

# Deep Forward Networks

## - Introduction

Wednesday  
8h00 – 8h45

- Continue from the neural network optimisation
  - backpropagation
  - minimize loss function
- Introduction of different layer type
  - conv layer
  - pooling
  - padding

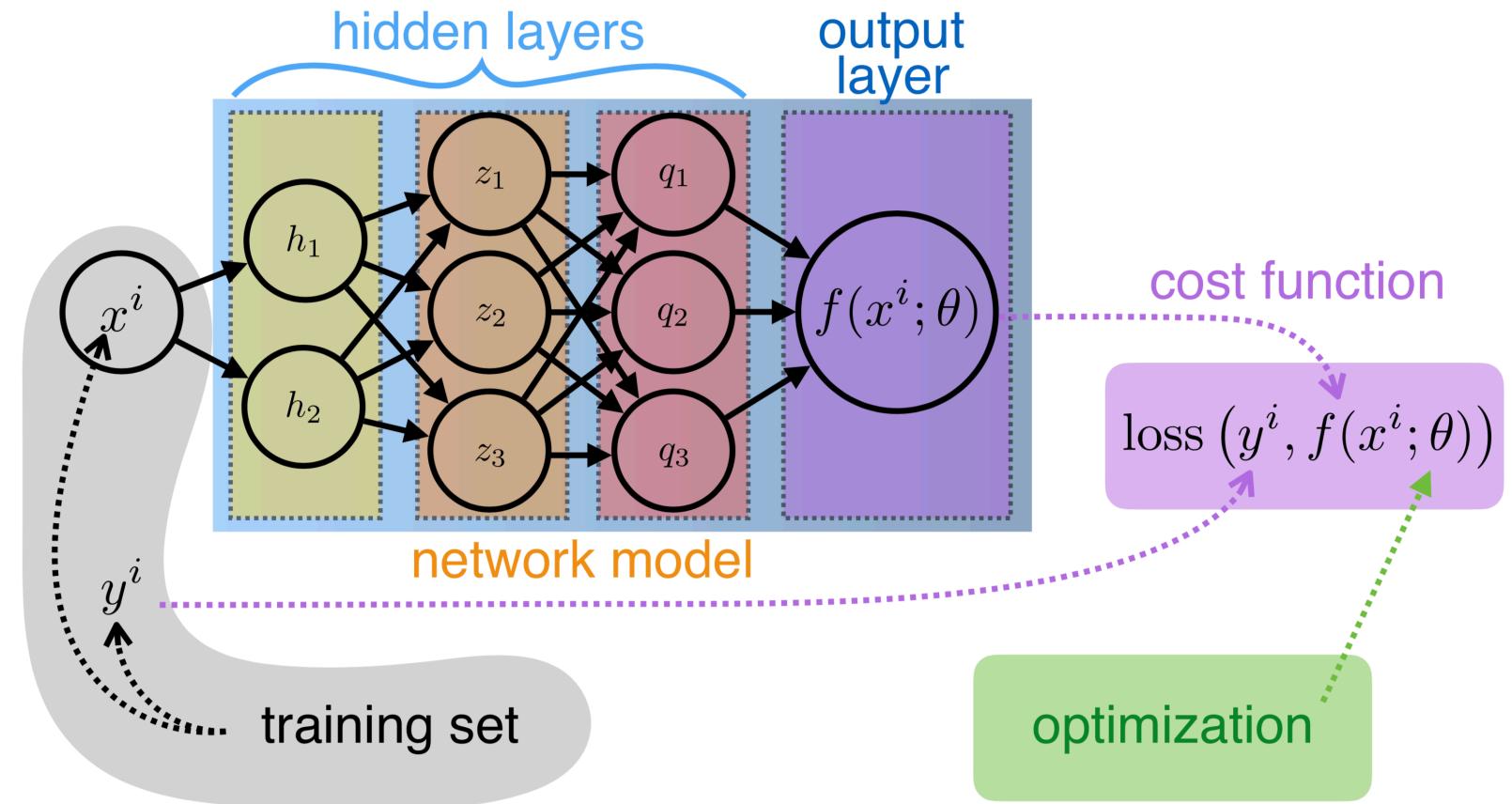
# Deploying a Neural Network

Given a task (in terms of **I/O mappings**), we need :

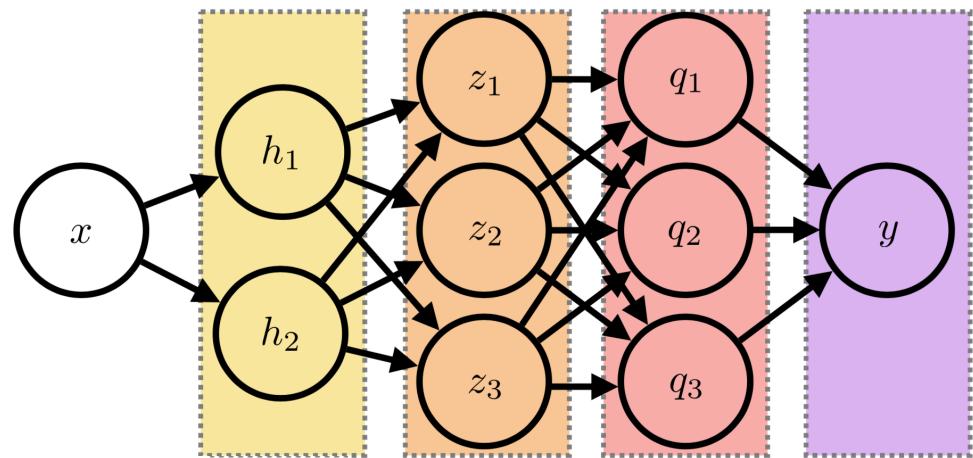
1) **Network model**

2) **Cost function**

3) **Optimization**



# Activation functions



$$y = f_4(q_1, q_2, q_3)$$

$$y = f_4(f_{3,1}(f_{2,1}(f_{1,1}(x), f_{1,2}(x)), \dots), \dots)$$

Hierarchical representation

$$\begin{aligned} h_1 &= f_{1,1}(x) \\ h_2 &= f_{1,2}(x) \end{aligned}$$

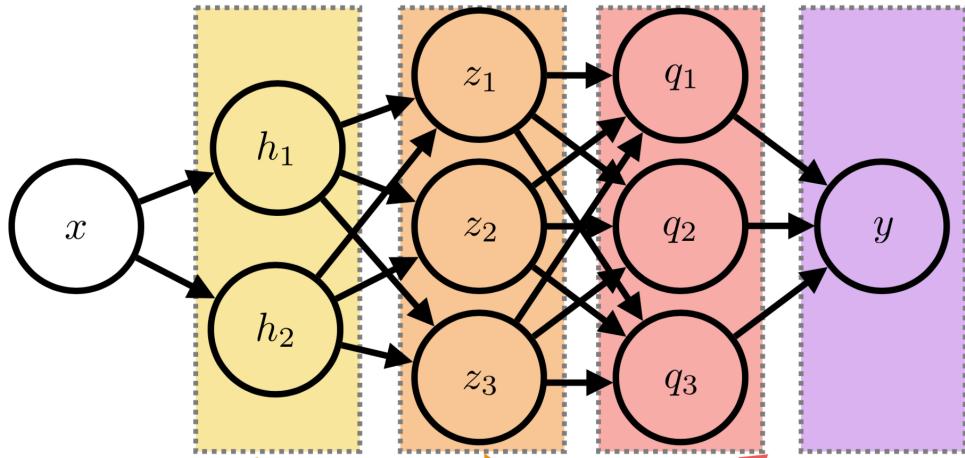
$$\begin{aligned} z_1 &= f_{2,1}(h_1, h_2) \\ z_2 &= f_{2,2}(h_1, h_2) \\ z_3 &= f_{2,3}(h_1, h_2) \end{aligned}$$

$$\begin{aligned} q_1 &= f_{3,1}(z_1, z_2, z_3) \\ q_2 &= f_{3,2}(z_1, z_2, z_3) \\ q_3 &= f_{3,3}(z_1, z_2, z_3) \end{aligned}$$

Bias (constant)

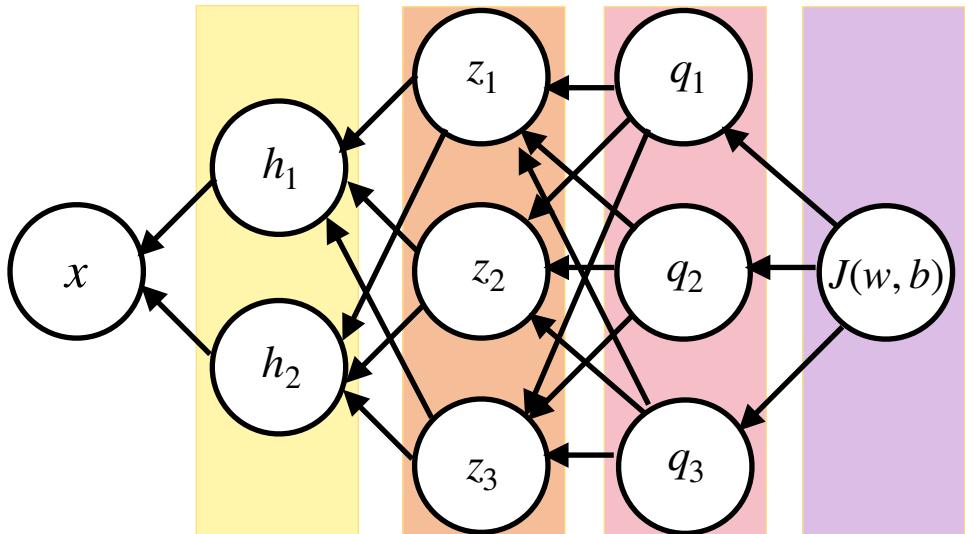
$$f_{2,2}(h_1, h_2) = w_1 h_1 + w_2 h_2 + b_{2,2}$$

# 3) Optimization



Forward propagation

$$\hat{y} = f_4(f_{3,1}(f_{2,1}(f_{1,1}(x), f_{1,2}(x)), \dots), \dots)$$



Back propagation

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

# Gradient descent

- Function :  $J(\theta)$
- Objective : minimize  $J(\theta)$
- Outline:
  - Start with an initial  $\theta$  (mostly we initialise with  $\theta = 0$  )
  - Keep updating  $\theta$  to reduce  $J(\theta)$  until we hopefully end up at a minimum

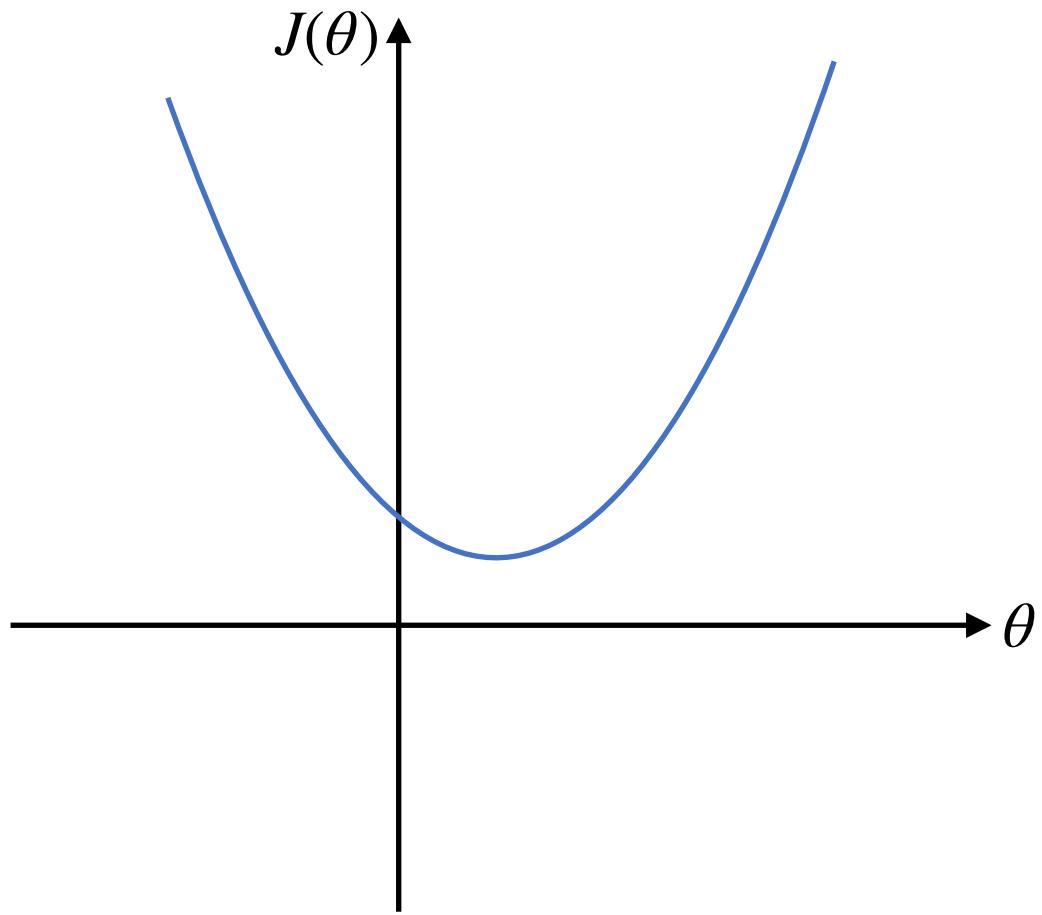
Criteria to update  $\theta$  :

$$\theta_{t+1} := \theta_t - \alpha \triangle J(\theta_t)$$

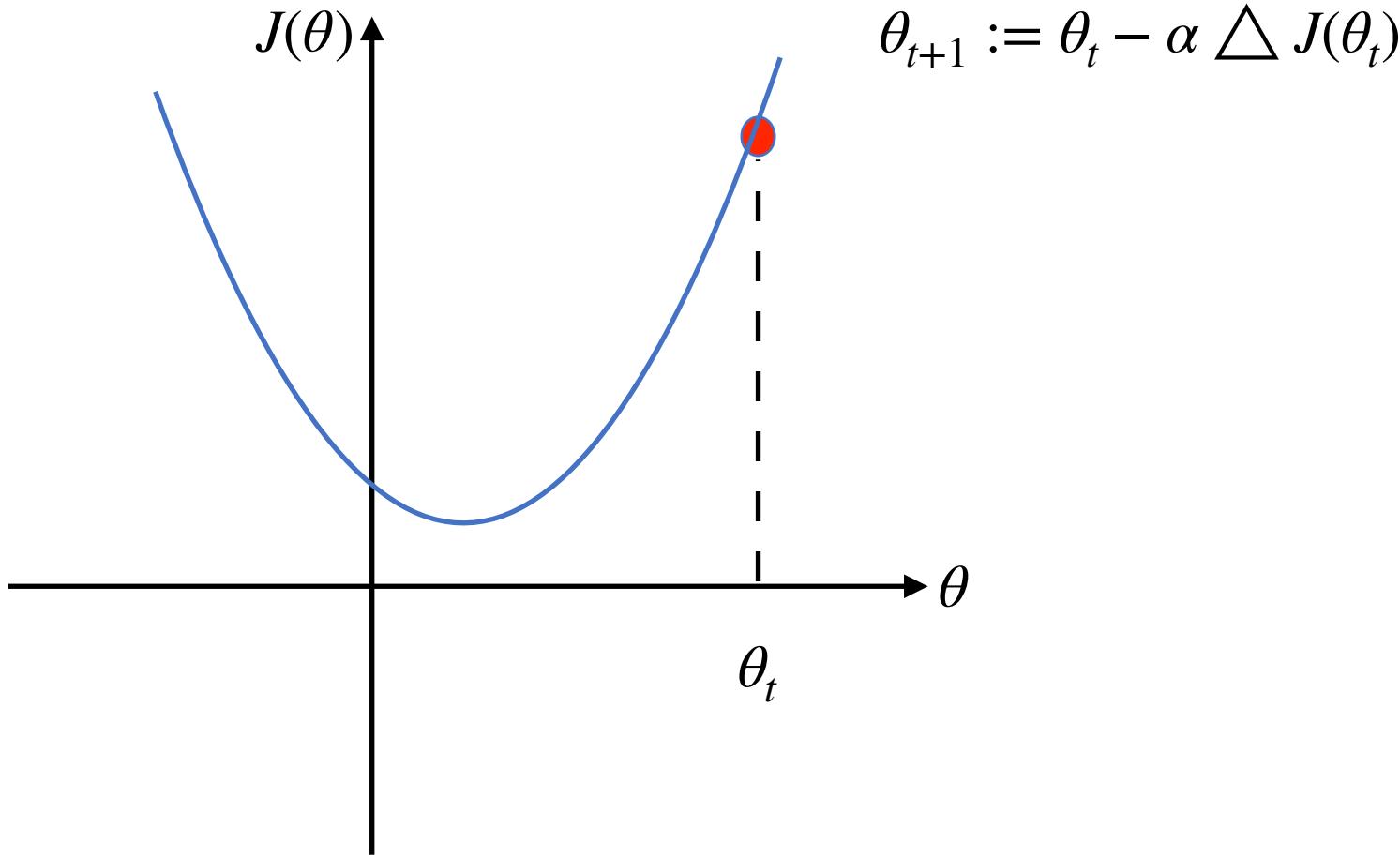
The diagram shows the update rule for gradient descent. A red arrow labeled "Learning rate (Positive constant)" points from below to the term  $\alpha$  in the equation. A blue arrow labeled "Gradient" points from below to the term  $\triangle J(\theta_t)$ .

The gradient vector can be interpreted as the "direction and rate of fastest increase".

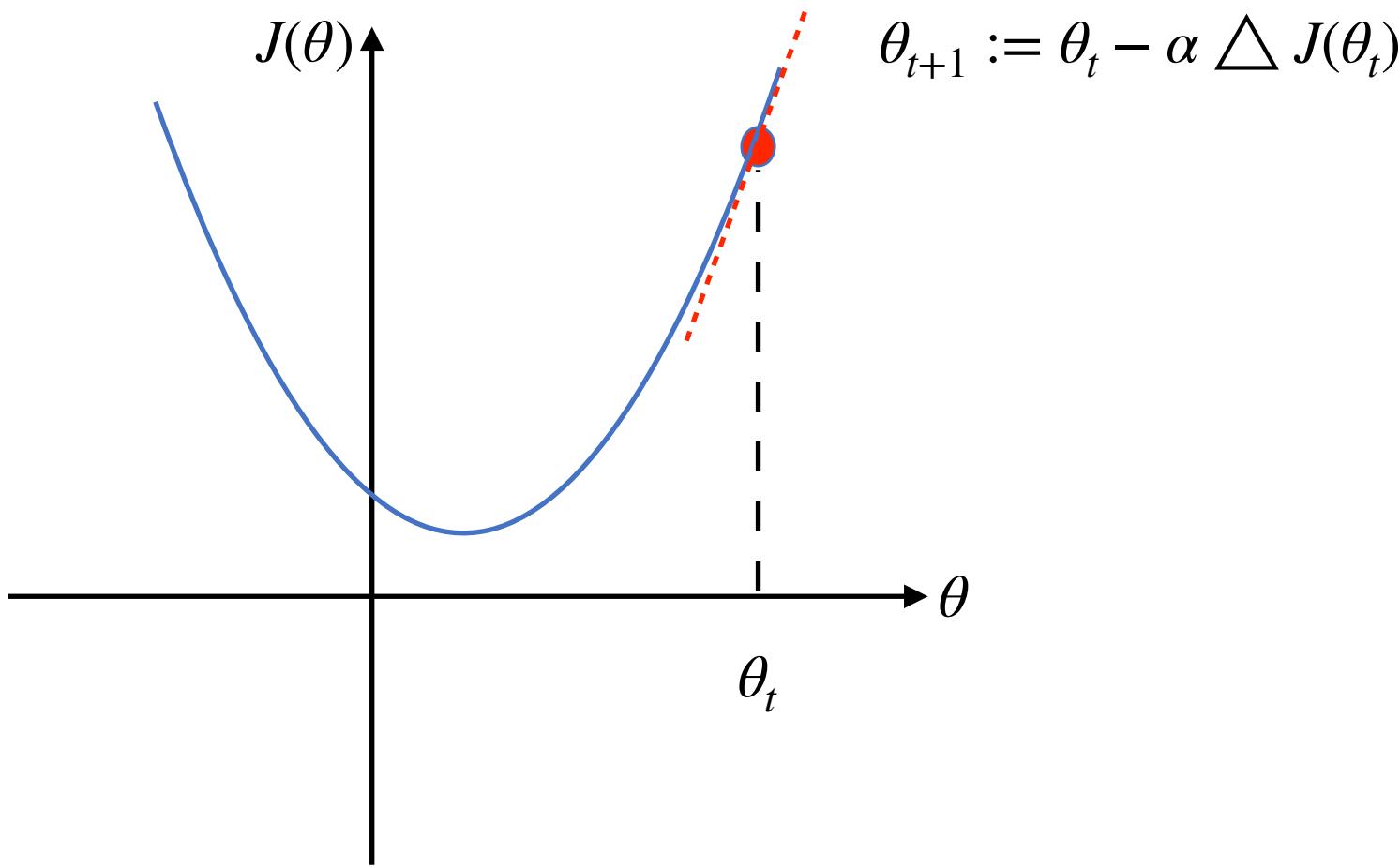
# Gradient descent



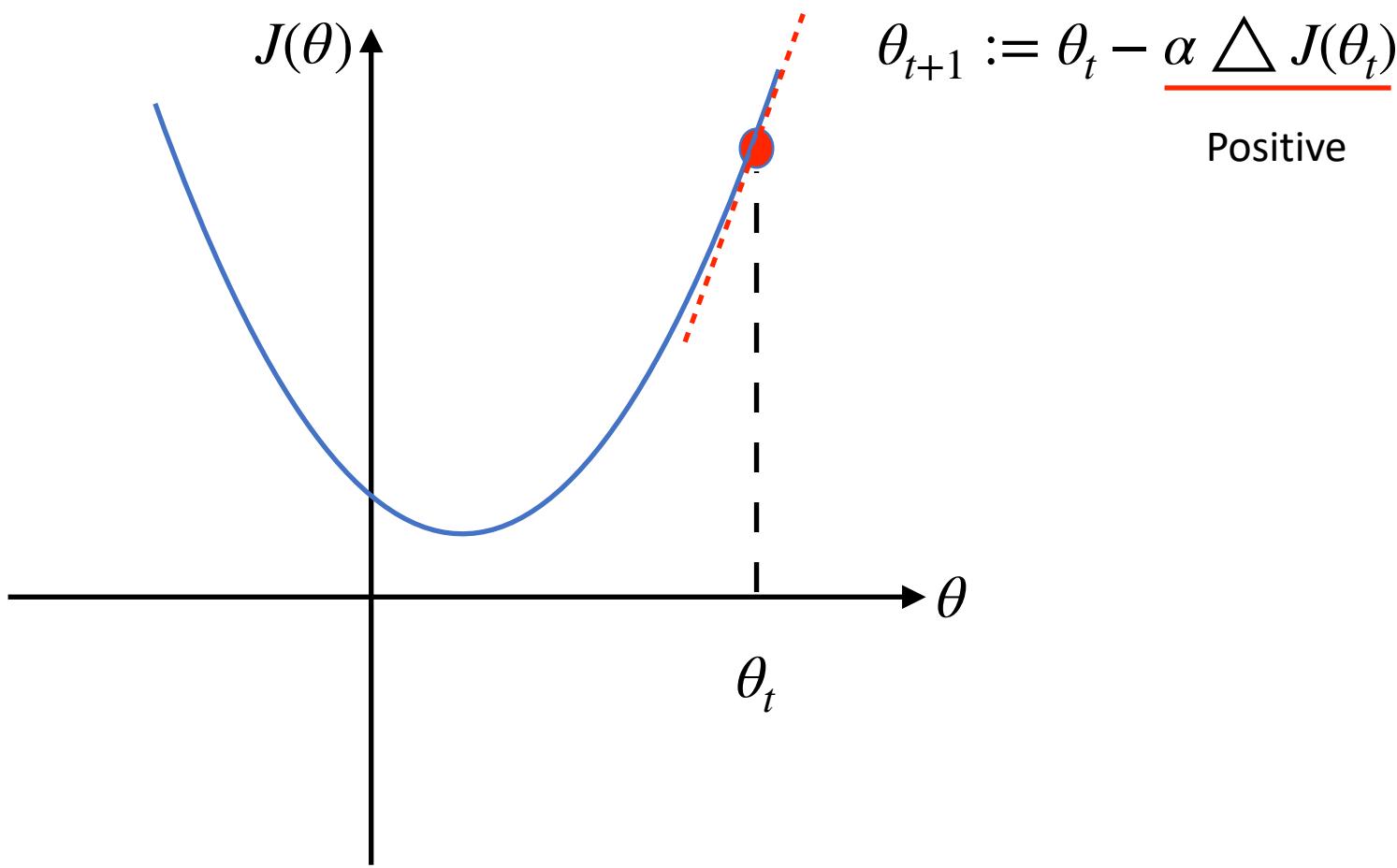
# Gradient descent



# Gradient descent



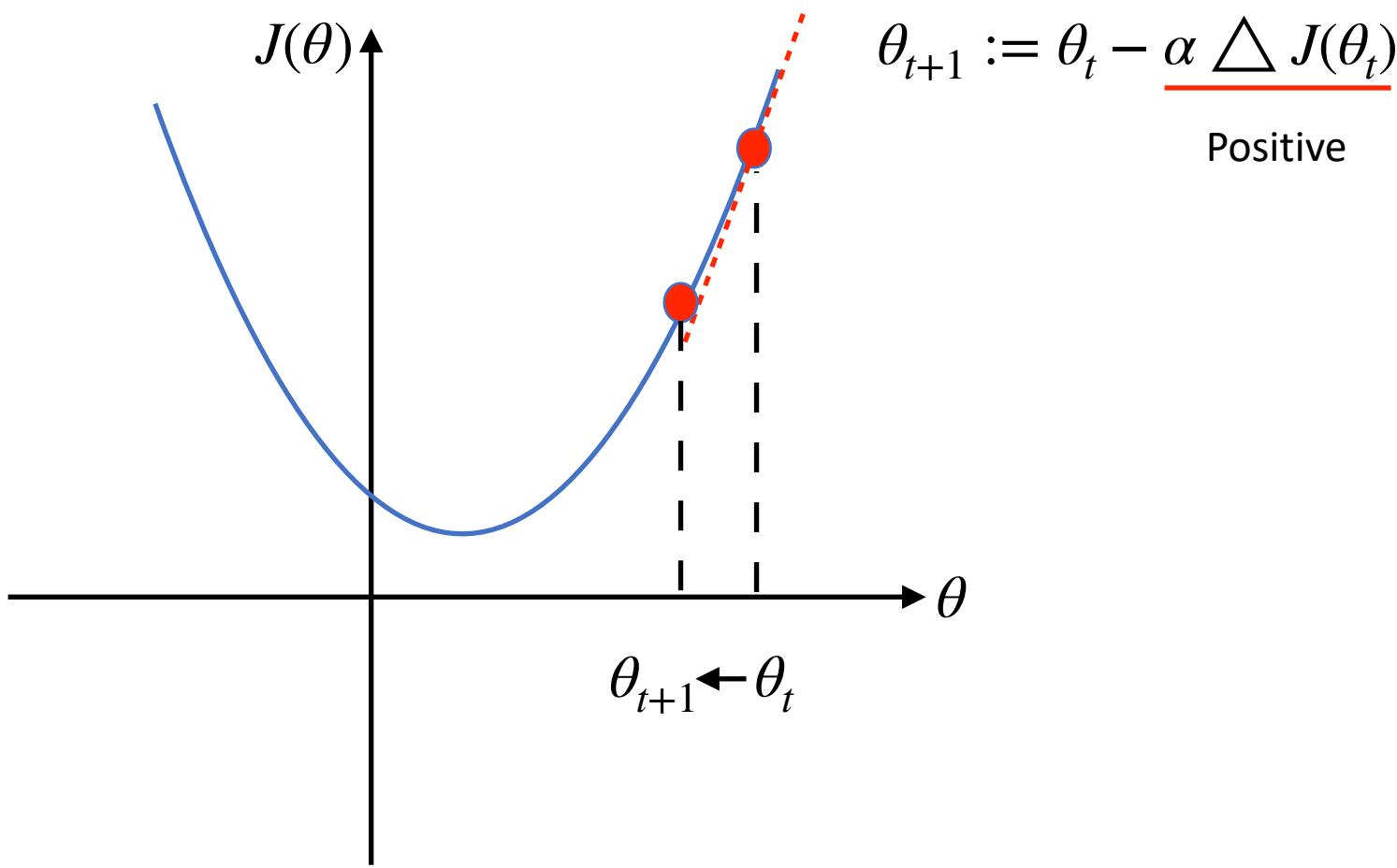
# Gradient descent



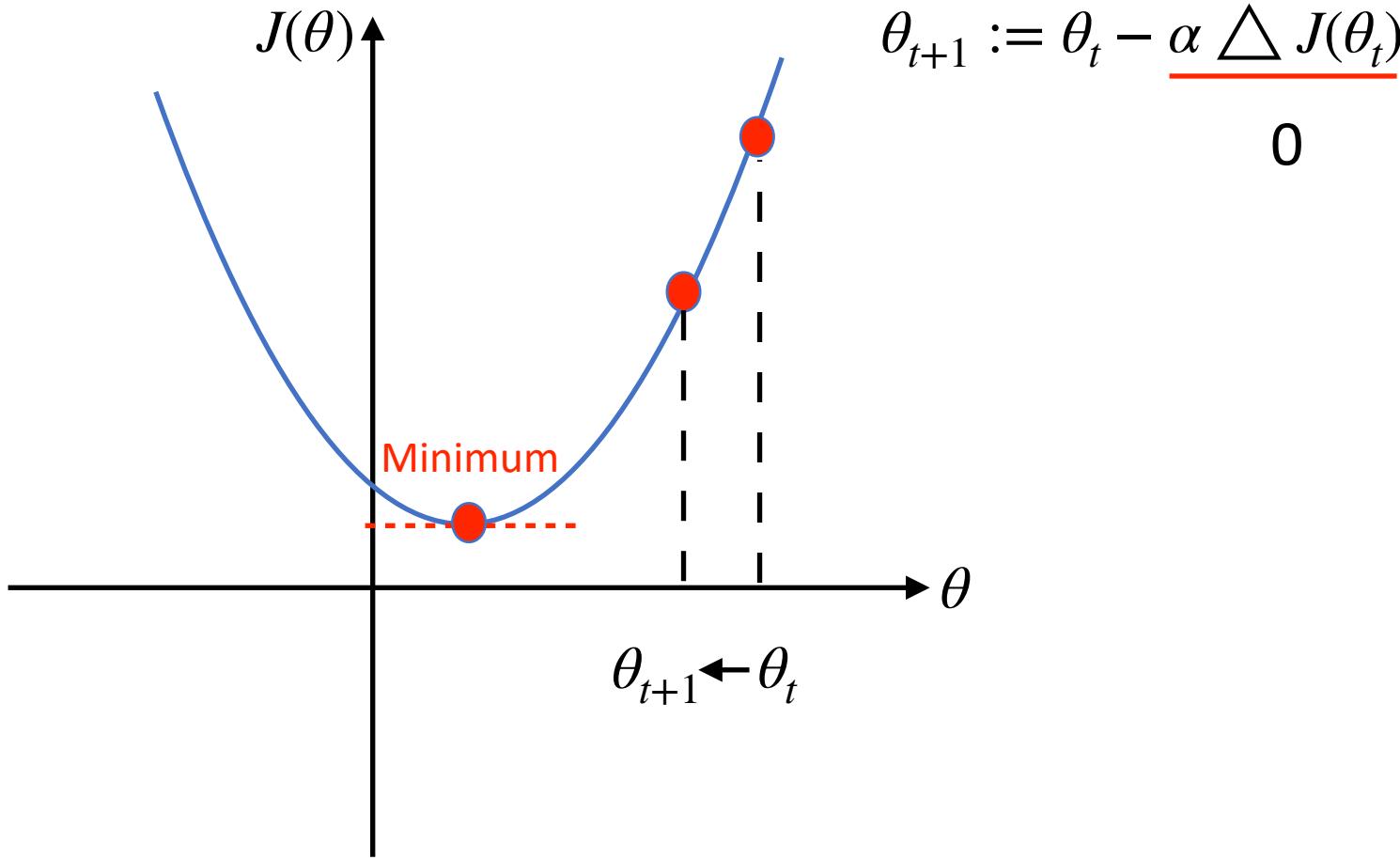
$$\theta_{t+1} := \theta_t - \alpha \triangle J(\theta_t)$$

Positive

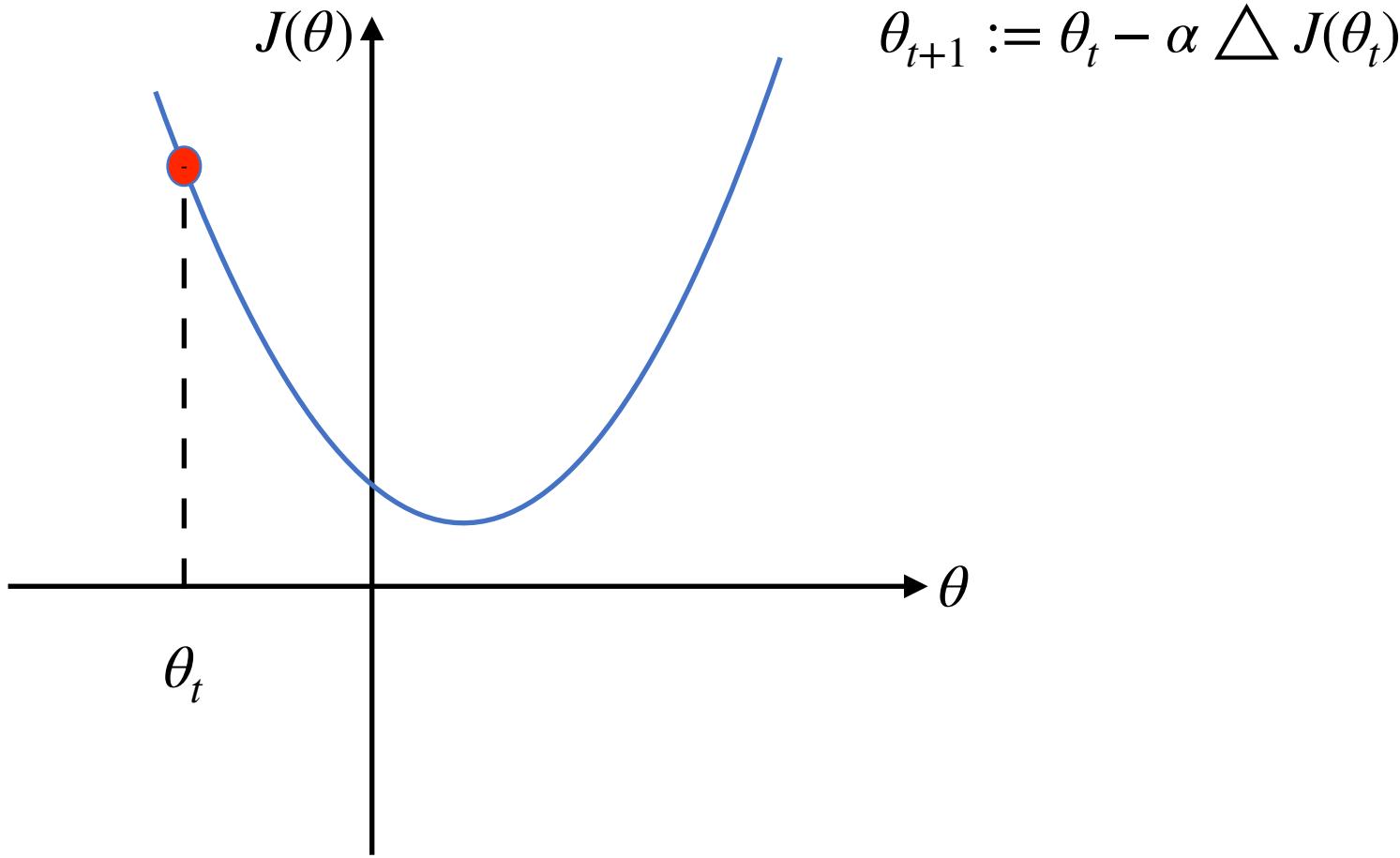
# Gradient descent



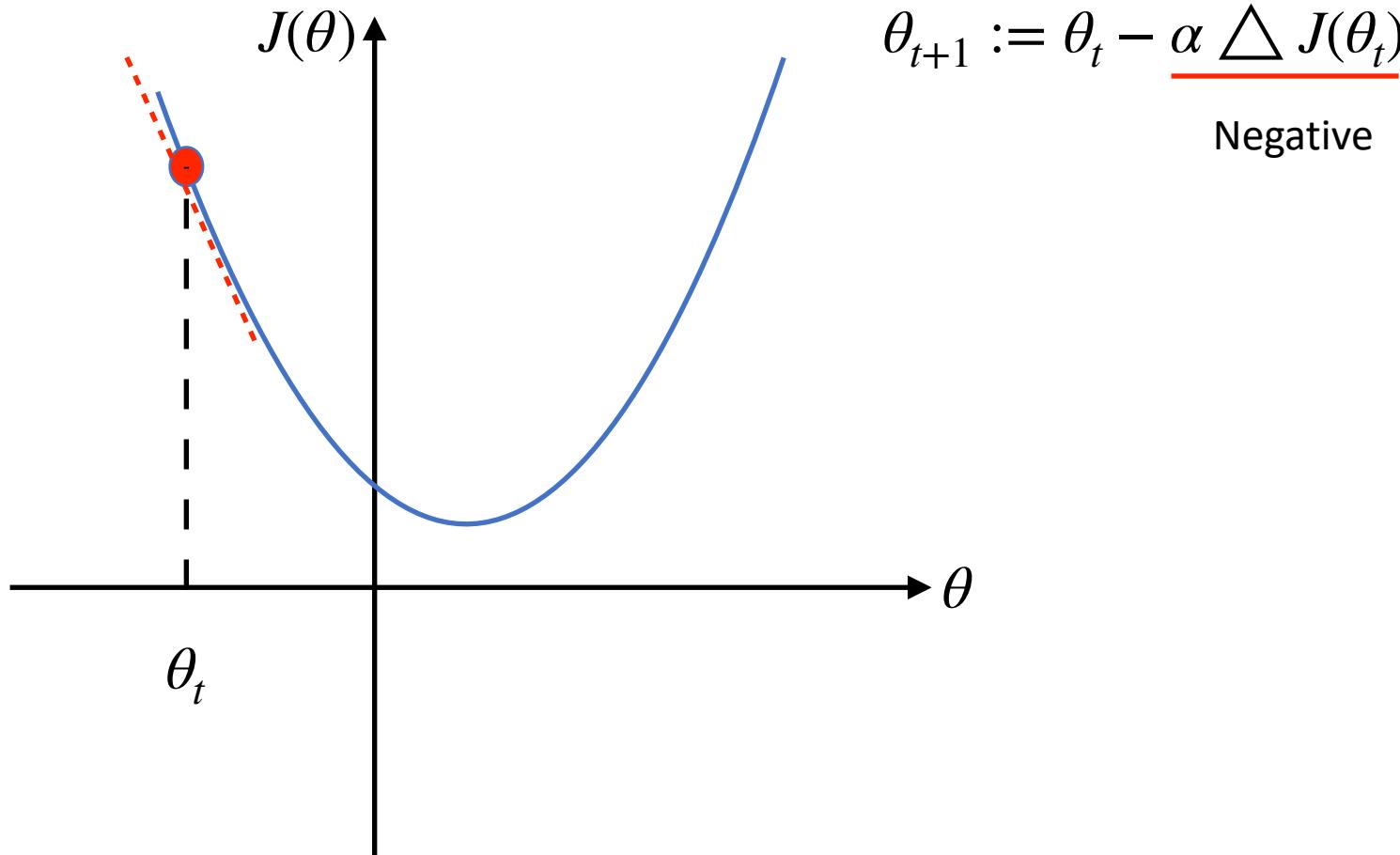
# Gradient descent



# Gradient descent



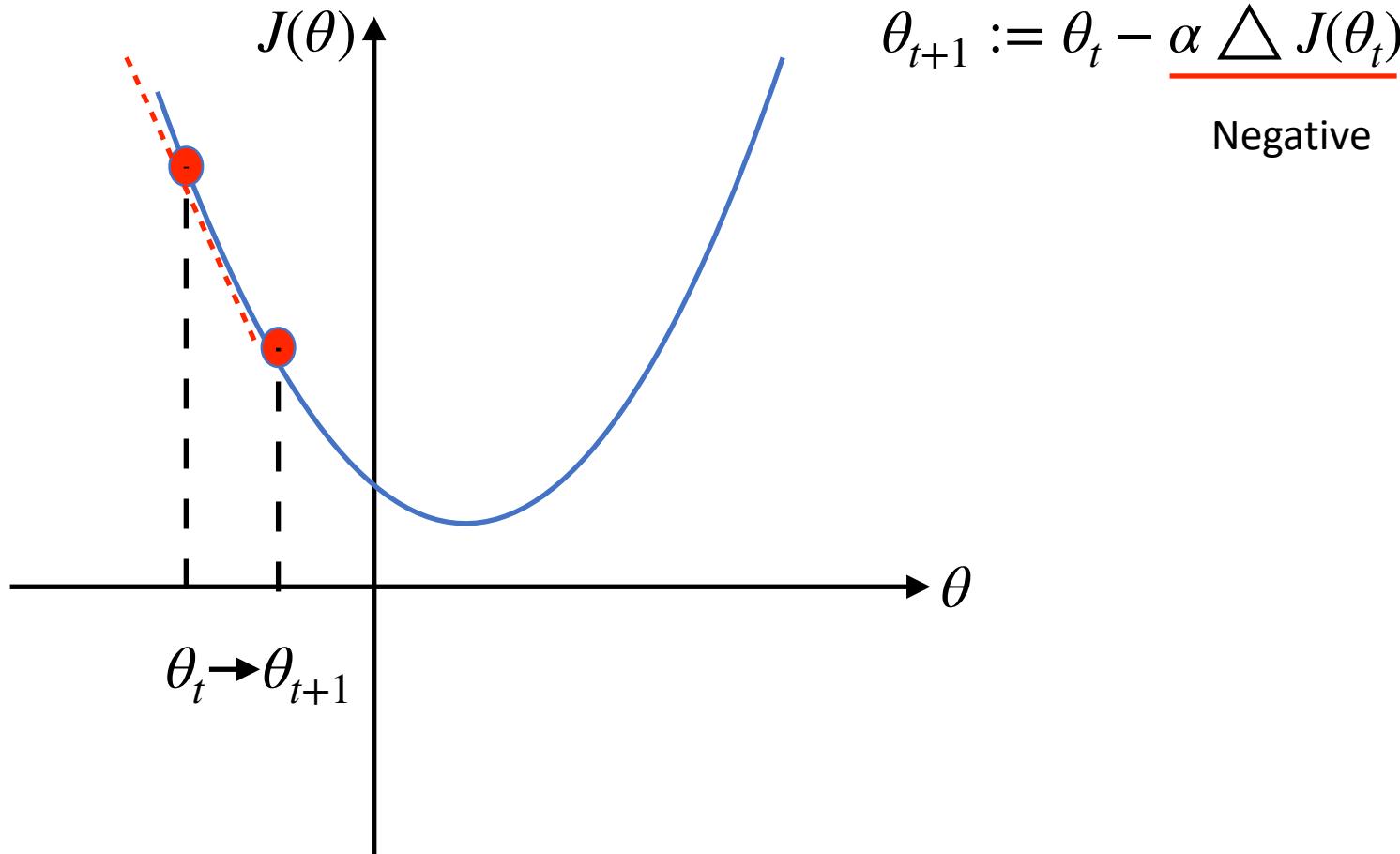
# Gradient descent



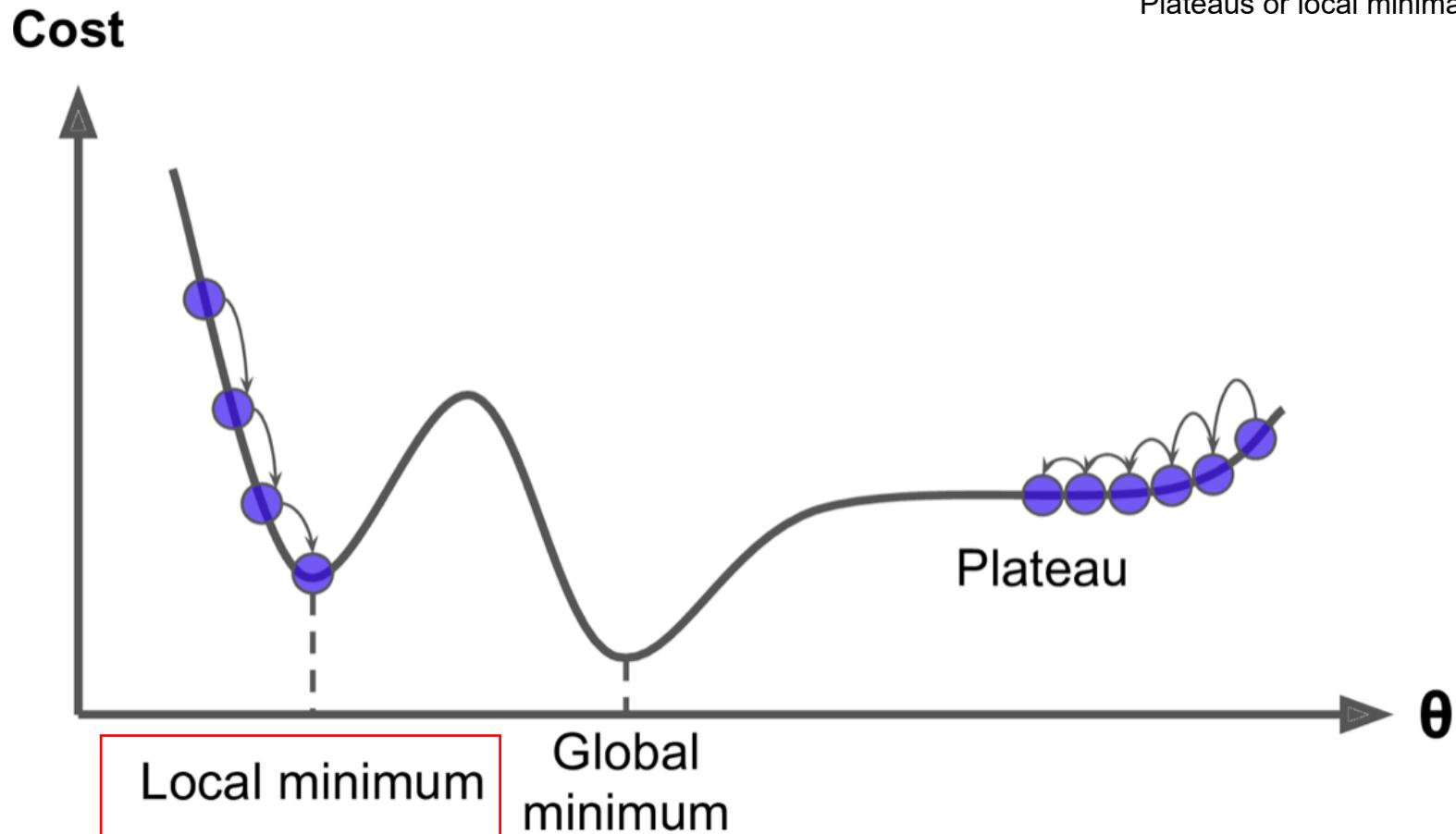
$$\theta_{t+1} := \theta_t - \underline{\alpha \triangle J(\theta_t)}$$

Negative

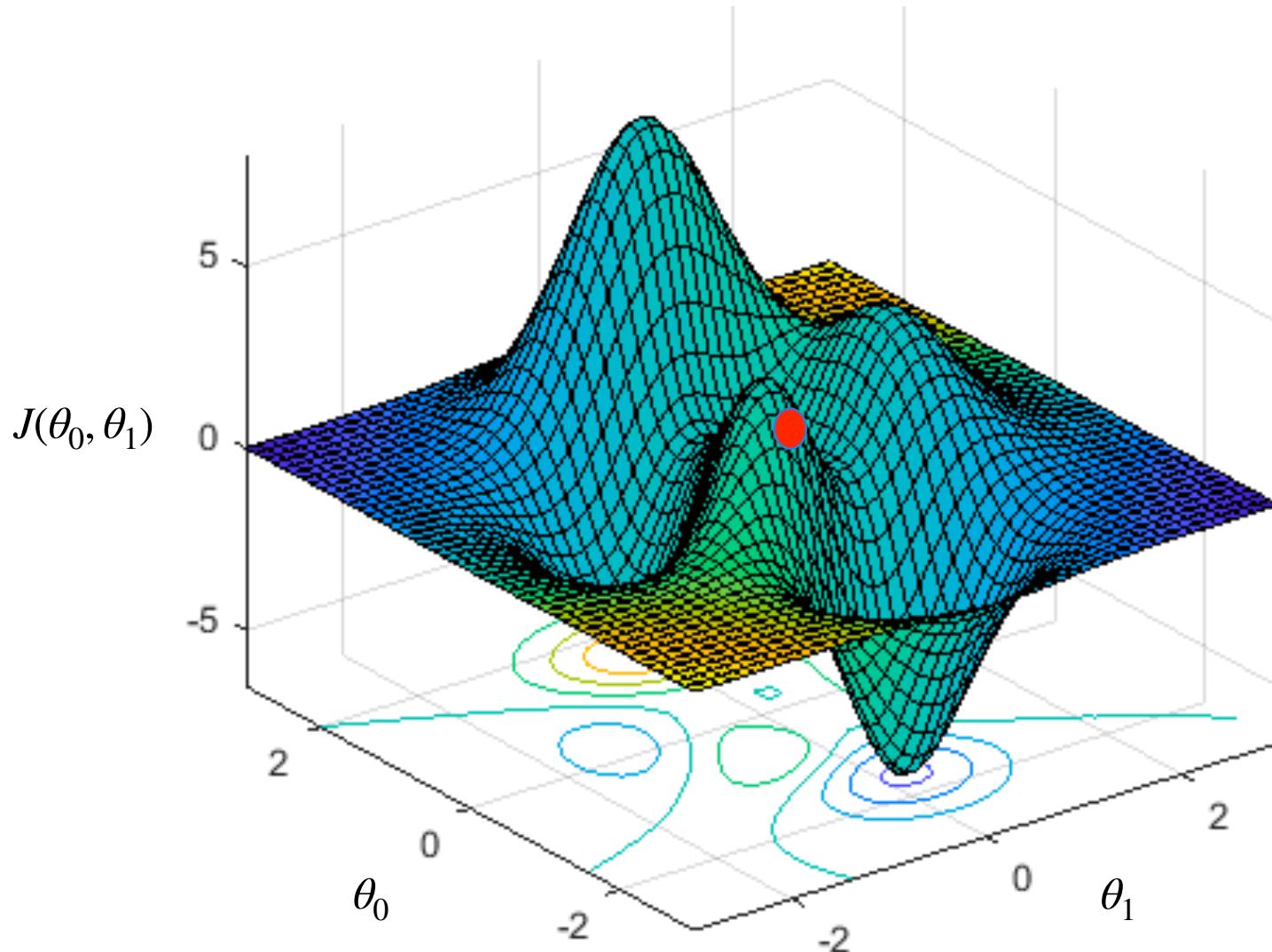
# Gradient descent



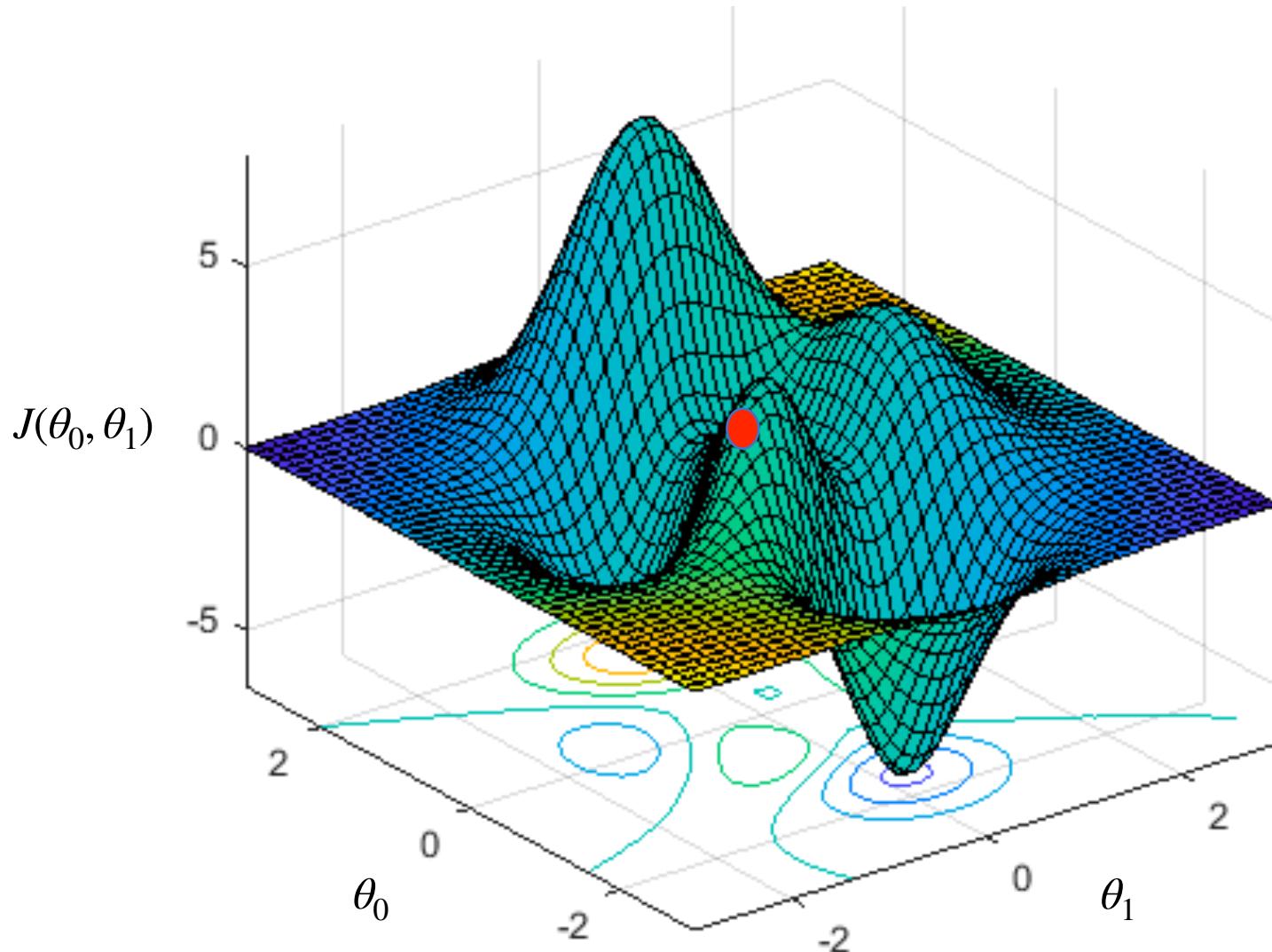
# Optimization pitfalls



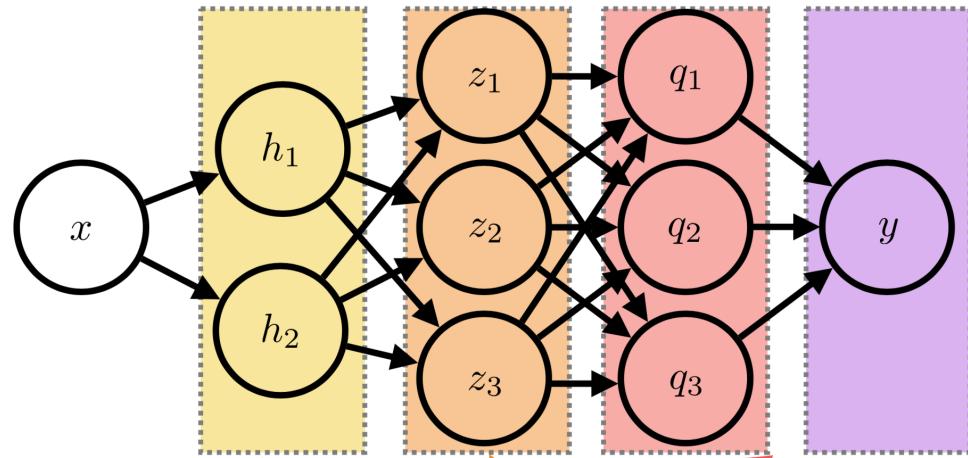
# Optimization pitfalls



# Optimization pitfalls

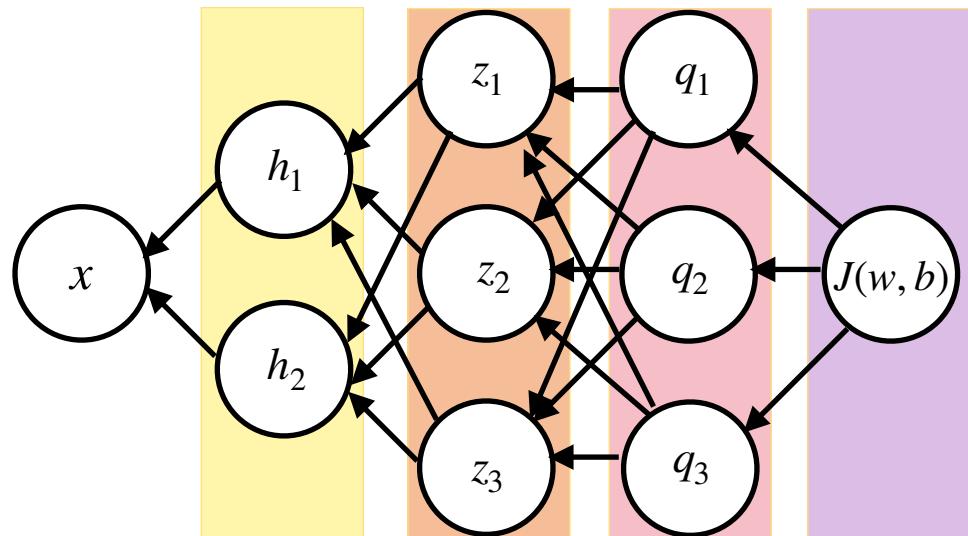


# Backpropagation



Forward propagation

$$\hat{y} = f_4(f_{3,1}(f_{2,1}(f_{1,1}(x), f_{1,2}(x)), \dots), \dots)$$



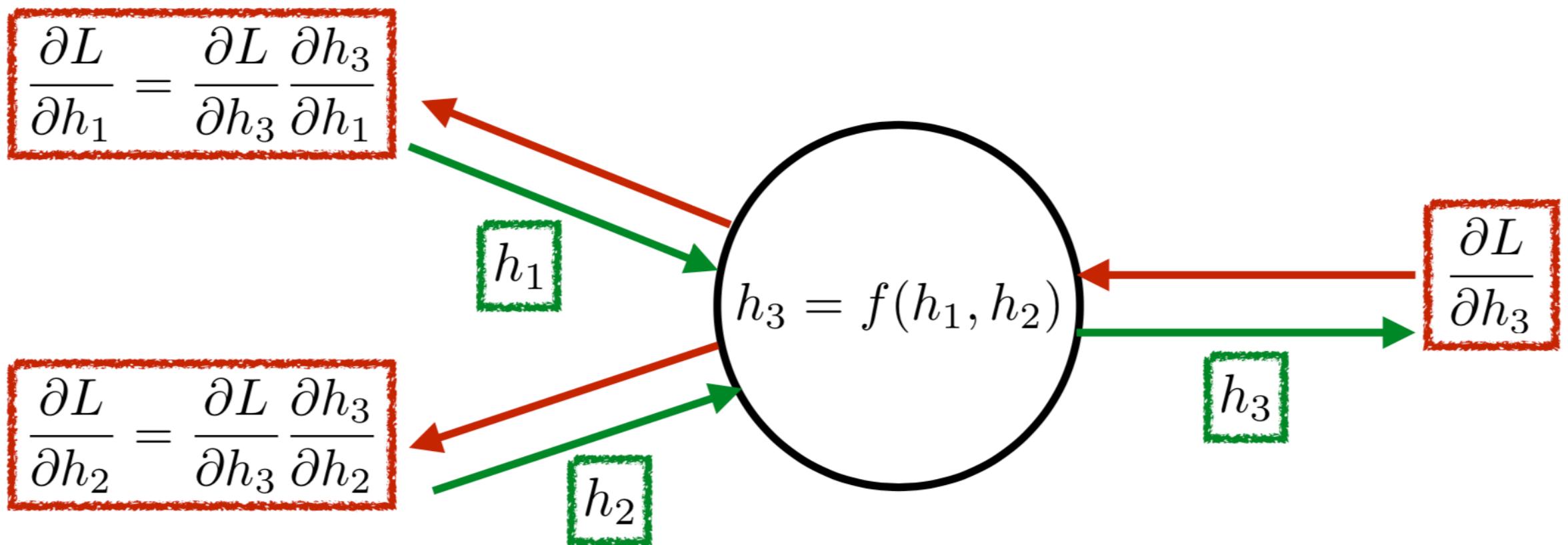
Back propagation

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

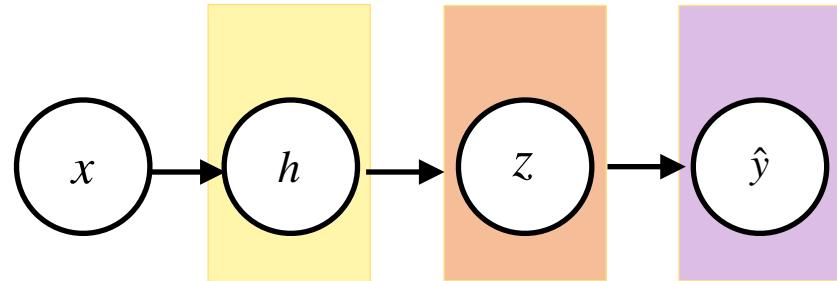
# Backpropagation

- Efficient implementation of the **chain-rule** to compute derivatives with respect to network weights

Calculating gradient and optimizing weights with chain rule.



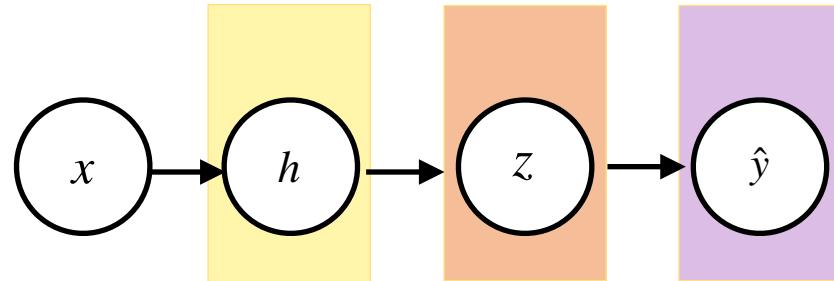
# Backpropagation



$$h = f_1(x) \quad z = f_2(h) \quad \hat{y} = f_3(z)$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

# Backpropagation

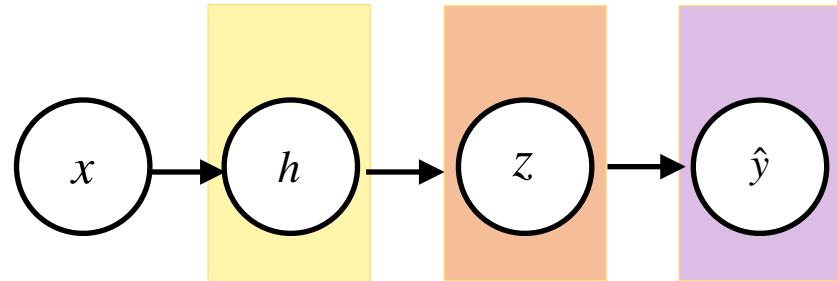


$$\begin{aligned} h &= f_1(x) & z &= f_2(h) & \hat{y} &= f_3(z) \\ &= w_1x + b_1 & = w_2x + b_2 & = w_3x + b_3 \end{aligned}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$f = wx + b$$

# Backpropagation



$$\begin{aligned} h &= f_1(x) & z &= f_2(h) & \hat{y} &= f_3(z) \\ &= w_1x + b_1 & = w_2x + b_2 & = w_3x + b_3 \end{aligned}$$

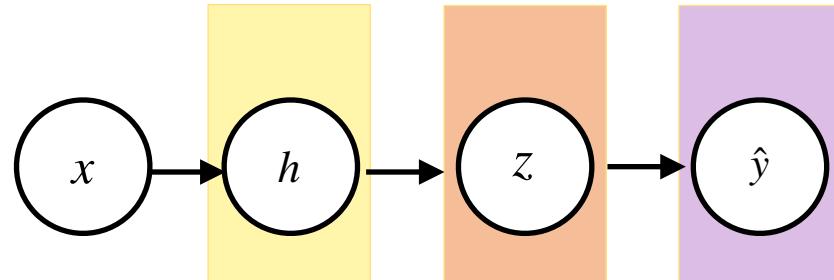
$$f = wx + b$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$J(w_1, b_1, w_2, b_2, w_3, b_3)$$

$$\Delta J(w_1, b_1, w_2, b_2, w_3, b_3) = \left( \frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial b_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial b_2}, \frac{\partial J}{\partial w_3}, \frac{\partial J}{\partial b_3} \right)$$

# Backpropagation



$$\begin{aligned} h &= f_1(x) & z &= f_2(h) & \hat{y} &= f_3(z) \\ &= w_1x + b_1 & = w_2x + b_2 & = w_3x + b_3 \end{aligned}$$

$$f = wx + b$$

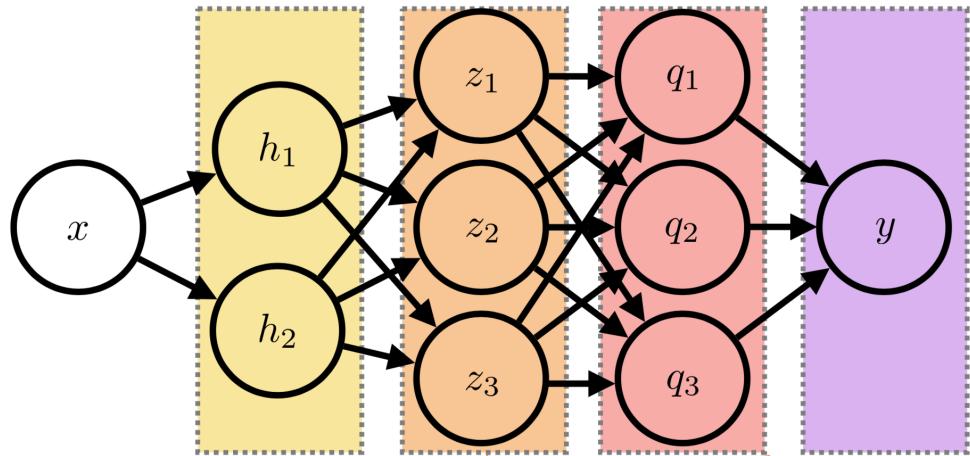
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$J(w_1, b_1, w_2, b_2, w_3, b_3)$$

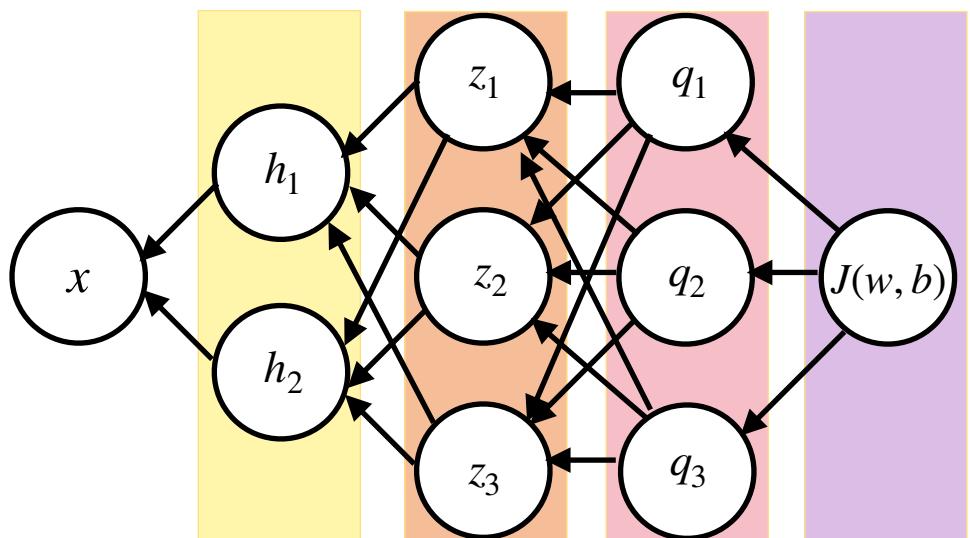
$$\Delta J(w_1, b_1, w_2, b_2, w_3, b_3) = \left( \frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial b_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial b_2}, \frac{\partial J}{\partial w_3}, \frac{\partial J}{\partial b_3} \right)$$

$$\frac{\partial J}{\partial w_1} = \left( \frac{\partial J}{\partial f_3} \frac{\partial f_3}{\partial f_2} \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial w_1} \right)$$

# Training iteration



1. Forward propagation, get output
2. Compute loss
3. Back propagation, update parameters
4. Forward propagation again with updated parameters, get new output
5. Repeat step 1-4 until converge



# Loss function choice

Loss function choice depends on what you want to achieve with the model.

- Choice determined by the **output representation**
  - Probability vector (**classification**) : Cross-entropy

$$\hat{y} = \sigma(w^\top h + b)$$

$$p(y|\hat{y}) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

$$L(\hat{y}, y) = -\log p(y|\hat{y}) = -\left( y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}) \right) \quad (\text{binary classification})$$

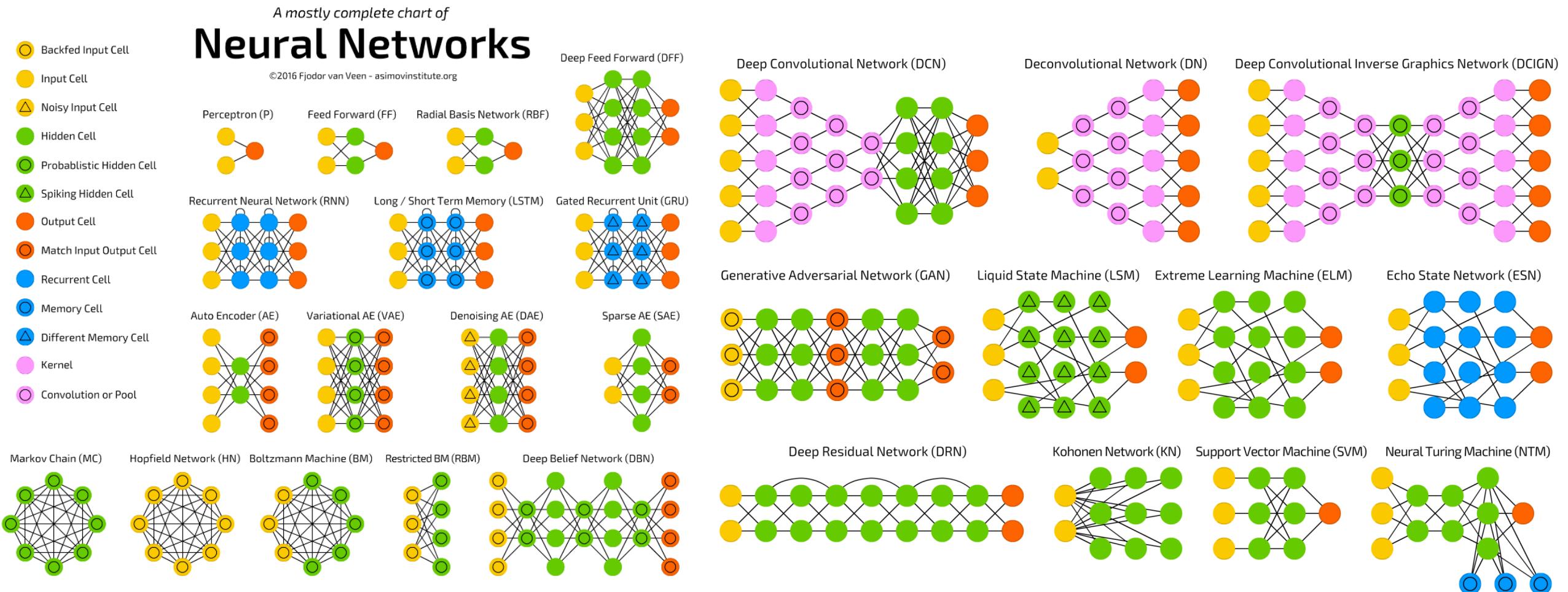
- Mean estimate (**regression**) : Mean Squared Error, L2 loss

$$\hat{y} = W^\top h + b$$

$$p(y|\hat{y}) = N(y; \hat{y})$$

$$L_2(\hat{y}, y) = -\log p(y|\hat{y}) = \sum_{i=0}^m (y^i - \hat{y}^i)^2$$

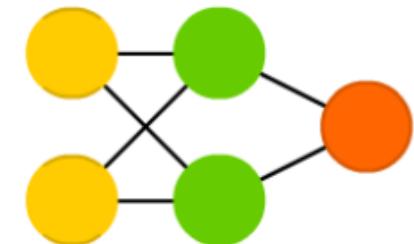
# 1) Network Model



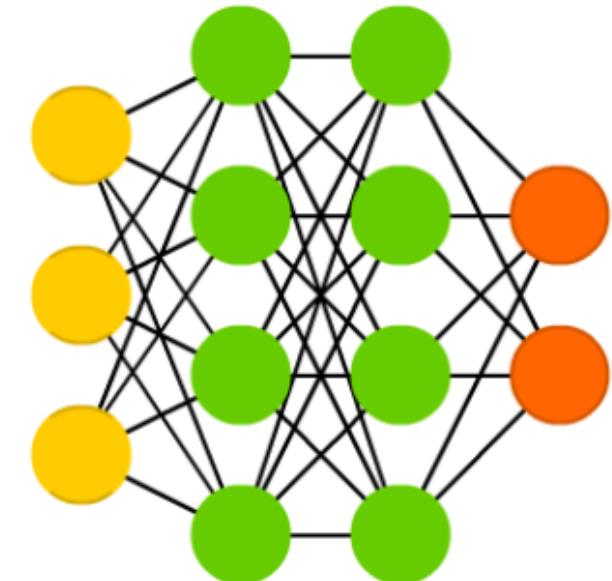
# (Deep) Feedforward NN (DFF)

- the **simplest type** of neural network
- All units are **fully connected**
- information flows from **input** to **output** layer **without back loops**
- The first single-neuron network was proposed already in 1958 by AI pioneer Frank Rosenblatt
- Deep for “more than 1 **hidden layer**”

Feed Forward (FF)



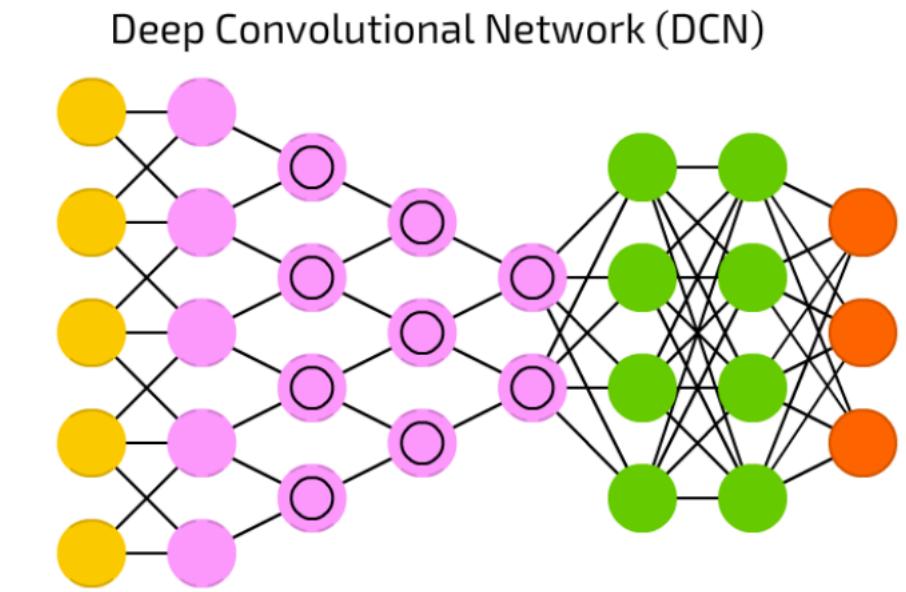
Deep Feed Forward (DFF)



1 hidden layer networks are simply called "feedforward NN", with more than one hidden layer we call it a "deep feedforward NN"

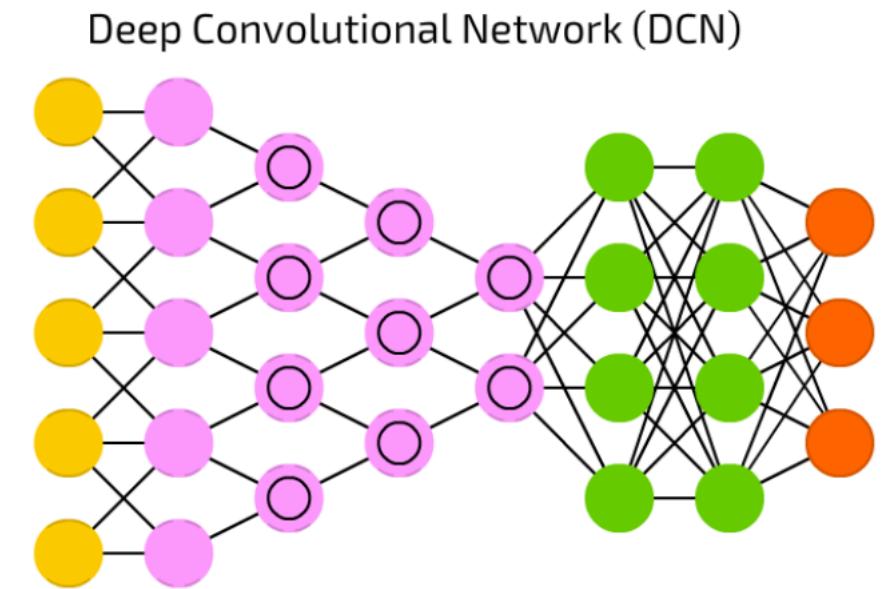
# Convolutional Neural Networks (CNN)

- inspired by the organization of the animal visual cortex
- Kernel and convolution or pool cells used to process and simplify input data
  - Weight sharing between local regions
- well suited for computer vision tasks
  - Image classification
  - Object detection



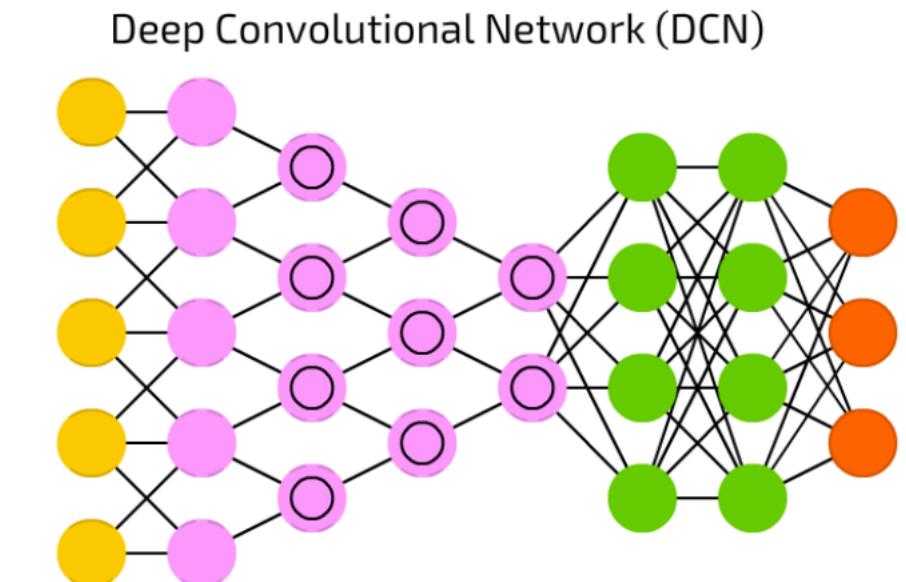
# Convolutional Neural Networks (CNN)

- inspired by the organization of the animal visual cortex
- Kernel and convolution or pool cells used to process and simplify input data
  - Weight sharing between *local regions*
- well suited for computer vision tasks
  - Image classification
  - Object detection



# Convolutional Neural Networks (CNN)

- inspired by the organization of the animal visual cortex
- Kernel and convolution or pool cells used to process and simplify input data
  - Weight sharing between *local regions*
- well suited for computer vision tasks
  - Image classification
  - Object detection



In deep learning, a convolutional neural network (CNN, or ConvNet) is a class of deep neural networks, most commonly applied to analyzing visual imagery. Multilayer perceptrons usually mean fully connected networks, that is, each neuron in one layer is connected to all neurons in the next layer. The "fully-connectedness" of these networks makes them prone to overfitting data. Typical ways of regularization include adding some form of magnitude measurement of weights to the loss function. CNNs take a different approach towards regularization: they take advantage of the hierarchical pattern in data and assemble more complex patterns using smaller and simpler patterns. Therefore, on the scale of connectedness and complexity, CNNs are on the lower extreme.

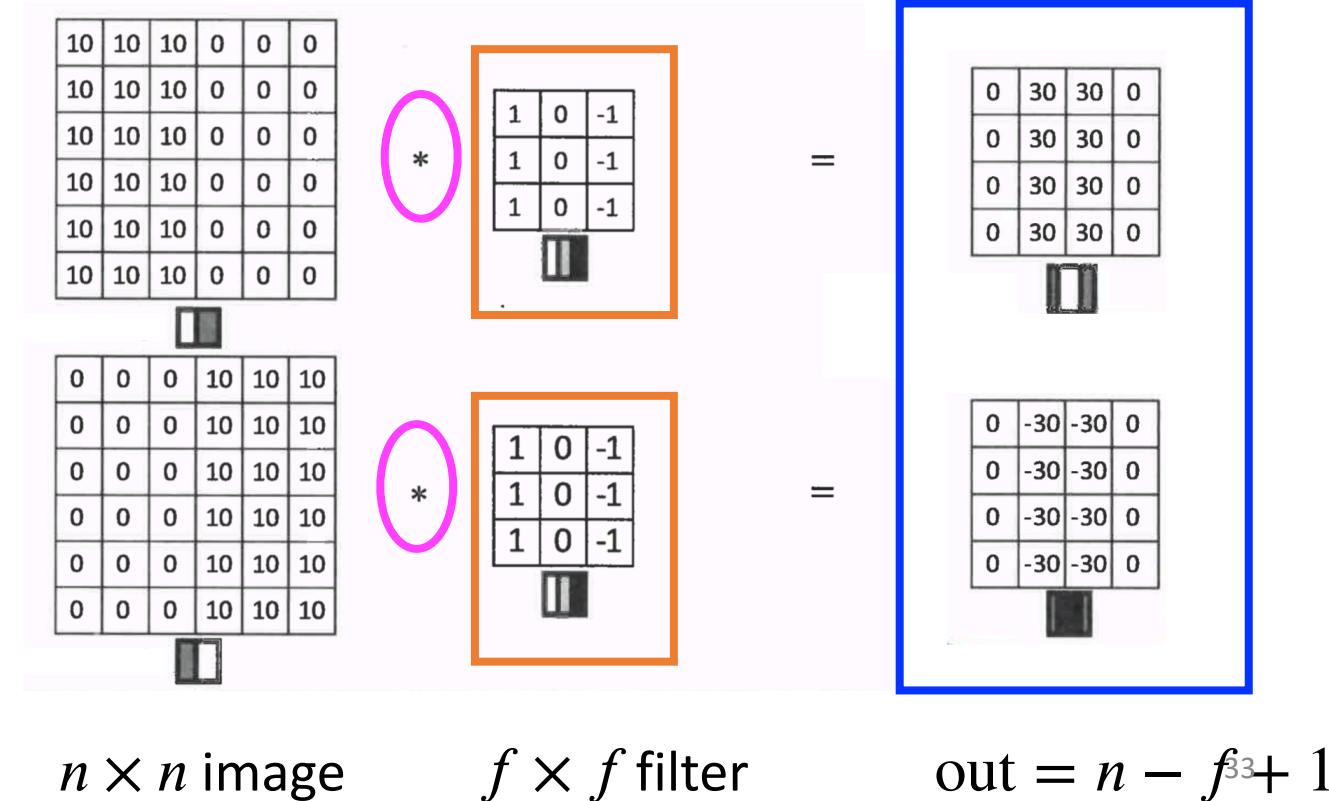
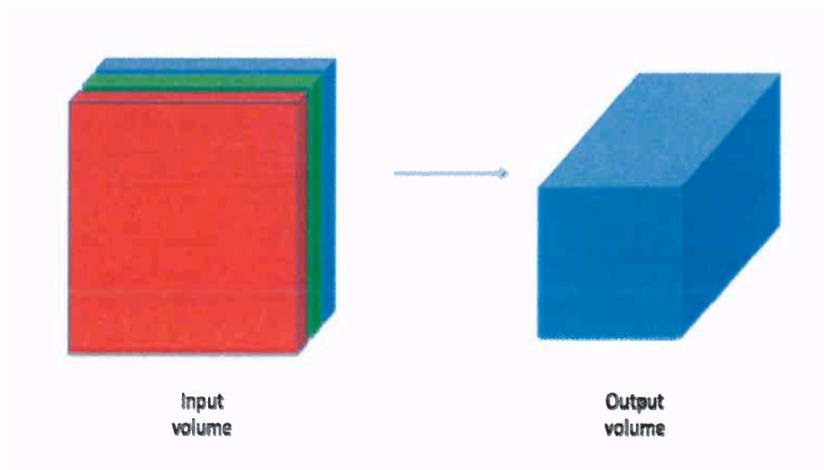
# Types of layer in CNN

- Convolution (CONV)
- Pooling (POOL)
- Fully connected (FC)
- *In feature learning tasks, usually multiple CONV layers followed by a POOL layer, and FC layers in the last few layers*

# Convolution Layer (CONV)

- **Convolution** transforms an input volume into an **output volume** of different size, also called feature map
- **Filter kernels** are used to detect features (for example, edge detection in 1<sup>st</sup> hidden layer)

The objective of the Convolution Operation is to extract the high-level features such as edges, from the input image.



# Padding

*Adds zeros around the border of an image*

This can be used if the dimension of the input data should remain the same.

0	0	0	0	0	0	0
0	60	113	56	139	85	0
0	73	121	54	84	128	0
0	131	99	70	129	127	0
0	80	57	115	69	134	0
0	104	126	123	95	130	0
0	0	0	0	0	0	0

Kernel

0	-1	0
-1	5	-1
0	-1	0

114				

# Padding

*Adds zeros around the border of an image*

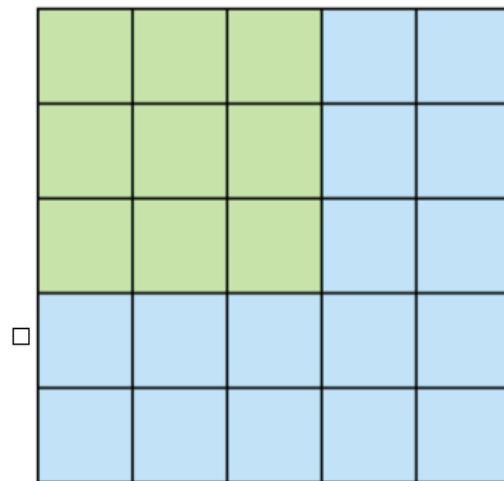
## Use cases :

- Keeps more information at the border of an image
- Allows to *use a CONV layer without shrinking* the height and width of the volumes (important for deeper networks)

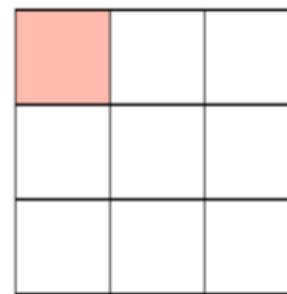
# Strided convolutions

The stride defines the step size of the kernel when traversing the image. While its default is usually 1, we can use a stride of 2 for downsampling an image similar to MaxPooling.

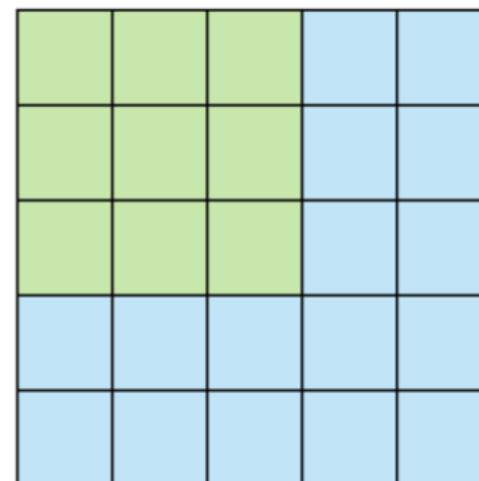
*By how much you move the filter*



Stride 1



Feature Map



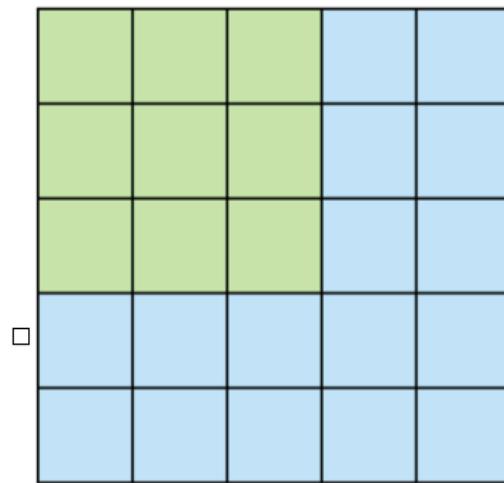
Stride 2



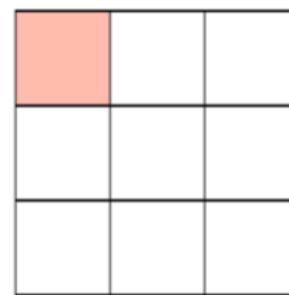
*Increasing stride from 1 to 2*

# Strided convolutions

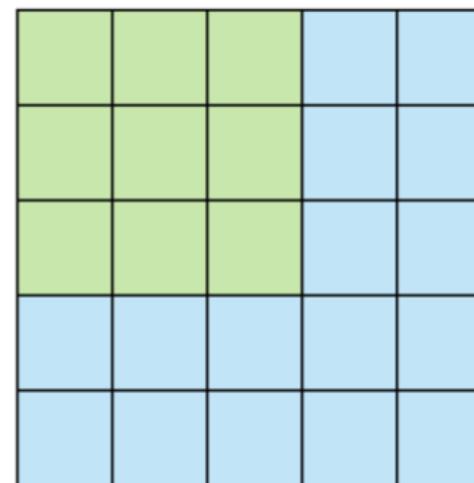
*By how much you move the filter*



Stride 1



Feature Map

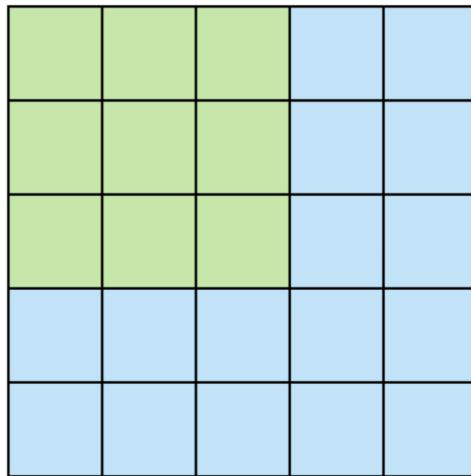


Stride 2

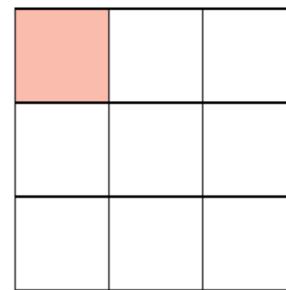
Increasing stride from 1 to 2

# Strided convolutions

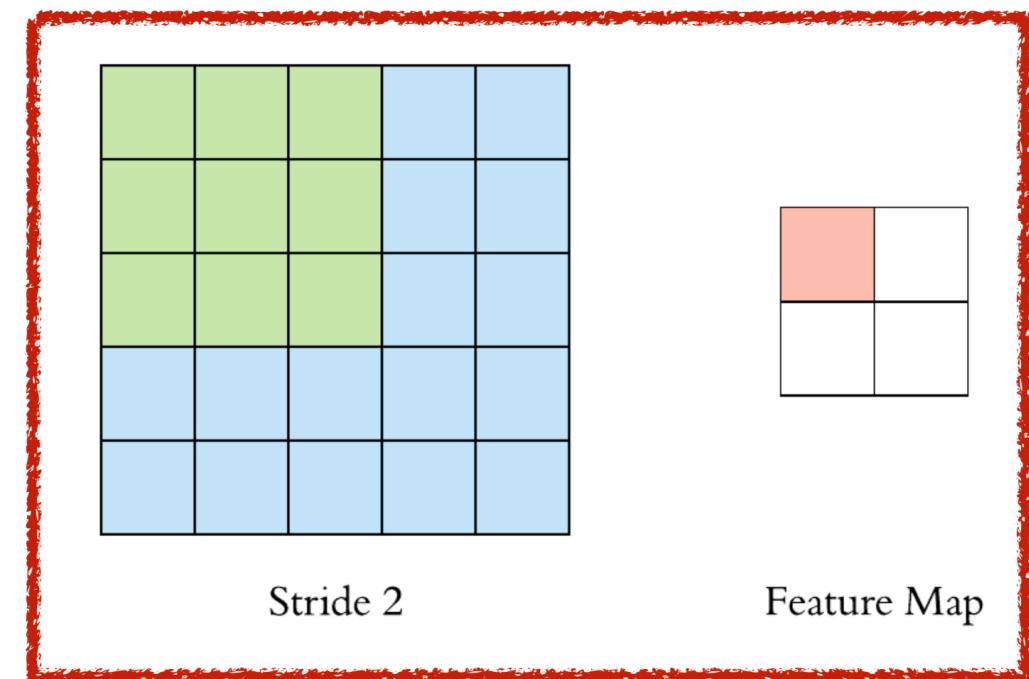
*By how much you move the filter*



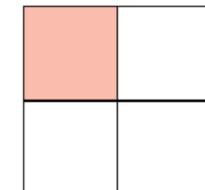
Stride 1



Feature Map



Stride 2

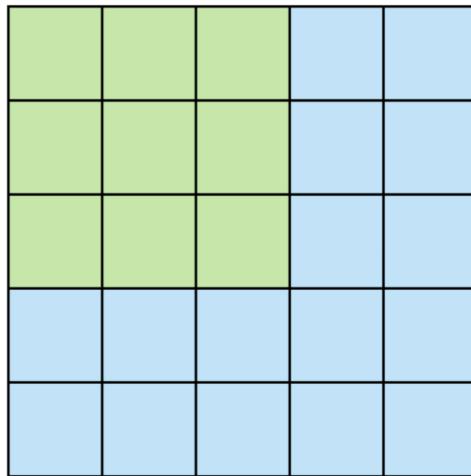


Feature Map

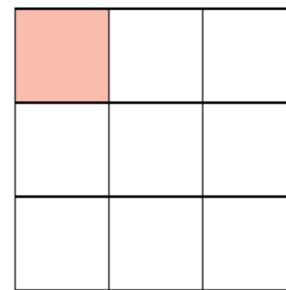
Increasing stride from 1 to 2

# Strided convolutions

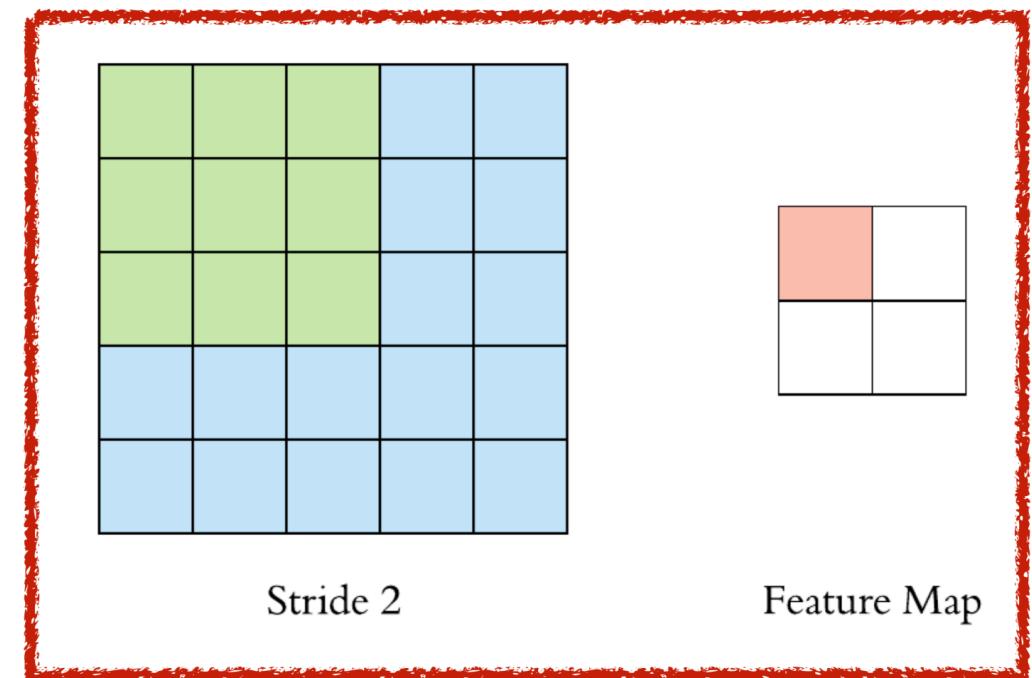
*By how much you move the filter*



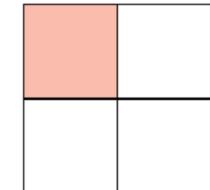
Stride 1



Feature Map



Stride 2



Feature Map

Increasing stride from 1 to 2

# Pooling Layer (POOL)

Similar to the Convolutional Layer, the Pooling layer is responsible for reducing the spatial size of the Convolved Feature. This is to decrease the computational power required to process the data through dimensionality reduction. Furthermore, it is useful for extracting dominant features which are rotational and positional invariant, thus maintaining the process of effectively training of the model.

- reduces the spatial dimension ( $n_H$  and  $n_W$ ) to **decrease resource (memory) load**

There are two types of Pooling: Max Pooling and Average Pooling. Max Pooling returns the maximum value from the portion of the image covered by the Kernel. On the other hand, Average Pooling returns the average of all the values from the portion of the image covered by the Kernel.

Max Pool

2	3	1	9
4	7	3	5
8	2	2	2
1	3	4	5



7	9
8	5

Max-Pool with a  
2 by 2 filter and  
stride 2.

*Get the max value*

Max pooling also performs as a noise suppressant --> discarding noisy activations --> usually performs better than average pooling just takes the average of all activations.

Average Pool

2	3	1	9
4	7	3	5
8	2	2	2
1	3	4	5

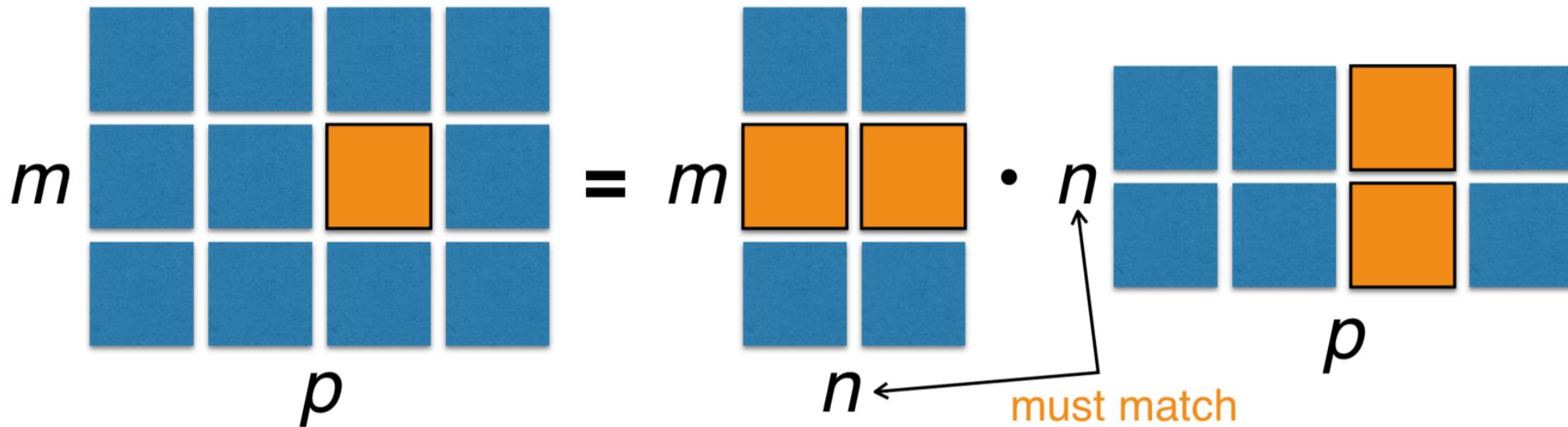


4	4.5
3.25	3.25

Average Pool with  
a 2 by 2 filter and  
stride 2.

*Get the average value*

# Fully Connected Layer (FC)



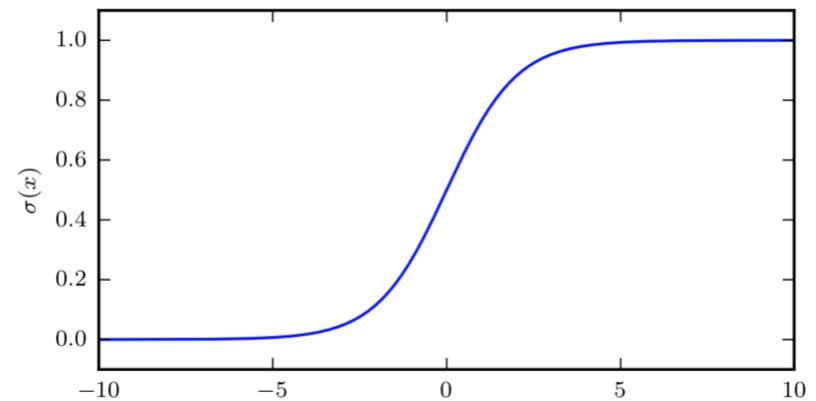
Adding a Fully-Connected layer is a (usually) cheap way of learning non-linear combinations of the high-level features as represented by the output of the convolutional layer. The Fully-Connected layer is learning a possibly non-linear function in that space.

**matrix product**

# Output layer : activation functions

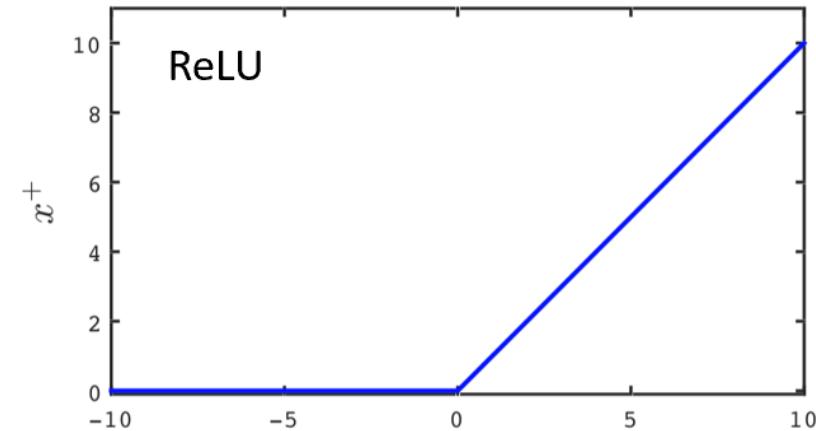
## 1) Classification : probability vector

- Sigmoid (binary class)
- Softmax (multiple class)

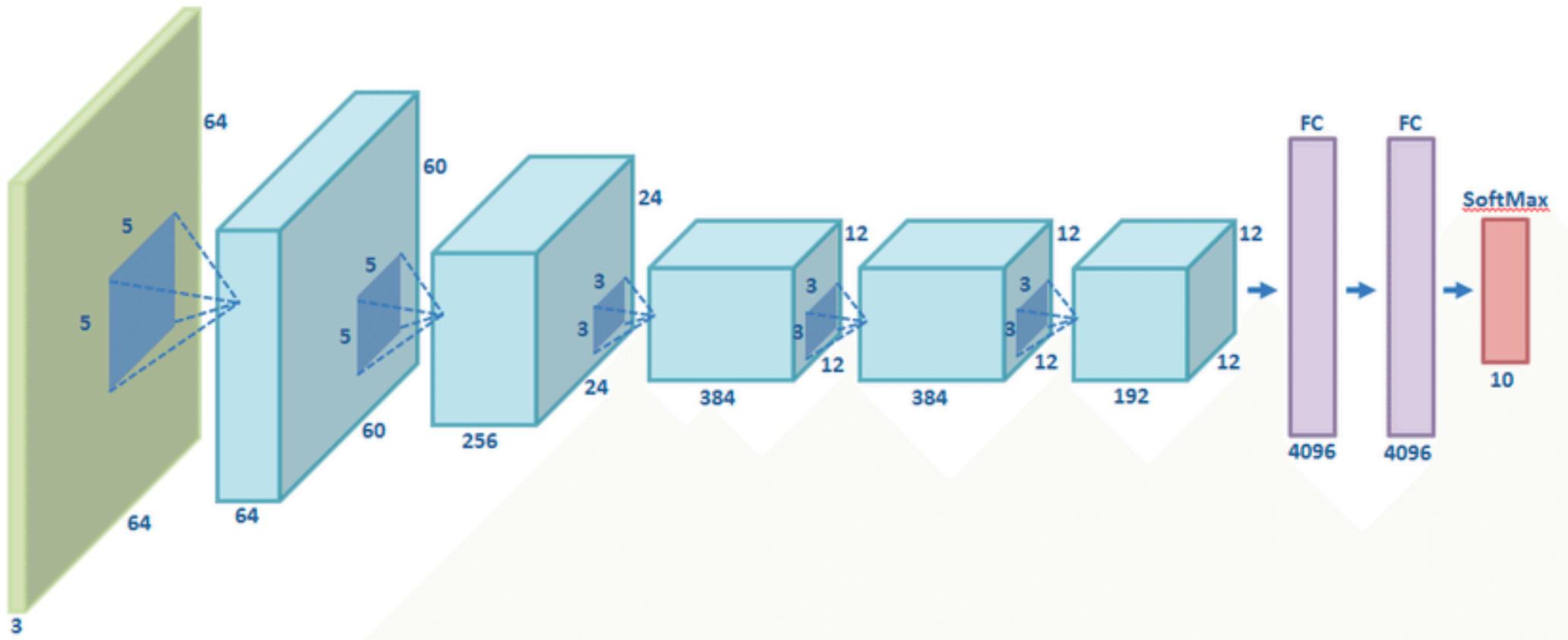


## 2) Regression : mean estimate

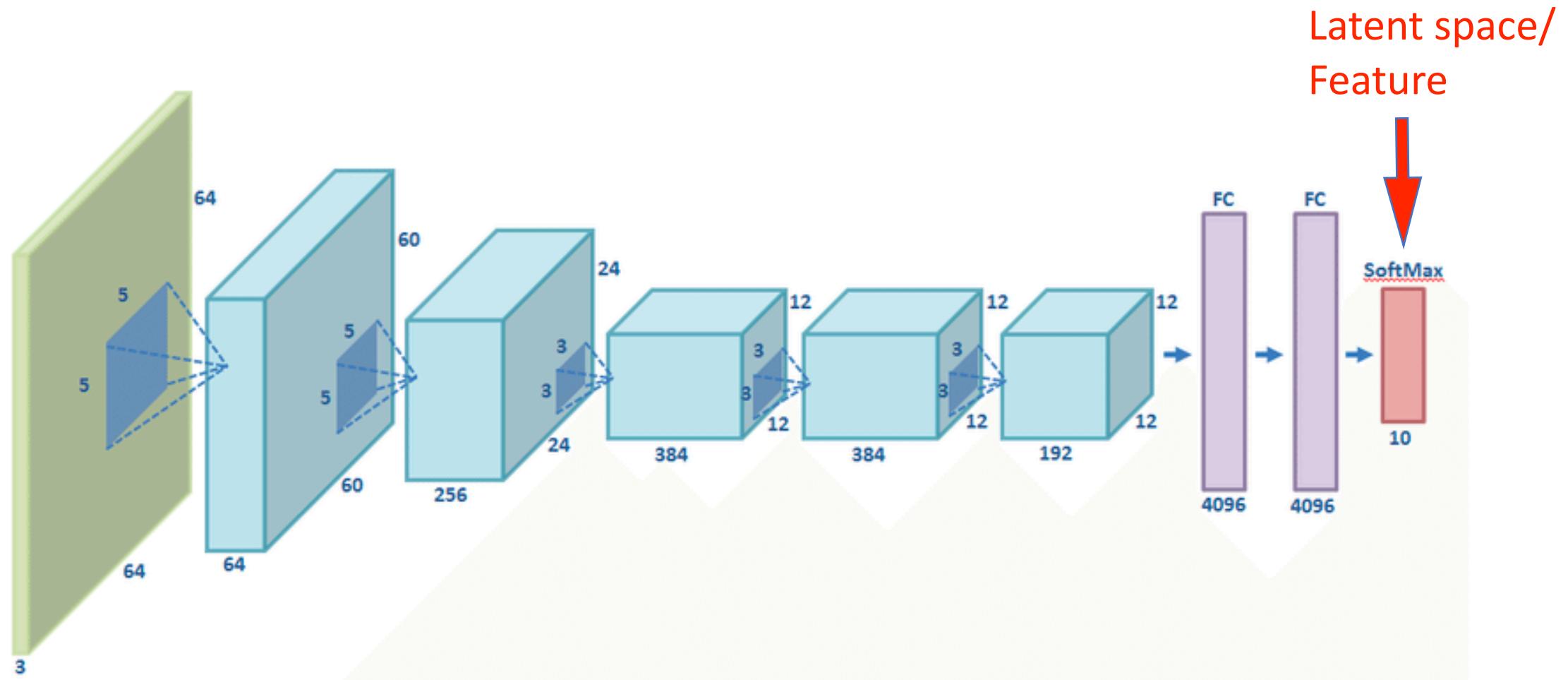
- ReLU
- Softplus
- Smoothed max
- Generalization of ReLU (leaky ReLU,...)



# Convolutional neural network



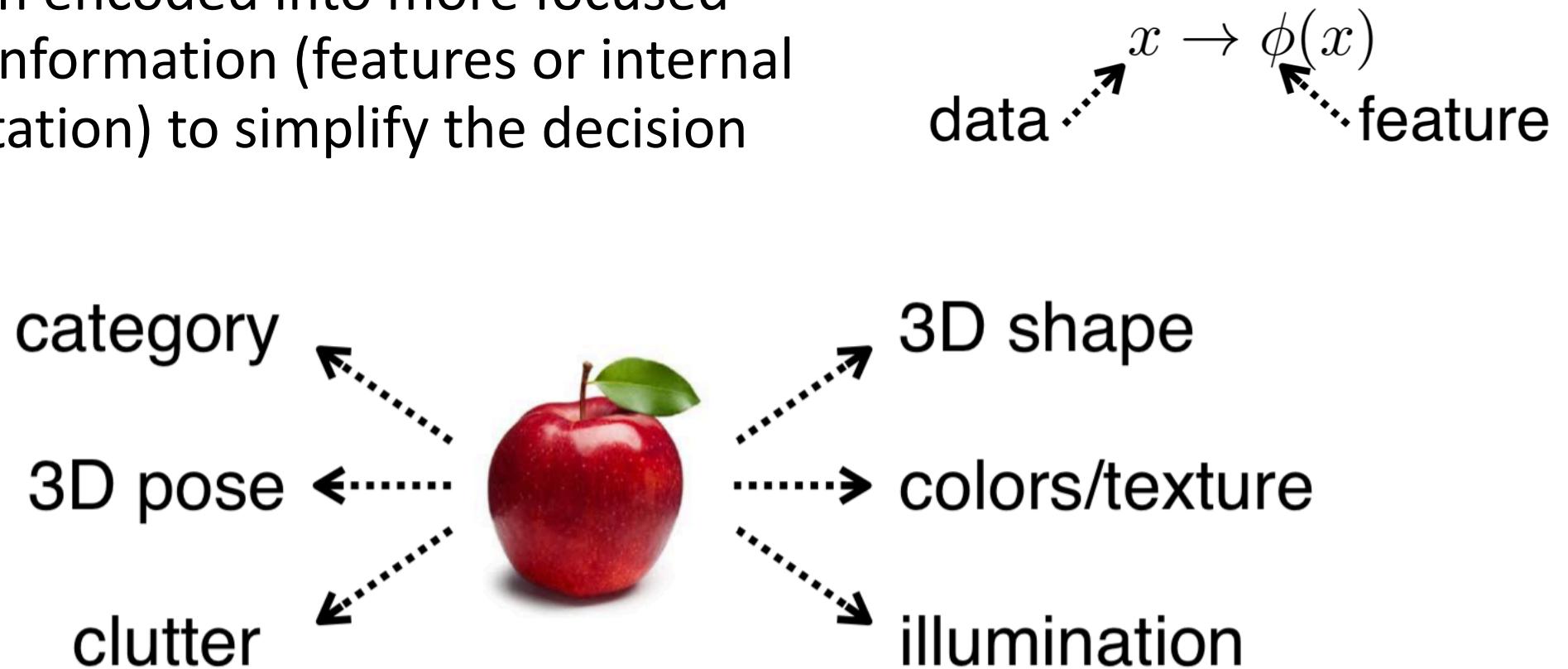
# Convolutional neural network



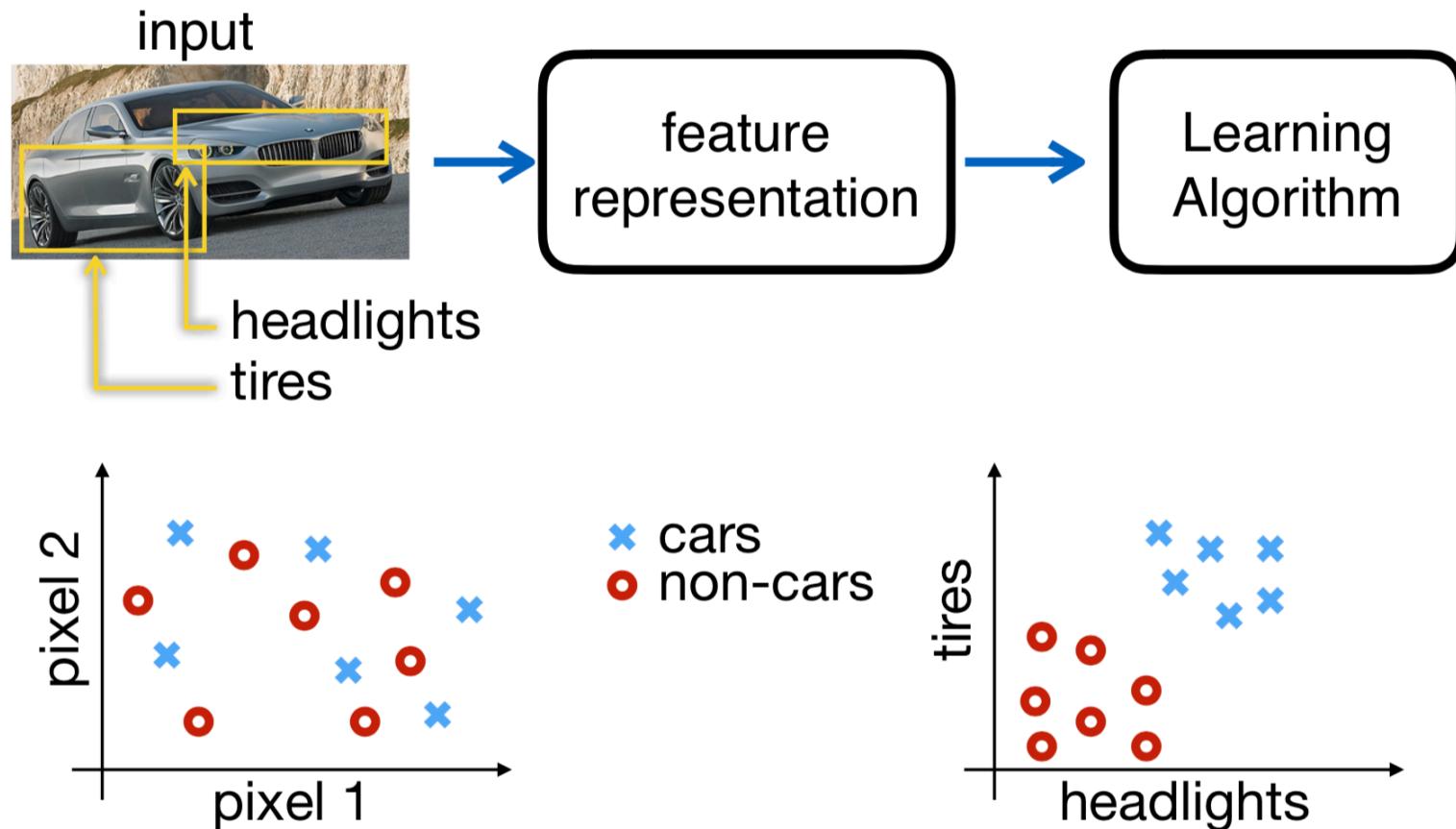
# Features

Features are high level data such as edges, color, shape, noise, etc. that can be detected with the right use of CNNs.

- Data often encoded into more focused relevant information (features or internal representation) to simplify the decision

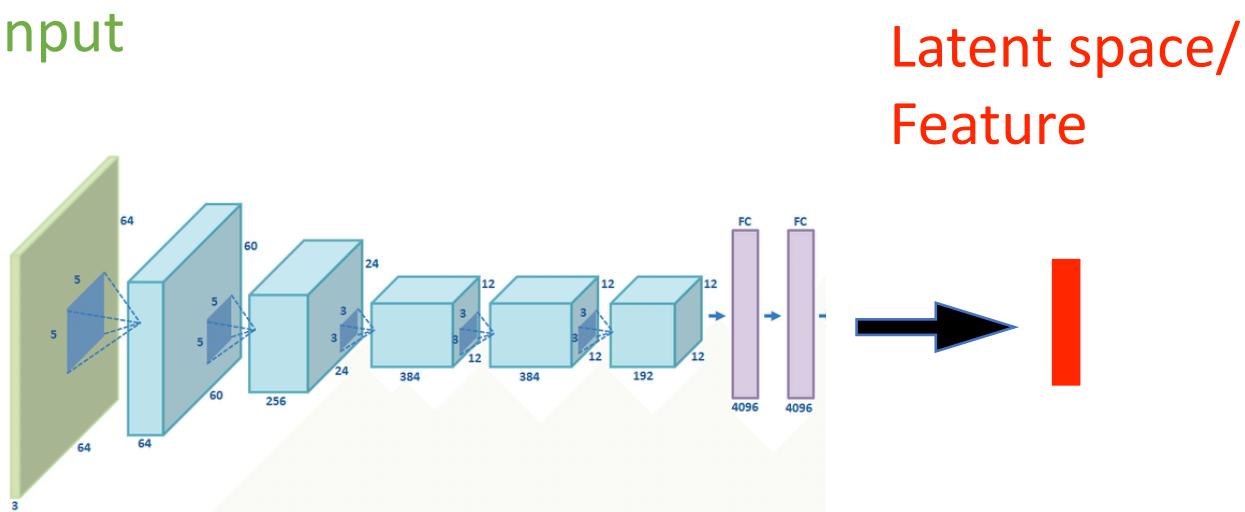


# Features Example :Image classification



# Convolutional neural network

Input

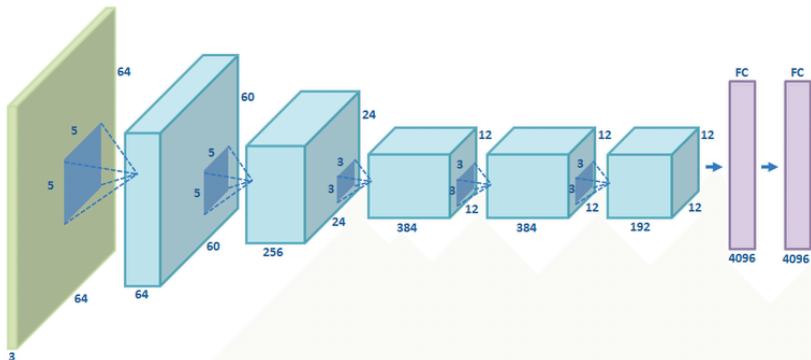


Latent space/  
Feature

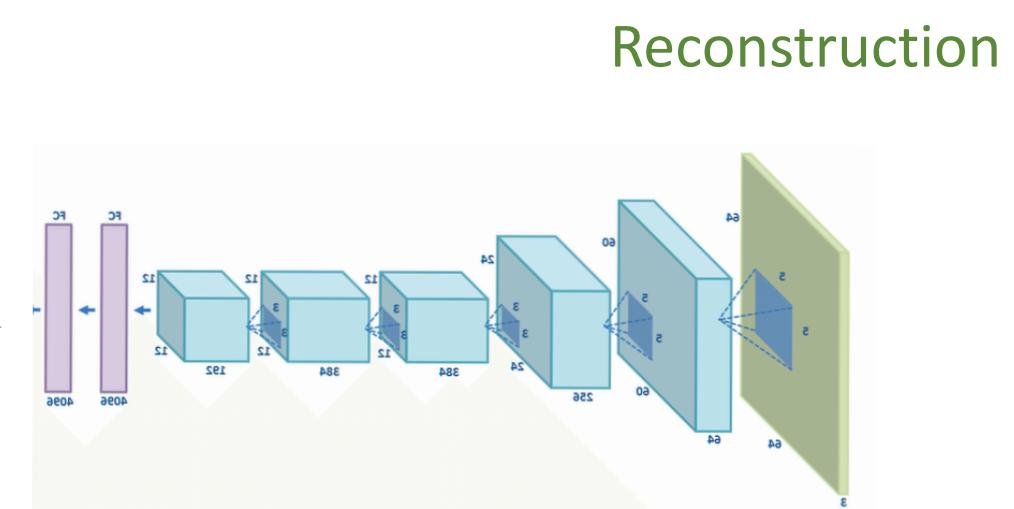
# AutoEncoder

Autoencoder is an unsupervised artificial neural network that learns how to efficiently compress and encode data then learns how to reconstruct the data back from the reduced encoded representation to a representation that is as close to the original input as possible.  
Autoencoder, by design, reduces data dimensions by learning how to ignore the noise in the data.

Input



Latent space/  
Feature



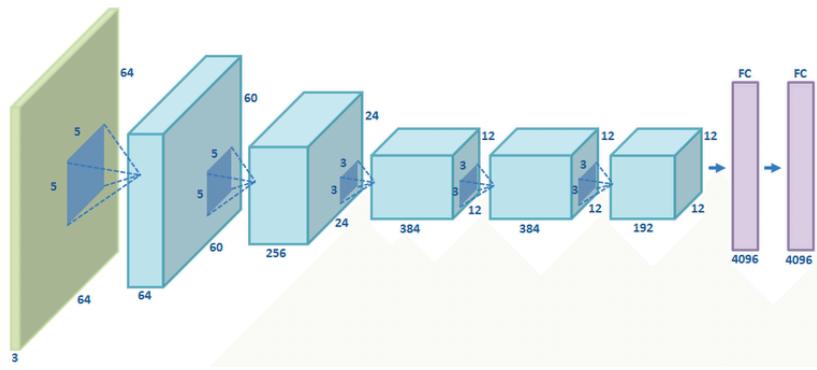
1- Encoder: In which the model learns how to reduce the input dimensions and compress the input data into an encoded representation.

2- Bottleneck: which is the layer that contains the compressed representation of the input data. This is the lowest possible dimensions of the input data

3- Decoder: In which the model learns how to reconstruct the data from the encoded representation to be as close to the original input as possible.

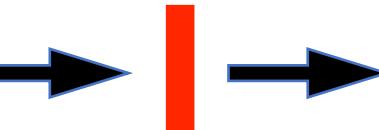
# AutoEncoder

Input

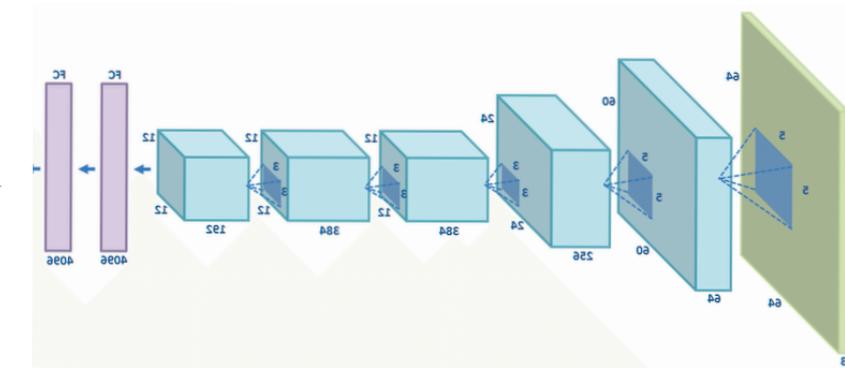


Encoder

Latent space/  
Feature

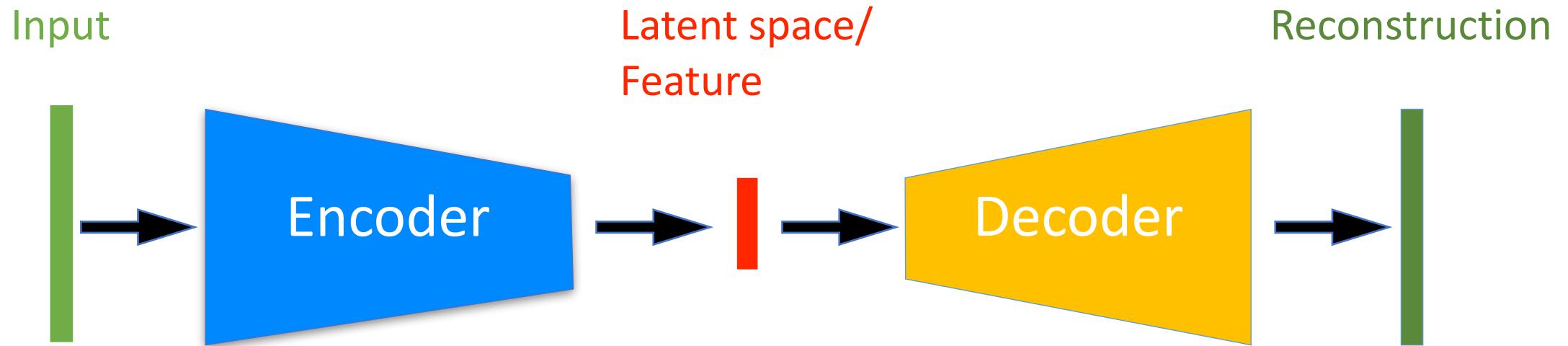


Reconstruction



Decoder

# AutoEncoder



Tutorial 1 : 11h-12h15:  
Introduction to Tensorflow

Tutorial 2: 17h-18h30  
Optimization in Tensorflow  
and NN introduction

