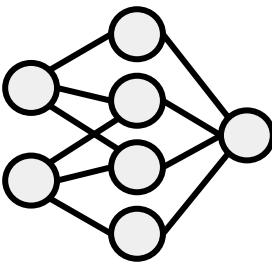


# Introduction to machine learning



# Look familiar?

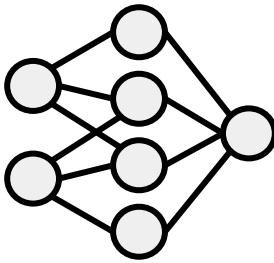


Figure from Sergey Ovchinnikov,  
Boston Protein Design and Modeling Club

# Let's start by deciphering this...

$\delta^{(2)}$

**“proportional error”**

**Layer 1**

$$\frac{\partial J(\theta)}{\partial \theta_{11}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right)$$

$$\frac{\partial J(\theta)}{\partial \theta_{12}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{12}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{12}^{(1)}} \right)$$

$$\frac{\partial J(\theta)}{\partial \theta_{21}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_2^{(2)}} \right) \left( \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \right) \left( \frac{\partial z_2^{(2)}}{\partial \theta_{21}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_2^{(2)}} \right) \left( \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \right) \left( \frac{\partial z_2^{(2)}}{\partial \theta_{21}^{(1)}} \right)$$

$$\frac{\partial J(\theta)}{\partial \theta_{22}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_2^{(2)}} \right) \left( \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \right) \left( \frac{\partial z_2^{(2)}}{\partial \theta_{22}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_2^{(2)}} \right) \left( \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \right) \left( \frac{\partial z_2^{(2)}}{\partial \theta_{22}^{(1)}} \right)$$

derivative chain for  $\delta_1^{(3)}$

derivative chain for  $\delta_2^{(3)}$

# Let's start by deciphering this...

Just kidding

Layer 1

$$\frac{\partial J(\theta)}{\partial \theta_{11}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right)$$
$$\frac{\partial J(\theta)}{\partial \theta_{12}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{12}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{12}^{(1)}} \right)$$
$$\frac{\partial J(\theta)}{\partial \theta_{21}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_2^{(2)}} \right) \left( \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \right) \left( \frac{\partial z_2^{(2)}}{\partial \theta_{21}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_2^{(2)}} \right) \left( \frac{\partial a_2^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{21}^{(1)}} \right)$$
$$\frac{\partial J(\theta)}{\partial \theta_{22}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_2^{(2)}} \right) \left( \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \right) \left( \frac{\partial z_2^{(2)}}{\partial \theta_{22}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_2^{(2)}} \right) \left( \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \right) \left( \frac{\partial z_2^{(2)}}{\partial \theta_{22}^{(1)}} \right)$$

derivative chain for  $\delta_1^{(3)}$

derivative chain for  $\delta_2^{(3)}$

$\delta^{(2)}$

“proportional error”

$f'(a^{(2)})$

input

$\delta^{(2)}$

“proportional error”

$f'(a^{(2)})$

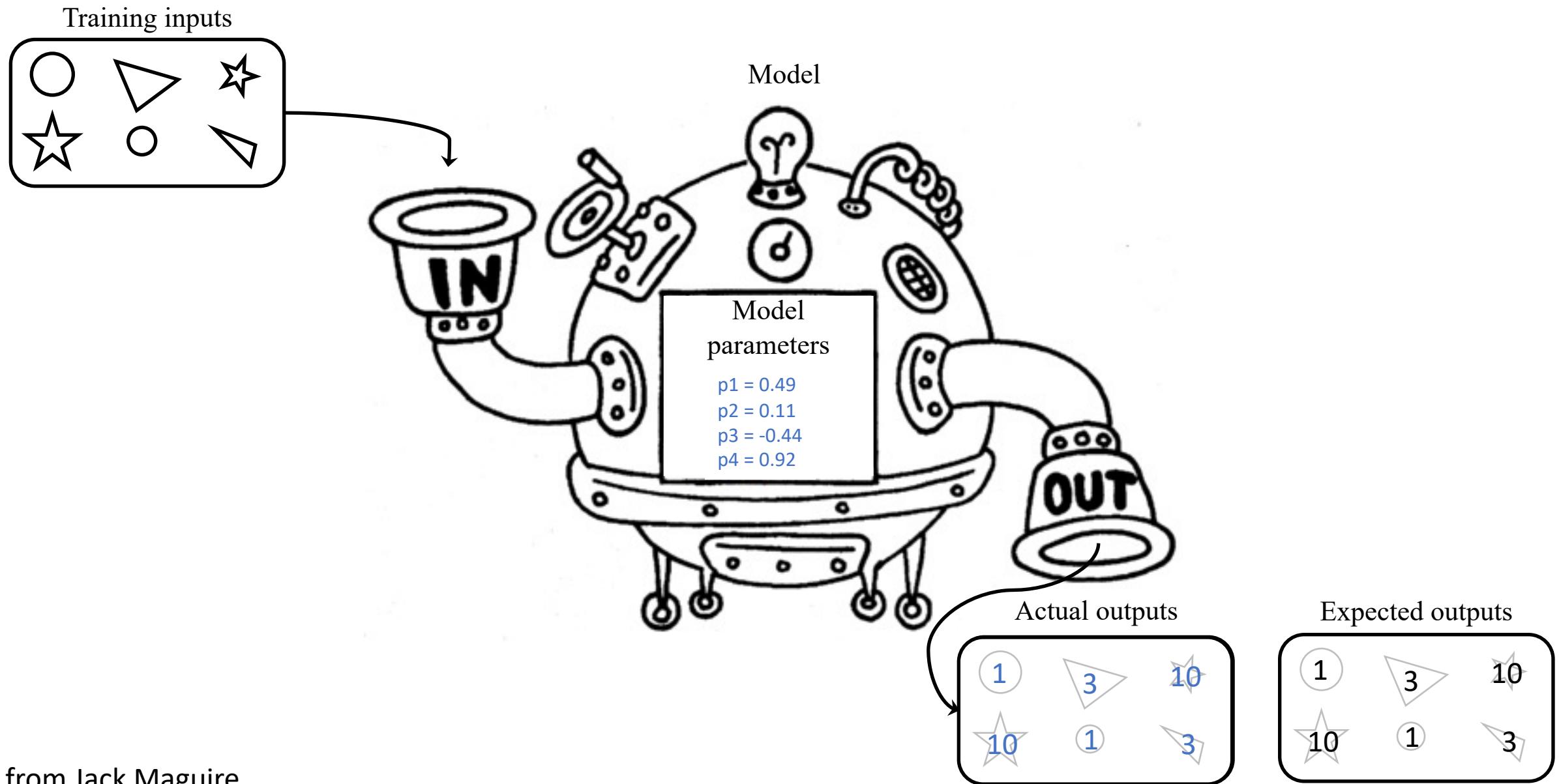
input

# First, we will learn machine learning basics

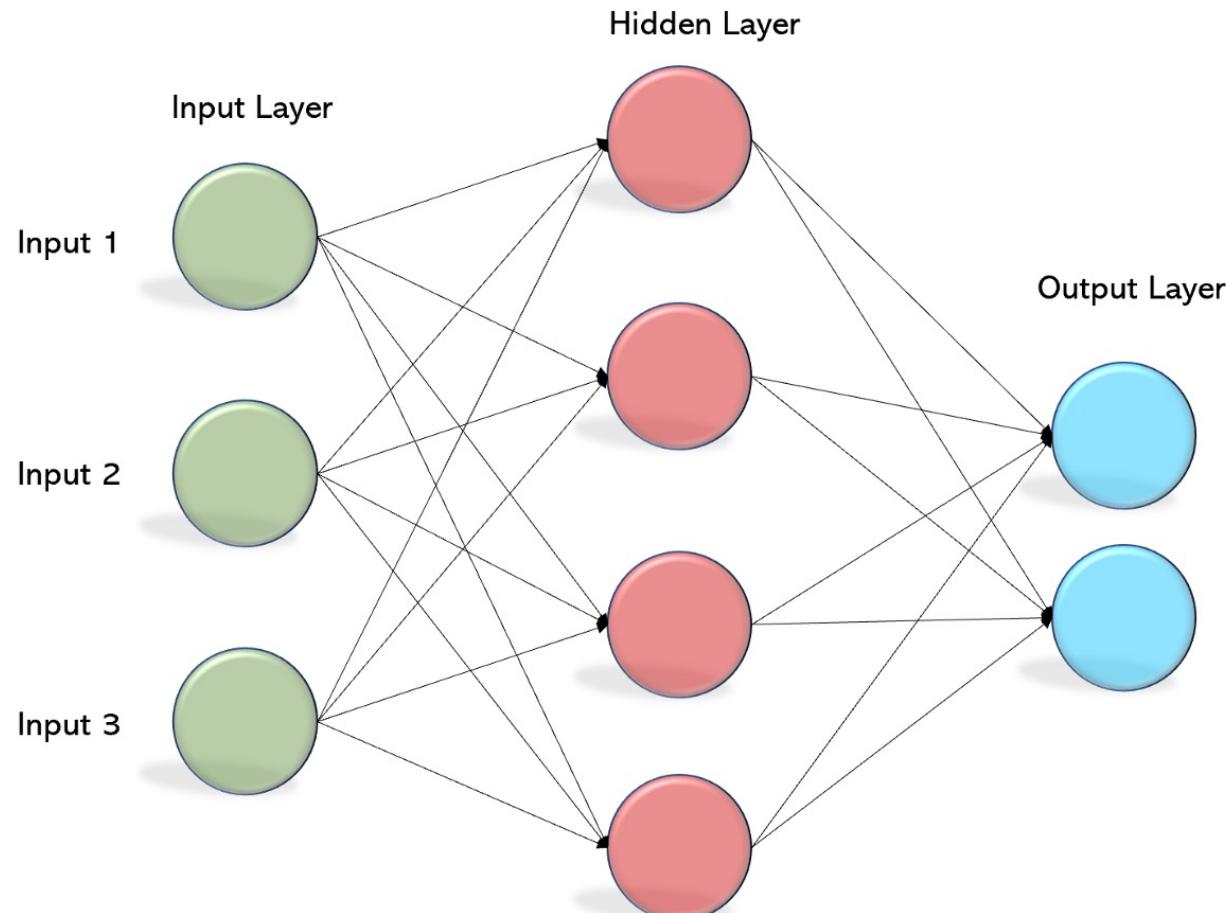
Main concepts that we will cover before structure prediction via ML

- Loss function (Mean Squared Error, Cross entropy)
- Optimization techniques (Gradient descent)
- Activation function (ReLU)
- Normalization
- Learning rate
- Training
- Backpropagation

# What is supervised machine learning?



# Basic Deep Neural Networks (DNNs) a.k.a Multilayer Perceptrons (MLPs)



# Loss functions

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$CE = - \sum_{i=1}^n Y_i \cdot \log \hat{Y}_i$$

# Loss functions

$$MSE = \frac{1}{n} \sum_{n=1}^n (Y_i - \hat{Y}_i)^2$$

$$CE = - \sum_{n=1}^n p(x) \log q(x)$$

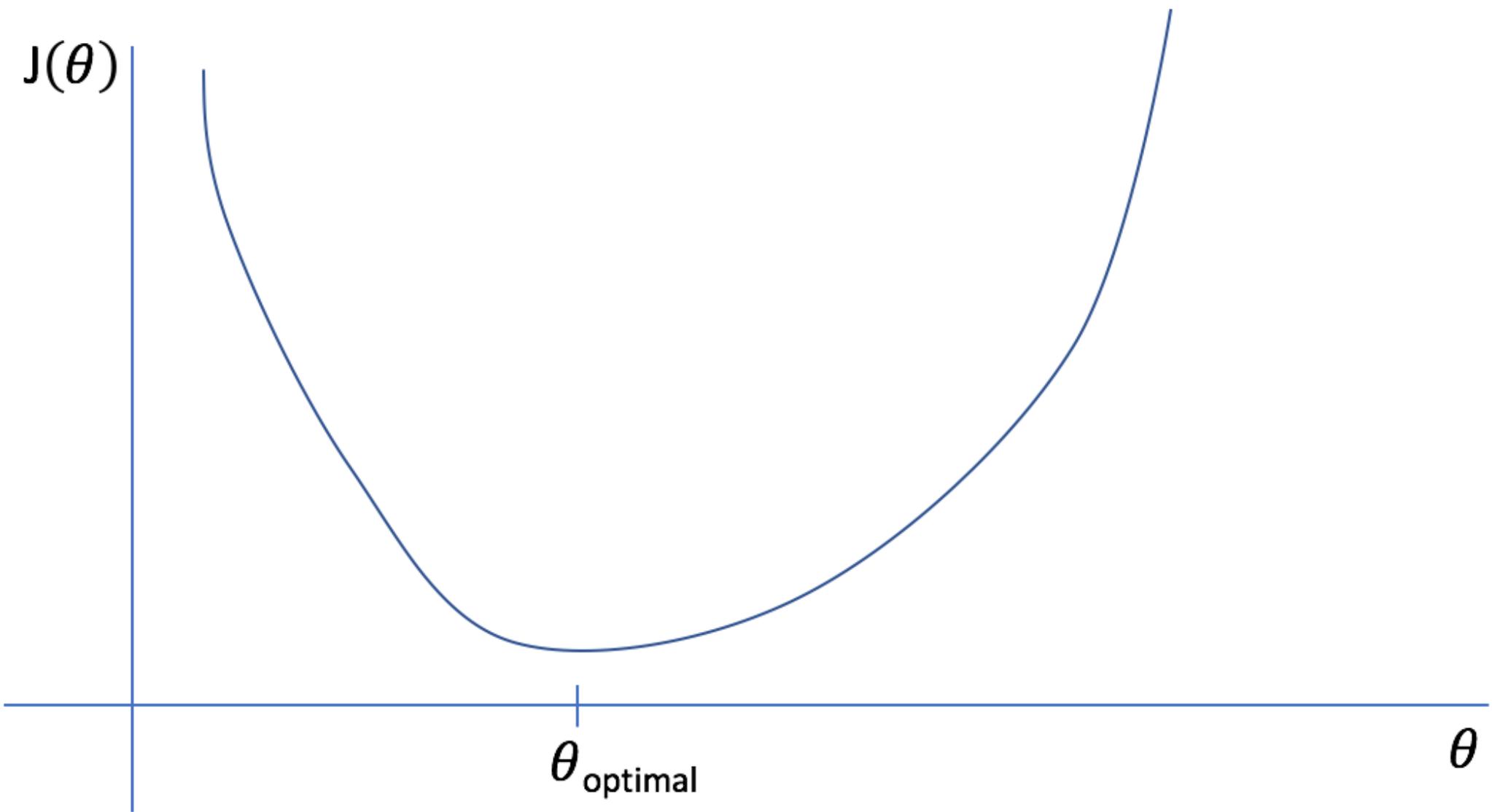
# Optimization via Gradient Descent

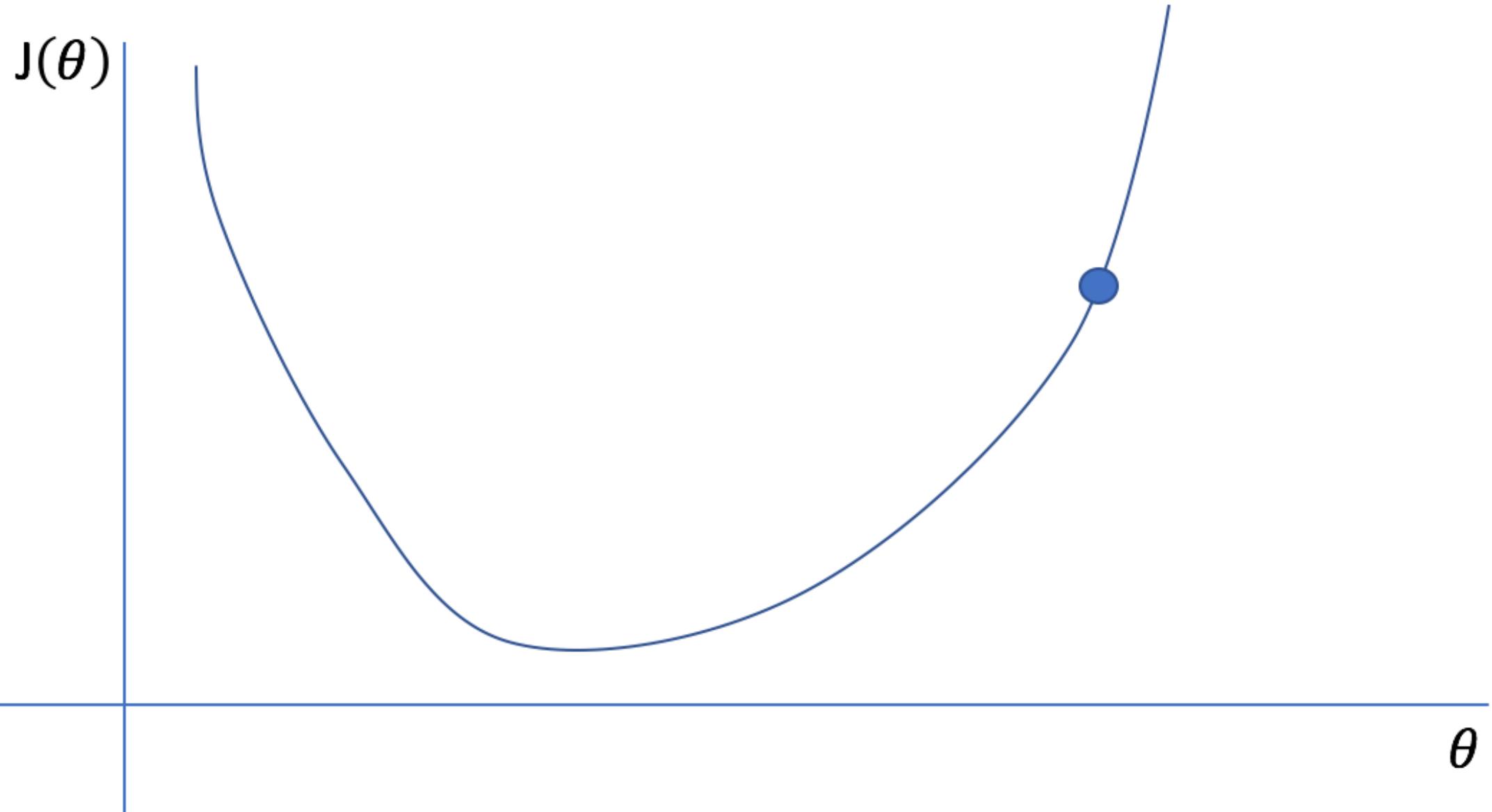
$$\theta_i := \theta_i + \Delta\theta_i$$

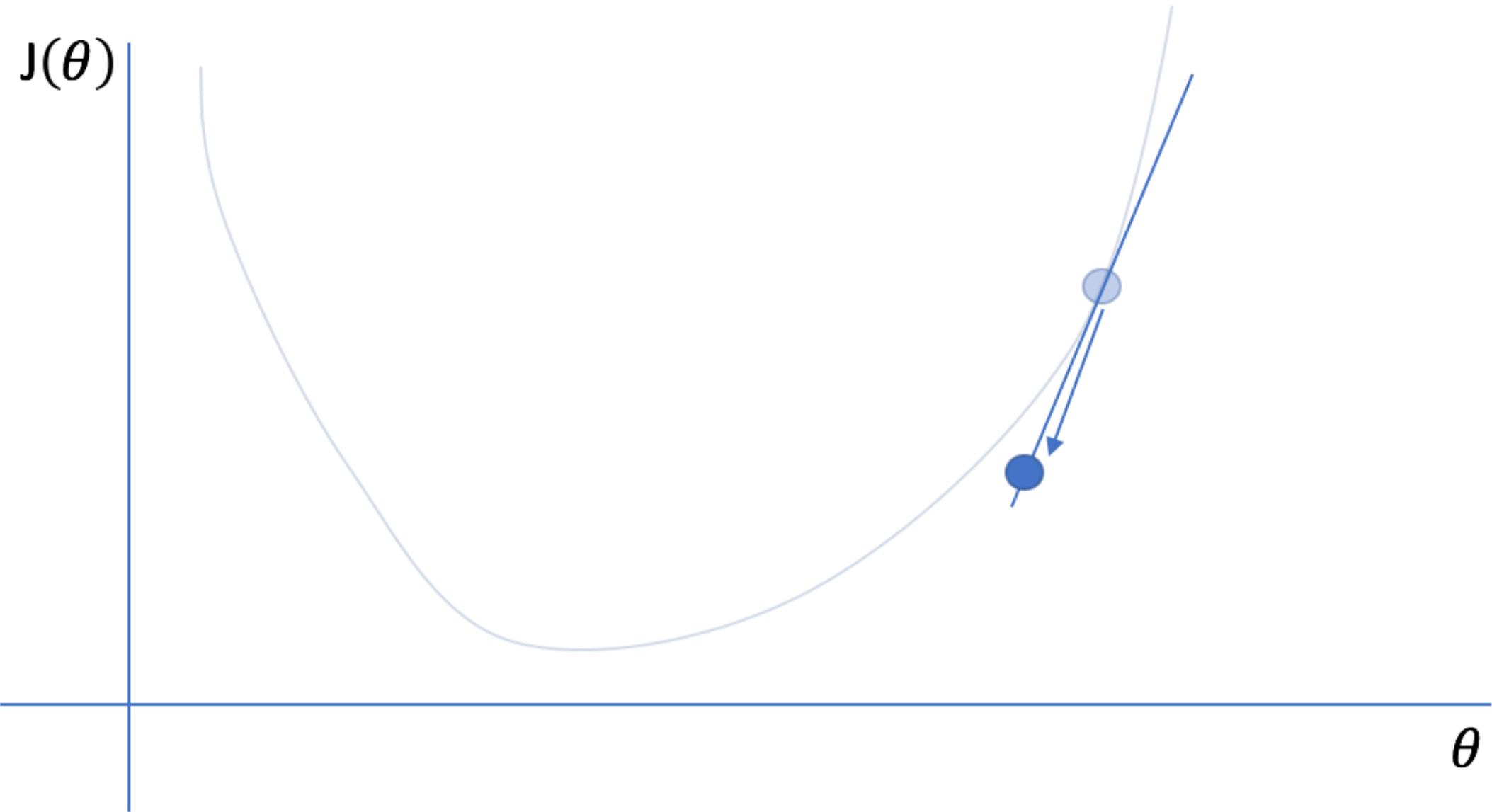
where

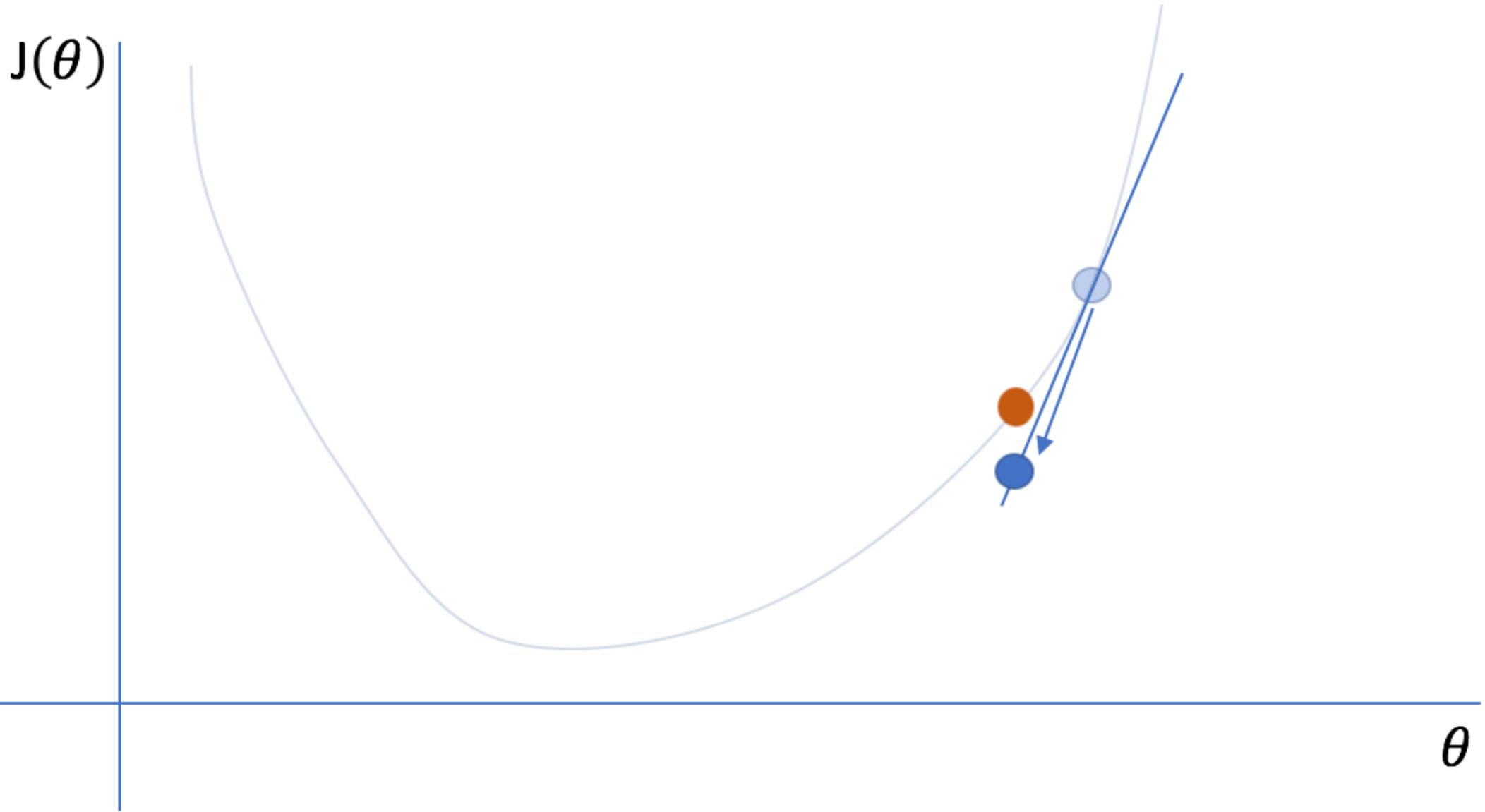
$\theta$  → parameter  
 $\alpha$  → learning rate

$$\Delta\theta_i = -\alpha \frac{\partial J(\theta)}{\partial \theta_i}$$





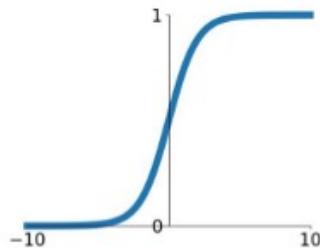




# Activation functions

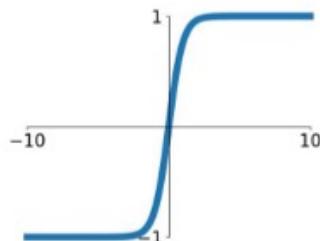
## Sigmoid

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



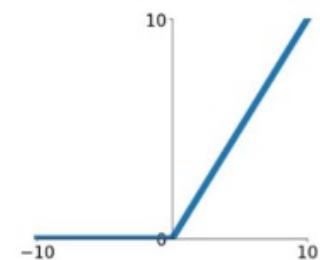
## tanh

$$\tanh(x)$$



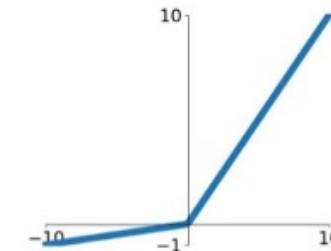
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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

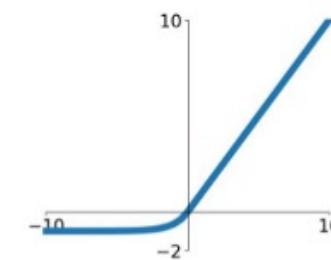


## Maxout

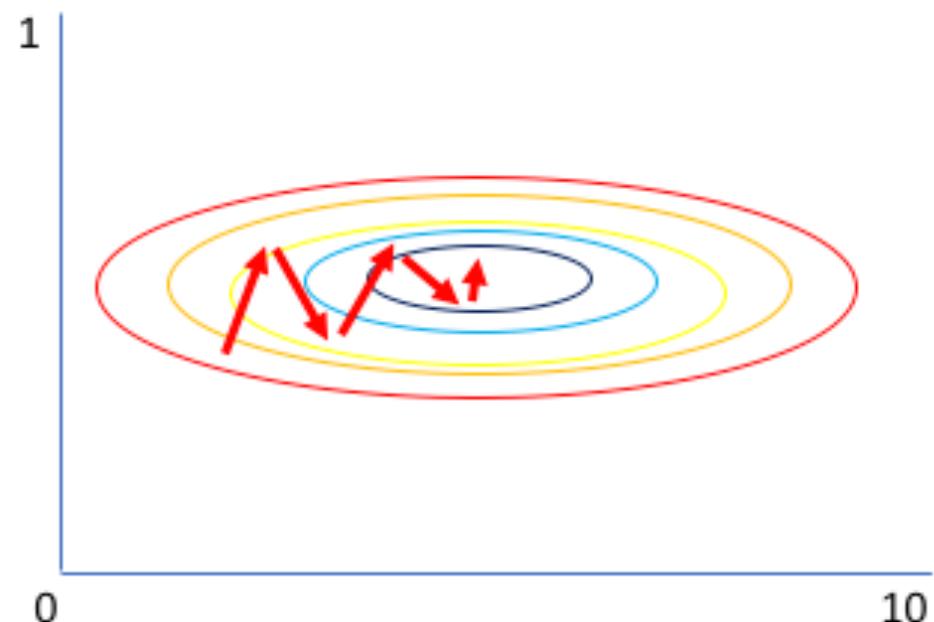
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

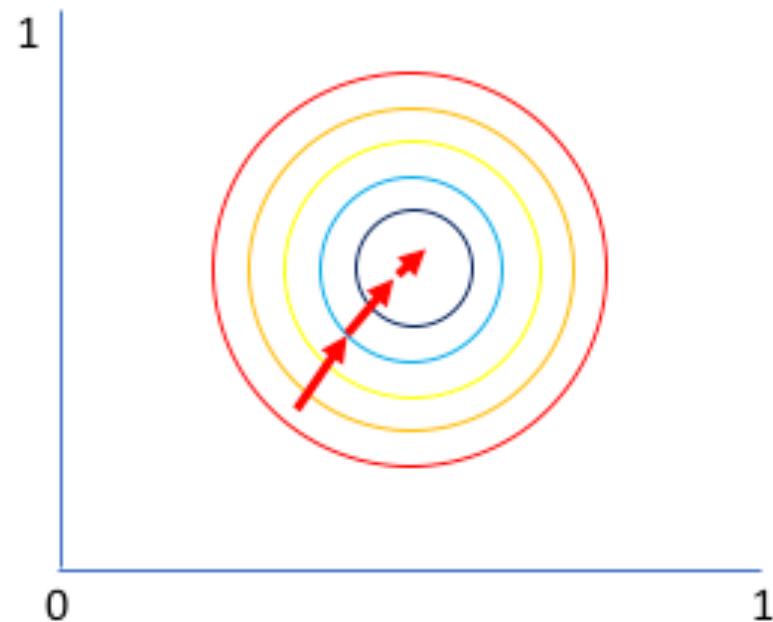
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Normalization

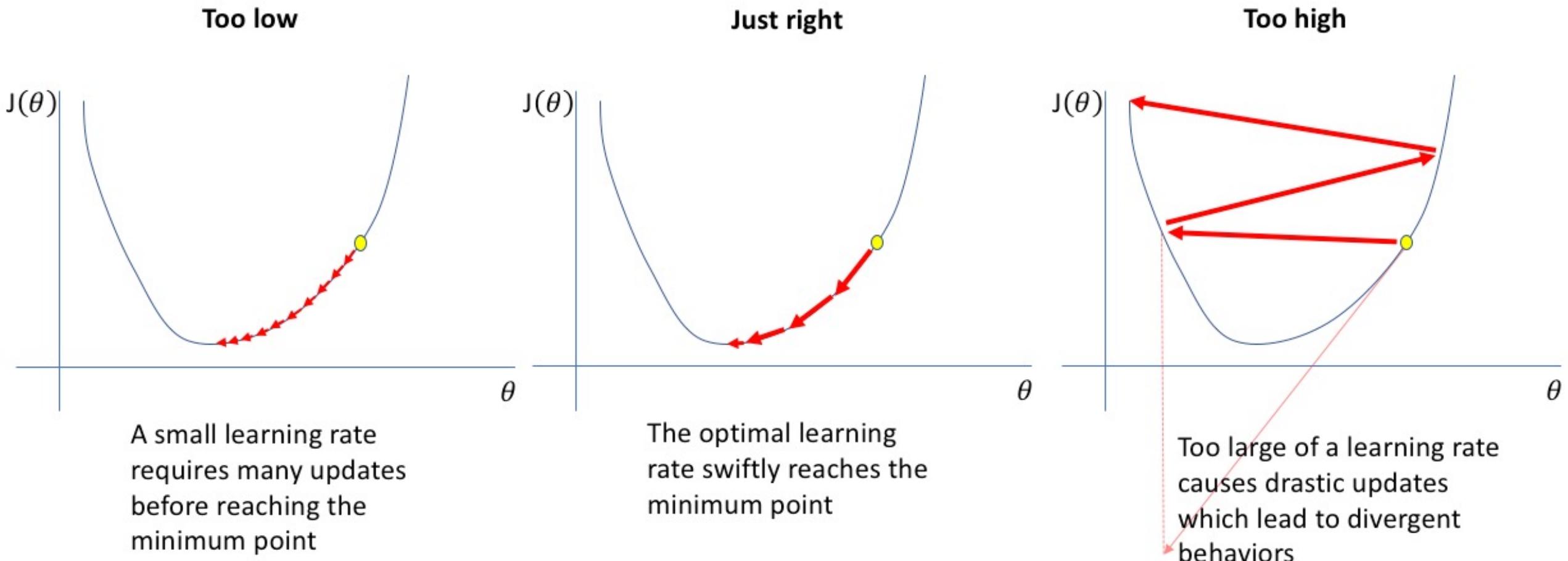


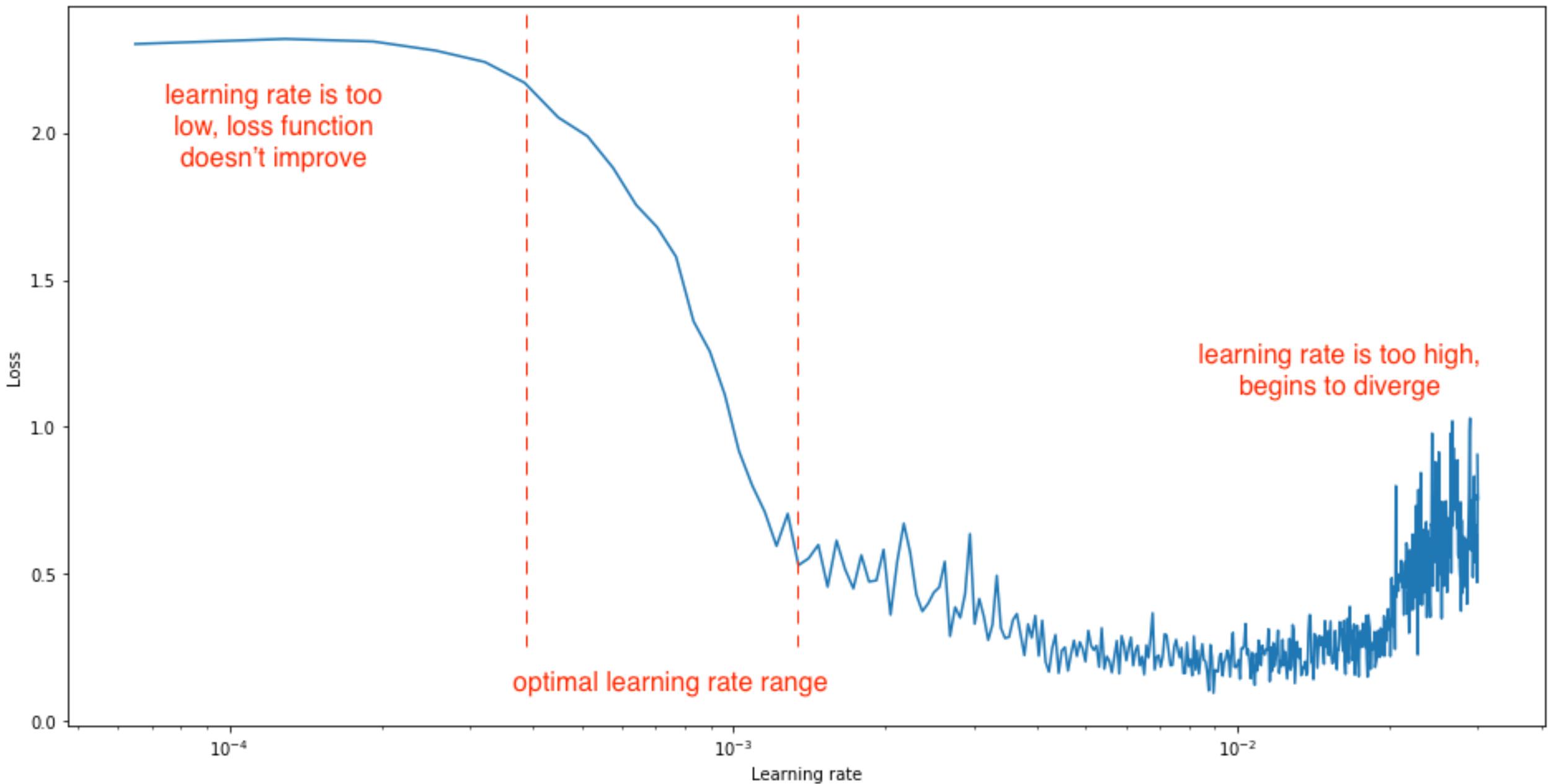
Gradient of larger parameter  
dominates the update



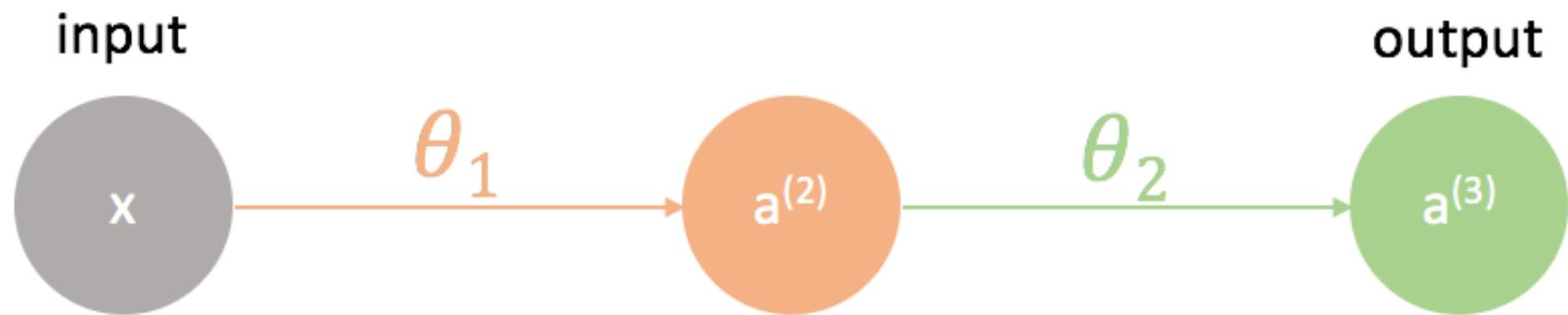
Both parameters can be  
updated in equal proportions

# Learning rate





# Training with backpropagation



$$a^{(2)} = g(\theta_1 x)$$

$$a^{(3)} = g(\theta_2 a^{(2)})$$

input neuron



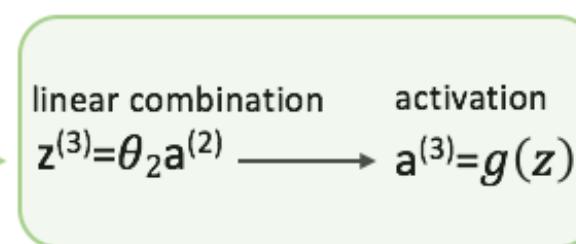
$\theta_1$



hidden neuron

$$\begin{array}{ccc} \text{linear combination} & & \text{activation} \\ z^{(2)} = \theta_1 x & \longrightarrow & a^{(2)} = g(z) \end{array}$$

output neuron



$\theta_2$



output:  $a^{(3)}$   
target:  $y$

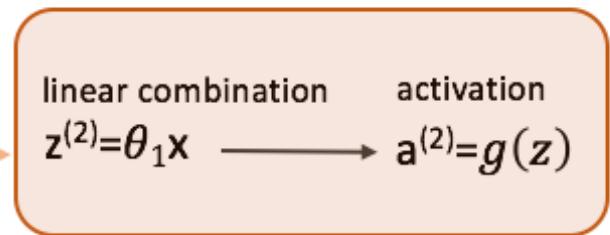
input neuron



$\theta_1$



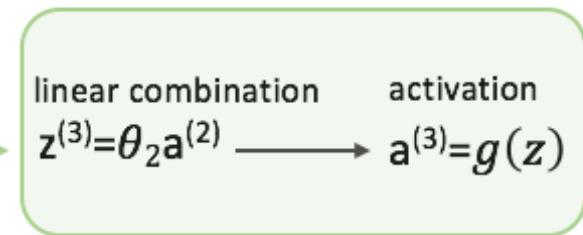
hidden neuron



$\theta_2$

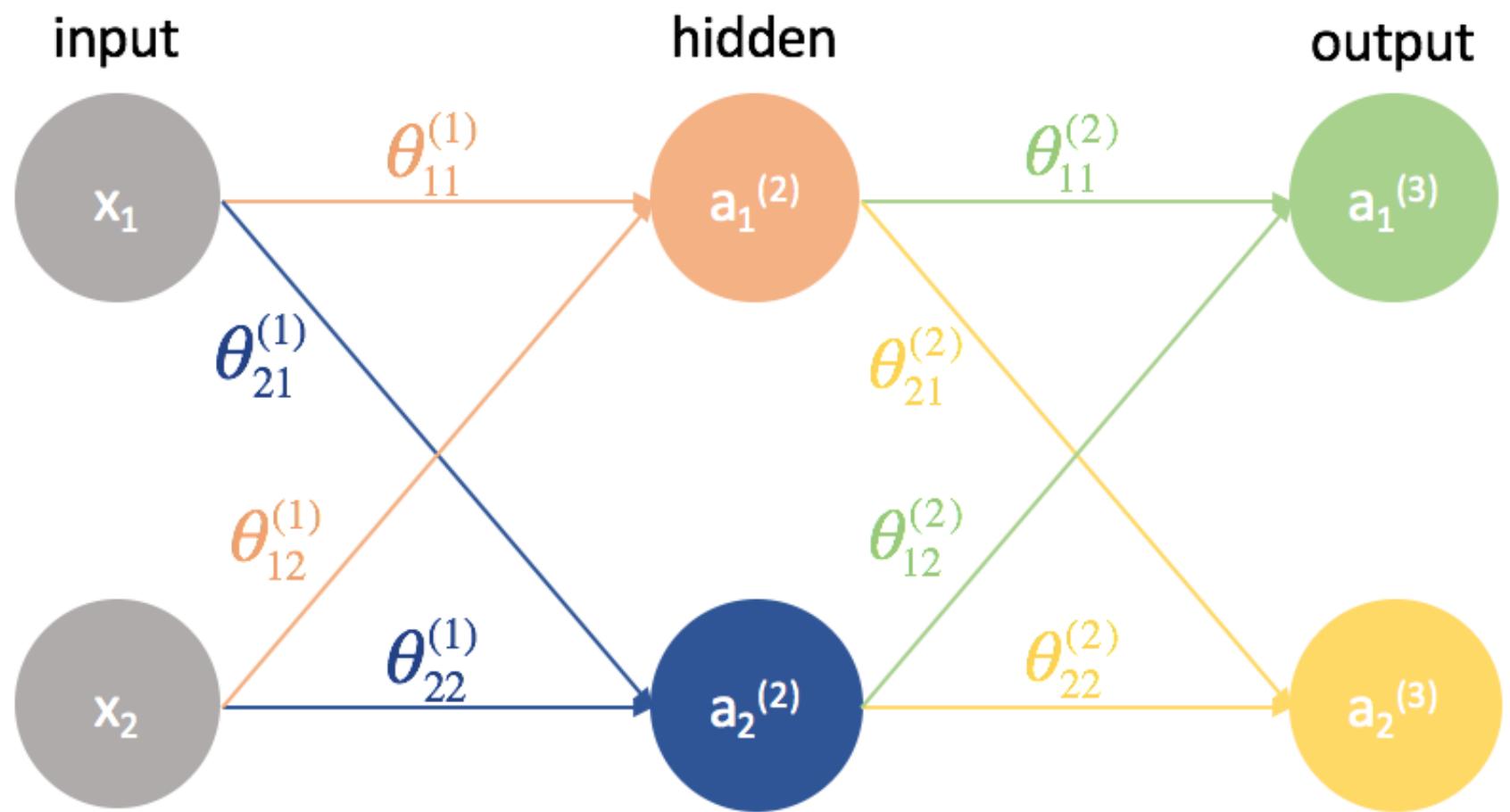


output neuron



performance

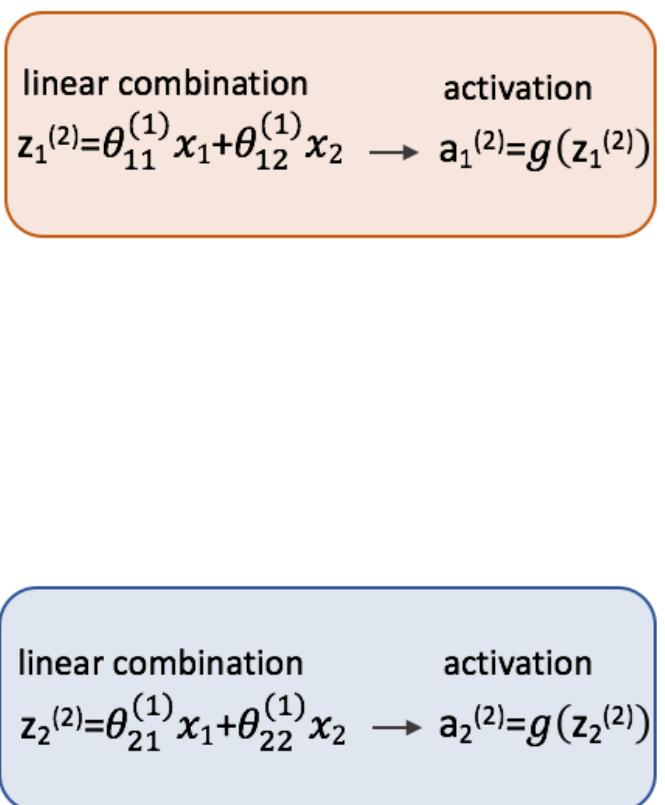
$$J(\theta) = \frac{1}{2} (y - a^{(3)})^2$$



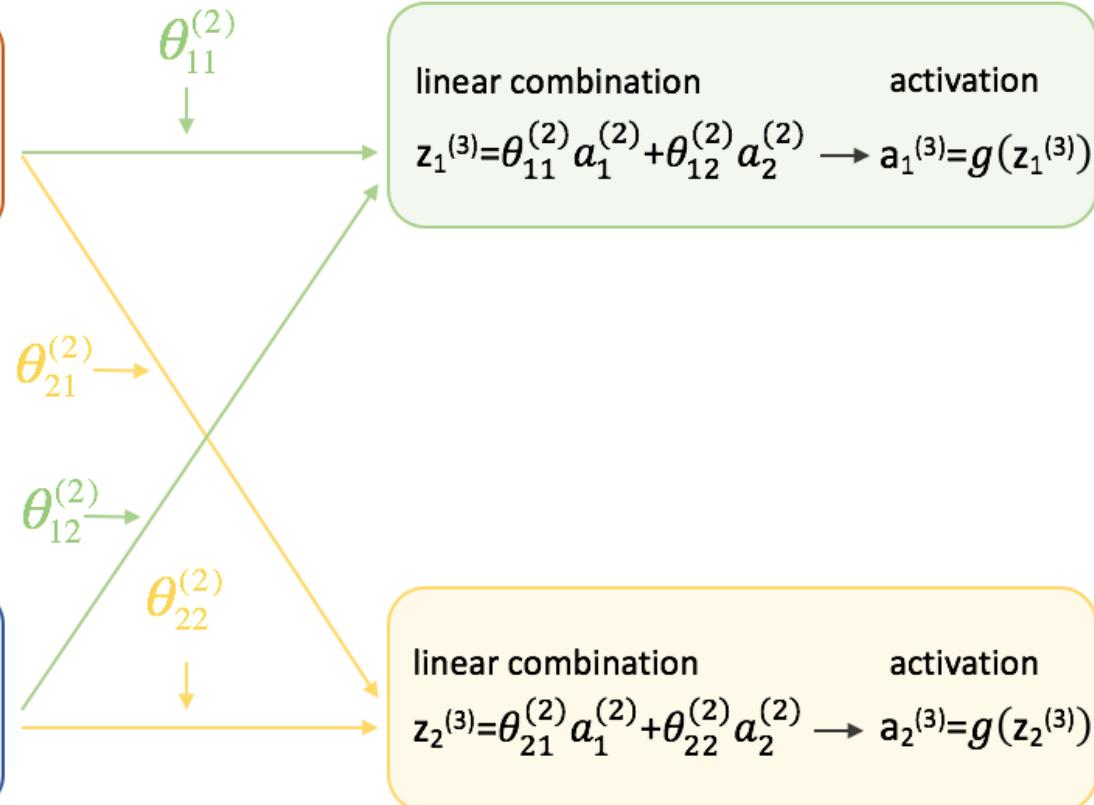
input layer



hidden layer

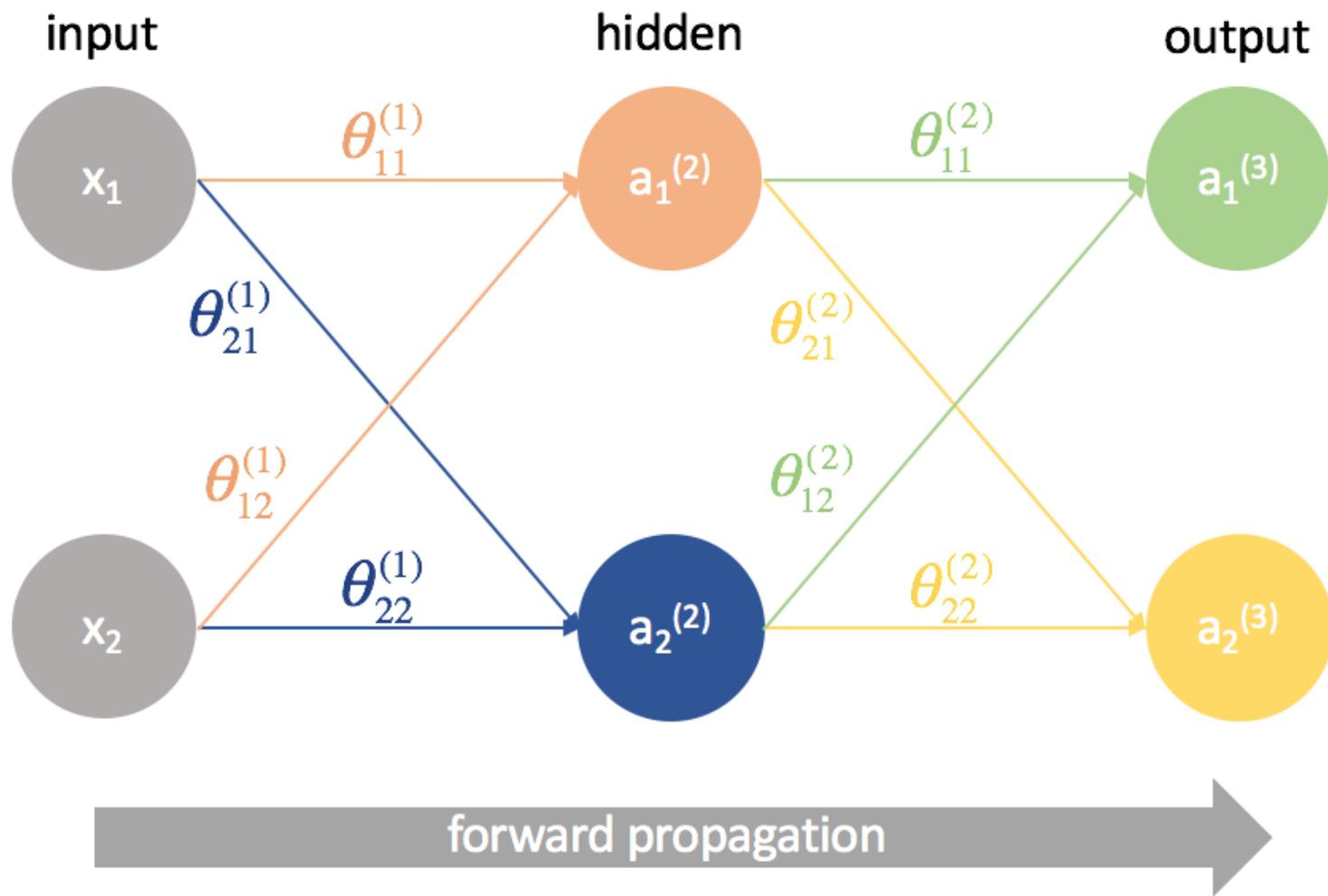


output layer

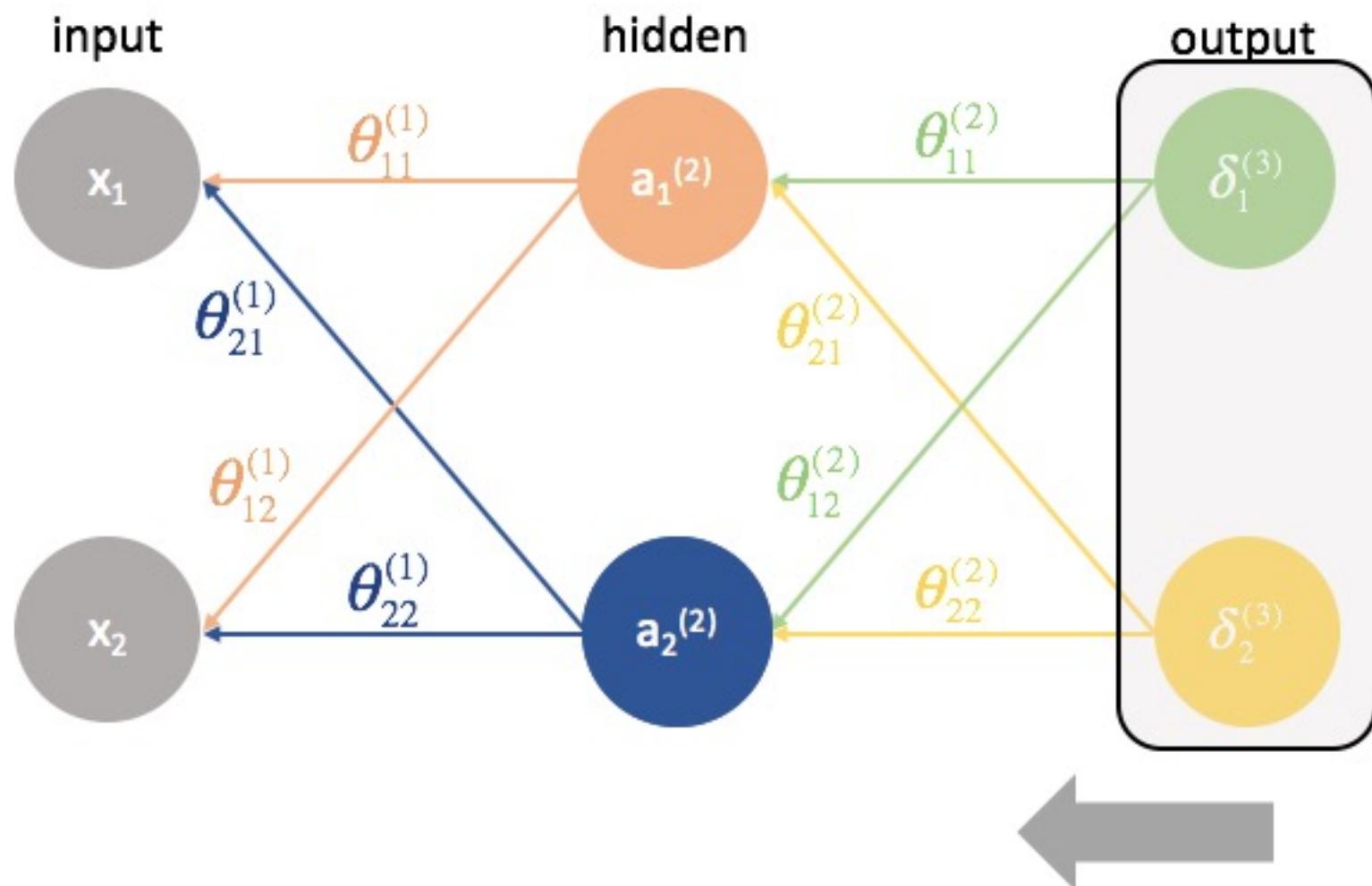


output:  $a_1^{(3)}$   
target:  $y_1$

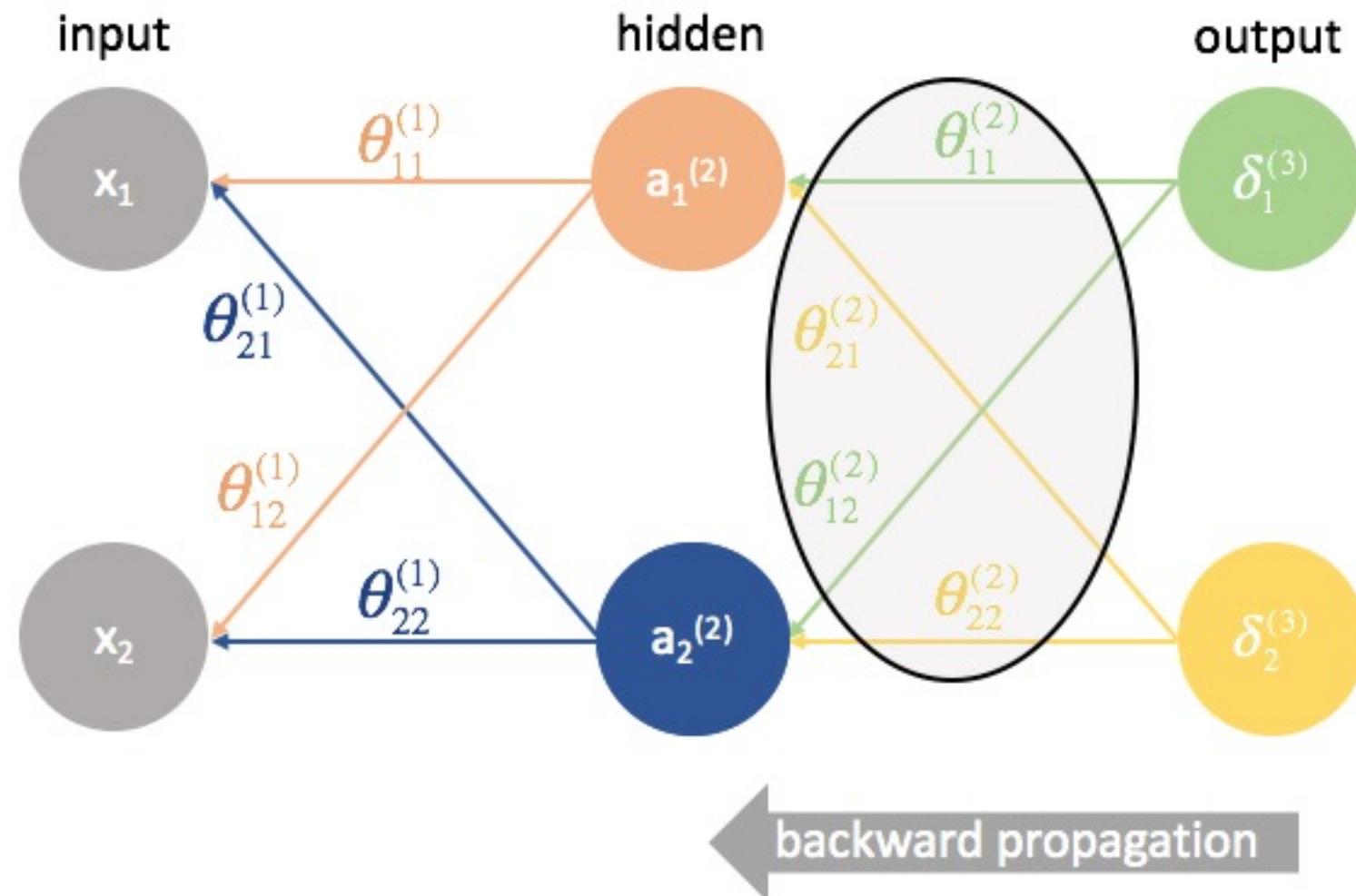
output:  $a_2^{(3)}$   
target:  $y_2$



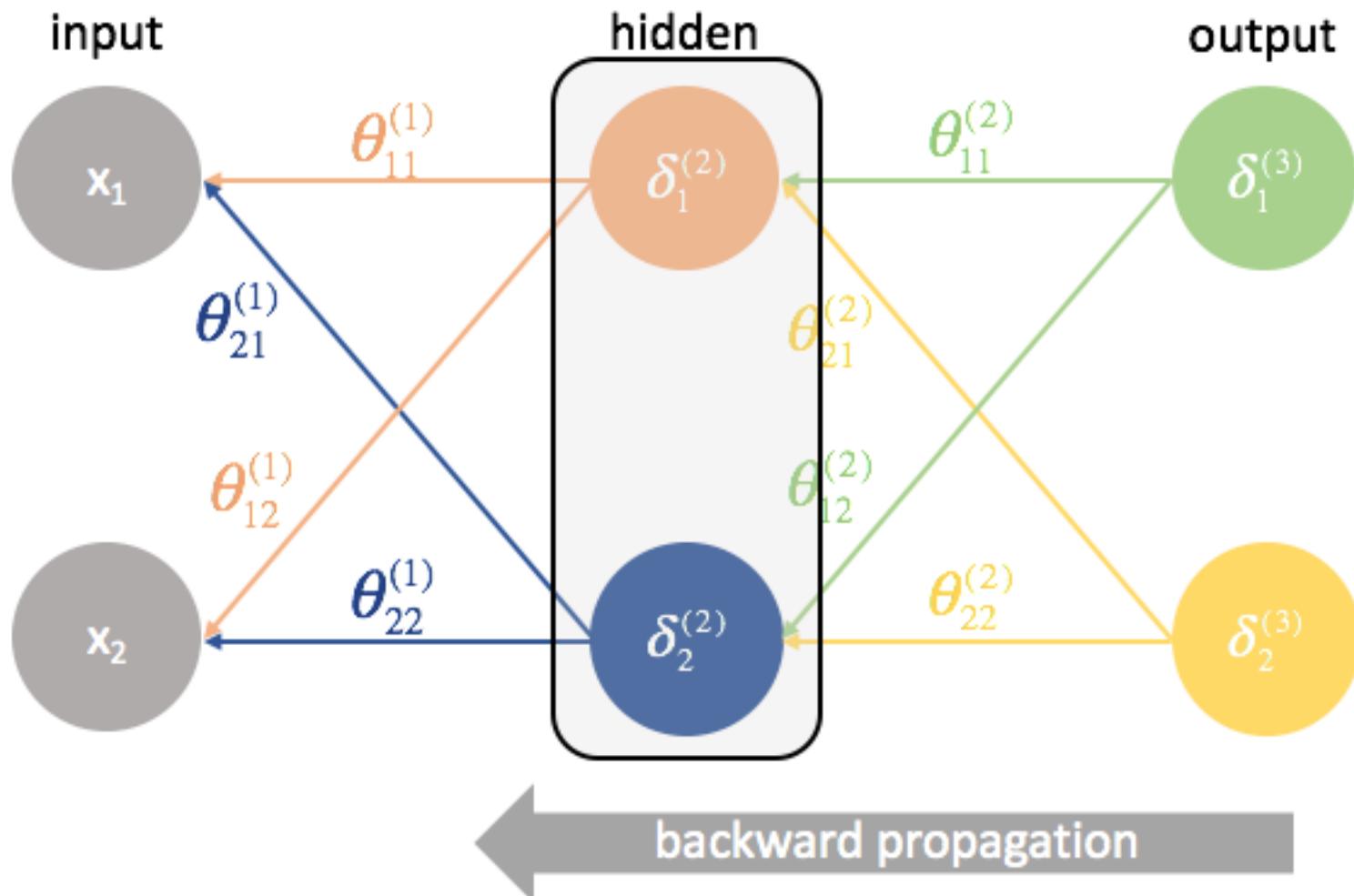
$$\delta^{(3)} = \frac{1}{m} (y - a^{(3)}) f'(a^{(3)})$$



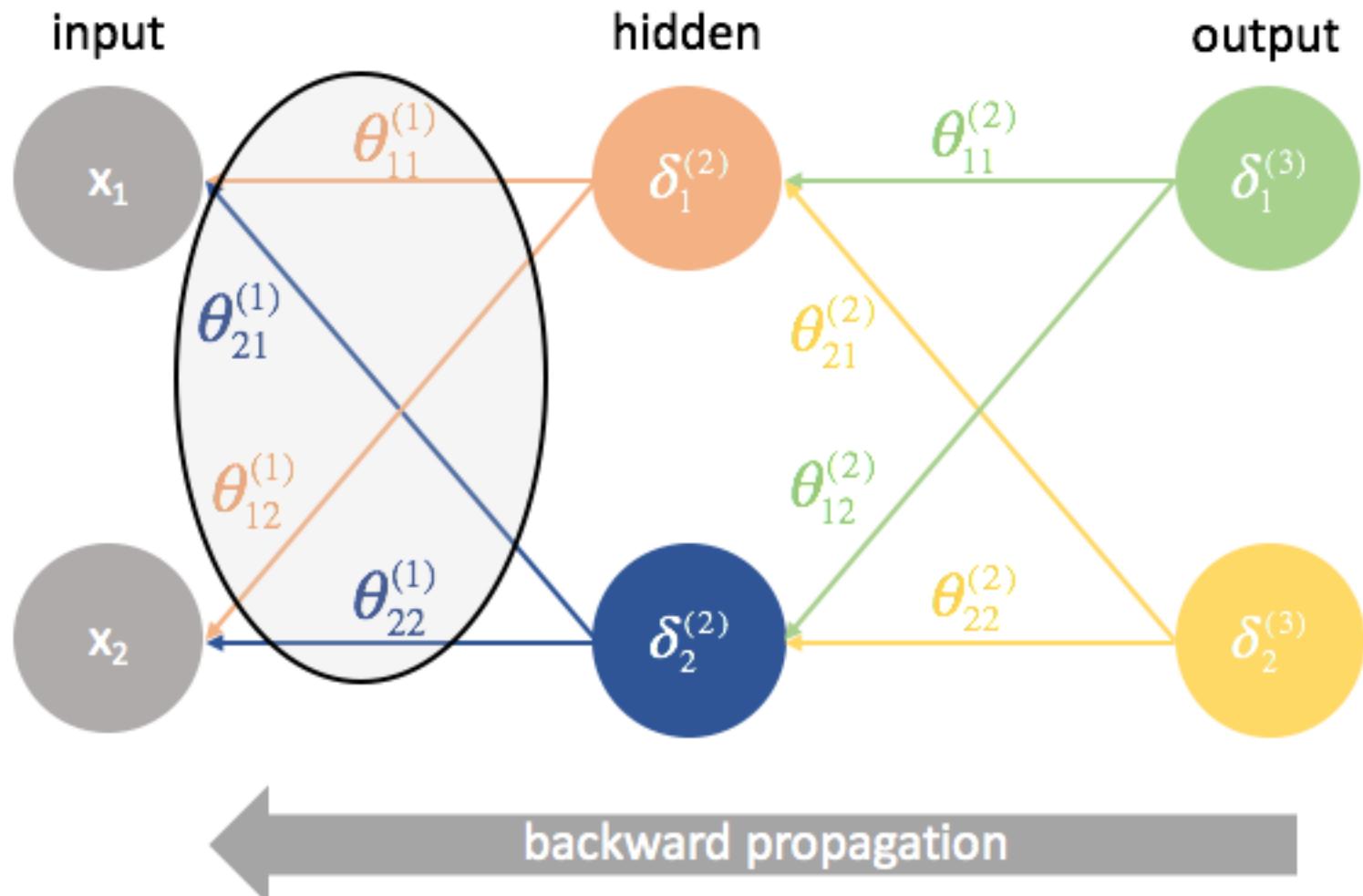
$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(2)}} = (\delta^{(3)})^T a^{(2)}$$



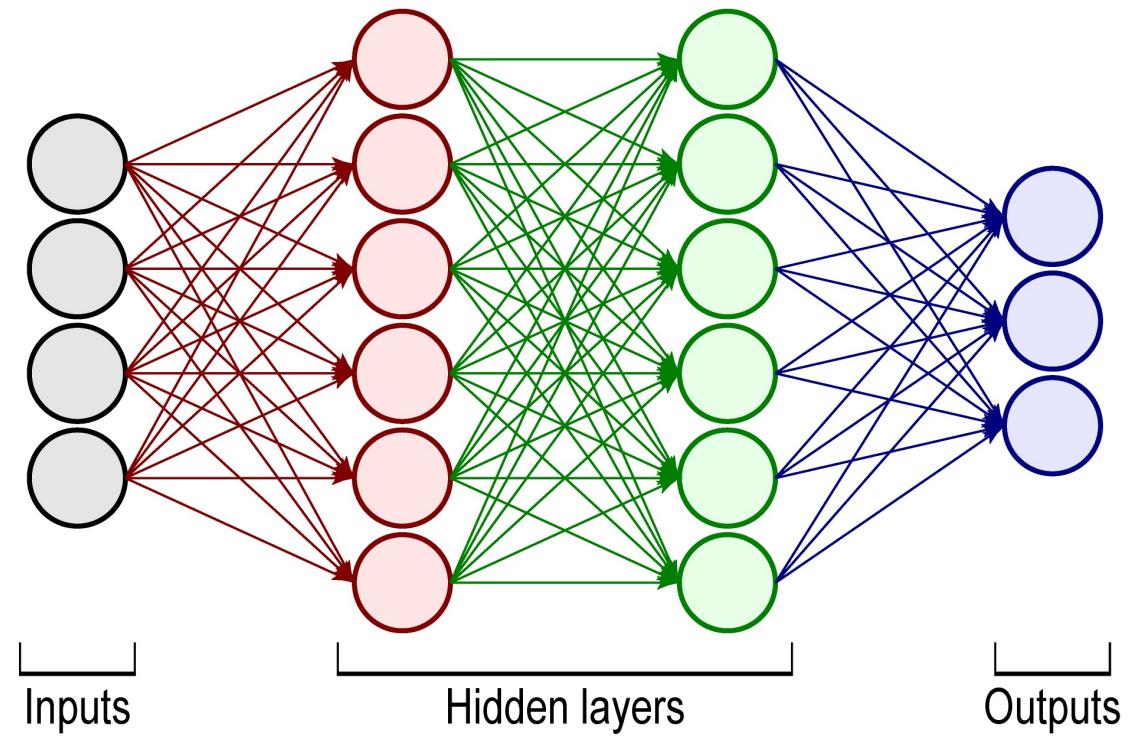
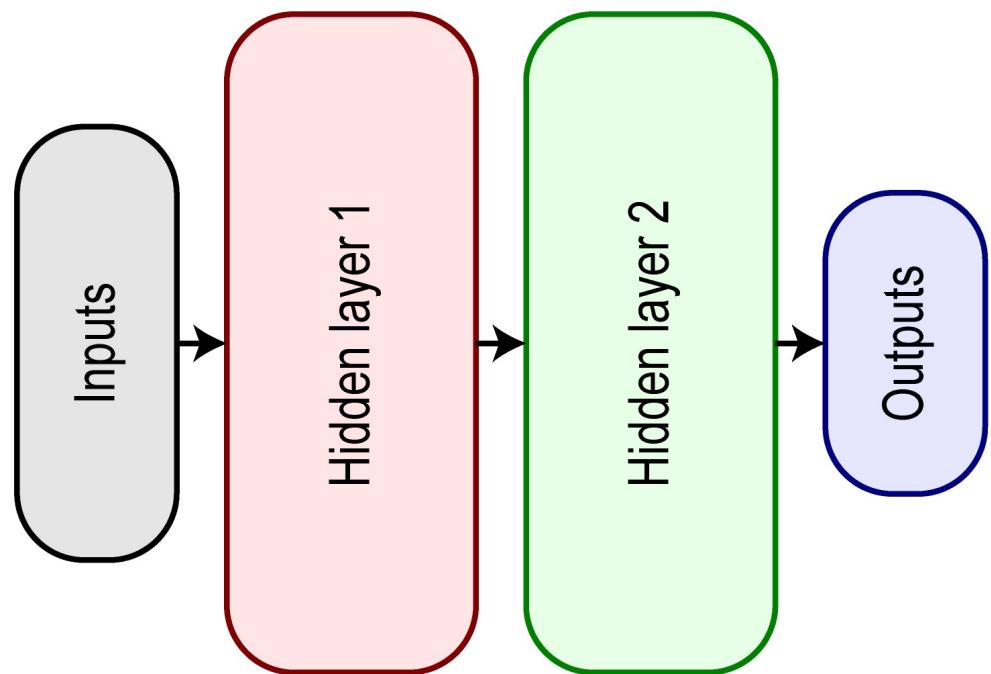
$$\delta^{(2)} = \delta^{(3)} \Theta^{(2)} f'(a^{(2)})$$



$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(1)}} = (\delta^{(2)})^T x$$

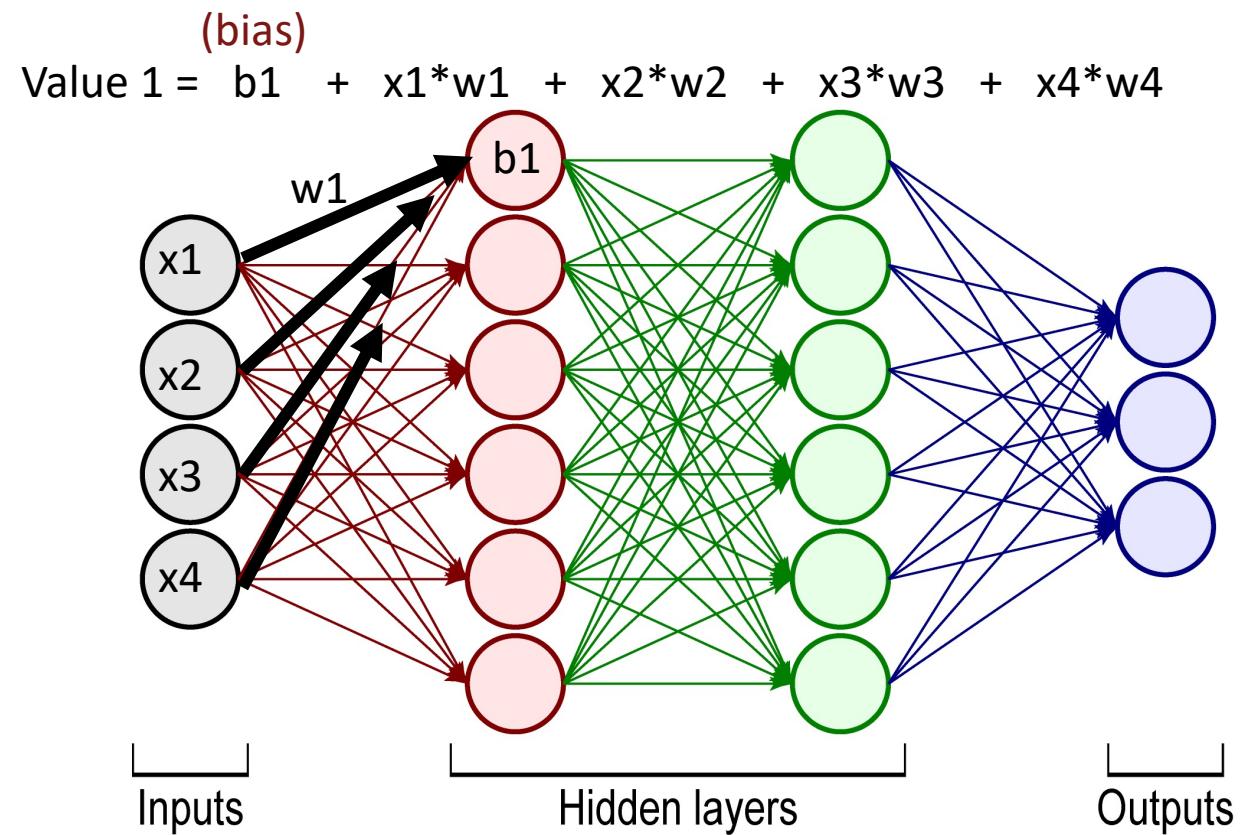
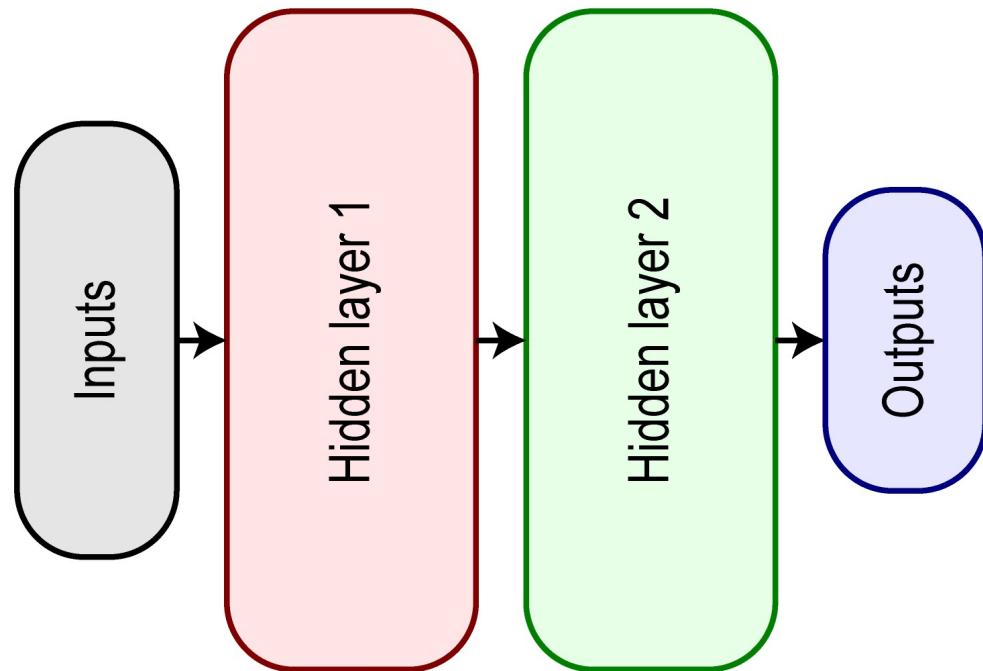


# Basic deep neural networks (a.k.a. multilayer perceptrons/MLPs)



$$f_o \left( f_b \left( f_a \left( [x_1 \ \dots \ x_4] \begin{bmatrix} w_{a,11} & \dots & w_{a,16} \\ \vdots & \ddots & \vdots \\ w_{a,41} & \dots & w_{a,46} \end{bmatrix} + [c_{a,1} \ \dots \ c_{a,6}] \right) \begin{bmatrix} w_{b,11} & \dots & w_{b,16} \\ \vdots & \ddots & \vdots \\ w_{b,61} & \dots & w_{b,66} \end{bmatrix} + [c_{b,1} \ \dots \ c_{b,6}] \right) \begin{bmatrix} w_{o,11} & \dots & w_{o,13} \\ \vdots & \ddots & \vdots \\ w_{o,61} & \dots & w_{o,63} \end{bmatrix} + [c_{o,1} \ c_{o,2} \ c_{o,3}] \right) = [y_1 \ y_2 \ y_3]$$

# Basic deep neural networks (a.k.a. multilayer perceptrons/MLPs)



$$f_o \left( f_b \left( f_a \left( [x_1 \dots x_4] \begin{bmatrix} w_{a,11} & \dots & w_{a,16} \\ \vdots & \ddots & \vdots \\ w_{a,41} & \dots & w_{a,46} \end{bmatrix} + [c_{a,1} \dots c_{a,6}] \right) \begin{bmatrix} w_{b,11} & \dots & w_{b,16} \\ \vdots & \ddots & \vdots \\ w_{b,61} & \dots & w_{b,66} \end{bmatrix} + [c_{b,1} \dots c_{b,6}] \right) \begin{bmatrix} w_{o,11} & \dots & w_{o,13} \\ \vdots & \ddots & \vdots \\ w_{o,61} & \dots & w_{o,63} \end{bmatrix} + [c_{o,1} \dots c_{o,3}] \right) = [y_1 \dots y_3]$$

# How do we figure out these weights?

- Two families:
  - Gradient-based minimization
    - This will work for us 99% of the time
    - Very similar to Rosetta's MinMover protocol
  - Derivative free optimization
    - “reinforcement learning”
    - AlphaGo, DeepStack, etc.
    - Learn to play video games, where “loss” is harder to calculate
    - Very similar to our docking, folding methods
      - Monte Carlo & genetic algorithms

# Summary

- Neural networks are chains of parameters
- If our system is differentiable, we can use gradient-based minimization to optimize these parameters
  - 99% use case
  - Can be tricky to ensure a healthy gradient
    - Normalize all input data
    - Use batch normalization for large models
  - Be aware of your learning rate!
- Otherwise, we must use derivative-free optimization
  - Less common
  - “reinforcement learning”
  - *Much* slower

# Resources

Math Overview: <https://www.jeremyjordan.me/neural-networks-training/>

Normalization: <https://www.jeremyjordan.me/batch-normalization/>

Learning Rate: <https://www.jeremyjordan.me/nn-learning-rate/>

## Interactive Demos

- <http://neuralnetworksanddeeplearning.com/chap4.html>
- <https://playground.tensorflow.org/>

It's time for a  
quick  
tensorflow  
playground  
NN demo



<https://playground.tensorflow.org/>

# Data tools and techniques

- **Basic Data Manipulation and Analysis**

Performing well-defined computations or asking well-defined questions (“queries”)

- **Data Mining**

Looking for patterns in data

- **Machine Learning**

Using data to build models and make predictions

- **Data Visualization**

Graphical depiction of data

- **Data Collection and Preparation**

# Regression

Using data to build models and make predictions

- Supervised
- Training data, each example:
  - Set of predictor values - “independent variables”
  - Numerical output value - “dependent variable”
- Model is function from predictors to output
  - Use model to predict output value for new predictor values
- Example
  - Predictors: mother height, father height, current age
  - Output: height

# Classification

Using data to build models and make predictions

- Supervised
- Training data, each example:
  - Set of feature values - numeric or categorical
  - Categorical output value - “label”
- Model is method from feature values to label
  - Use model to predict label for new feature values
- Example
  - Feature values: age, gender, income, profession
  - Label: buyer, non-buyer

# Other examples

## Medical diagnosis

- **Feature values:** age, gender, history, symptom1-severity, symptom2-severity, test-result1, test-result2
- **Label:** disease

## Email spam detection

- **Feature values:** sender-domain, length, #images, keyword<sub>1</sub>, keyword<sub>2</sub>, ..., keyword<sub>n</sub>
- **Label:** spam or not-spam

## Credit card fraud detection

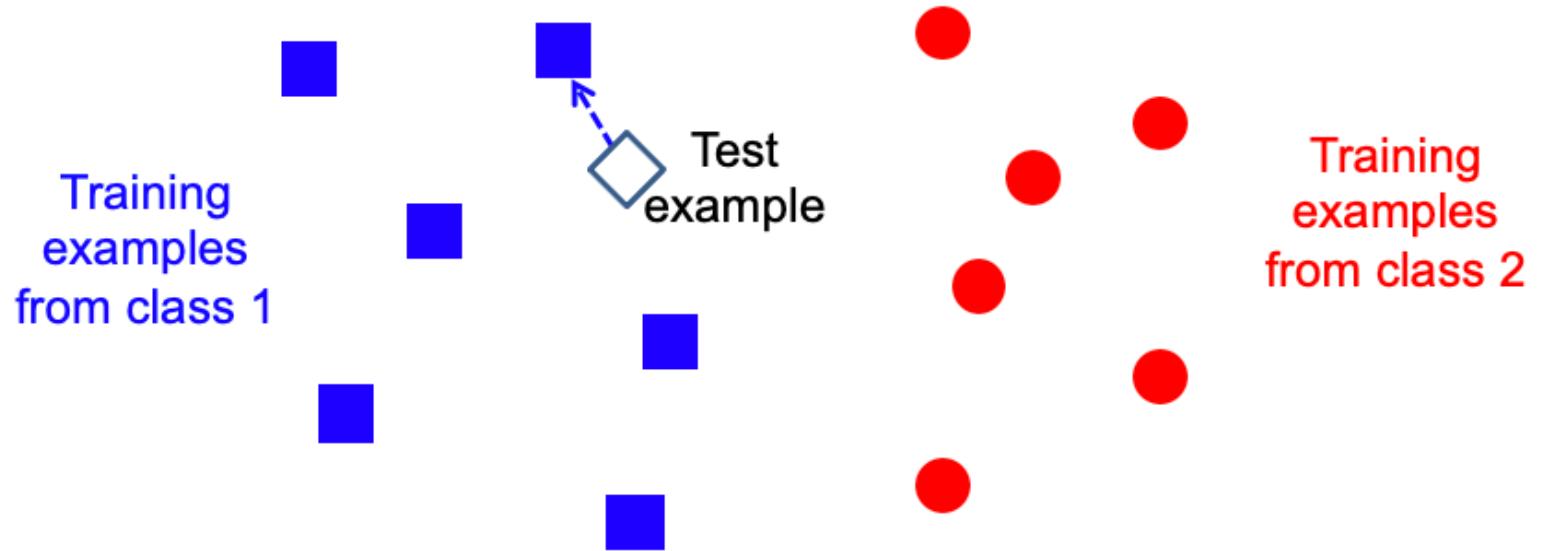
- **Feature values:** user, location, item, price
- **Label:** fraud or okay

# Algorithms for classification

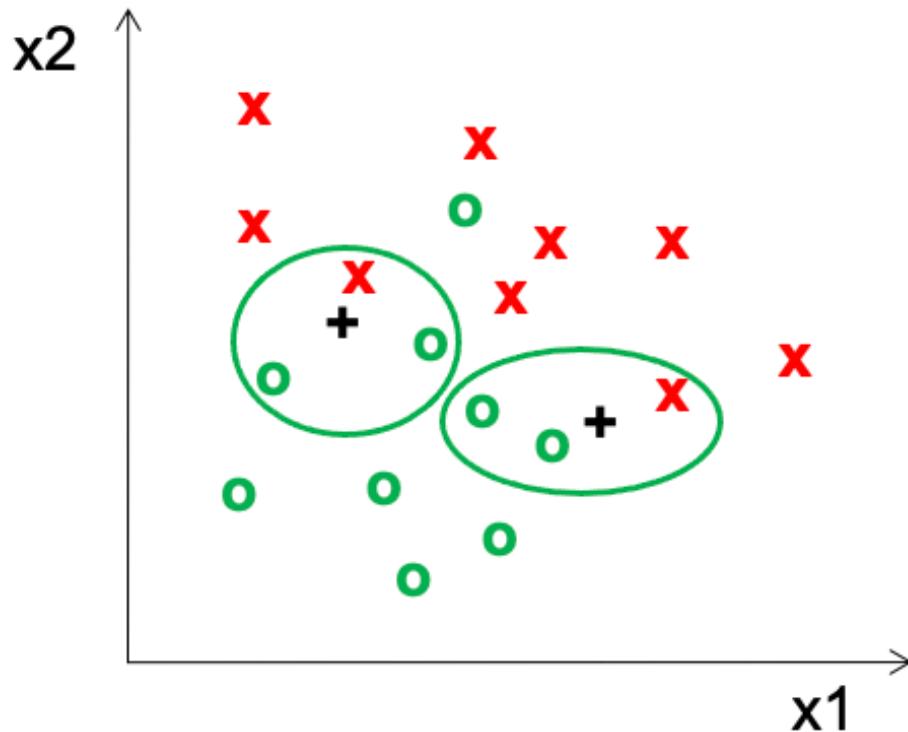
Despite similarity of problem statement to regression, non-numerical nature of classification leads to completely different approaches

- K-nearest neighbors
- Decision trees
- Naïve Bayes
- Deep neural networks
- ... and others

# K-Nearest Neighbors (KNN)



# K-Nearest Neighbors (KNN)



# K-Nearest Neighbors (KNN)

For any pair of data items  $i_1$  and  $i_2$ , from their feature values compute  $distance(i_1, i_2)$

Example:

Features - gender, profession, age, income, postal-code

$person_1 = (\text{male}, \text{teacher}, 47, \$25K, 94305)$

$person_2 = (\text{female}, \text{teacher}, 43, \$28K, 94309)$

$distance(person_1, person_2)$

$distance()$  can be defined as inverse of  $similarity()$

# K-Nearest Neighbors (KNN)

Features - gender, profession, age, income, postal-code

$\text{person}_1 = (\text{male}, \text{teacher}, 47, \$25\text{K}, 94305)$  buyer

$\text{person}_2 = (\text{female}, \text{teacher}, 43, \$28\text{K}, 94309)$  non-buyer

Remember training data has labels

# K-Nearest Neighbors (KNN)

Features - gender, profession, age, income, postal-code

$\text{person}_1 = (\text{male}, \text{teacher}, 47, \$25\text{K}, 94305)$  buyer

$\text{person}_2 = (\text{female}, \text{teacher}, 43, \$28\text{K}, 94309)$  non-buyer

Remember training data has labels

To classify a new item  $i$ : In the labeled data find  
the K closest items to  $i$ , assign most frequent label

$\text{person}_3 = (\text{female}, \text{doctor}, 40, \$40\text{K}, 95123)$

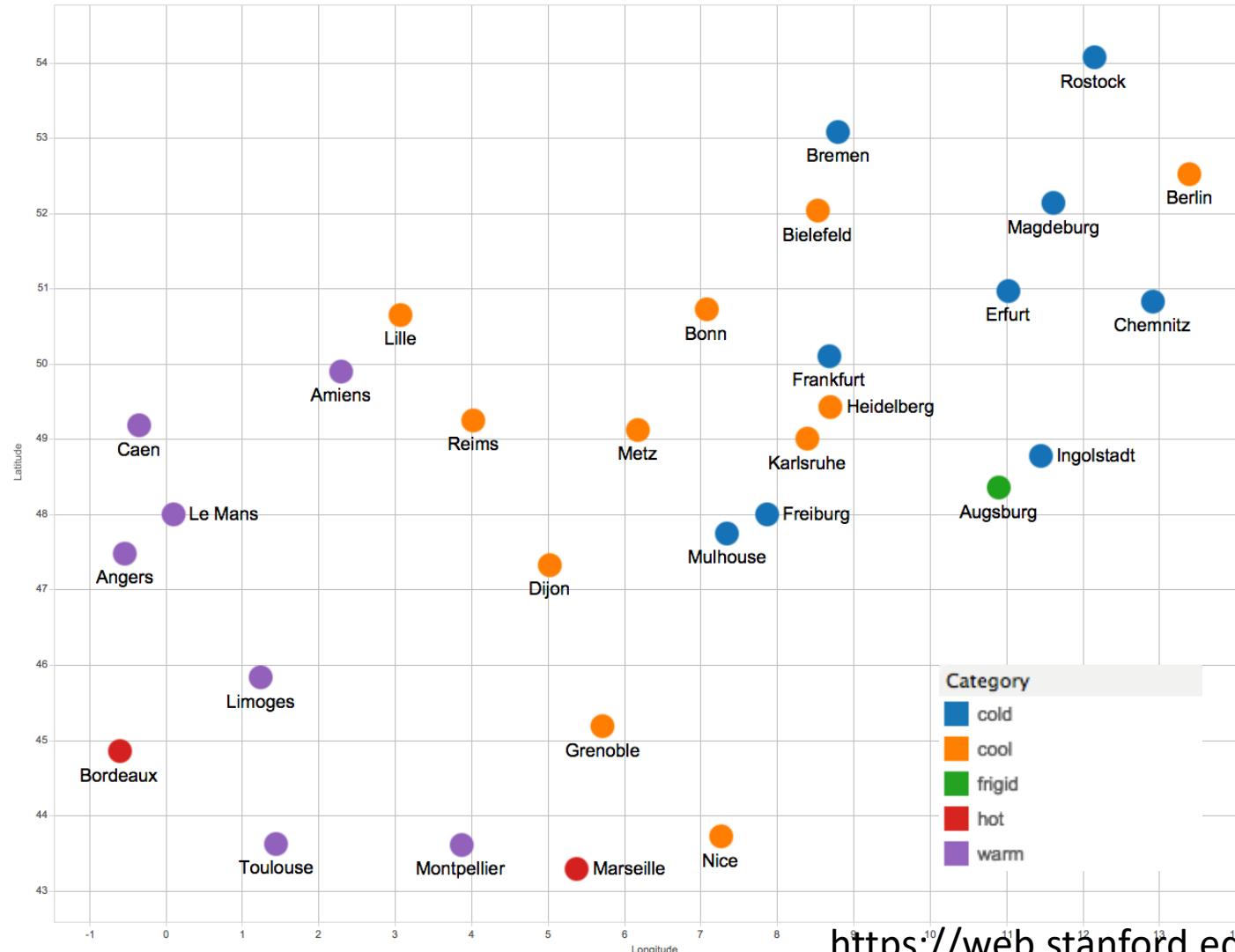
# KNN example

- City temperatures - France and Germany
- Features: longitude, latitude
- Distance is Euclidean distance
  - $distance([o_1, a_1], [o_2, a_2]) = \sqrt{(o_1 - o_2)^2 + (a_1 - a_2)^2}$
  - = actual distance in x-y plane
- Labels: frigid, cold, cool, warm, hot

Nice (7.27, 43.72) cool  
Toulouse (1.45, 43.62) warm  
Frankfurt (8.68, 50.1) cold  
.....

Predict temperature category from longitude and latitude

# KNN example



# KNN summary

To classify a new item  $i$ : find K closest items to  $i$  in the labeled data, assign most frequent label

- No hidden complicated math!
- Once distance function is defined, rest is easy
- Though not necessarily efficient
  - Real examples often have thousands of features
    - Medical diagnosis: symptoms (yes/no), test results
    - Email spam detection: words (frequency)

Database of labeled items might be enormous

# Regression using KNN

Features - gender, profession, age, income, postal-code

$\text{person}_1 = (\text{male}, \text{teacher}, 47, \$25K, 94305)$  \$250

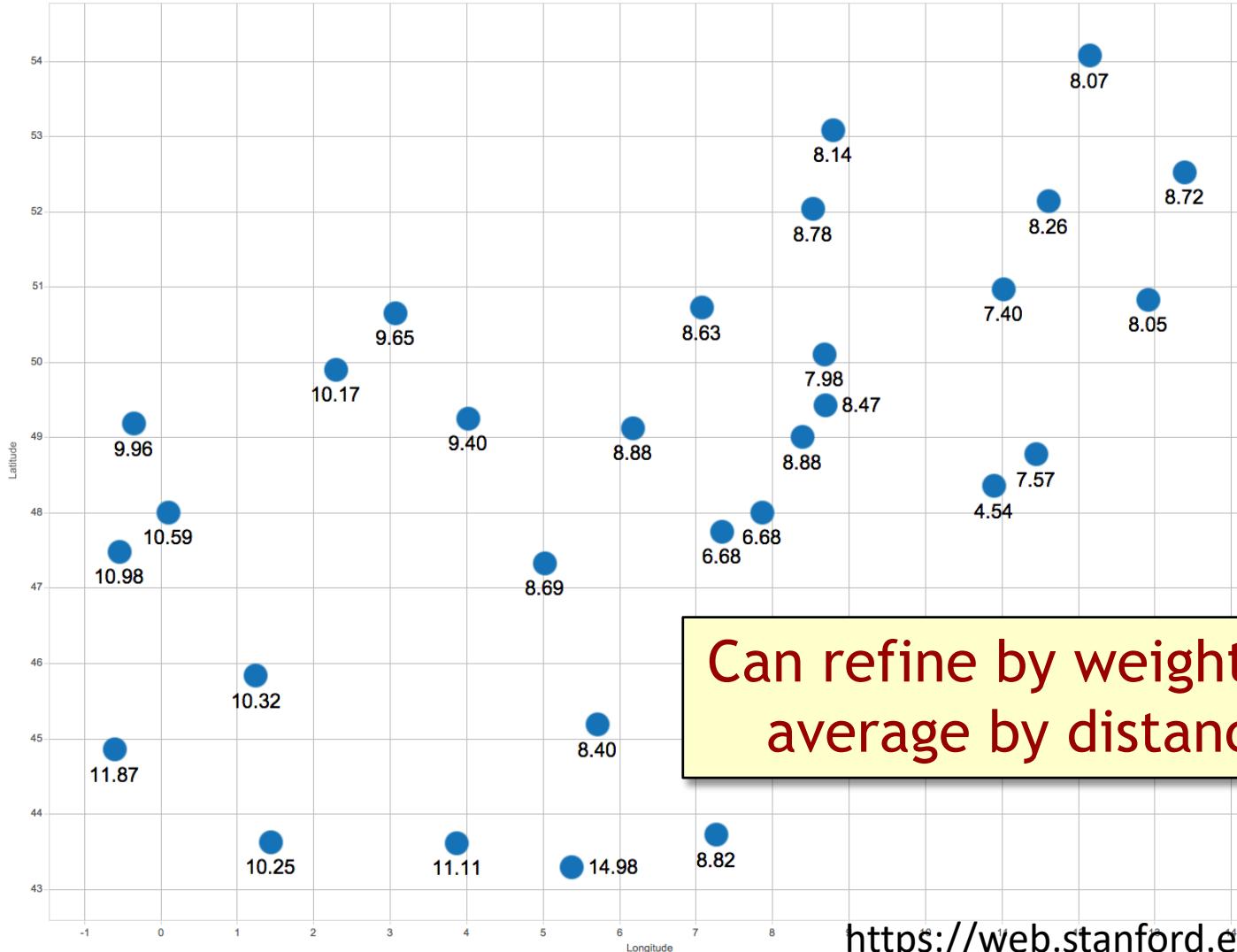
$\text{person}_2 = (\text{female}, \text{teacher}, 43, \$28K, 94309)$  \$100

Remember training data has labels

To classify a new item  $i$ , find K closest items to  $i$   
in the labeled data, assign average value of labels

$\text{person}_3 = (\text{female}, \text{doctor}, 40, \$40K, 95123)$

# Regression using KNN



Can refine by weighting  
average by distance

# Decision trees

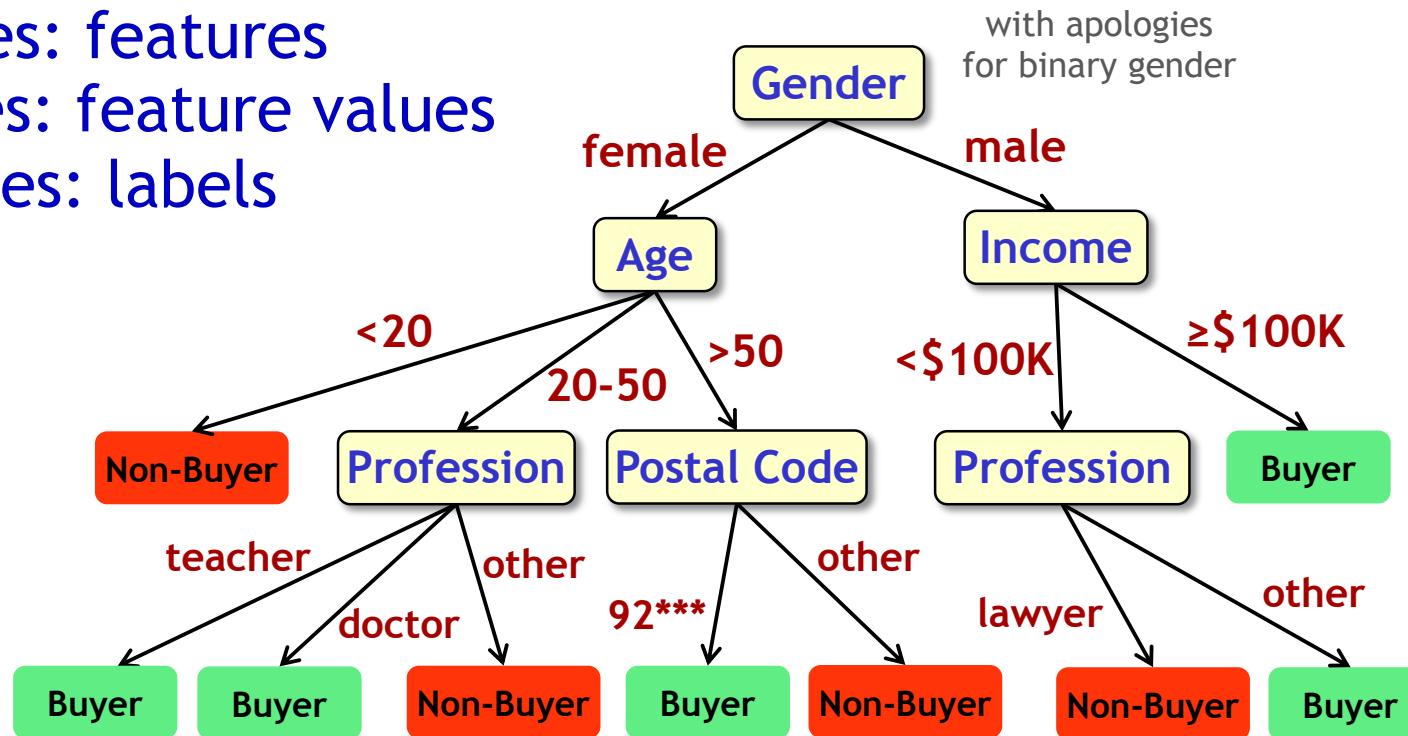
- Use the training data to construct a decision tree
- Use the decision tree to classify new data

# Decision trees

Nodes: features

Edges: feature values

Leaves: labels



New data item to classify:  
Navigate tree based on feature values

# Decision trees

Primary challenge is building good decision trees from training data

- Which features and feature values to use at each choice point
- HUGE number of possible trees even with small number of features and values

Common approach: “forest” of many trees, combine the results

- Still impossible to consider all trees

# Deep Neural Networks (DNNs)

## Neural Networks

- Machine learning method modeled loosely after connected neurons in brain
- Invented decades ago but not successful
- Recent resurgence enabled by:
  - Powerful computing that allows for many layers (making the network “deep”)
  - Massive data for effective training

# Deep Neural Networks (DNNs)

= Deep Learning

- Huge breakthrough in effectiveness and reach of machine learning
- Accurate predictions across many domains
- Big plus: Automatically identifies features in unstructured data  
(e.g., images, videos, text)

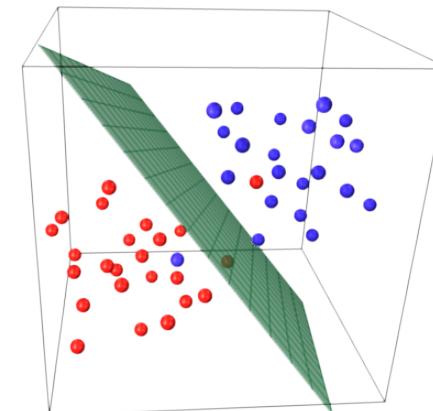
# Other common methods

## Logistic Regression

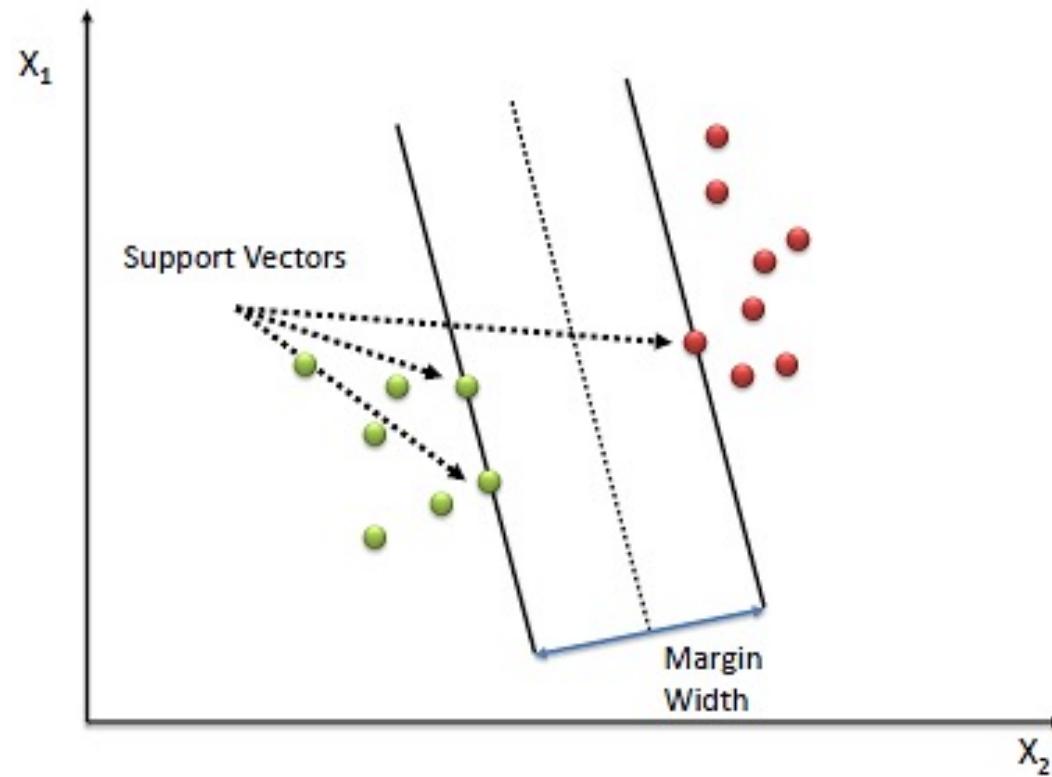
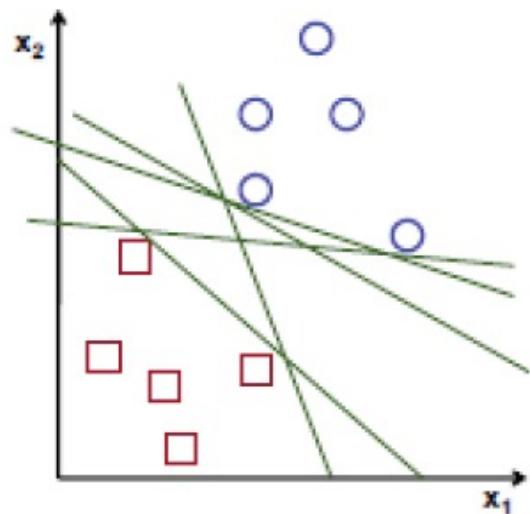
- Typically for two labels only (“binary classifier”)
- Recall regression model is function  $f$  from predictor values to numeric output value
- Labels  $L_1 + L_2$ , from training data obtain function:  
 $f(\text{feature-values}) = \text{probability of item having label } L_1$

## Support Vector Machine

- Also for binary classification
- Features = multidimensional space
- From training data SVM finds hyper-plane that best divides space according to labels



# Support vector machines (SVMs)



# Slide Credit

Deniz Akpinaroglu