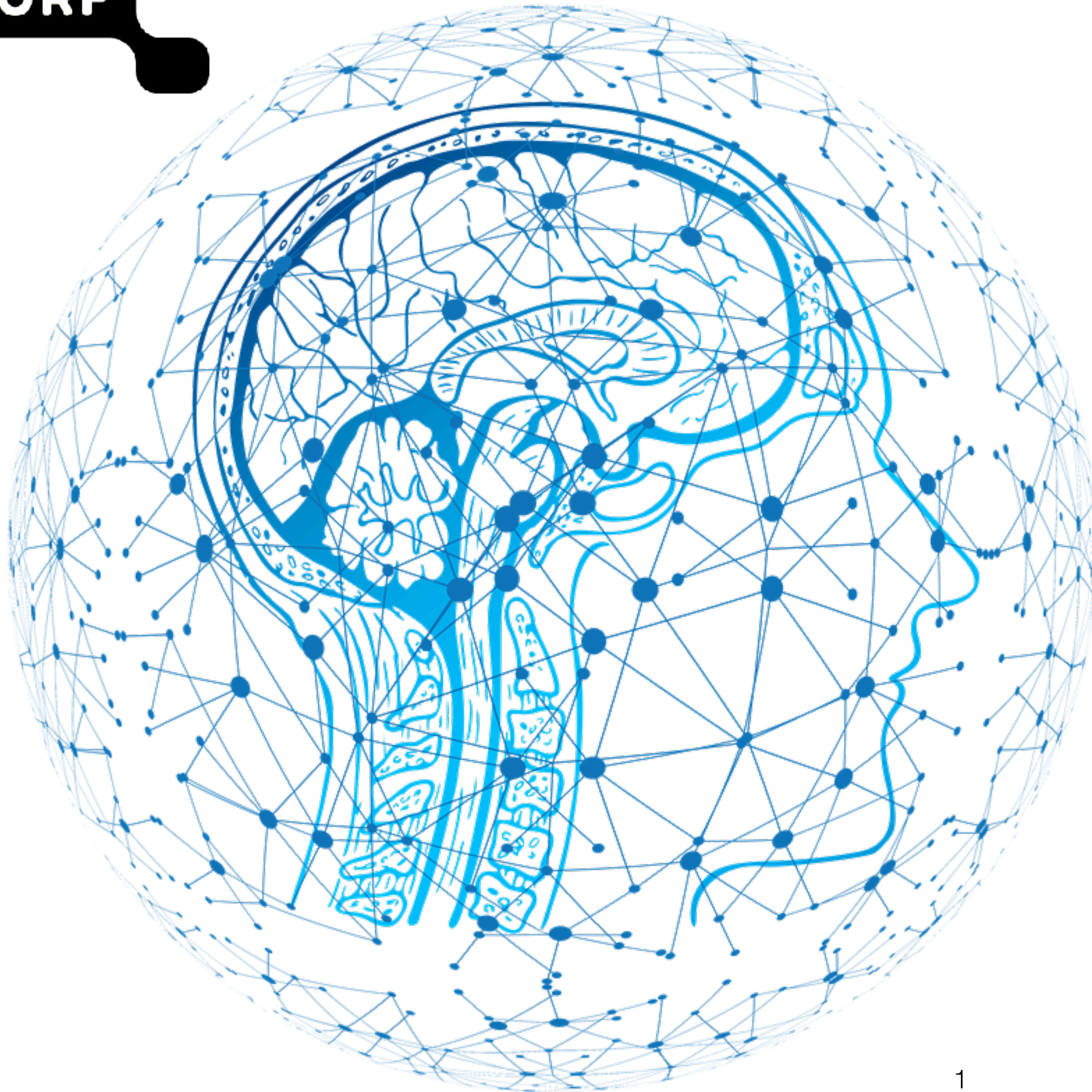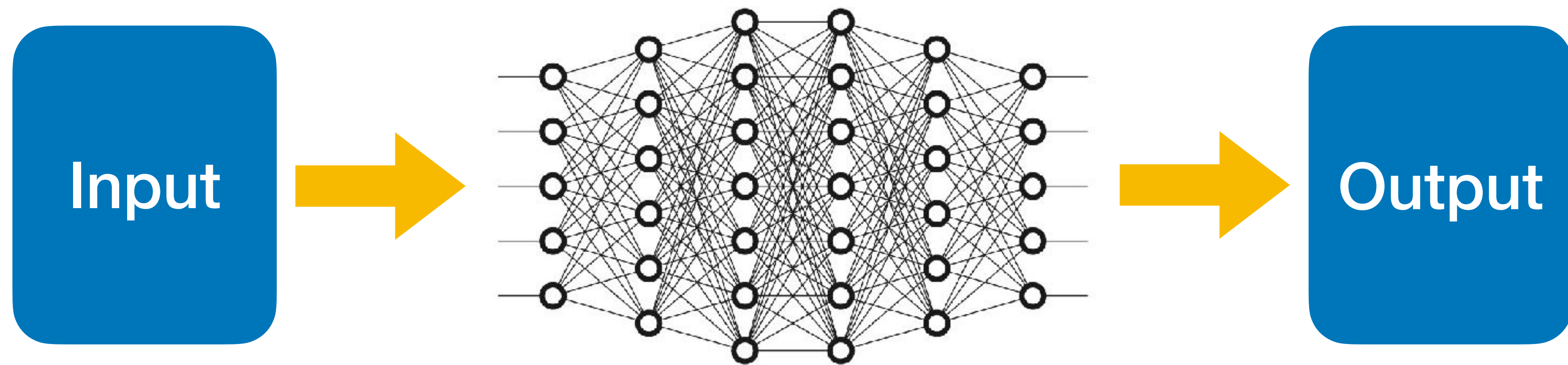# Introduction to Deep Learning

## Multi-layer Perceptron

*Maxwell Cai*
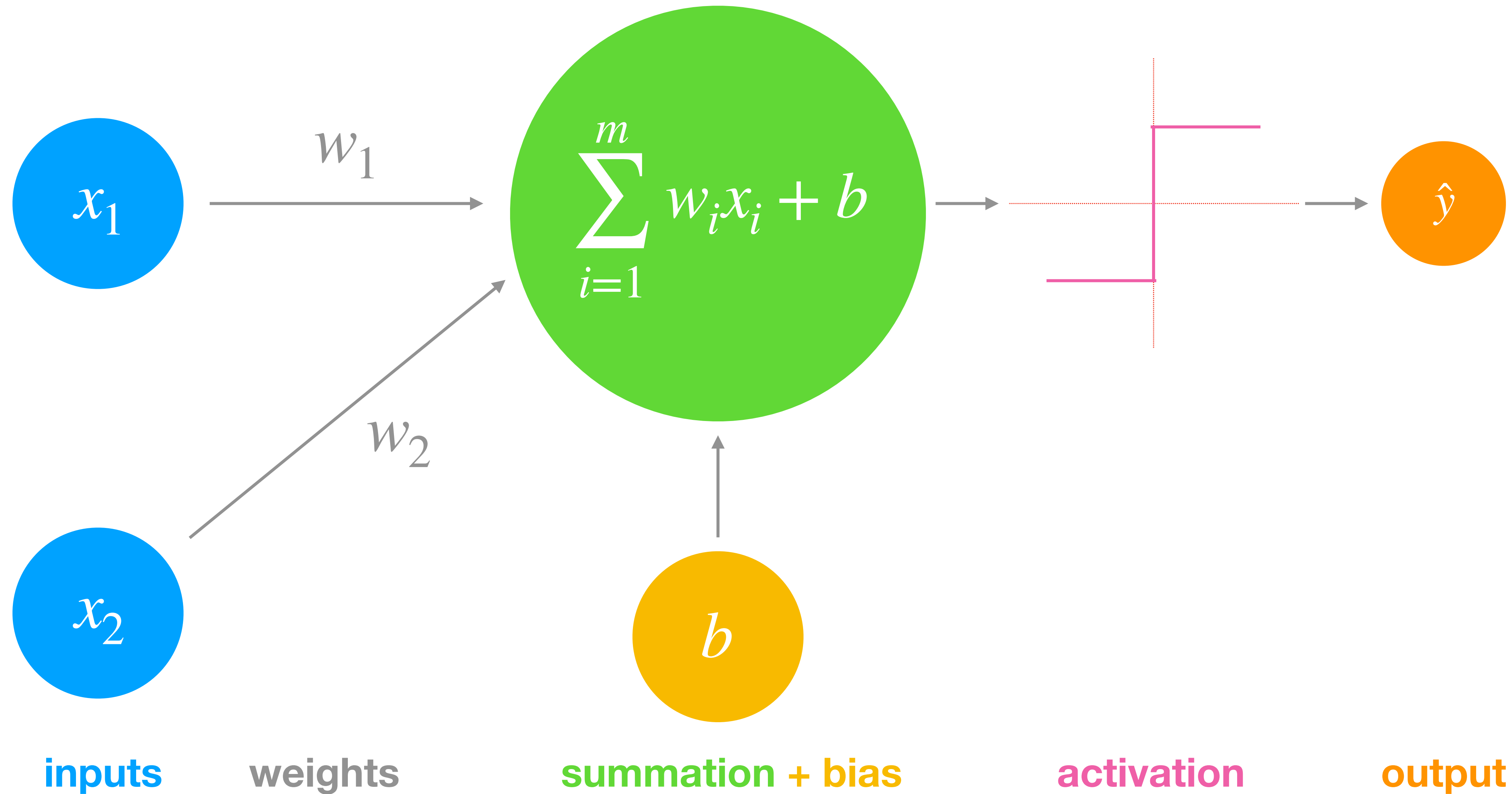
maxwellcai.com

# How does a neural network work?



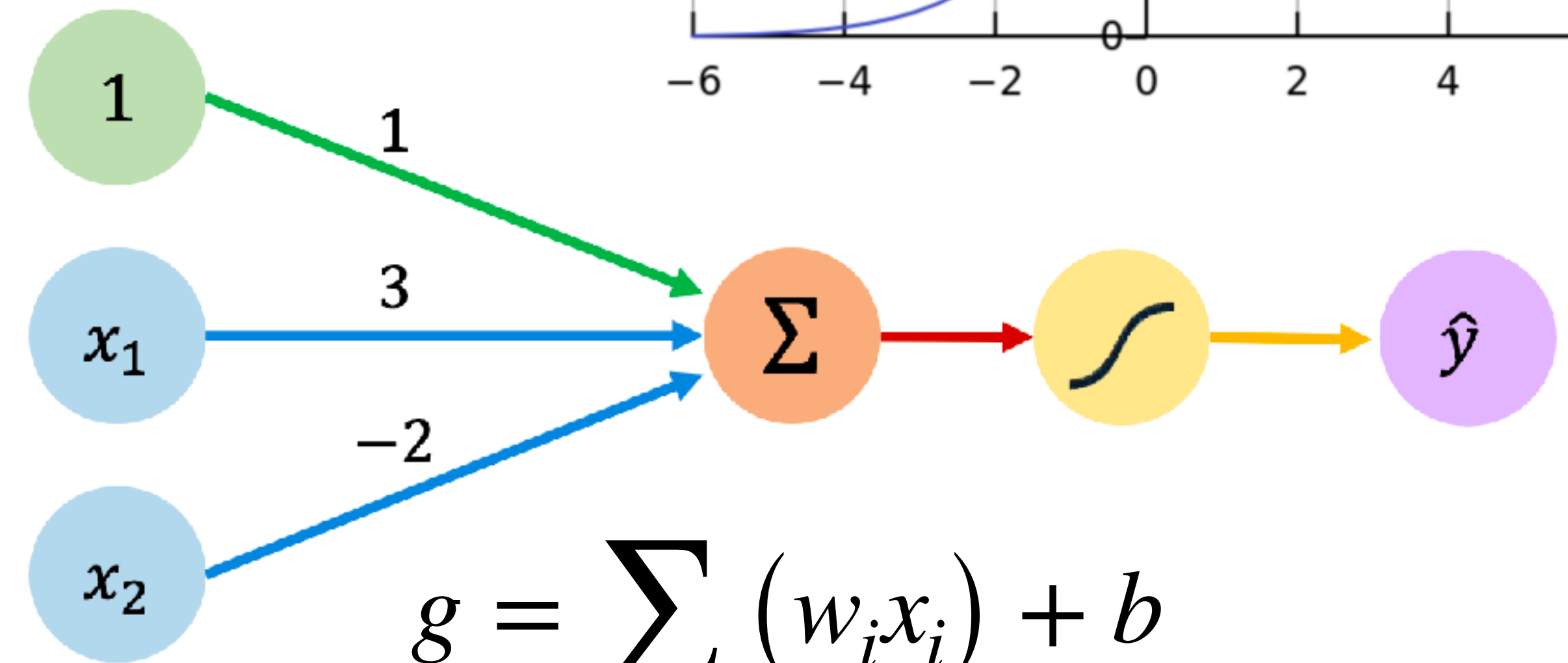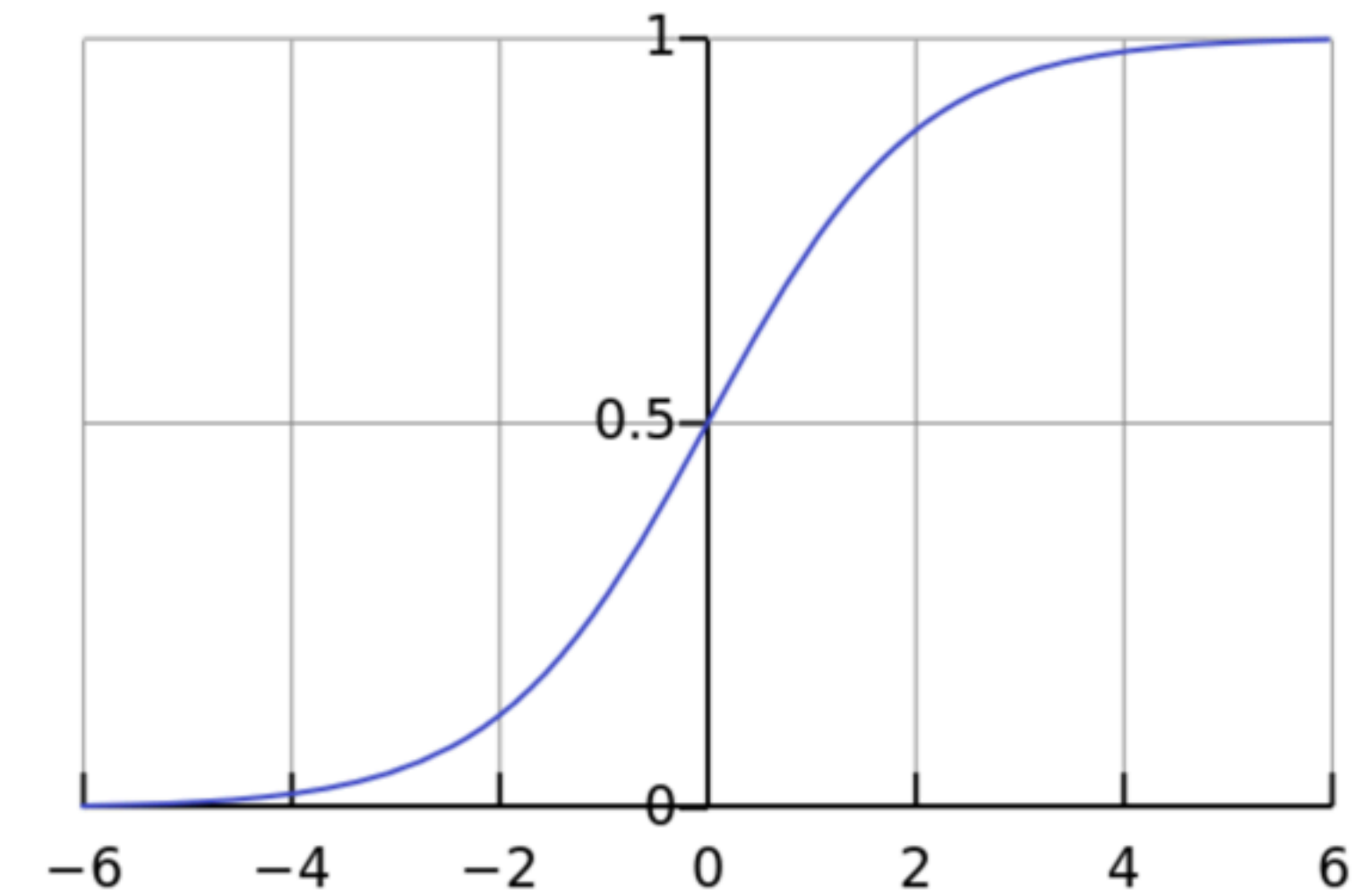$$\hat{y} = f_{\mathrm{NN}}(x_1, x_2, \ldots, x_n)$$

# Prediction



$$\sum_{i=1}^{m} w_i x_i + b$$

$x_1$   $w_1$   $x_2$   $w_2$   $b$   $\hat{y}$

**inputs**    **weights**    **summation + bias**    **activation**    **output**

3

$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

$$g = \sum \left( w_i x_i \right) + b$$

# Live demo: Multi-layer Perceptrons



*https://playground.tensorflow.org*

# Why multiple layers?

**Example**: house price prediction model (designed by humans)



# beds

Size

Zip code

Wealth

Family size

Walkability

School quality

Price

**Input layer**     **Hidden layer**     **Output layer**

6

# Why multiple layers?

**Example**: house price prediction model (designed by machines)



# beds

Size

Zip code

Wealth

?

?

?

Price

*During the **optimization** process, the machine learns to **encode** a presentation that **maps** the input to the corresponding output.*

**Input layer**     **Hidden layer**     **Output layer**

# Multi-Layer Perception



**MNIST dataset**
(**M**odified **N**ational **I**nstitute for **S**tandards and **T**echnology)

Deep neural networks learn hierarchical feature representations

*Source: 3blue1brown*

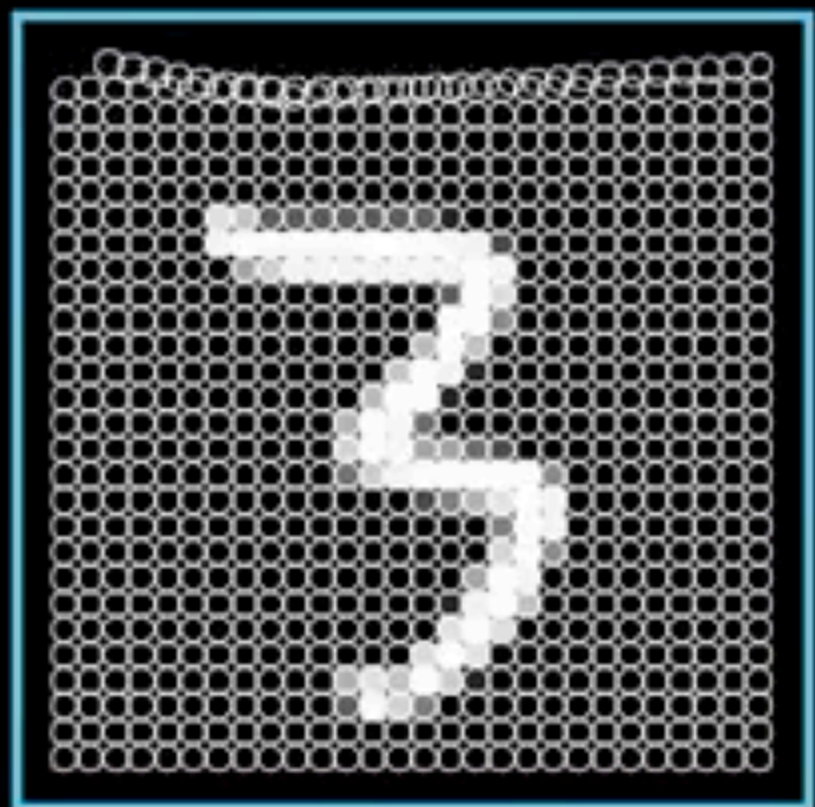What's the "cost" of this difference?

Utter trash

Training: backward propagation

Animation: 3blue1brown

**Training: backward propagation**

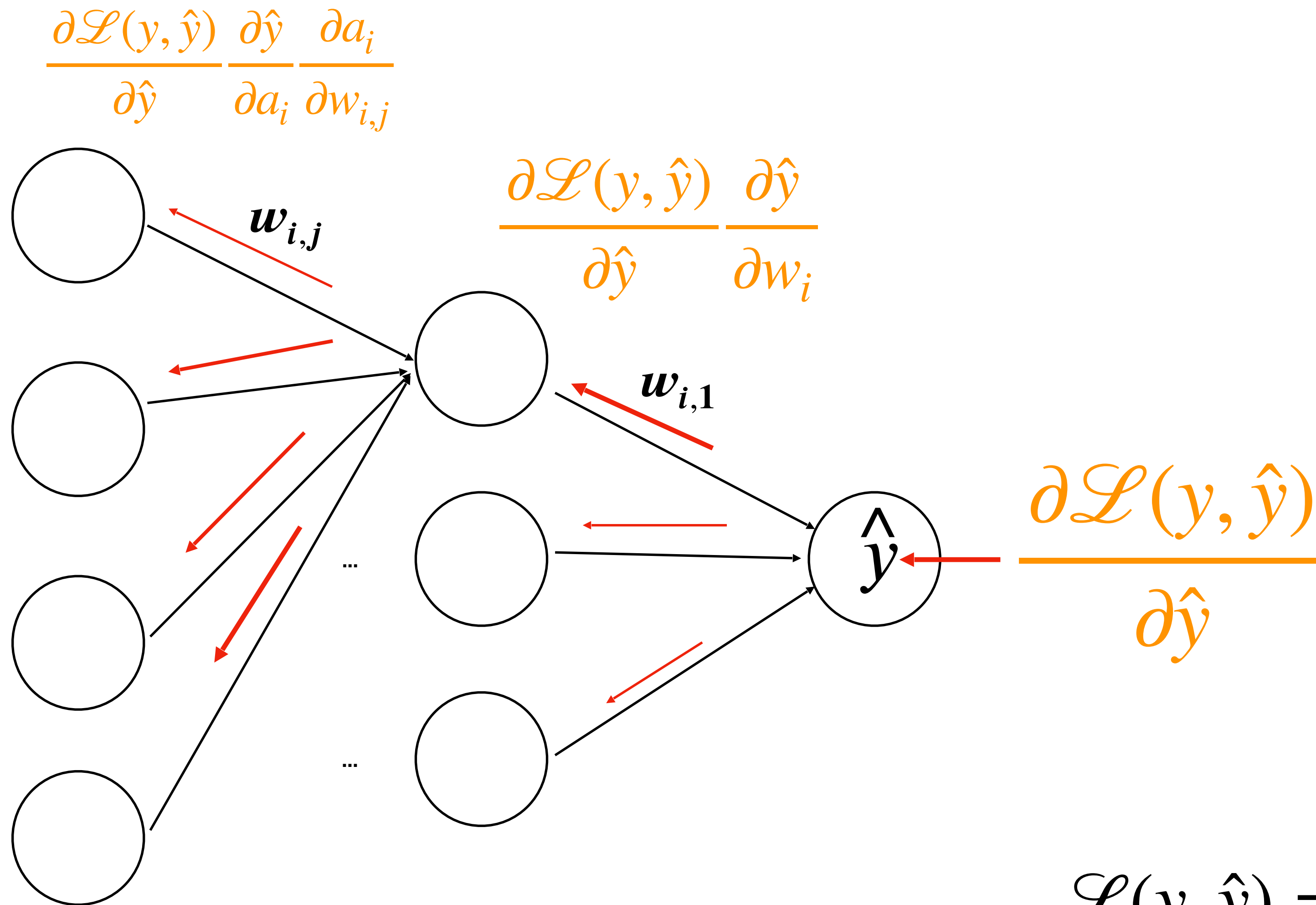Animation: _3blue1brown_

Cost of
one example

784

0
1
2
3
4
5
6
7
8
9

0
1
2
3
4
5
6
7
8
9

**Training: backward propagation**

$$\frac{\partial \mathscr{L}(y, \hat{y})}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_i} \frac{\partial a_i}{\partial w_{i,j}}$$

$$\frac{\partial \mathscr{L}(y, \hat{y})}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_i}$$

$$w_{i,j}$$

$$w_{i,1}$$

$$\hat{y}$$

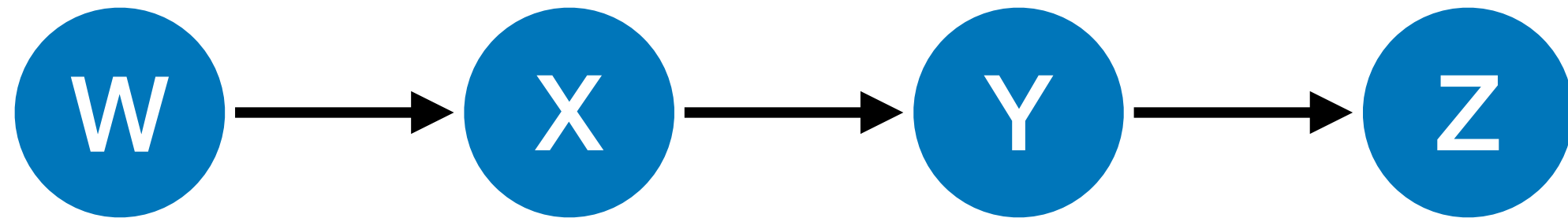$$\frac{\partial \mathscr{L}(y, \hat{y})}{\partial \hat{y}}$$

- 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**.

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

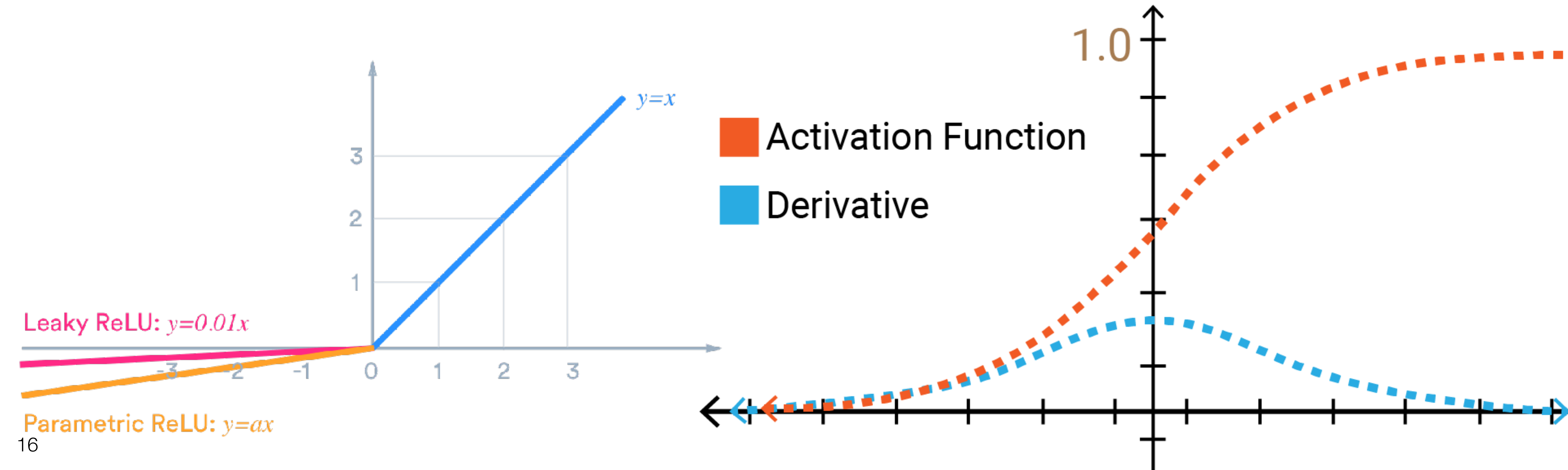# Neural Network: the deeper, the better?
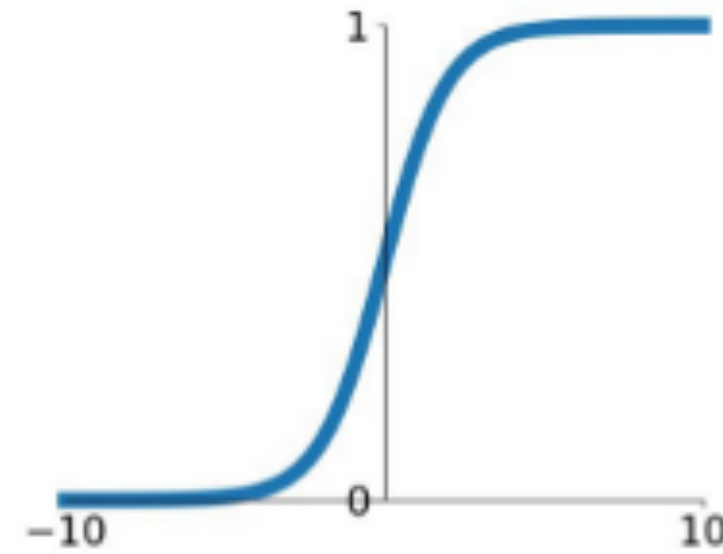
*Not really.*

# The vanishing gradient problem

W → X → Y → Z

**Chain rule**

$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y}\frac{\partial y}{\partial x}\frac{\partial x}{\partial w}$$
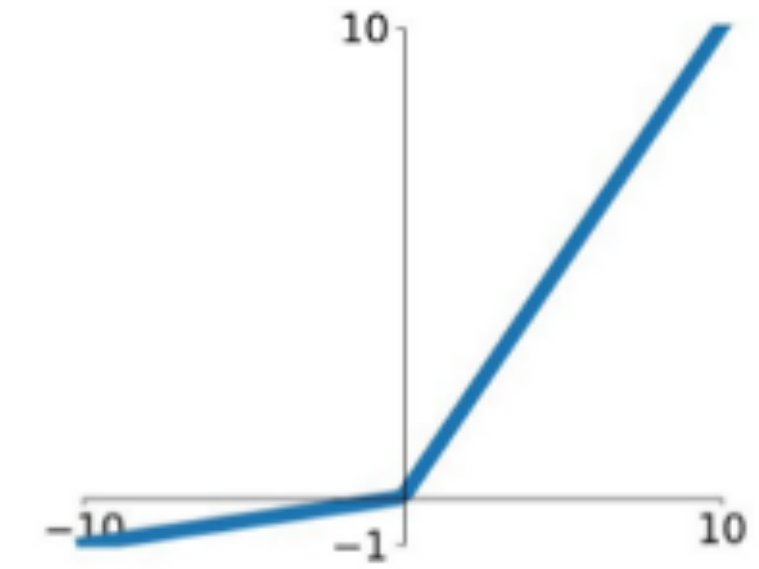
Leaky ReLU: *y=0.01x*

Parametric ReLU: *y=ax*

$y=x$

3
2
1
-3  -2  -1  0  1  2  3

1.0

■ Activation Function

■ Derivative

**Sigmoid**

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

**tanh**

$$\tanh(x)$$

**ReLU**

$$\max(0, x)$$

**Leaky ReLU**

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

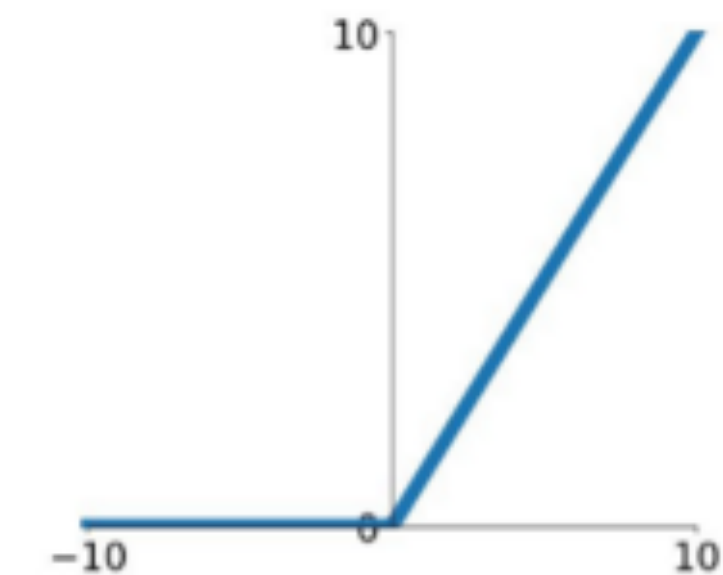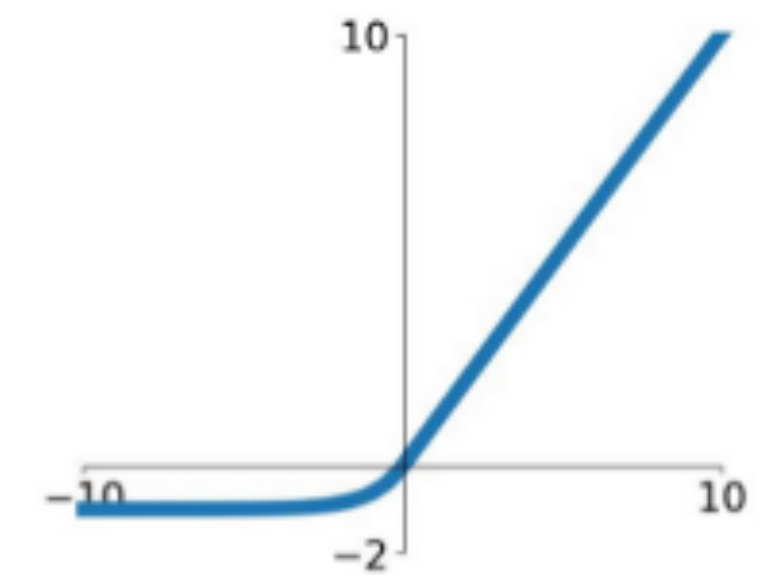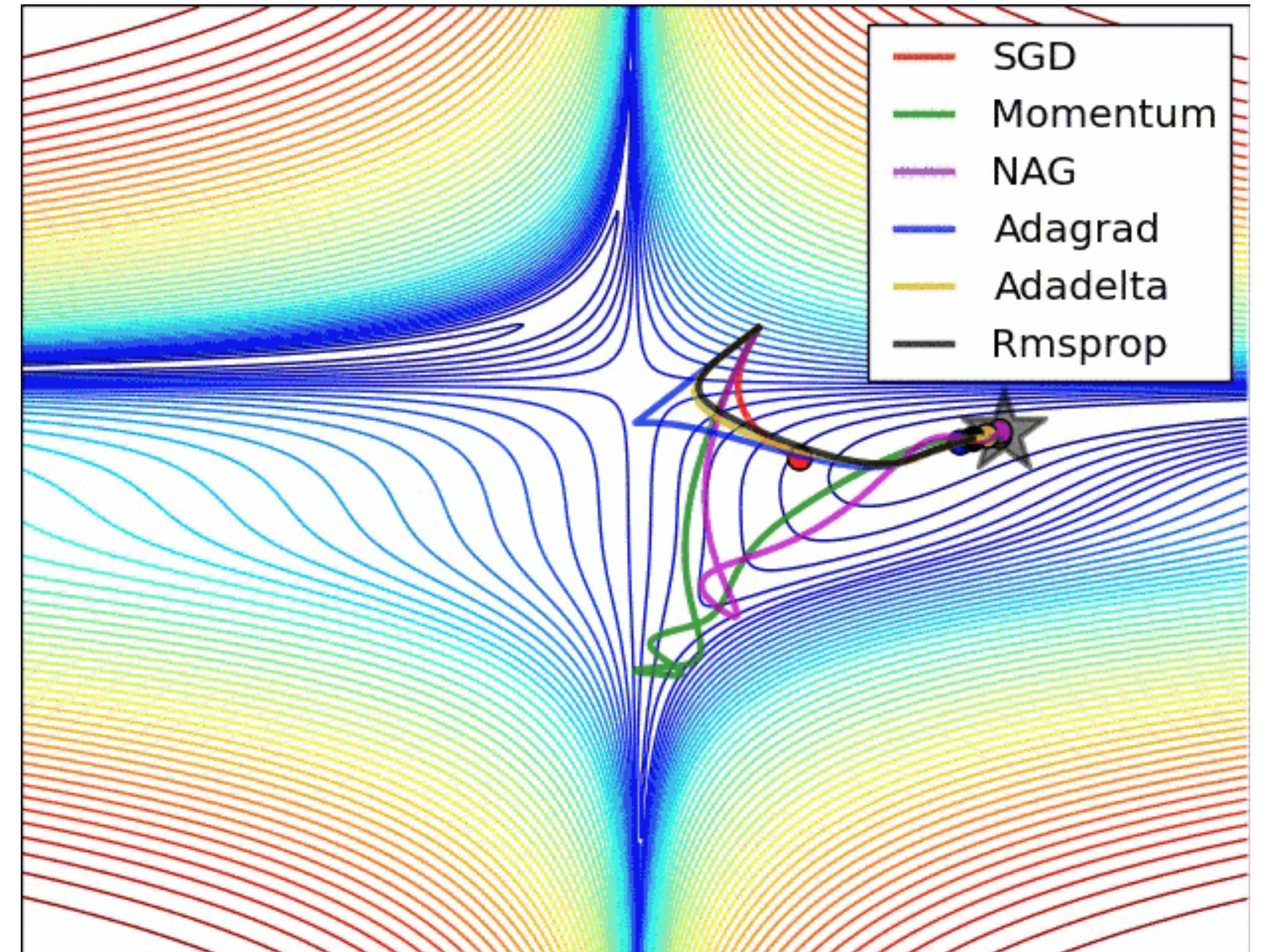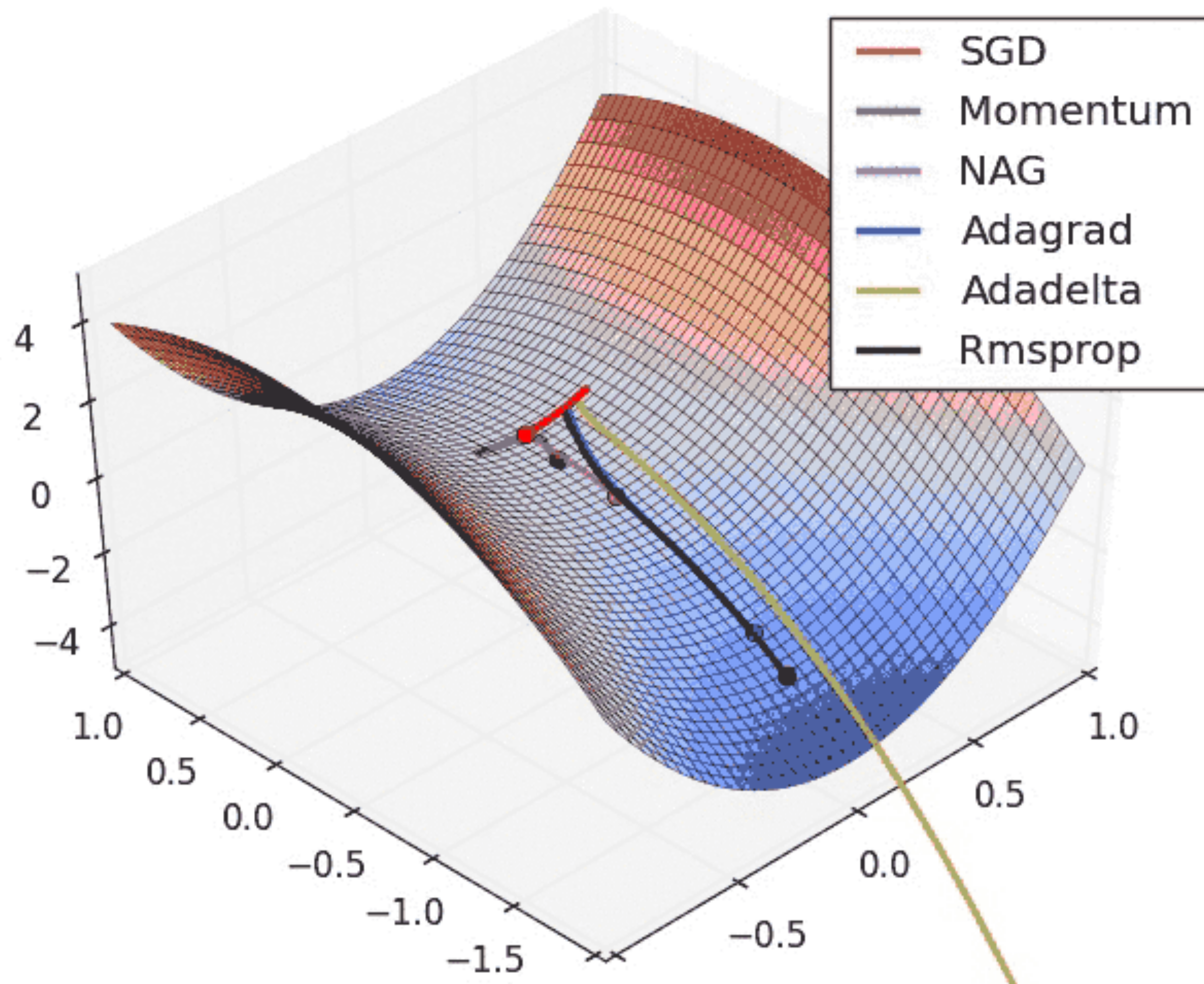**Maxout**

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

**ELU**

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

**But one can design their own activation functions!**

*Source: medium.com/analytics-vidhya*

# Common practice for loss functions
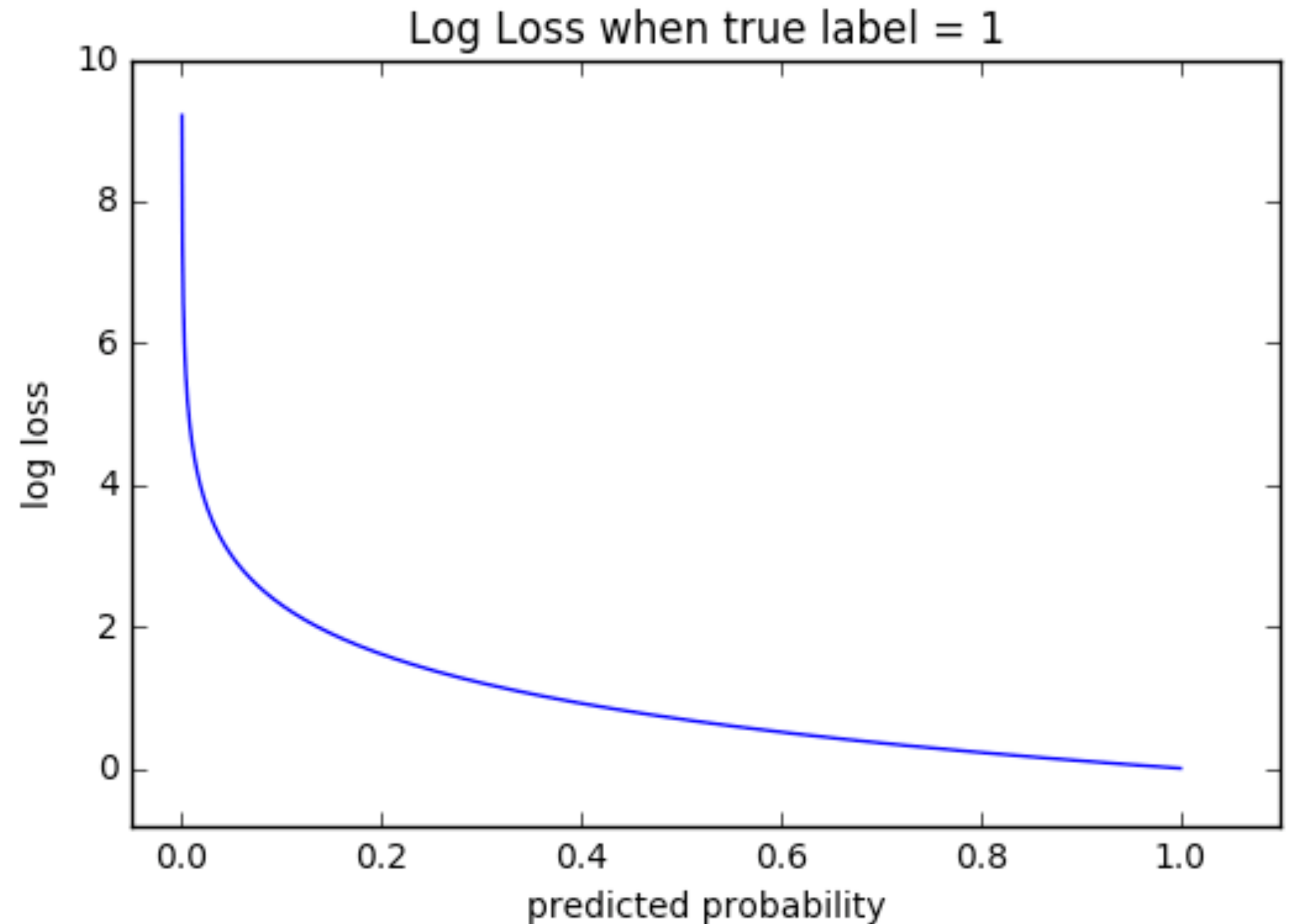
**Regression**
- Mean squared error
- Mean squared logarithmic error
- Mean absolute error

**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



**Cross-entropy loss outputs a log probability**

# 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.
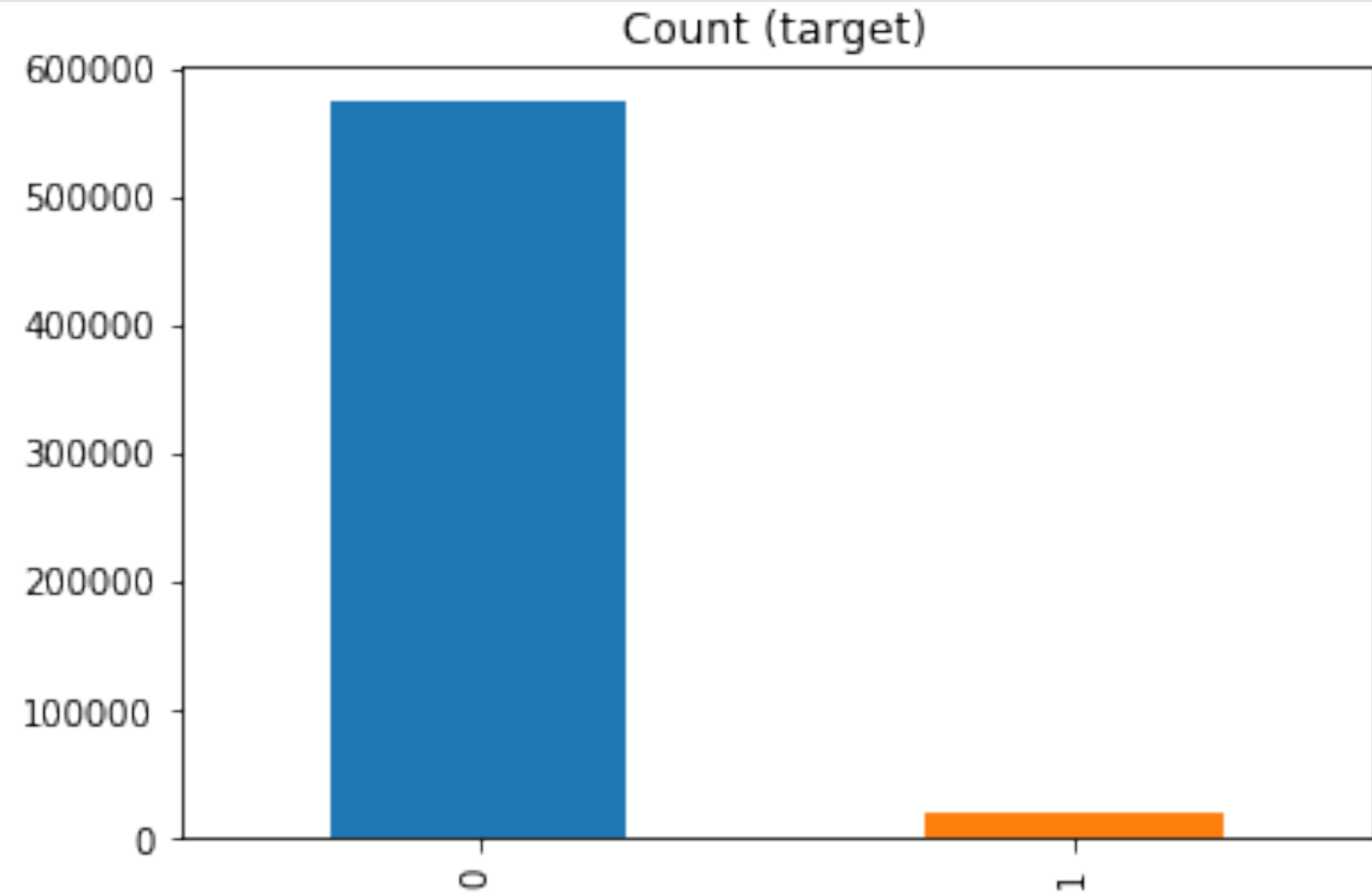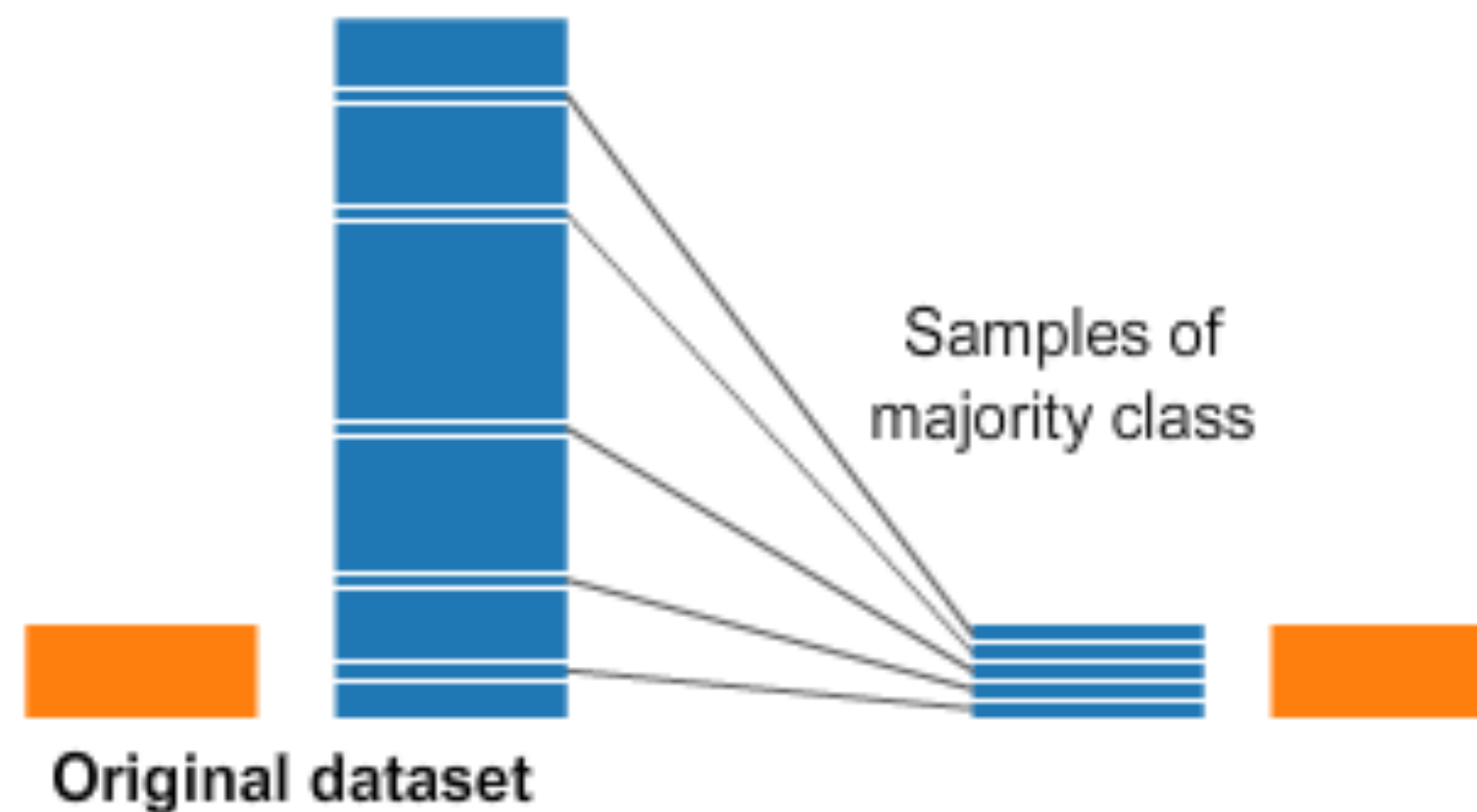
# Model Evaluation

# Balanced/Imbalanced training set



Count (target)

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

**Undersampling**

Samples of majority class

Original dataset

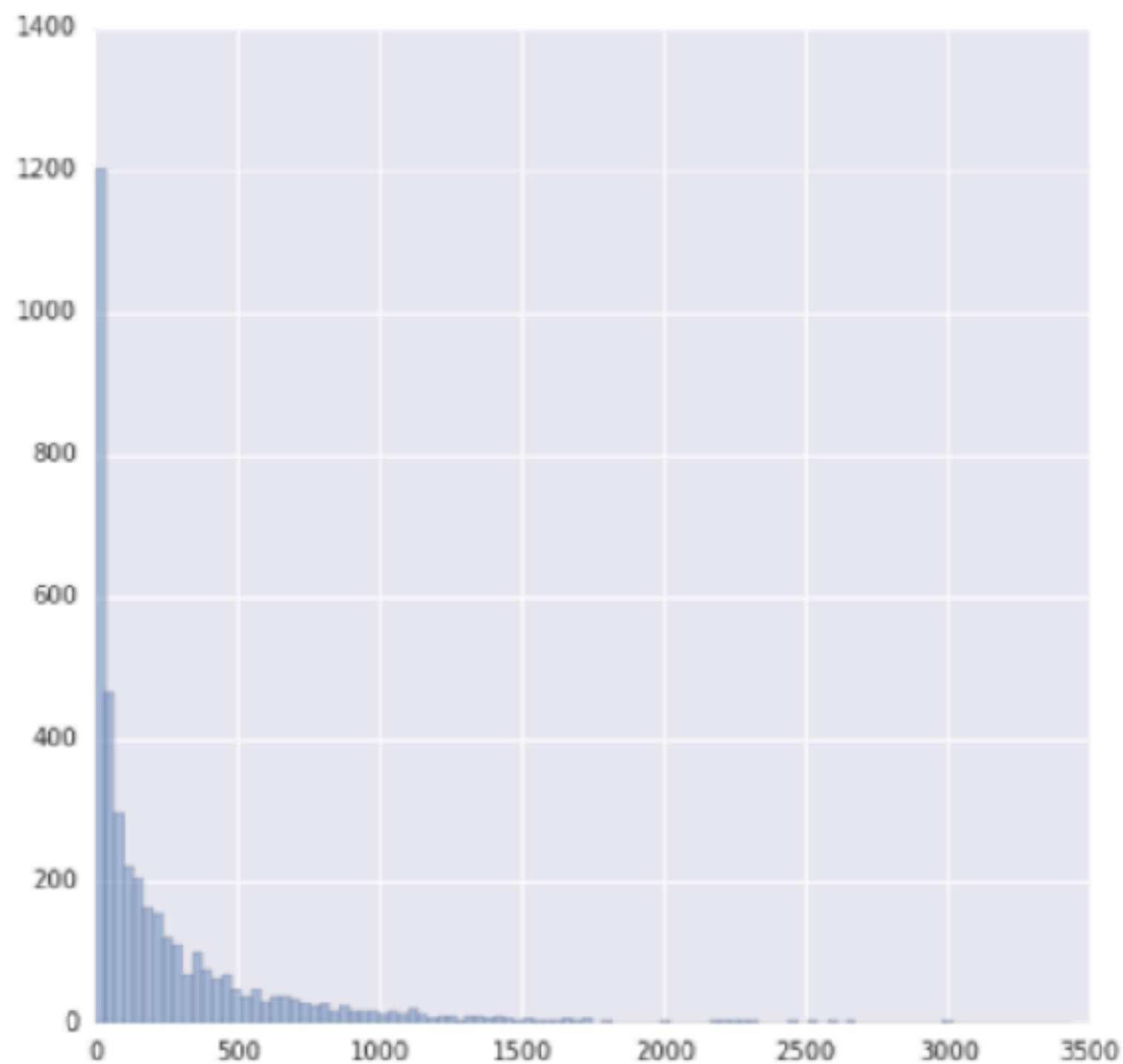**Oversampling**

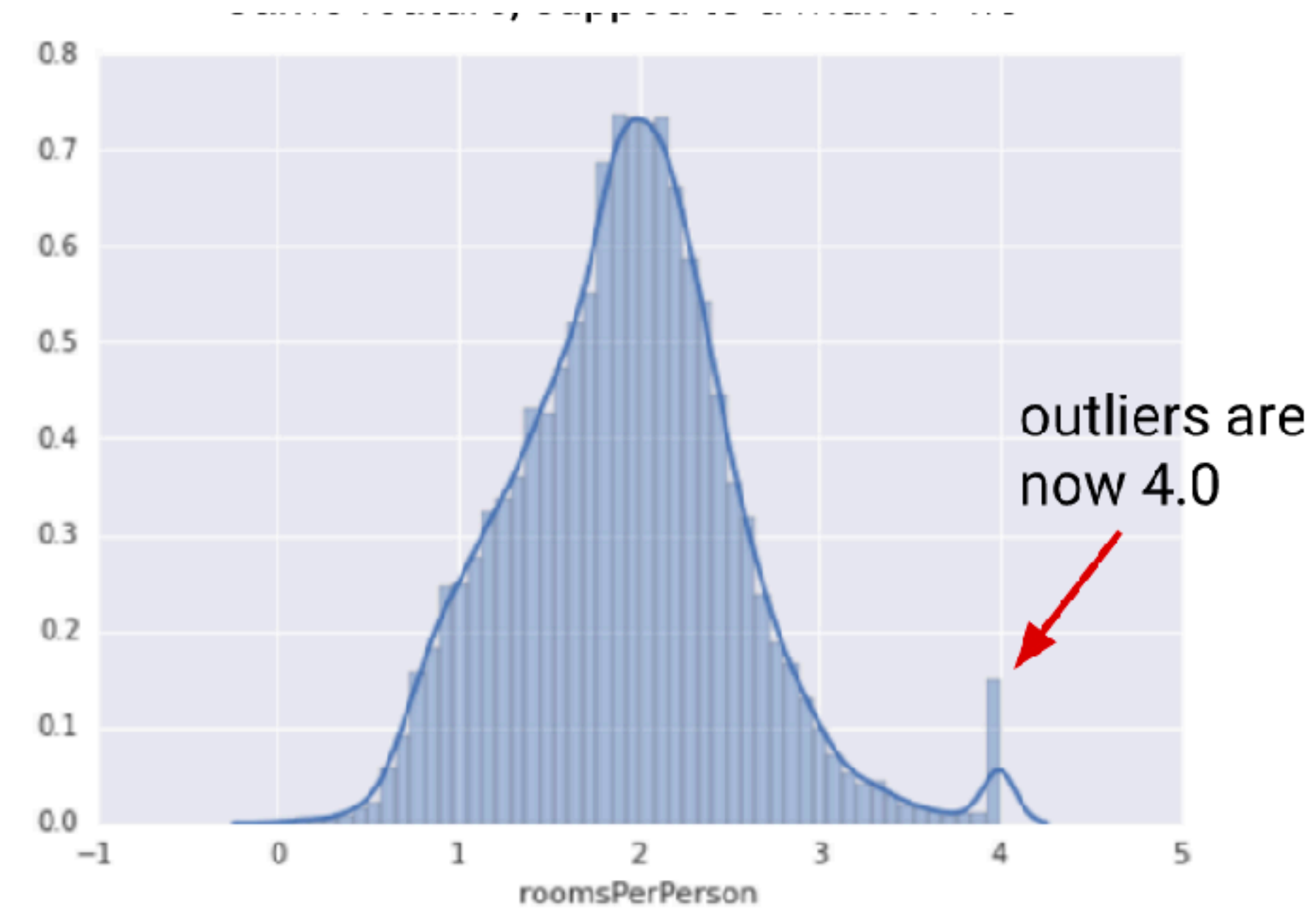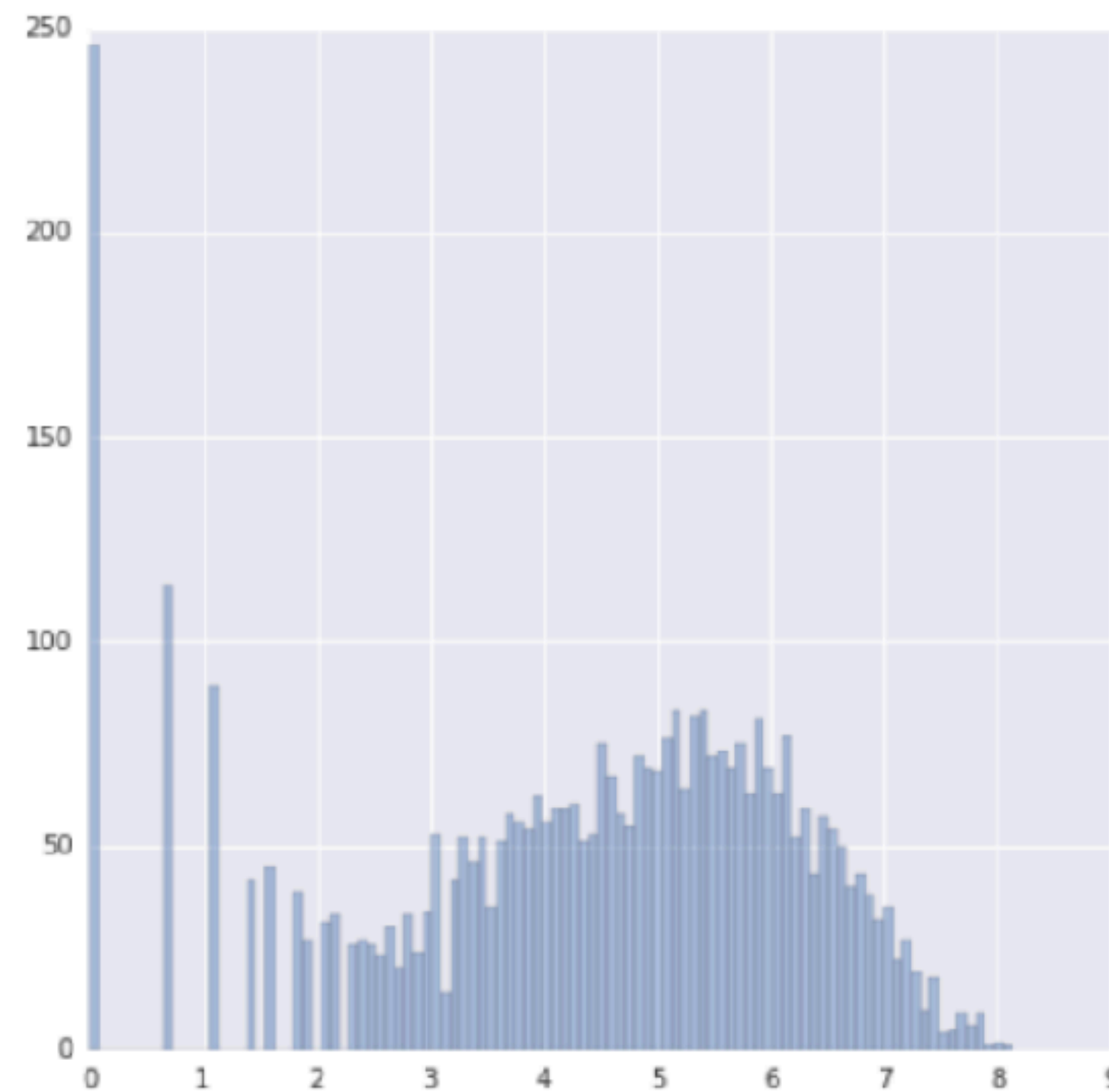Copies of the minority class

Original dataset

# Data Normalization

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



Ratings per movie



Log ratings per movie

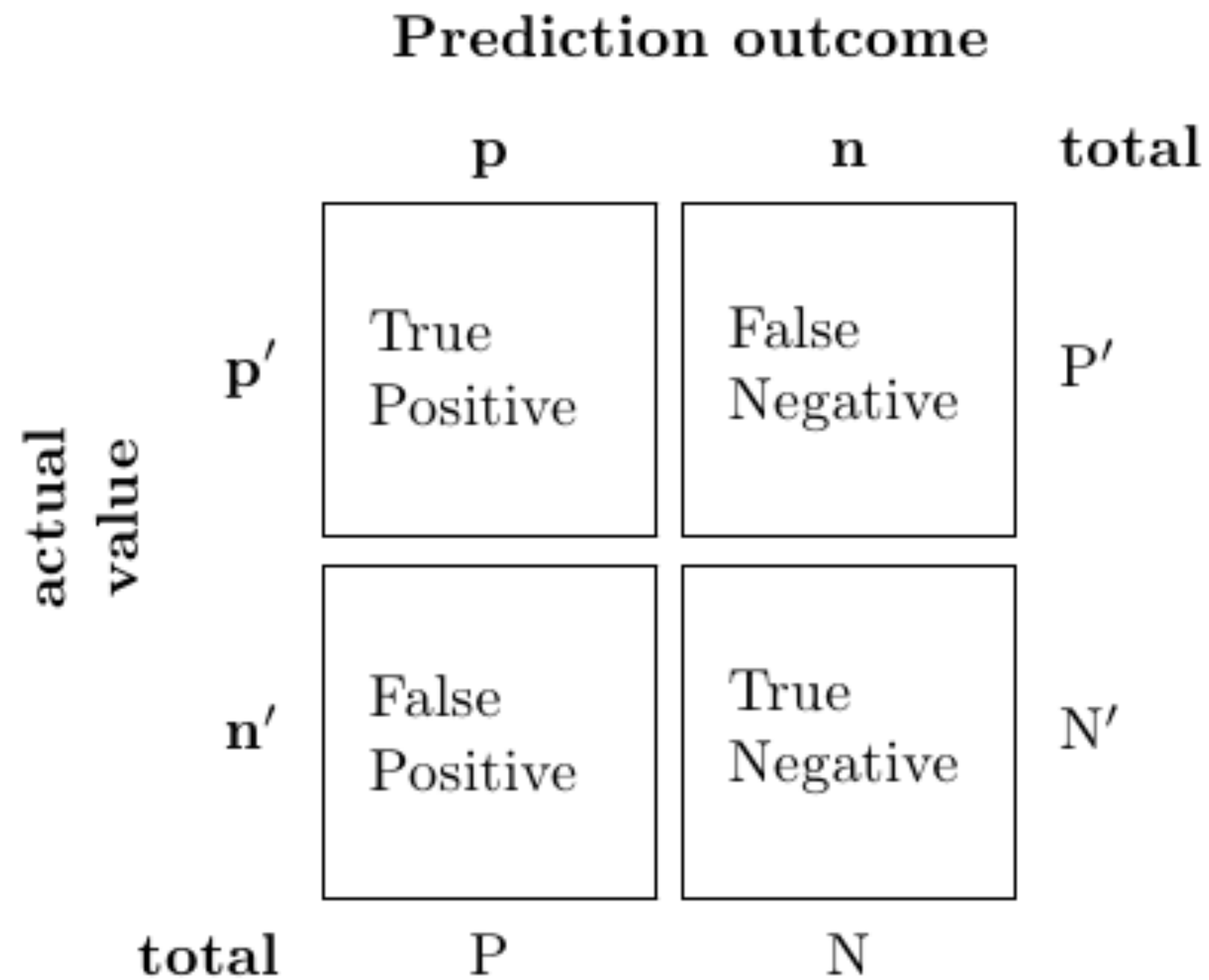outliers are now 4.0

**Further reading: <u>sklearn scalers</u>**

*Source: developers.google.com*

# Dataset splitting

## Training set

**Model training**



**Best model parameters**
Weights

**70%**

## Validation set

**Model selection**



**Best hyperparameters**
Learning rate
#neurons
#layers

**20%**

## Test set

**Model testing**



**Final model**
Accuracy
Sensitivity/specificity
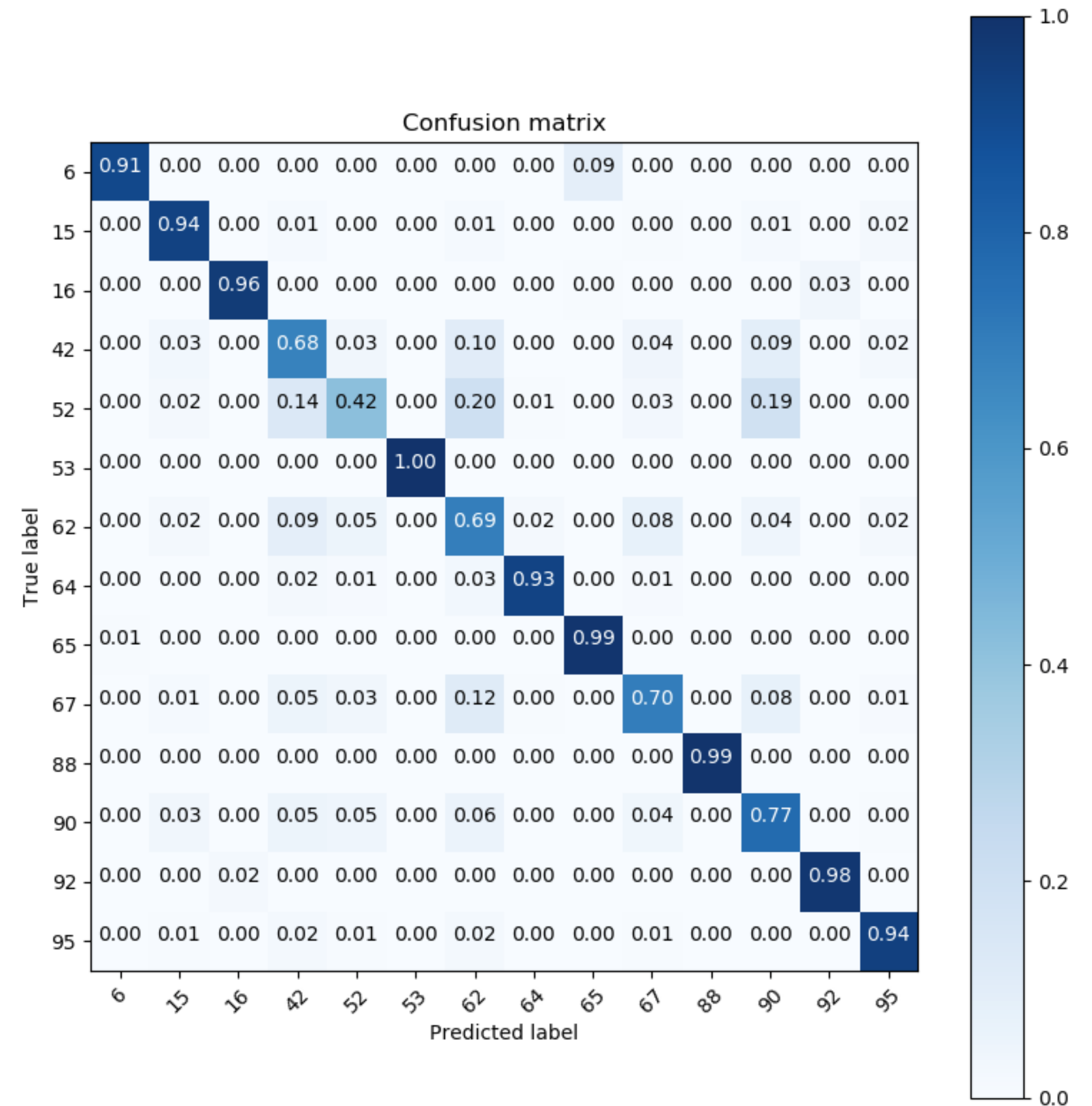
**10%**

# Confusion Matrix



Accuracy = (TP + TN) / (TP + FN + FP + TN)
Precision (p) = TP / (TP + FP)
Recall (r) = TP / (TP + FN)

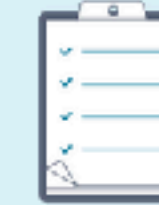$$F_1 = \frac{2}{r^{-1} + p^{-1}}$$

# General Workflow of ML/DL

Source: hackernoon.com

# Open Datasets

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



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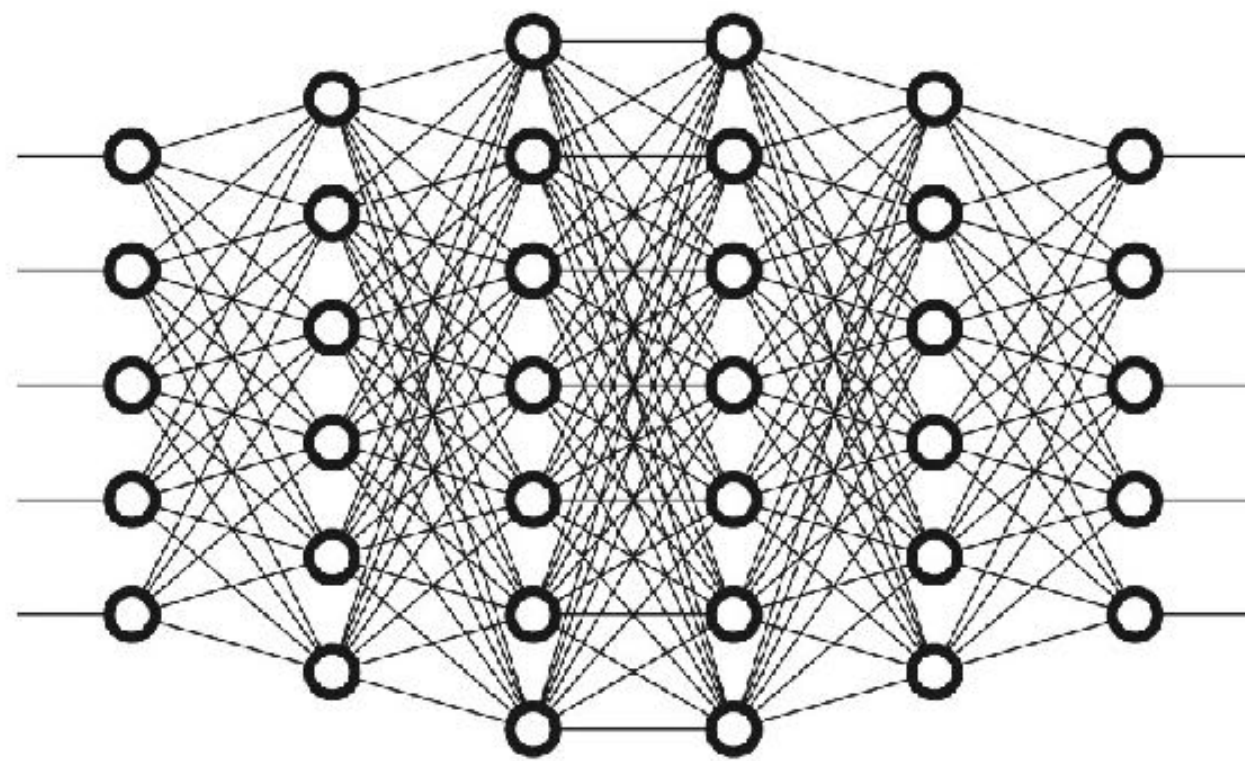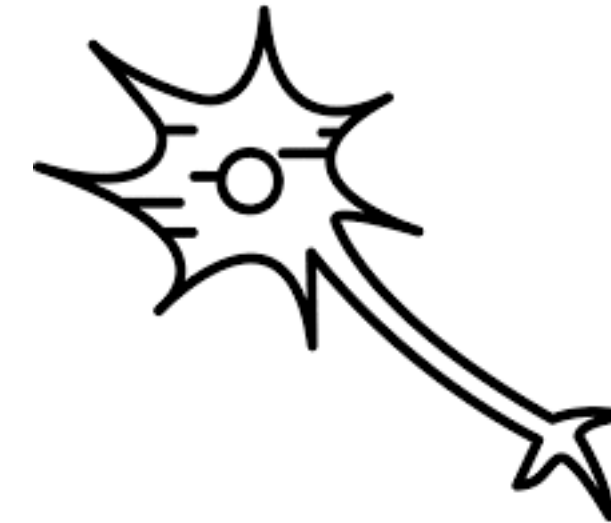https://github.com/awesomedata/awesome-public-datasets

# 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