

Shallow to Deep Neural Networks

Computer Vision and Artificial Intelligence

Dr Varun Ojha

varun.ojha@ncl.ac.uk

School of Computing
Newcastle University

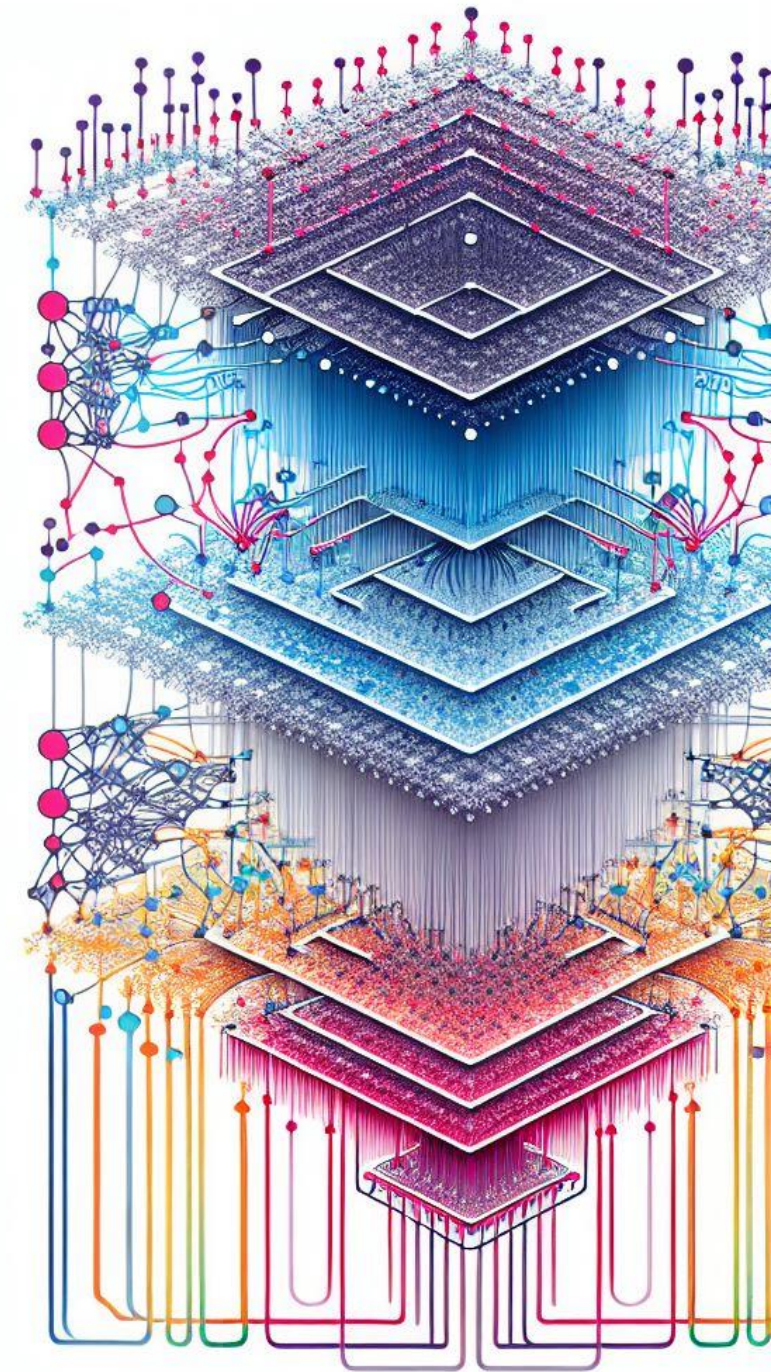
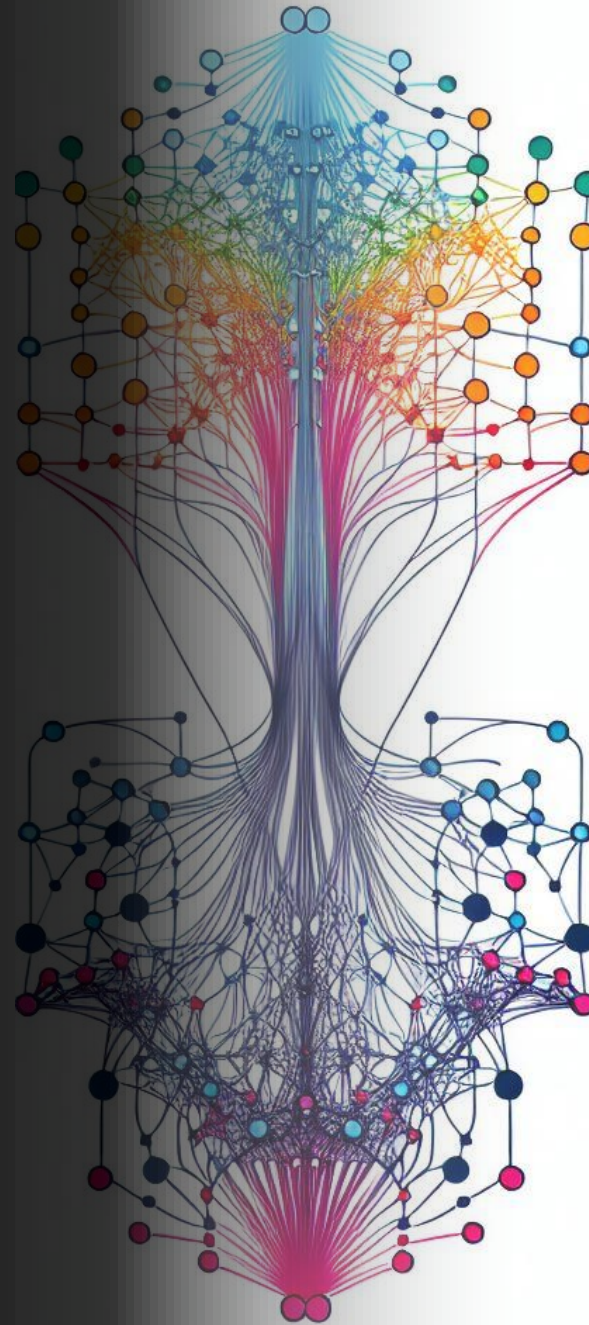


Image by Dall-E 3

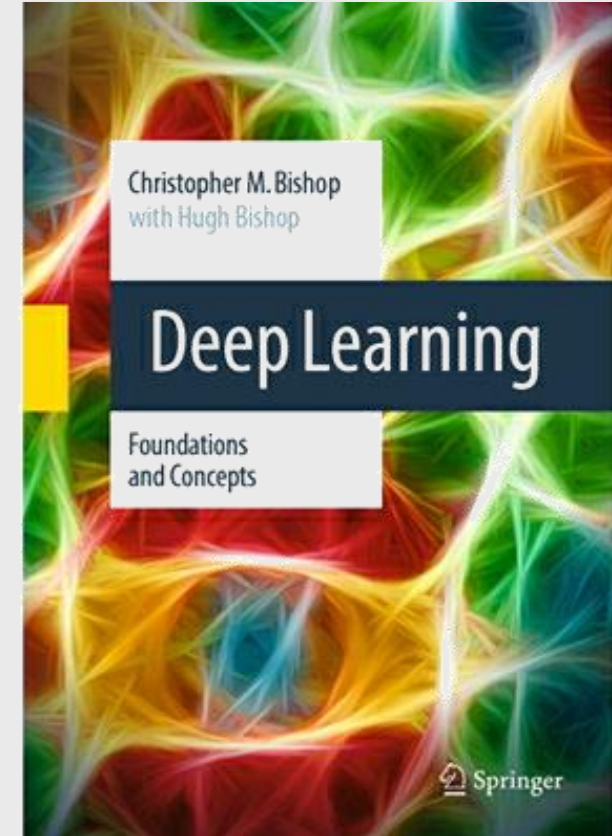
Learning objectives (neural networks)

By the end of this week, you will be able to

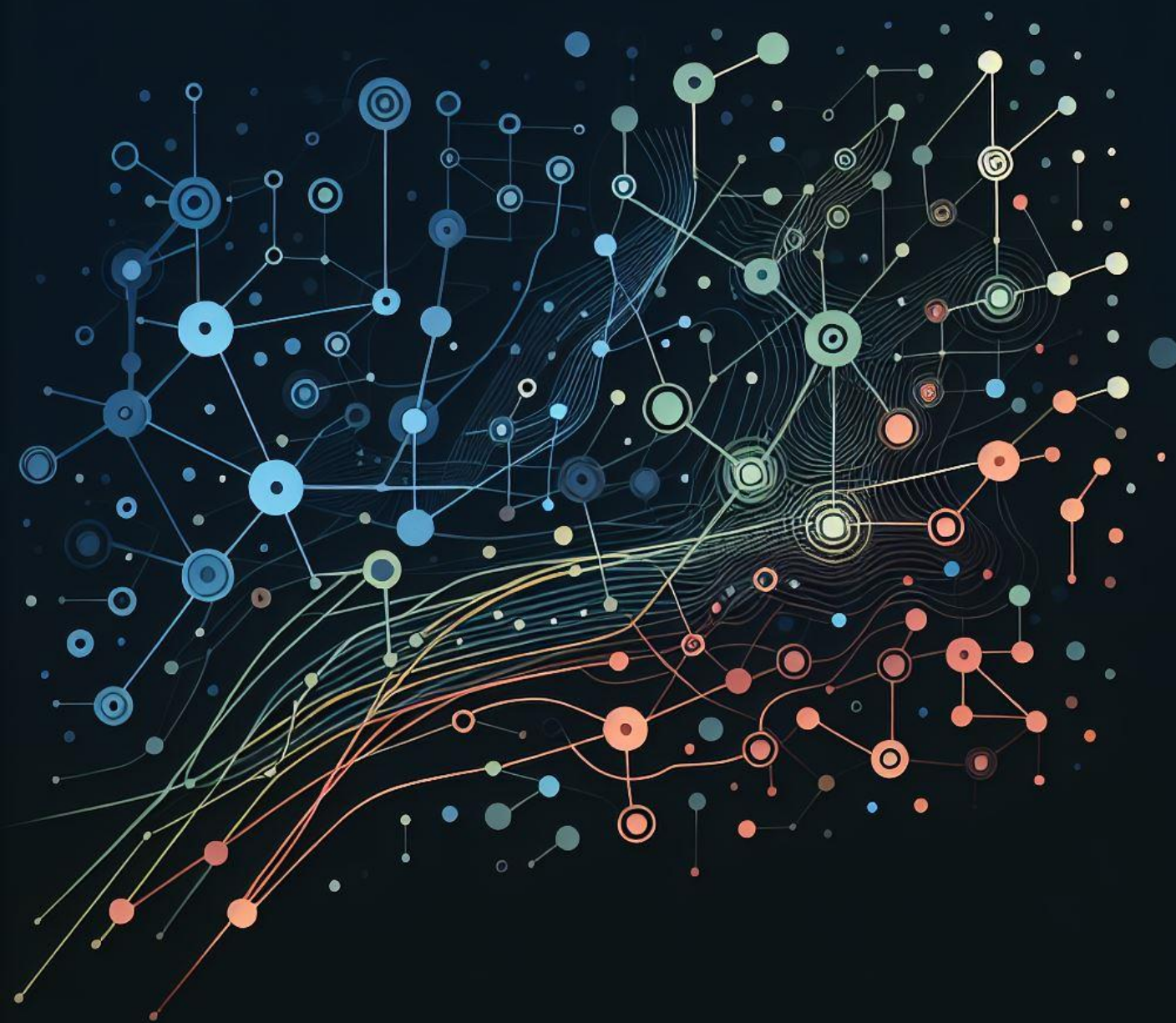
- Learn ‘concepts of learning’ in Neural Networks
- Understand gradient descent and backpropagation algorithms
- Distinguish shallow and deep neural network architectures
- Apply and evaluate neural networks for a pattern recognition (image classification) problem – *in practical session*

Content (neural networks)

- **Part 1: Introduction**
- **Part 2: Fundamental Theory**
 - Supervised Learning problem
 - Design of learning process
 - Gradient descent optimization
- **Part 3: Neural Network Architectures**
 - Neural Network Components
 - The Backpropagation algorithm
 - Deep Neural Networks
- **Part 4: Practical Exercise**



I recommend *Deep Learning* by C Bishop
<https://www.bishopbook.com/>



Part 1 Introduction

Intrinsic Intelligence? *Inside a baby's mind*

Experiment:
Warneken & Tomasello (2006)



Video Source:
<https://www.youtube.com/watch?v=cUWllxpUfM0>
(Accessed on 21 Feb 2021)

© Warneken & Tomasello

Causal understanding of water displacement by a crow

Experiment: Sarah et al. (2014), Auckland and Cambridge
Video Source: <https://www.youtube.com/watch?v=ZerUbHmuY04>



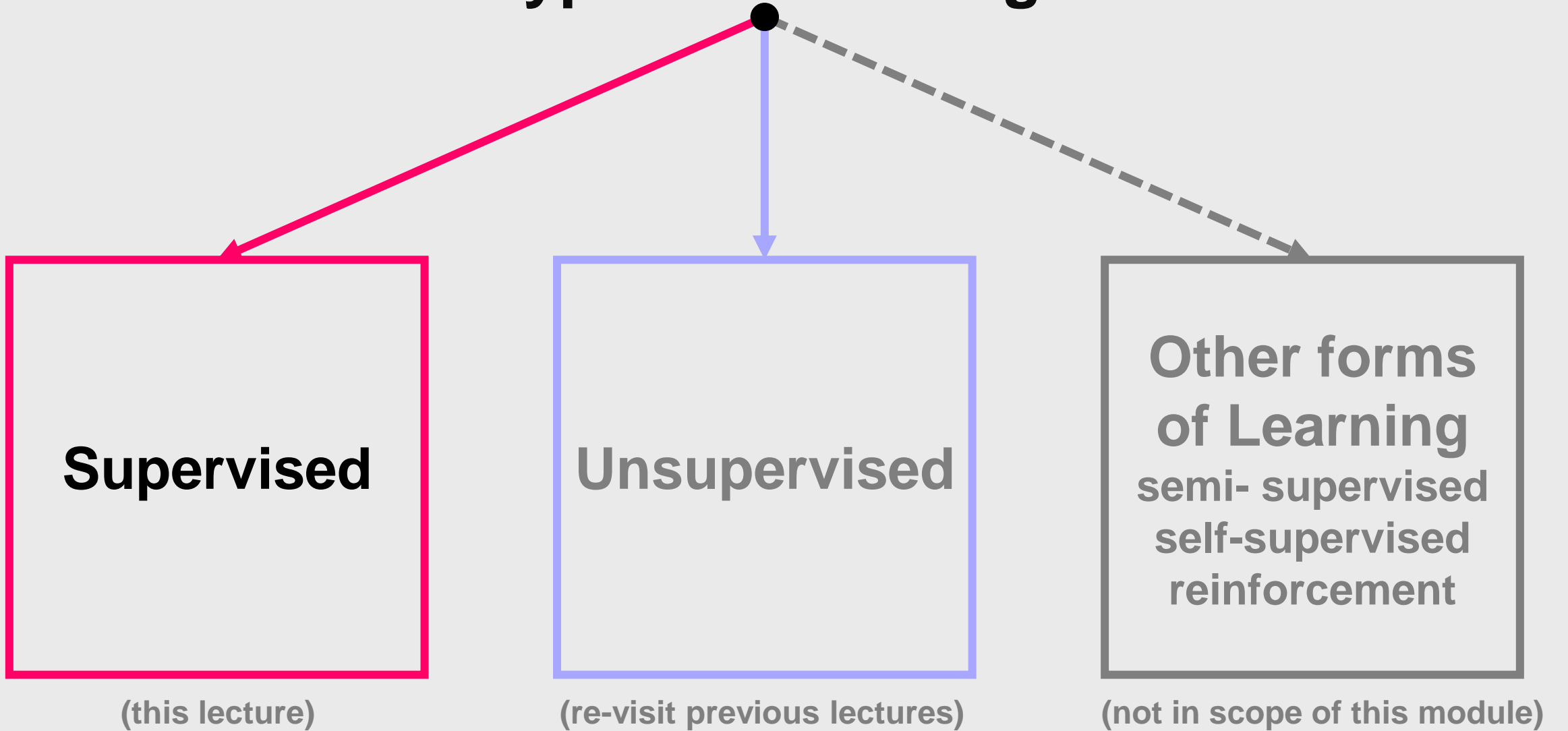
Learning by example

Learning / Training

Video Source:
<https://www.youtube.com/watch?v=Ak7bPuR2rDw>
(Accessed on 21 Feb 2021)

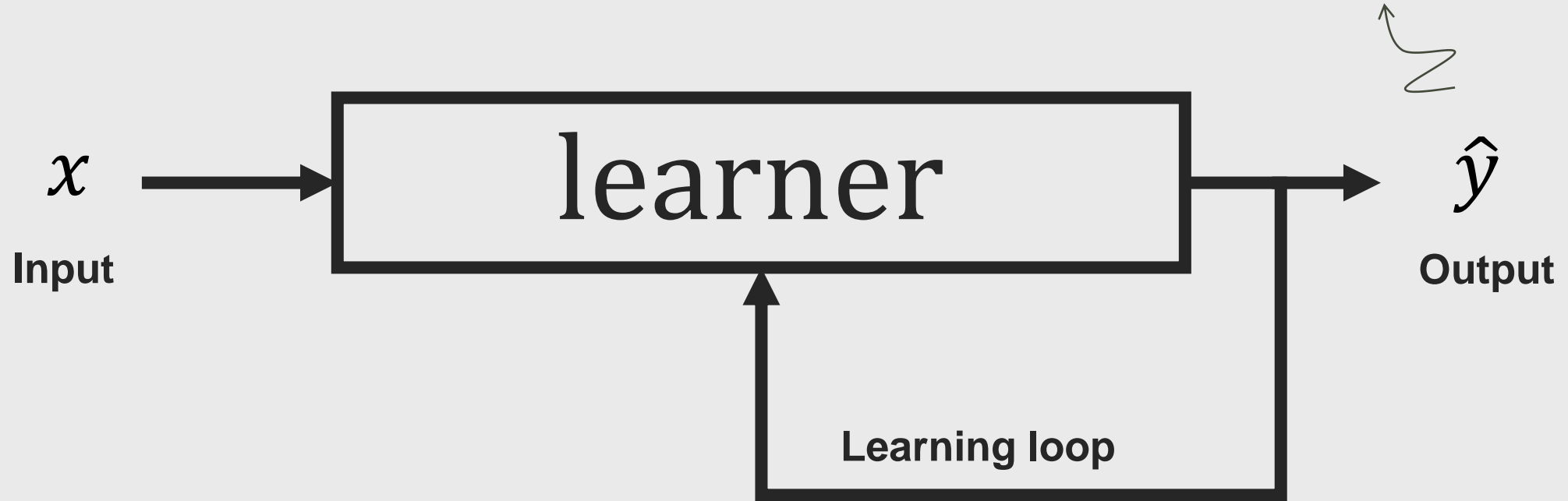


Types of Learning



Unsupervised

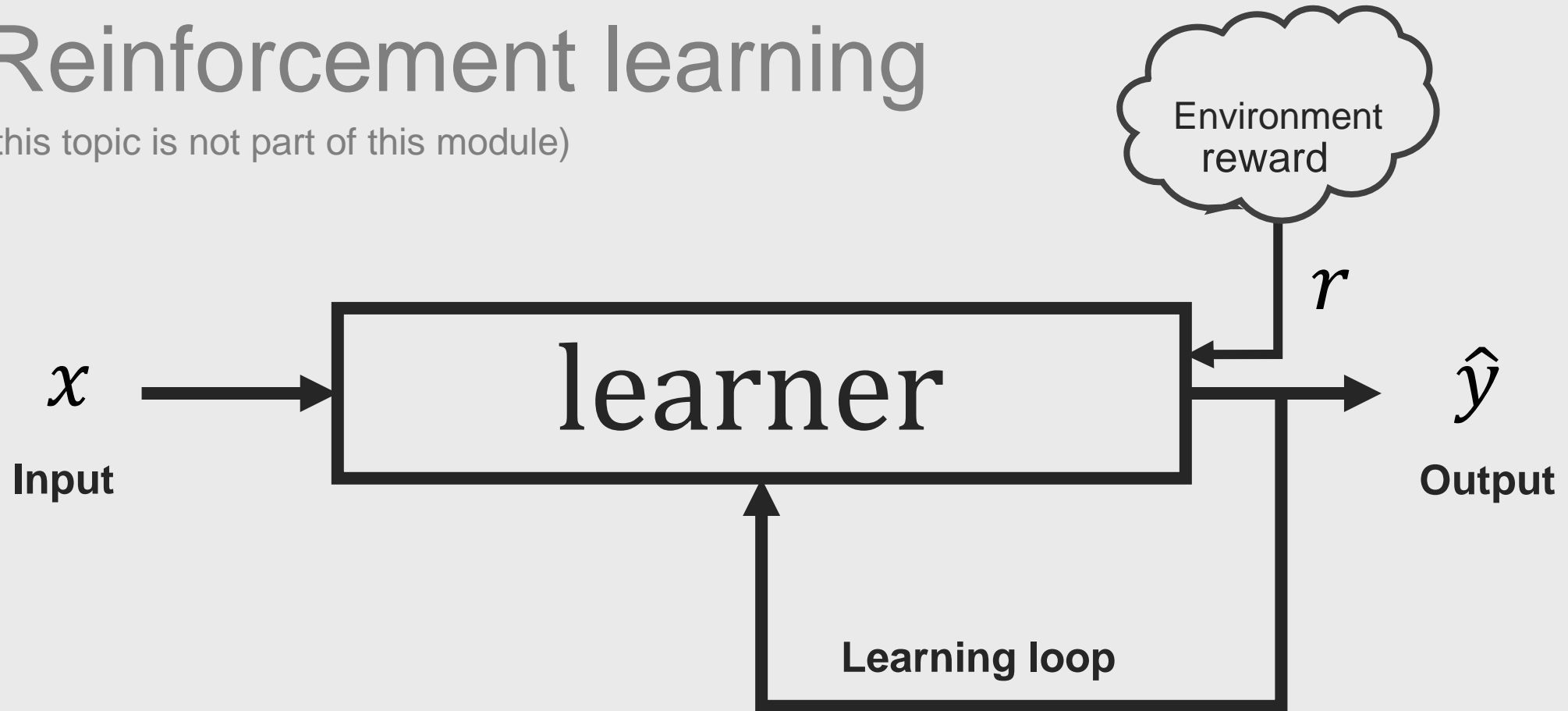
Categorise input features
or learns the input
features representation
or learns the structure of
the input



Re-visit previous lectures: Clustering (K-Means, DB Scan); Dimensionality reduction; Anomaly detection, Variational Autoencoder (VAE; *is not a part of this module*):

Reinforcement learning

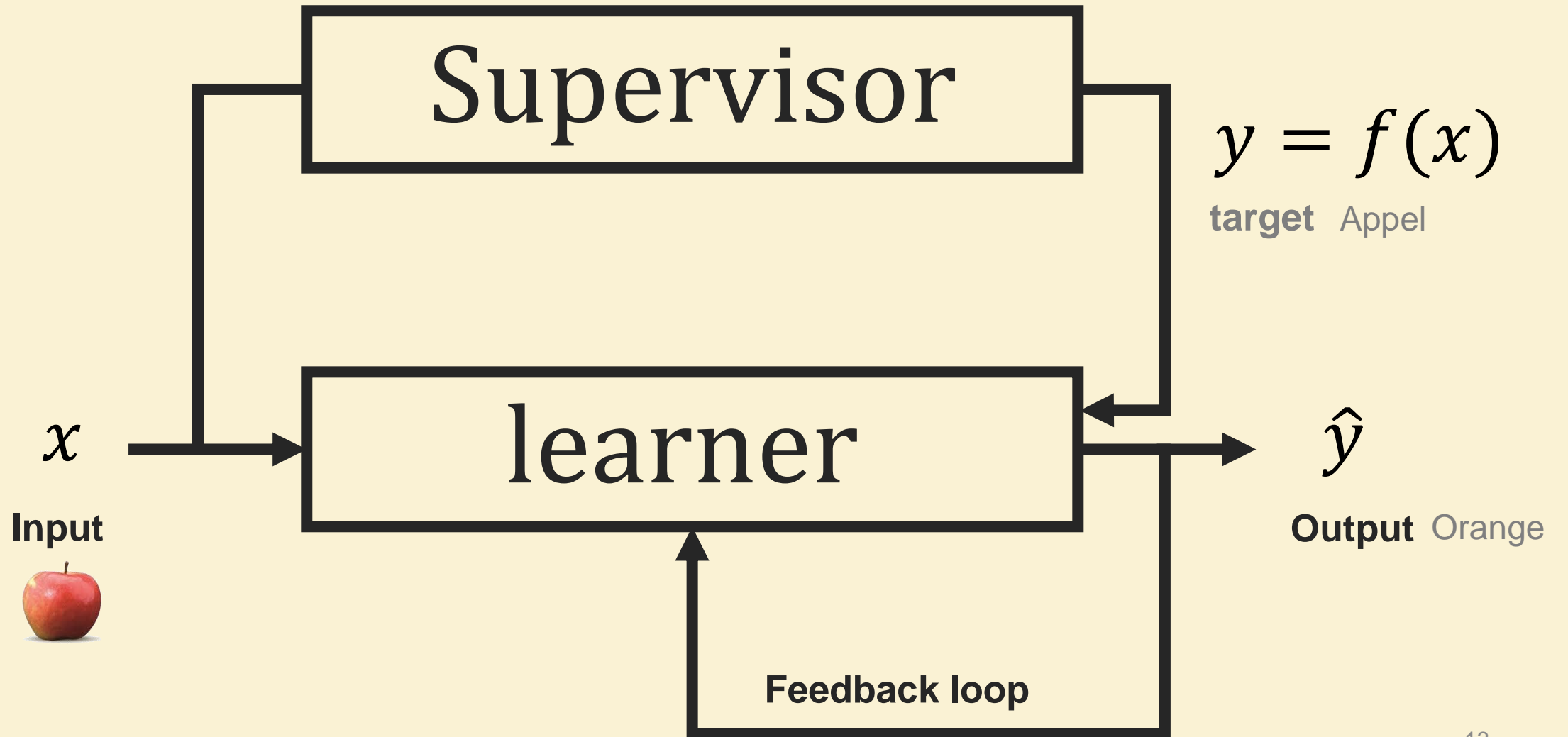
(this topic is not part of this module)



Example are game playing (most popular games are 'Atari' and 'Hide and Seek' where reinforcement learning (RL) is used, or RL is heavily used in robotics for learning control and actions)

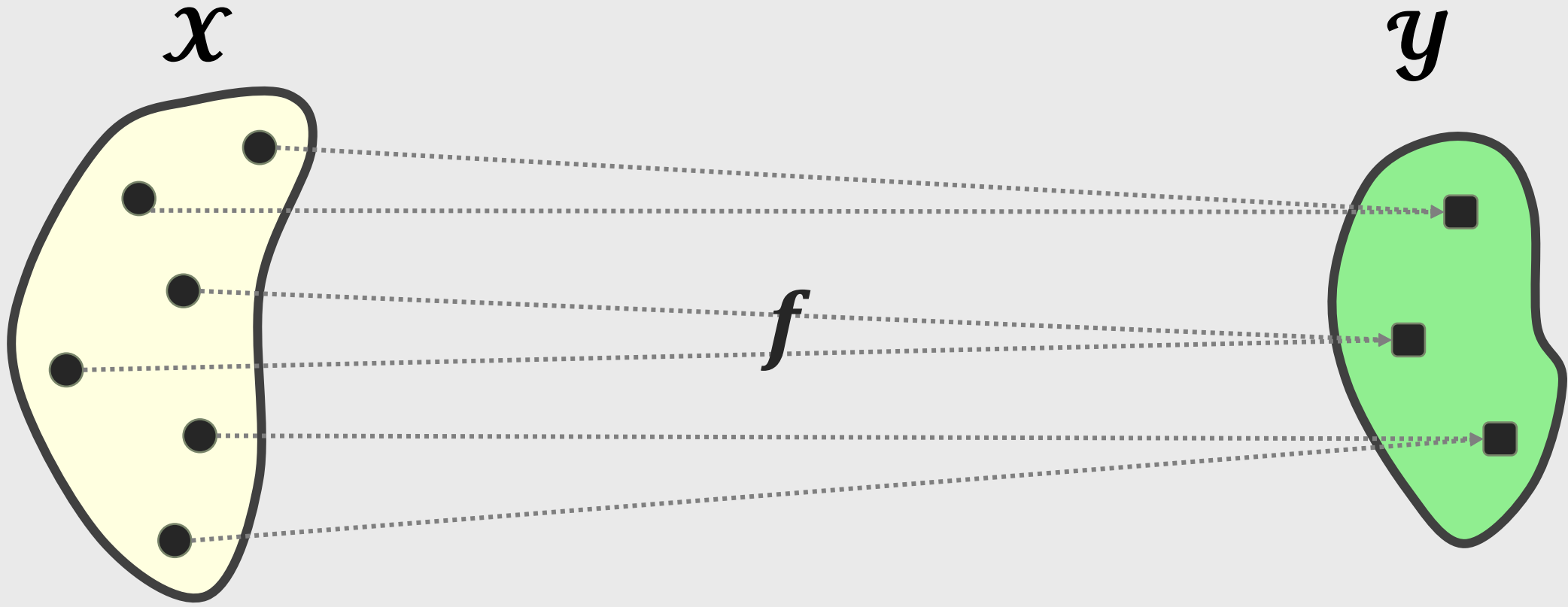
Explore (not part of this module): Reinforcement Learning: An Introduction by Richard S Sutton

Supervised



Learning $f : \mathcal{X} \rightarrow \mathcal{Y}$

Supervised learning is a mapping f of inputs \mathcal{X} to outputs \mathcal{Y}

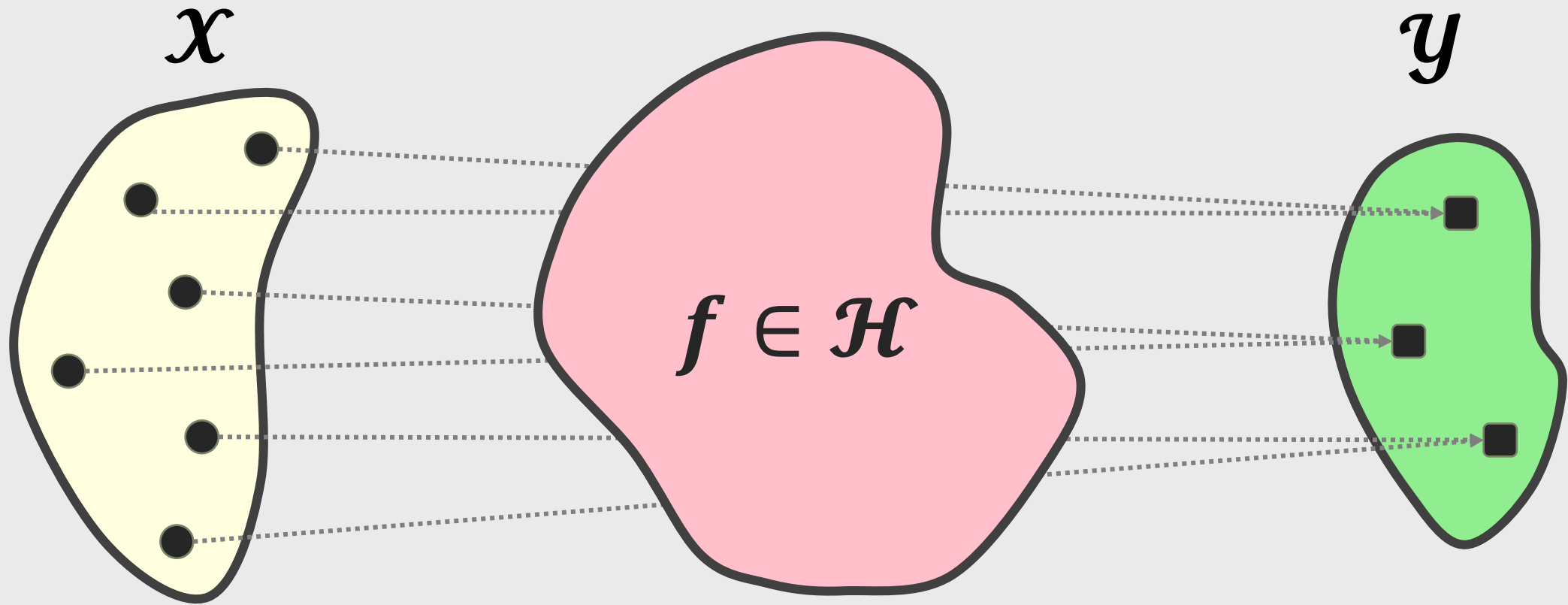


Inputs $\mathbf{X} \in$ Input space \mathcal{X}

outputs $\mathbf{y} \in$ output space \mathcal{Y}

Learning $f : X \rightarrow y$

We need to [find the unknown target function \$f\$](#) that maps x to y



Inputs $X \in$ Input space \mathcal{X}

model space \mathcal{H}

outputs $y \in$ output space \mathcal{Y}

Learning: $g(\mathbf{X}) \sim f(\mathbf{X})$

We need to search a function $g(X)$ that can approximate $f(X)$

Example Training Task: AND Logic Problem

	Colour	Shape	Fruit Name
	x_1	x_2	y
Fruit \mathbf{x}_1 :	0	0	0
\mathbf{x}_2 :	0	1	0
\mathbf{x}_3 :	1	0	0
\mathbf{x}_4 :	1	1	1

Input $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4)^T$, $\mathbf{x}_i = (x_{i1}, x_{i2})$,
Output $y = \{0,1\}$

Number of Inputs $d = 2$

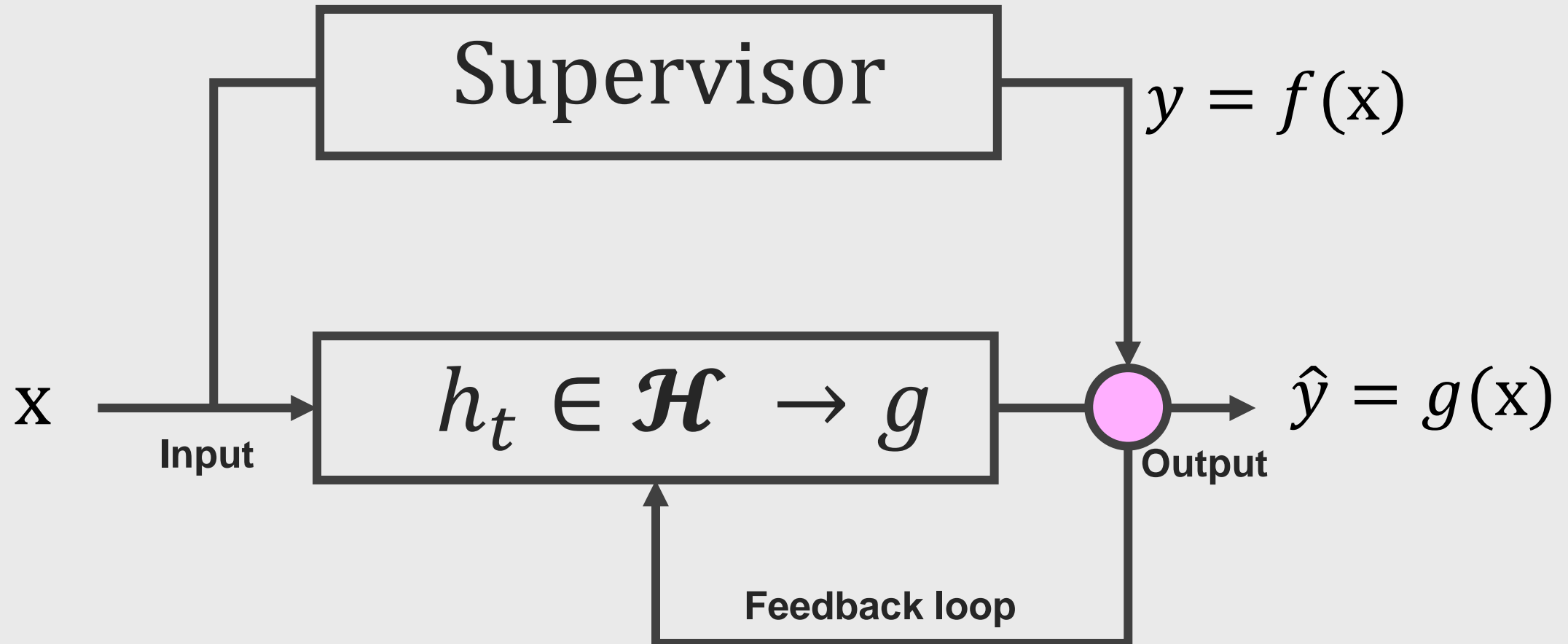
Each input x takes values either 0 or 1

Input-space $\mathcal{X} = 2^d = 2^2 = 4$

Number of outputs 1

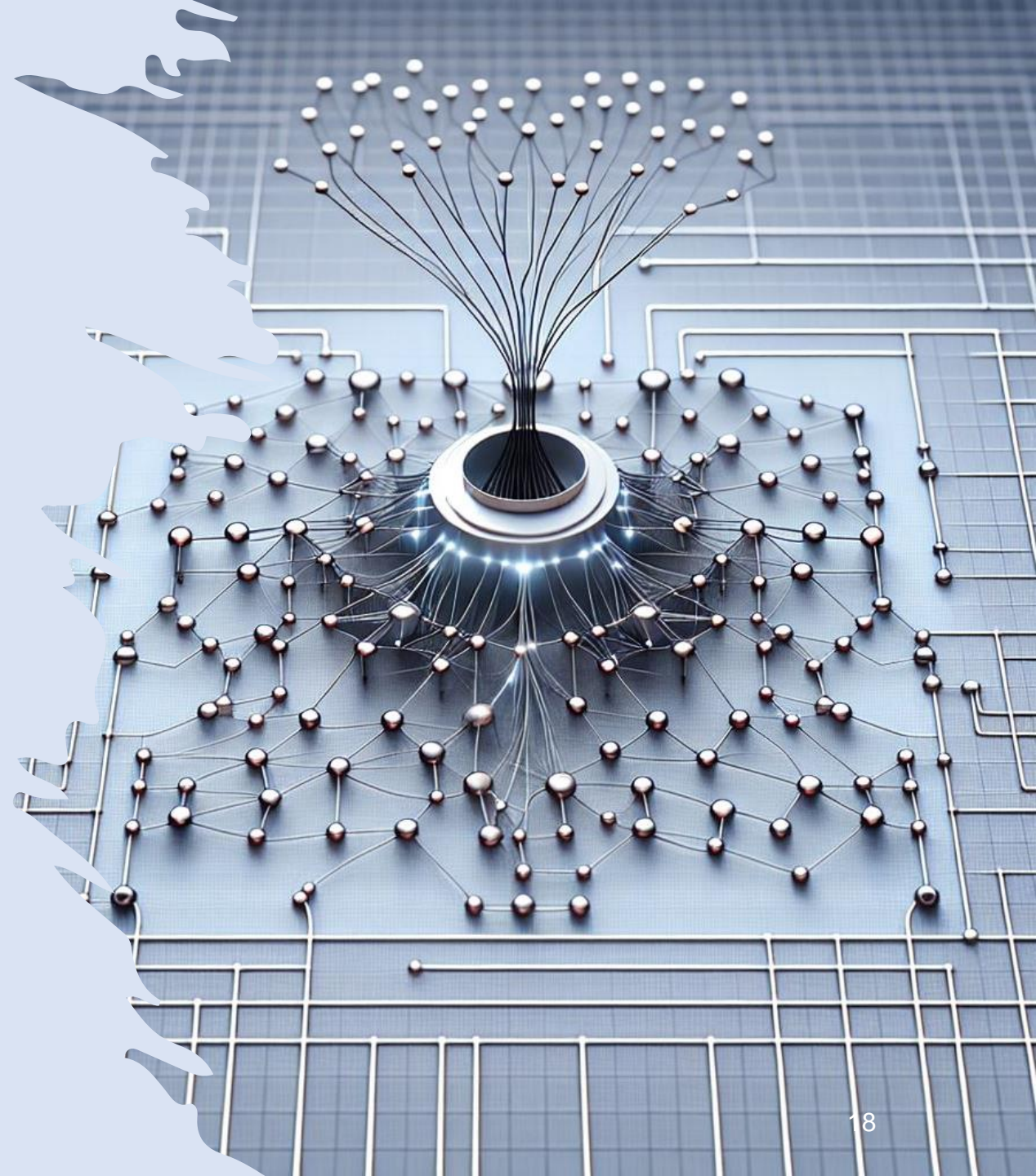
Output y takes 2 options from $\{0,1\}$

Search a function $g: X \rightarrow y$ that approximates $f(x)$



Part 2

Learning Theory



Requirements of Learning

Learning needs to

Represent a model (use a neural network architecture, deep neural networks)

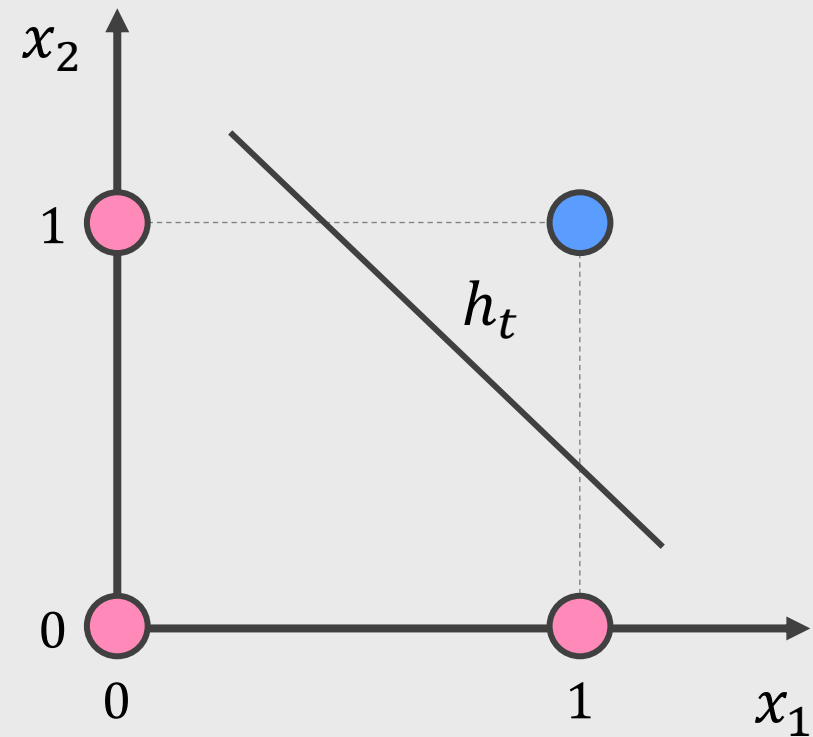
Evaluate the model (use a loss/cost function, e.g., Cross Entropy or MSE)

Optimize the model (use an optimizer, e.g., backpropagation – Adam or SGD)

Represent a model

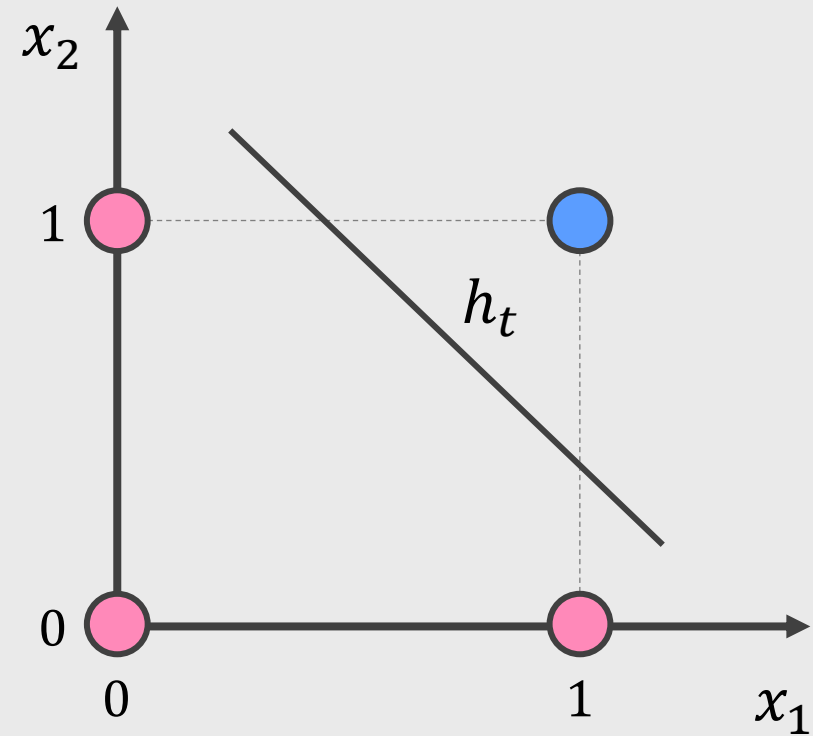
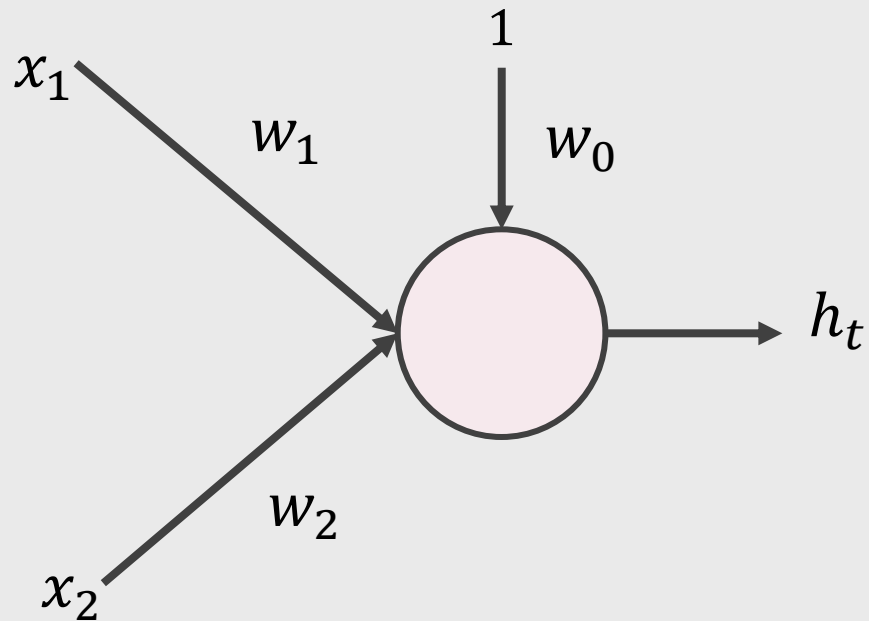
A line separating data can be considered as a model

	x_1	x_2	y
1:	0	0	0
2:	0	1	0
3:	1	0	0
4:	1	1	1



Represent a model

A line separating data can be considered a model which equivalent to a single neuron or a perceptron



Represent a model $h_t \in H$

A line separating data can be considered a model which equivalent to a single neuron or a perceptron

Perceptron is a simple linear combination of inputs, which is written as:

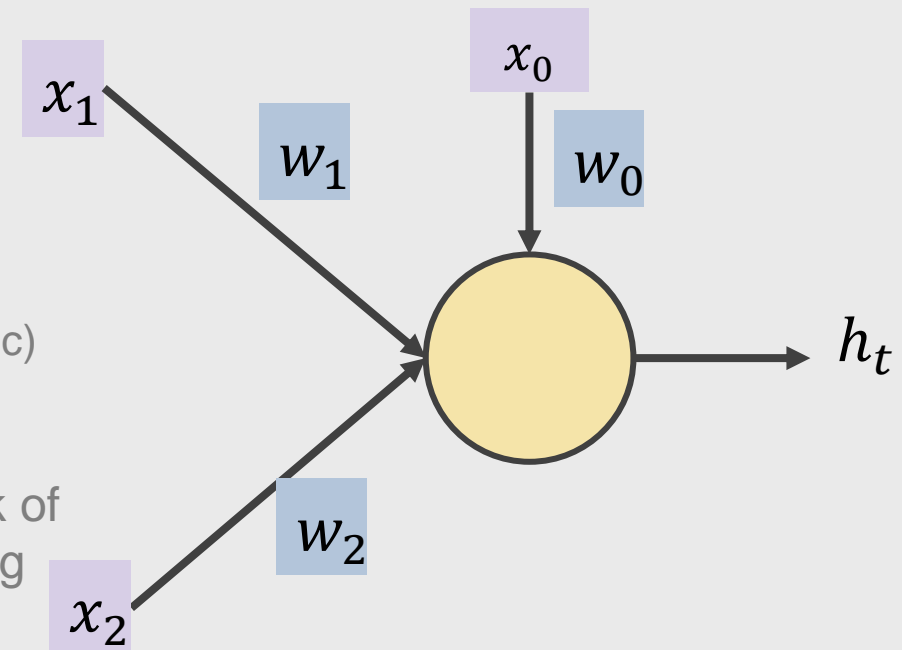
$$h_t = g(x) = \sum_{i=1}^d w_i x_i \geq x_0 w_0 ,$$

This equation is also equivalent to linear regression ($y = mx + c$)

where w_0 is a threshold.

The model h_t has the weights w_i and the threshold w_0 as its **trainable parameters**.

Real bottleneck of
Deep Learning



Model h_t as a **perceptron / single neuron**.

Represent a model $h_t \in H$

A line separating data can be considered as a model

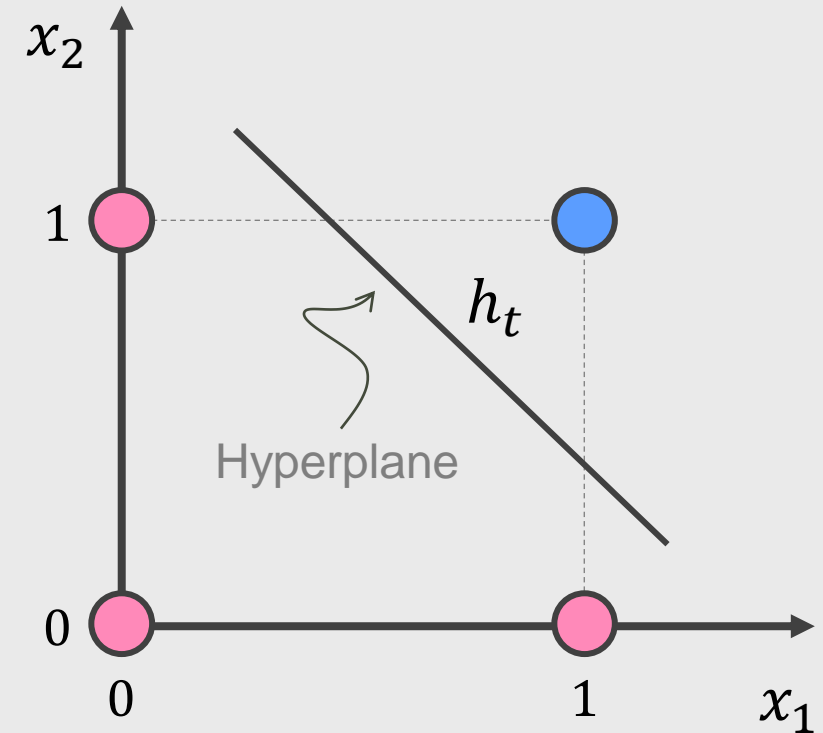
A model h_t as a perceptron.

$$\sum_{i=1}^d w_i x_i \geq x_0 w_0$$

$$\sum_{i=1}^d w_i x_i - x_0 w_0 = 0 \quad \text{This equation is also called a hyperplane}$$

For an artificial input (also called bias) $x_0 = 1$ we have:

$$\sum_{i=0}^d w_i x_i = 0 \quad \text{This is an equation of a single neuron}$$



hyperplane as decision boundary

Which model $h_t \in H$ to pick?

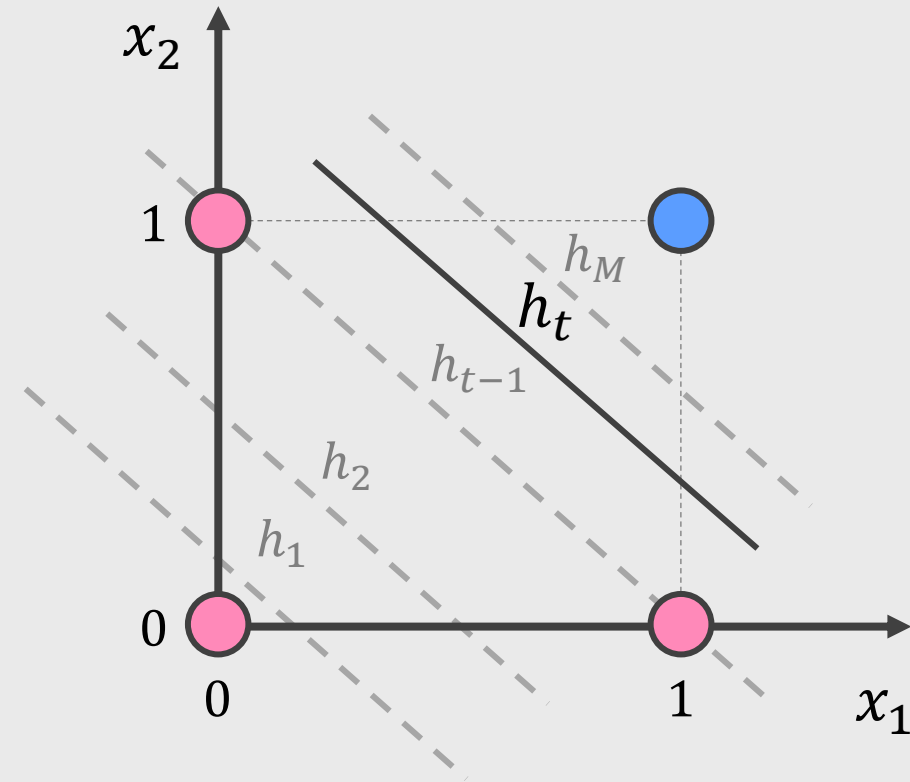
How to **evaluate** a model: compute cost of choosing a model

	x_1	x_2	$y = f(\mathbf{x})$
1:	0	0	0
2:	0	1	0
3:	1	0	0
4:	1	1	1

\mathcal{D}

Cost function such as the error rate:

$$E(h_t(\mathcal{D})) = \frac{1}{N} \sum_{j=1}^N (g_{\mathbf{w}}(\mathbf{x}_j) \neq f(\mathbf{x}_j))$$



Optimise model $h_t \in H$ by minimizing error

How to optimize a model: compute cost of adjust the model weights

Function g of the model has parameter \mathbf{w} :

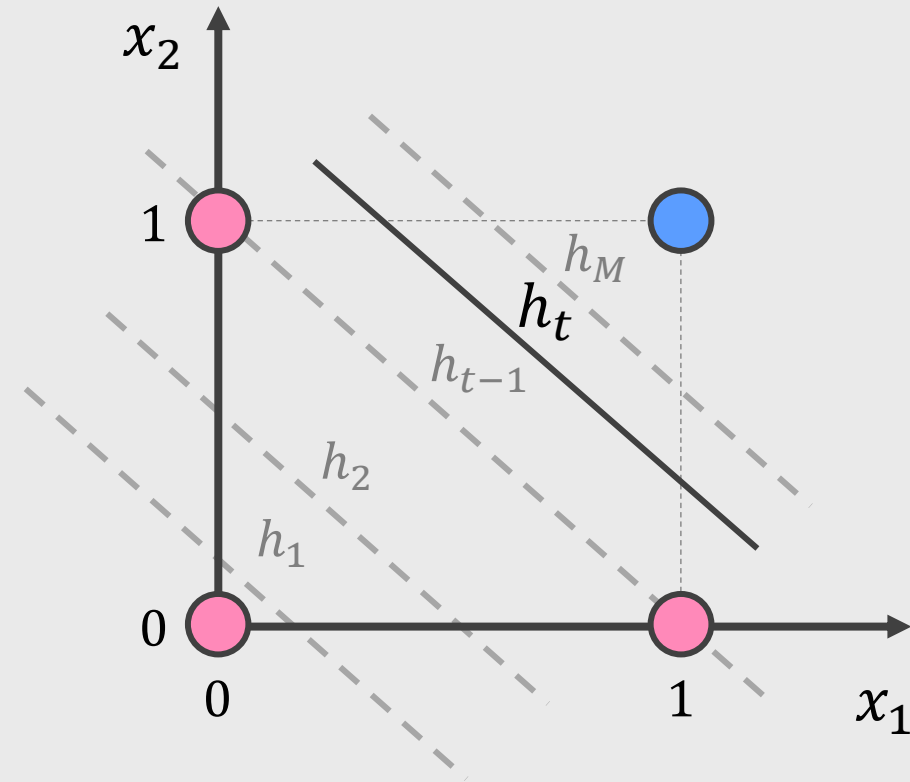
$$\hat{y} = g_{\mathbf{w}}(\mathbf{x}) = \sum_{i=0}^d w_i x_i = 0$$

Simple algorithm:

Repeat parameter \mathbf{w} update for $t = 2, 3, \dots, M$ as:

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \hat{y} \mathbf{x},$$

Until the error rate $E(h_t(\mathcal{D}))$ is acceptable or close to zero.



Does error $E(h_t(\mathcal{D}))$ minimization work?

9:37 AM

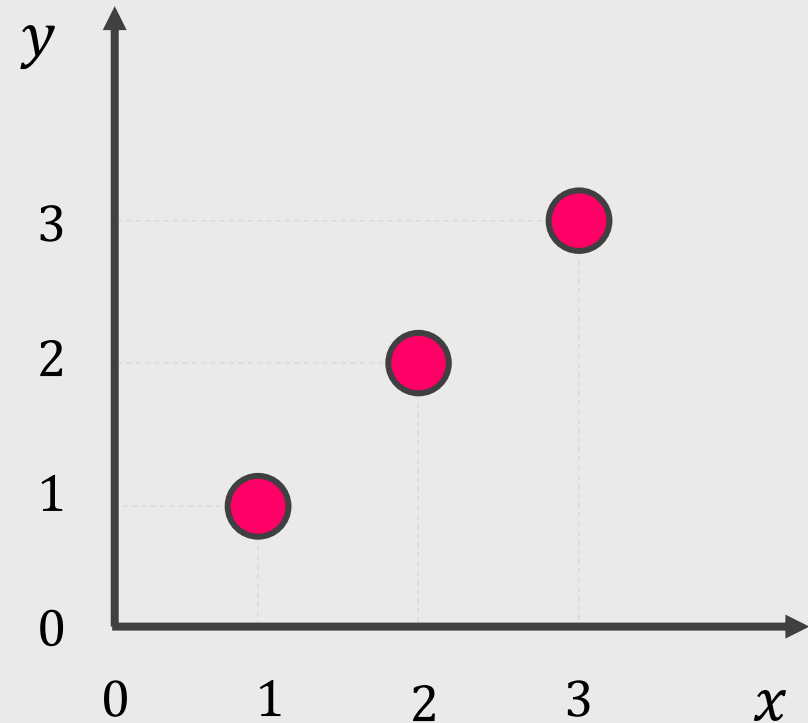
Let's see an example (house price):

	$x = \text{area}(m^2)$	$y = \text{price (in £)}$
1:	1000	100K
2:	2000	200K
3:	3000	300K

Now, the **cost function** is a squared error:

$$E(h_t(\mathbf{x})) = \frac{1}{2N} \sum_{j=1}^N (g_{\mathbf{w}}(\mathbf{x}_j) - f(\mathbf{x}_j))^2$$

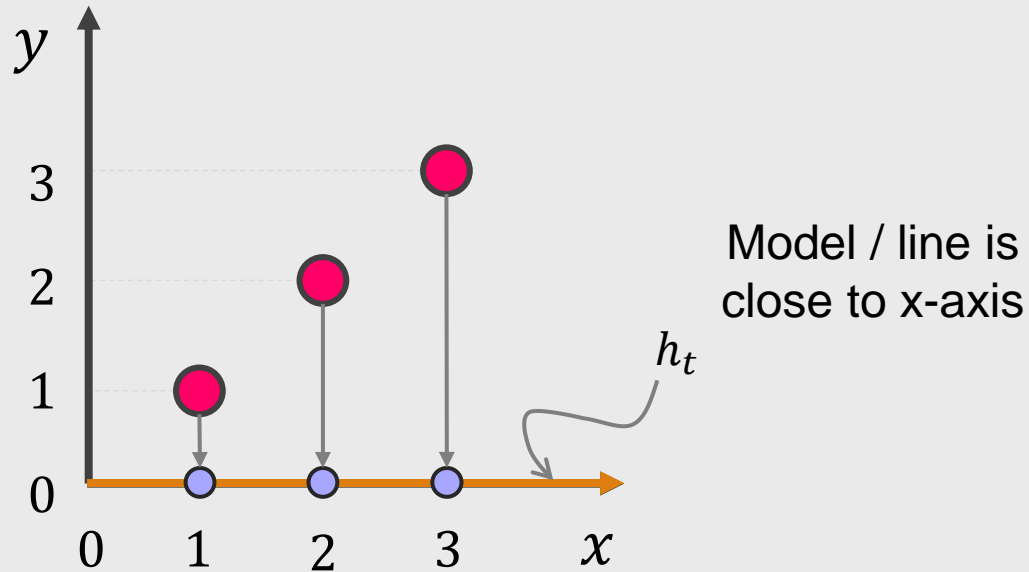
Note that y and x values are simplified to 1, 2 and 3



Does error $E(h_t(\mathcal{D}))$ minimization work?

9:37 AM

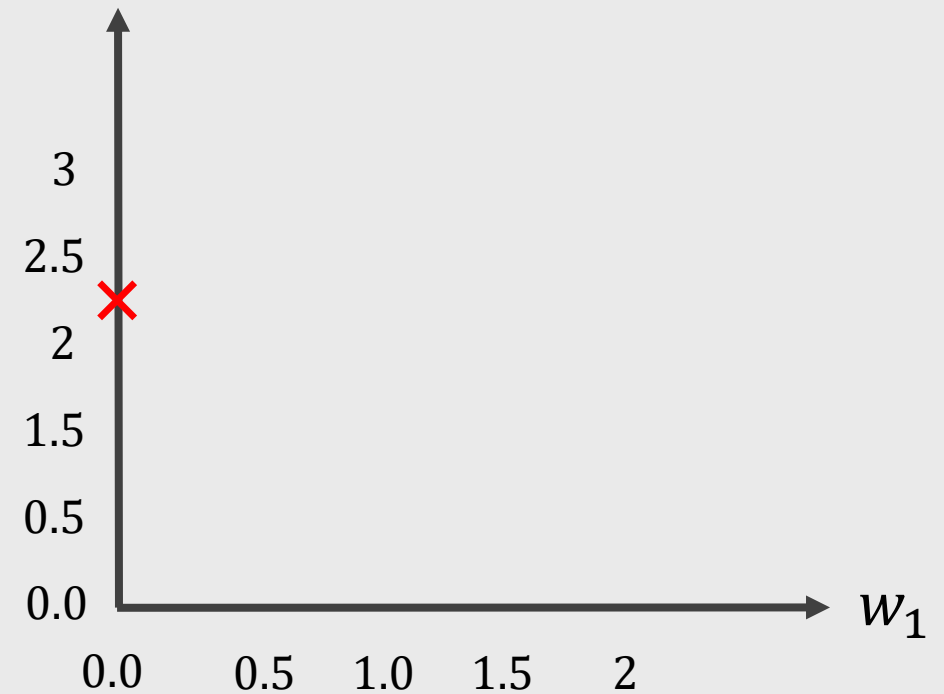
Note that y and x values are simplified to 1, 2 and 3



Model h_t for $w_0 = 0$ and $w_1 = 0.0$:

$$g_w(x_i) = w_0 + w_1 x_i \text{ for } i = 1, 2, 3$$

$$E(g_{w_1}(x)) = \frac{1}{2N} \sum_{j=1}^N (f(\mathbf{x}_j) - g_w(\mathbf{x}_j))^2$$



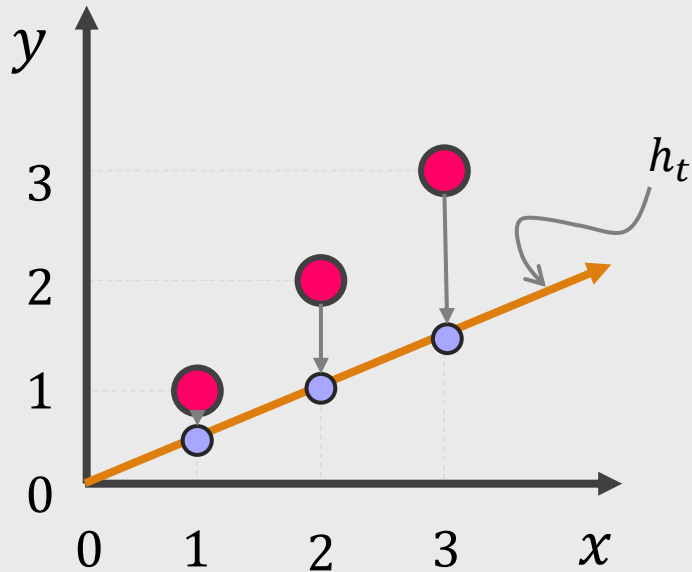
Error $E(w_1)$ for $w_0 = 0$ and $w_1 = 0.0$:

$$E(g_w(\mathbf{x})) = \frac{(1-0)^2 + (2-0)^2 + (3-0)^2}{2*3} = 2.33$$

Does error $E(h_t(\mathcal{D}))$ minimization work?

9:37 AM

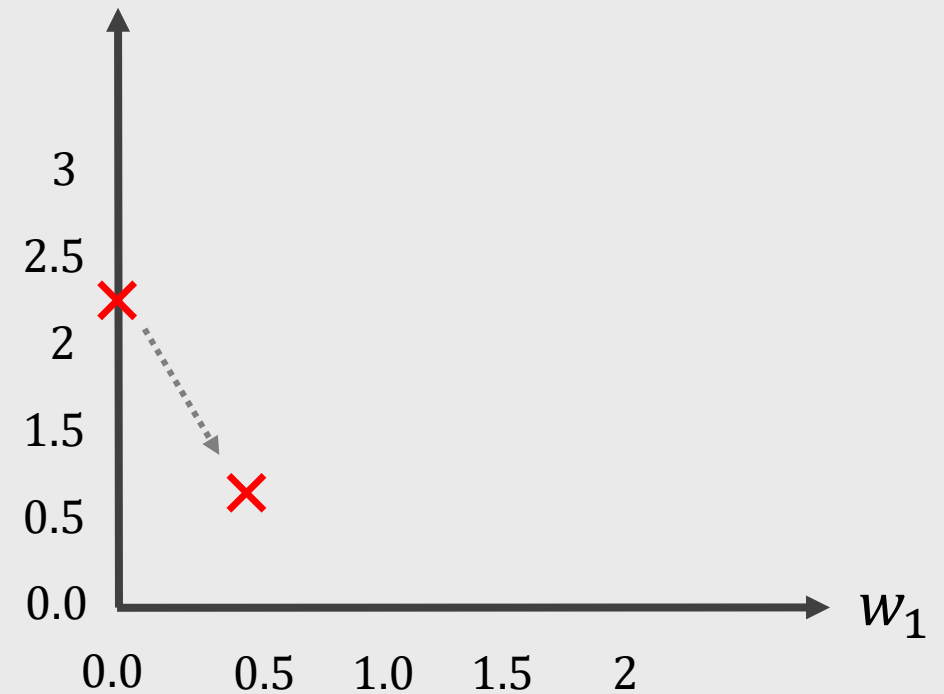
Note that y and x values are simplified to 1, 2 and 3



Model h_t for $w_0 = 0$ and $w_1 = 0.5$:

$$g_w(x_i) = w_0 + w_1 x_i \text{ for } i = 1, 2, 3$$

$$E(g_{w_1}(x)) = \frac{1}{2N} \sum_{j=1}^N (f(\mathbf{x}_j) - g_w(\mathbf{x}_j))^2$$



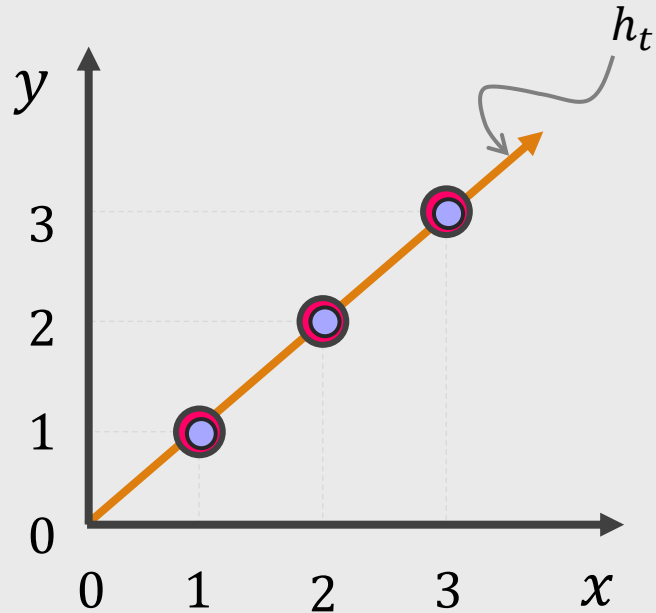
Error $E(w_1)$ for $w_0 = 0$ and $w_1 = 0.5$:

$$E(g_w(\mathbf{x})) = \frac{(1-0.5)^2 + (2-1)^2 + (3-1.5)^2}{2 \cdot 3} = 0.58$$

Does error $E(h_t(\mathcal{D}))$ minimization work?

9:37 AM

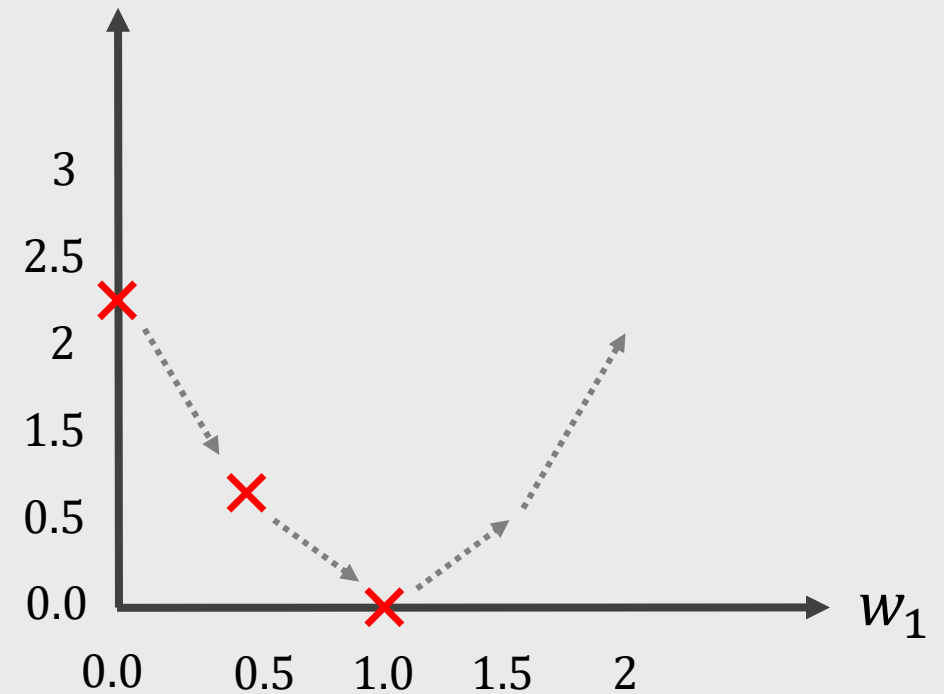
Note that y and x values are simplified to 1, 2 and 3



Model h_t for $w_0 = 0$ and $w_1 = 1$:

$$g_w(x_i) = w_0 + w_1 x_i \text{ for } i = 1, 2, 3$$

$$E(g_{w_1}(x)) = \frac{1}{2N} \sum_{j=1}^N (f(\mathbf{x}_j) - g_w(\mathbf{x}_j))^2$$

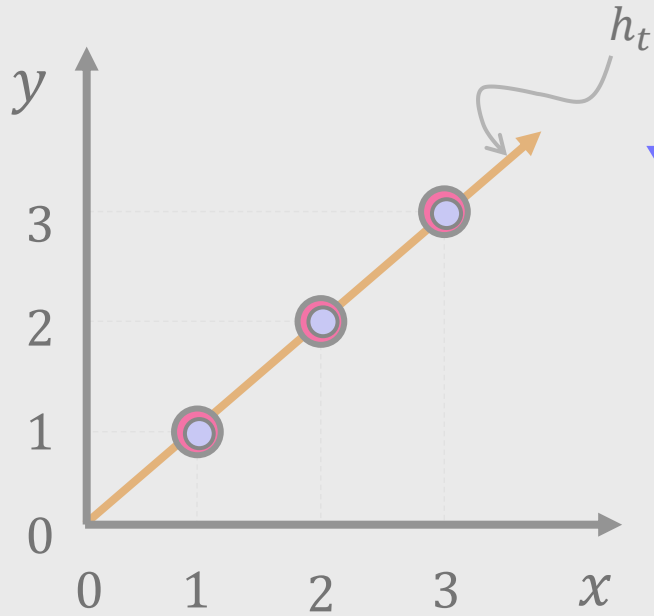


Error $E(w_1)$ for $w_0 = 0$ and $w_1 = 1$:

$$E(g_w(\mathbf{x})) = \frac{(1-1)^2 + (2-2)^2 + (3-3)^2}{2 \cdot 3} = 0.0$$

Does error $E(h_t(\mathcal{D}))$ minimization work?

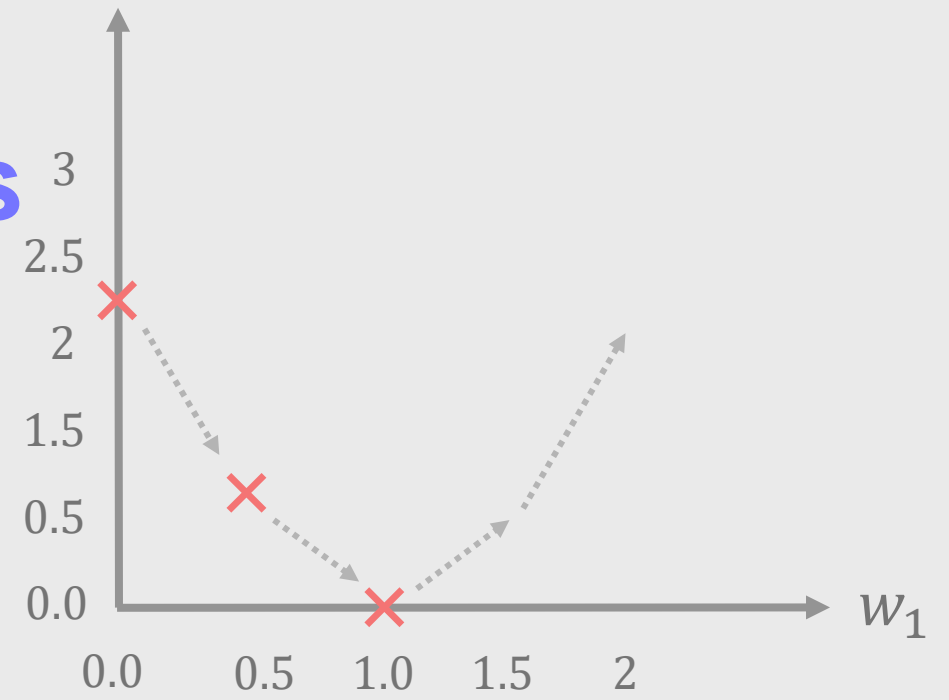
Note that y and x values are simplified to 1, 2 and 3



Yes, it does

But how much weight to change in each step?

$$E(g_{w_1}(x)) = \frac{1}{2N} \sum_{j=1}^N (f(\mathbf{x}_j) - g_w(\mathbf{x}_j))^2$$



Model h_t for $w_0 = 0$ and $w_1 = 1$:

$$g_w(x_i) = w_0 + w_1 x_i \text{ for } i = 1, 2, 3$$

Error $E(w_1)$ for $w_0 = 0$ and $w_1 = 1$:

$$E(g_w(\mathbf{x})) = \frac{(1-1)^2 + (2-2)^2 + (3-3)^2}{2 \cdot 3} = 0.0$$

Optimizer: Gradient Descent

Function g of the model has parameter \mathbf{w} :

$$g_{\mathbf{w}}(\mathbf{x}) = \sum_{i=0}^d w_i x_i = 0$$

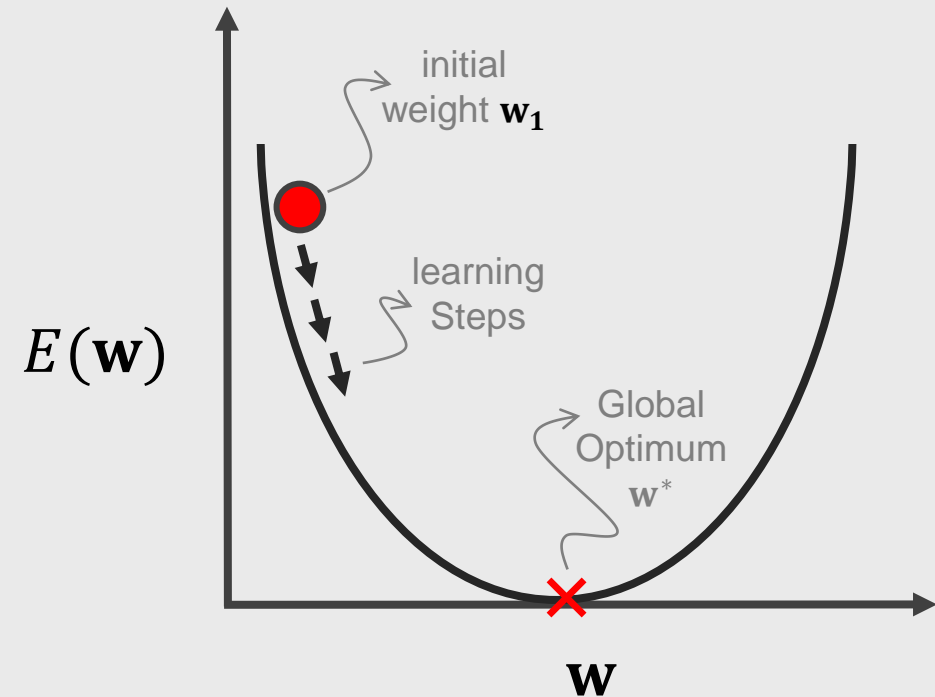
Gradient Descent Algorithm:

Repeat parameter \mathbf{w} update for $t = 2, 3, \dots, M$.

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \eta \frac{\partial E(g_{\mathbf{w}}(\mathbf{x}))}{\partial \mathbf{w}_t} \mathbf{x} \text{ for learning rate } \eta$$

$\mathbf{w}_t = \mathbf{w}_{t-1} + \hat{y} \mathbf{x}$

Until error rate $E(g_{\mathbf{w}}(\mathbf{x}))$ is acceptable or goes to zero.



Optimizer: Gradient Descent

Function g of the model has parameter \mathbf{w} :

$$g_{\mathbf{w}}(\mathbf{x}) = \sum_{i=0}^d w_i x_i = 0$$

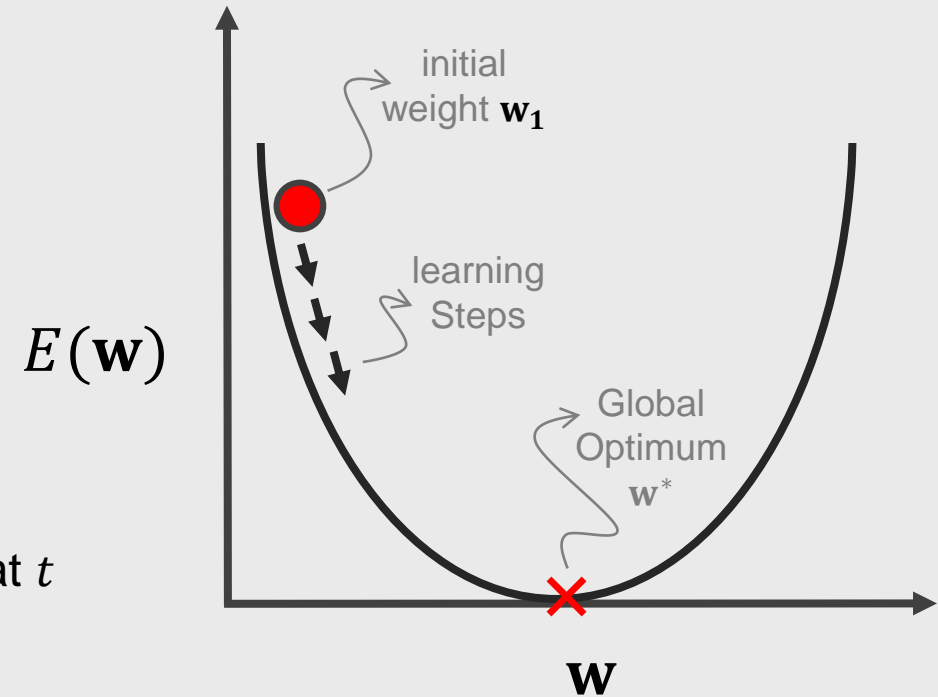
Gradient Descent Algorithm:

Repeat parameter \mathbf{w} update for $t = 2, 3, \dots, M$.

$\mathbf{w}_t = \mathbf{w}_{t-1} + \Delta \mathbf{w}_t$, where Δ is weight change (step) at t

$\nearrow \eta \frac{\partial E(g_{\mathbf{w}}(\mathbf{x}))}{\partial \mathbf{w}_t}$

Until error rate $E(g_{\mathbf{w}}(\mathbf{x}))$ is acceptable or goes to zero



Versions of Gradient Descent

Stochastic Gradient Descent

$t = 0$

\mathbf{w} initial weights

for t in epochs **do**

$\mathcal{D} \leftarrow \text{shuffle}(\mathcal{D})$

for $\mathbf{x}_j \in \mathcal{D}$ **do** // for each sample

$\nabla \mathbf{w}_j = \partial E(g_{\mathbf{w}_t}(\mathbf{x}_j)) / (\partial \mathbf{w}_t)$ // gradient of error *with respect to* weight \mathbf{w}_j

$$\mathbf{w}_j = \mathbf{w}_{j-1} + \eta \nabla \mathbf{w}_j \mathbf{x}_j$$

$t = t + 1$

Batch Gradient Descent

$t = 0$

\mathbf{w} initial weights

for t in epochs **do**

$\mathcal{D} \leftarrow \text{shuffle}(\mathcal{D})$

for $\mathbf{x}_j \in \mathcal{D}$ **do** // for each sample

$\nabla \mathbf{w} = \nabla \mathbf{w} + \partial E(g_{\mathbf{w}}(\mathbf{x}_j)) / (\partial \mathbf{w}) \mathbf{x}_j$ // gradient of error *with respect to* weight \mathbf{w}_j

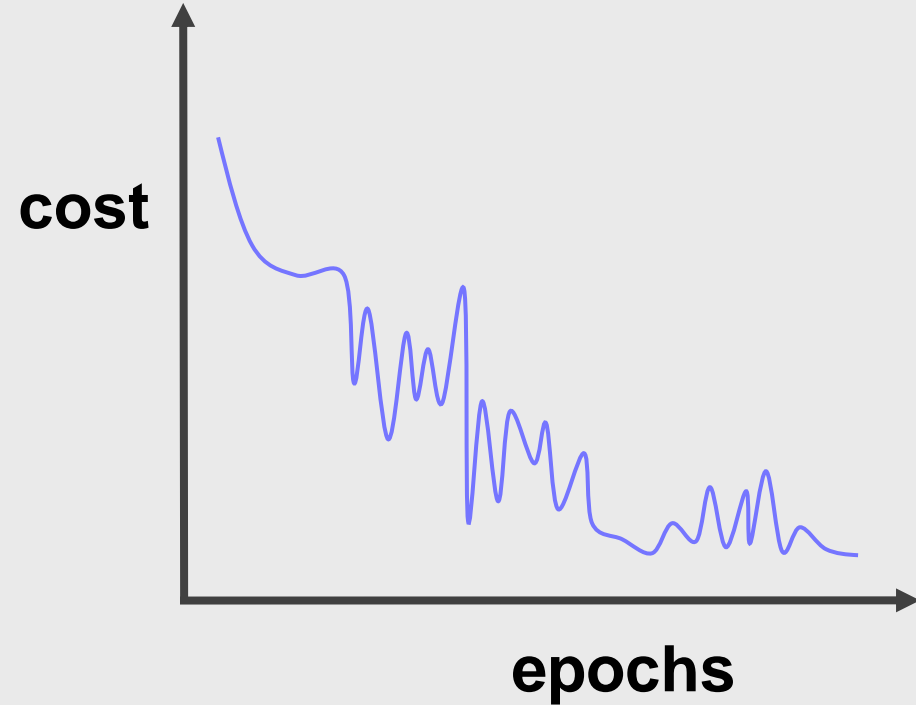
$$\mathbf{w}_t = \mathbf{w}_{t-1} + \eta \frac{\nabla \mathbf{w}}{|\mathcal{D}|}$$

$t = t + 1$

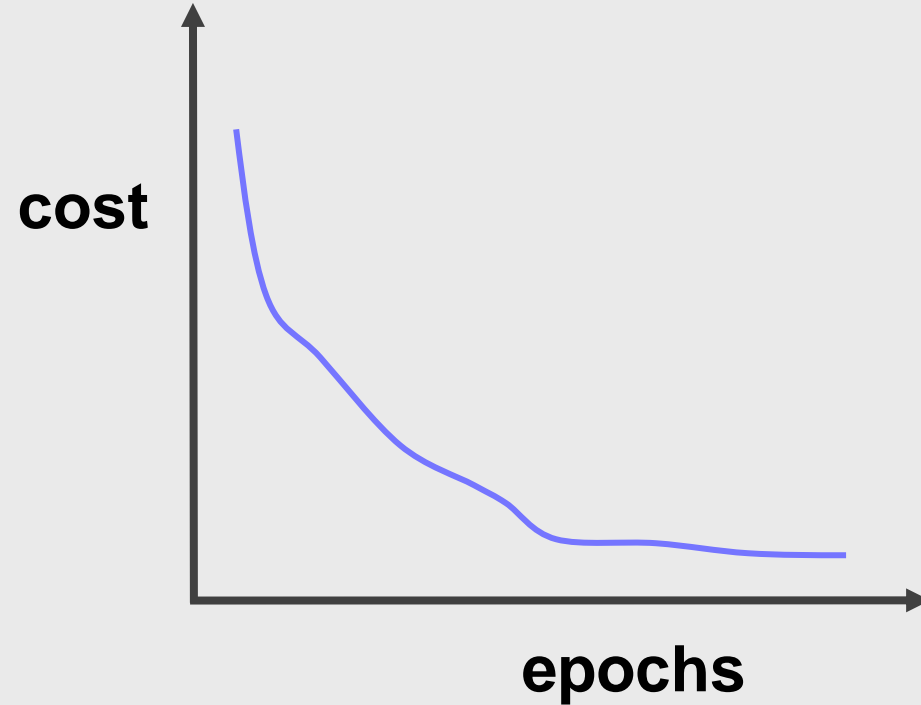
Gradient Descent: Versions

9:37 AM

Stochastic Gradient Descent



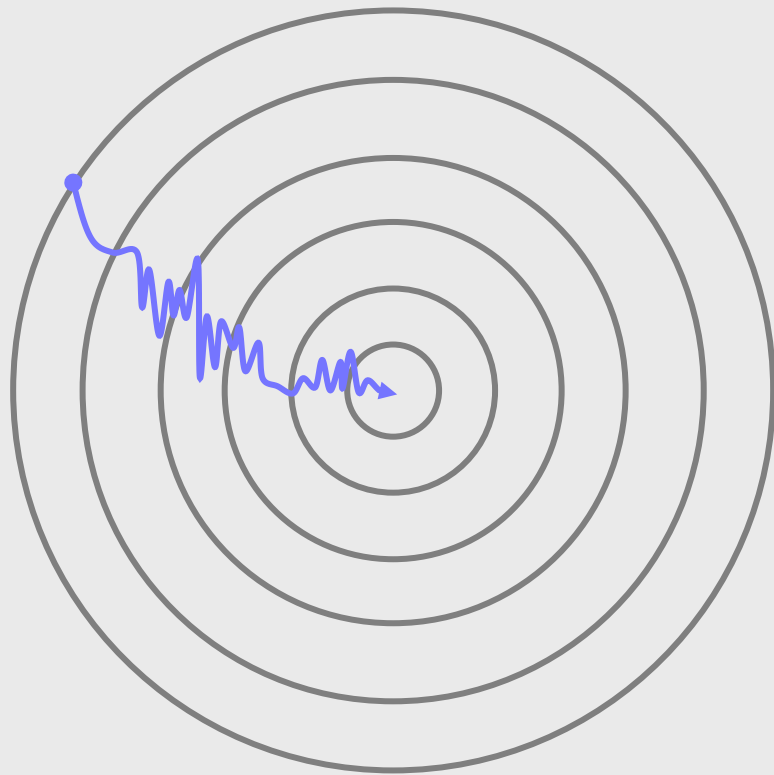
Batch Gradient Descent



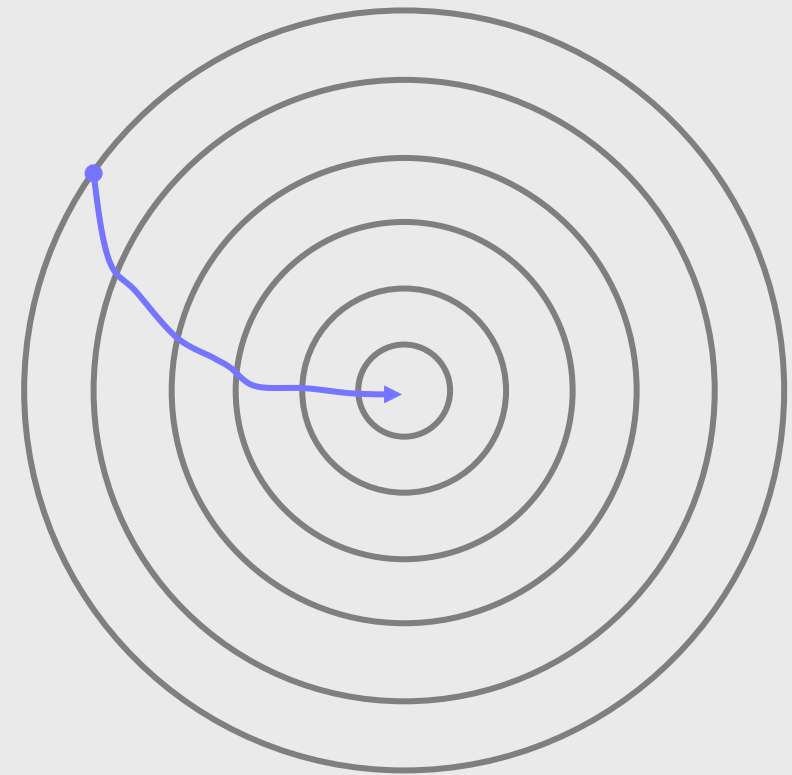
Gradient Descent: Versions

9:37 AM

Stochastic Gradient Descent



Batch Gradient Descent



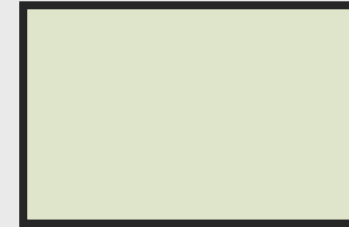
Training Method

(from previous lectures)

9:37 AM



Training Set

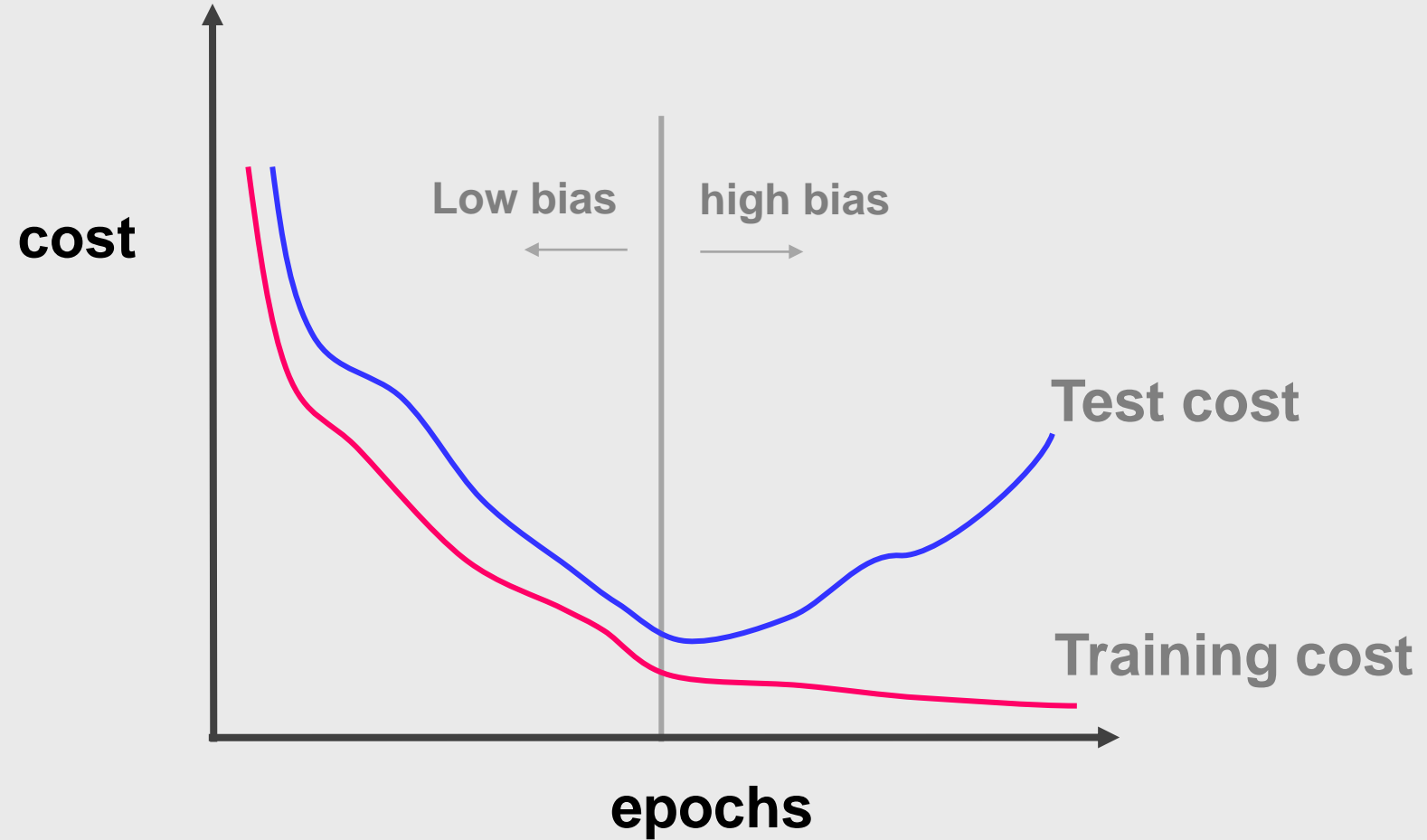


Test Set

Training Method

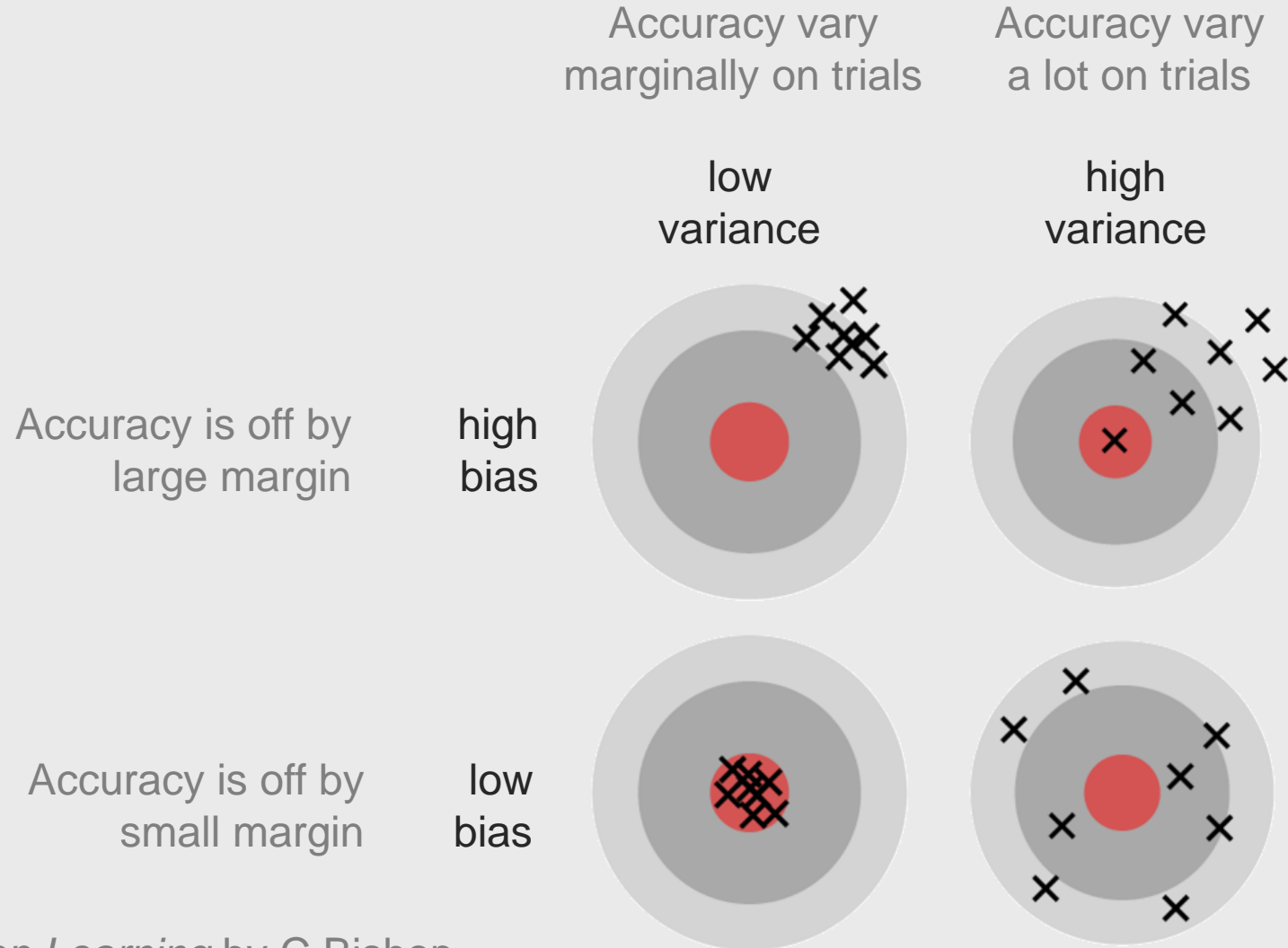
(from previous lectures)

9:37 AM



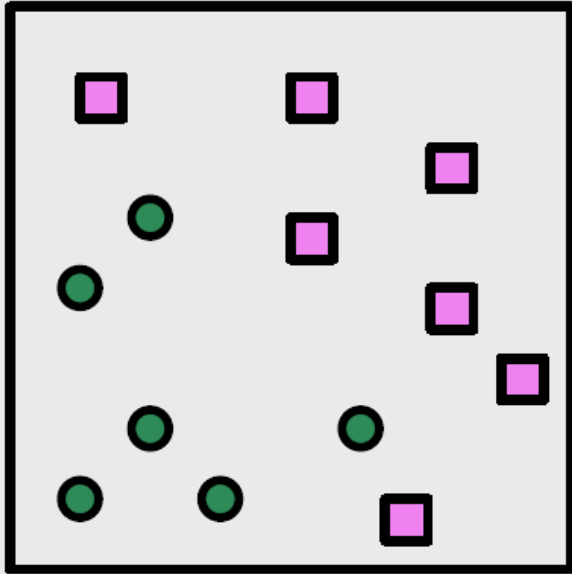
Bias-Variance Issue

(from previous lectures)

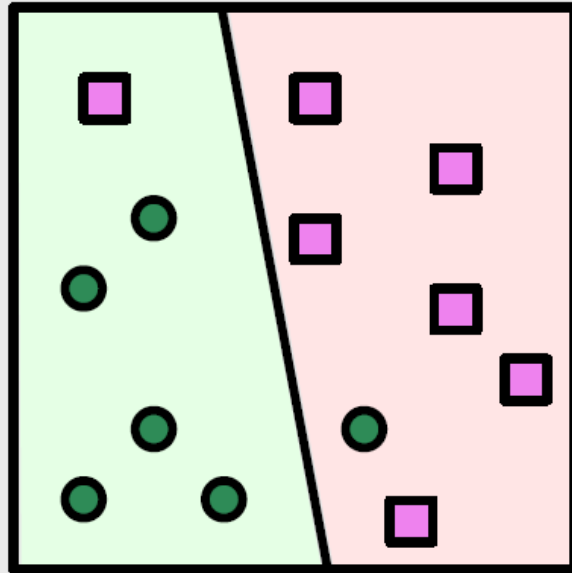


Is the chosen model good?

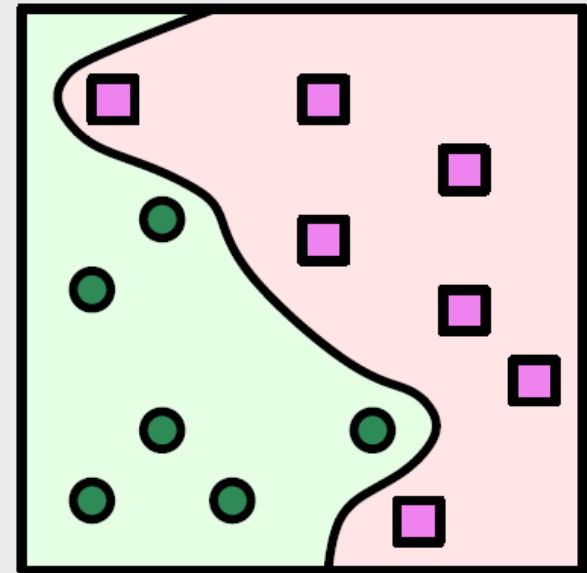
(from previous lectures)



Training data



Underfit



Overfit

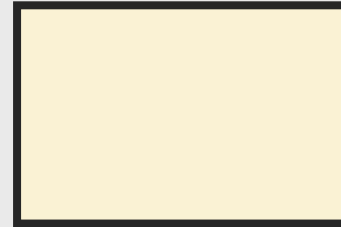
Avoid Overfitting

(from previous lectures)

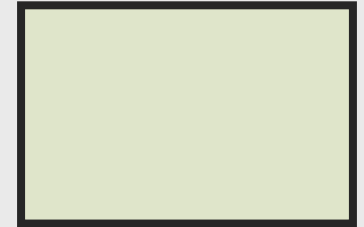
9:37 AM



Training Set



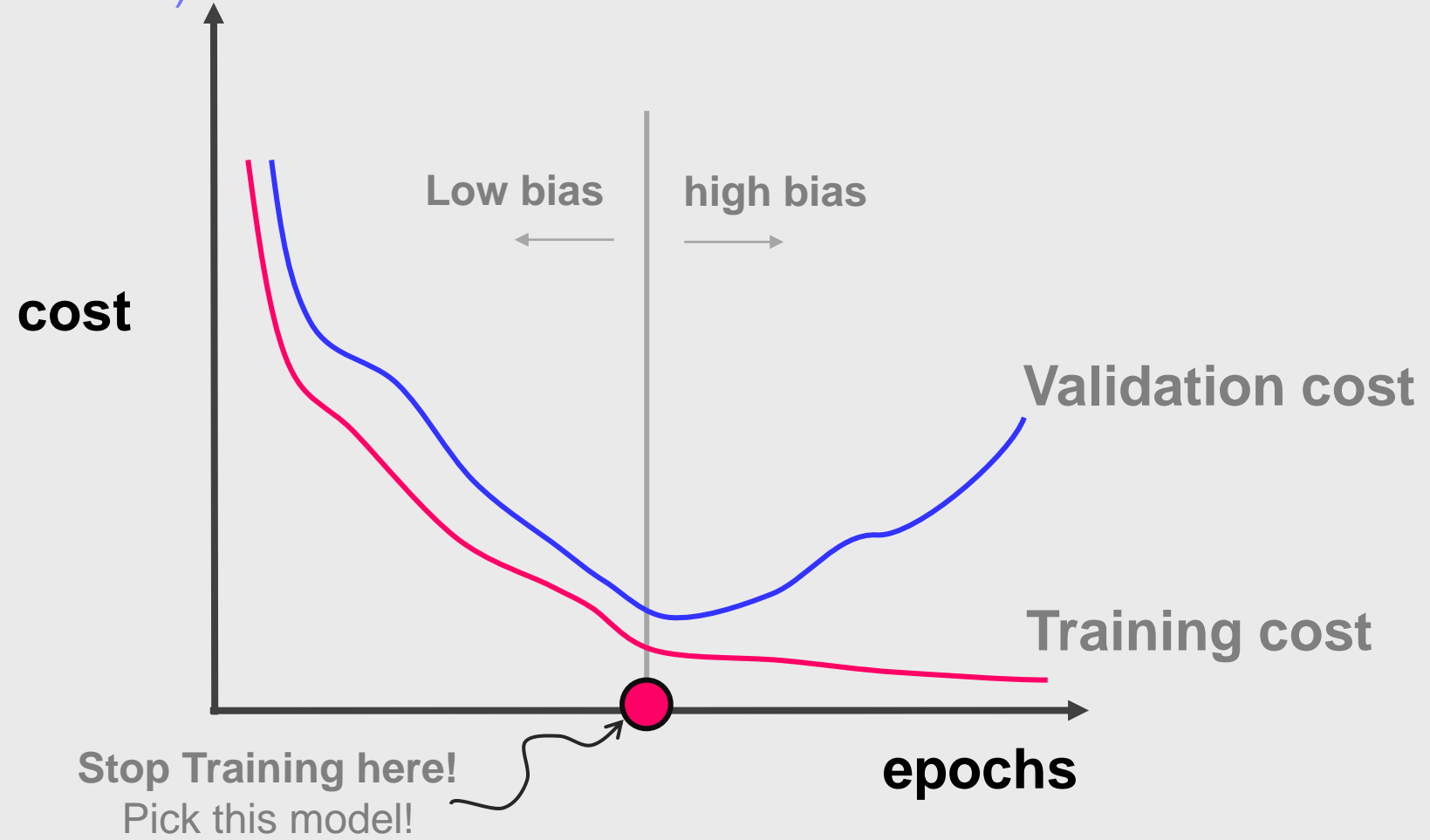
Validation Set



Test Set

Avoid Overfitting

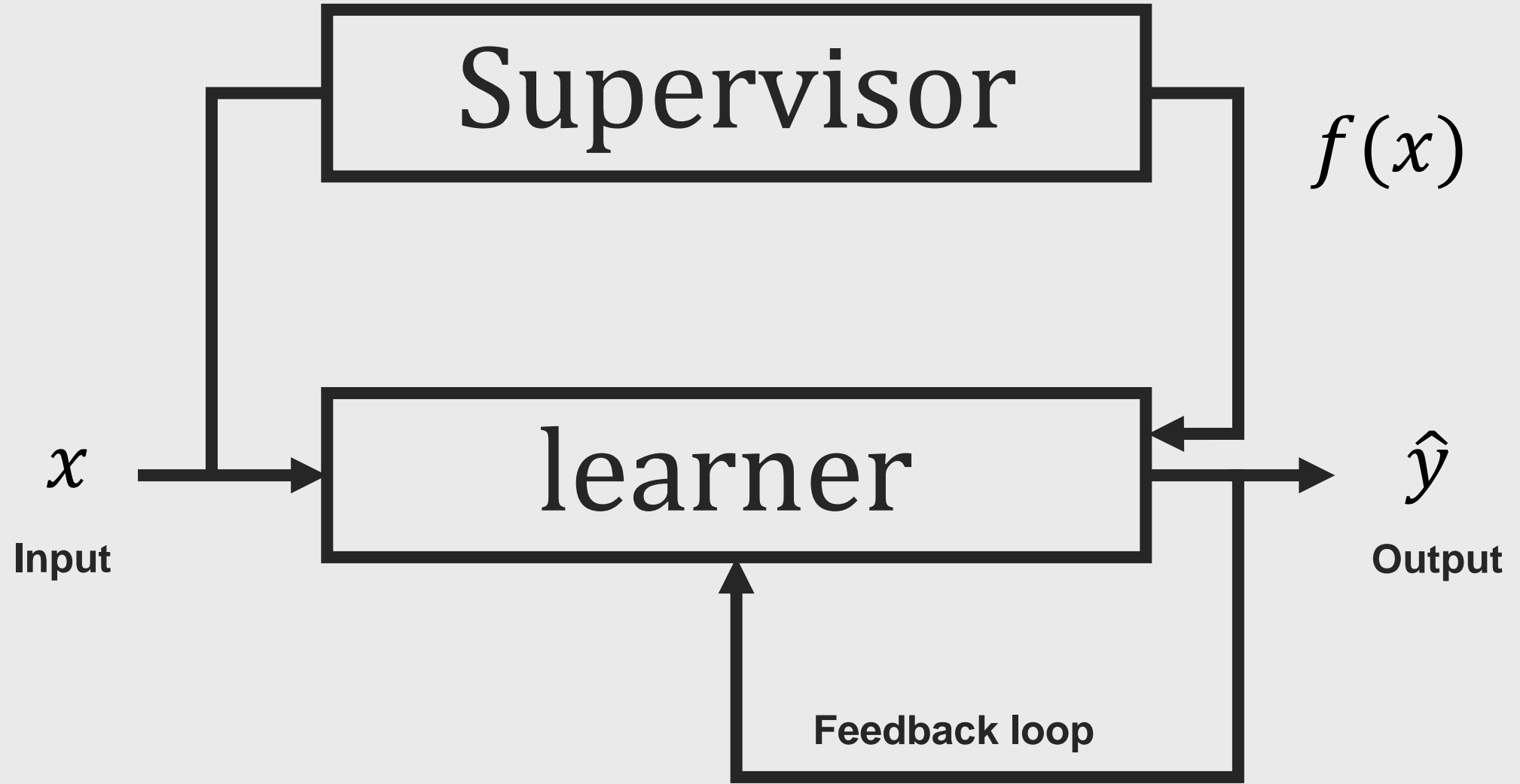
(from previous lecturers)



Part 3

Neural Network Architectures







Regression and Classification

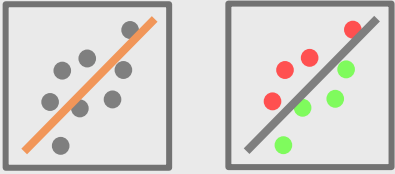
Class/Target attribute

	#	Inputs Attributes (Independent)		Target/Class/Output Attributes (Dependent)
		A1	A2	A3
Records	Ex. 0	A1 ₀	A2 ₀	A3 ₀
	Ex. 1	A1 ₁	A2 ₁	A3 ₁
	Ex. 2	A1 ₂	A2 ₂	A3 ₂
	Ex. 3	A1 ₃	A2 ₃	A3 ₃
	Ex. 4	A1 ₄	A2 ₄	A3 ₄
	Ex. 5	A1 ₅	A2 ₅	A3 ₅
	Ex. 6	A1 ₆	A2 ₆	A3 ₆
	Ex. 7	A1 ₇	A2 ₇	A3 ₇
	Ex. 8	A1 ₈	A2 ₈	A3 ₈
	Ex. 9	A1 ₉	A2 ₉	A3 ₉

Target (Class) Attributes (A3)

Regression
 Continuous
 (Numerical)
 labeled data

Classification
 Discrete
 (Categorical)
 labeled data



Tasks: Regression and Classification

Continuous labeled data

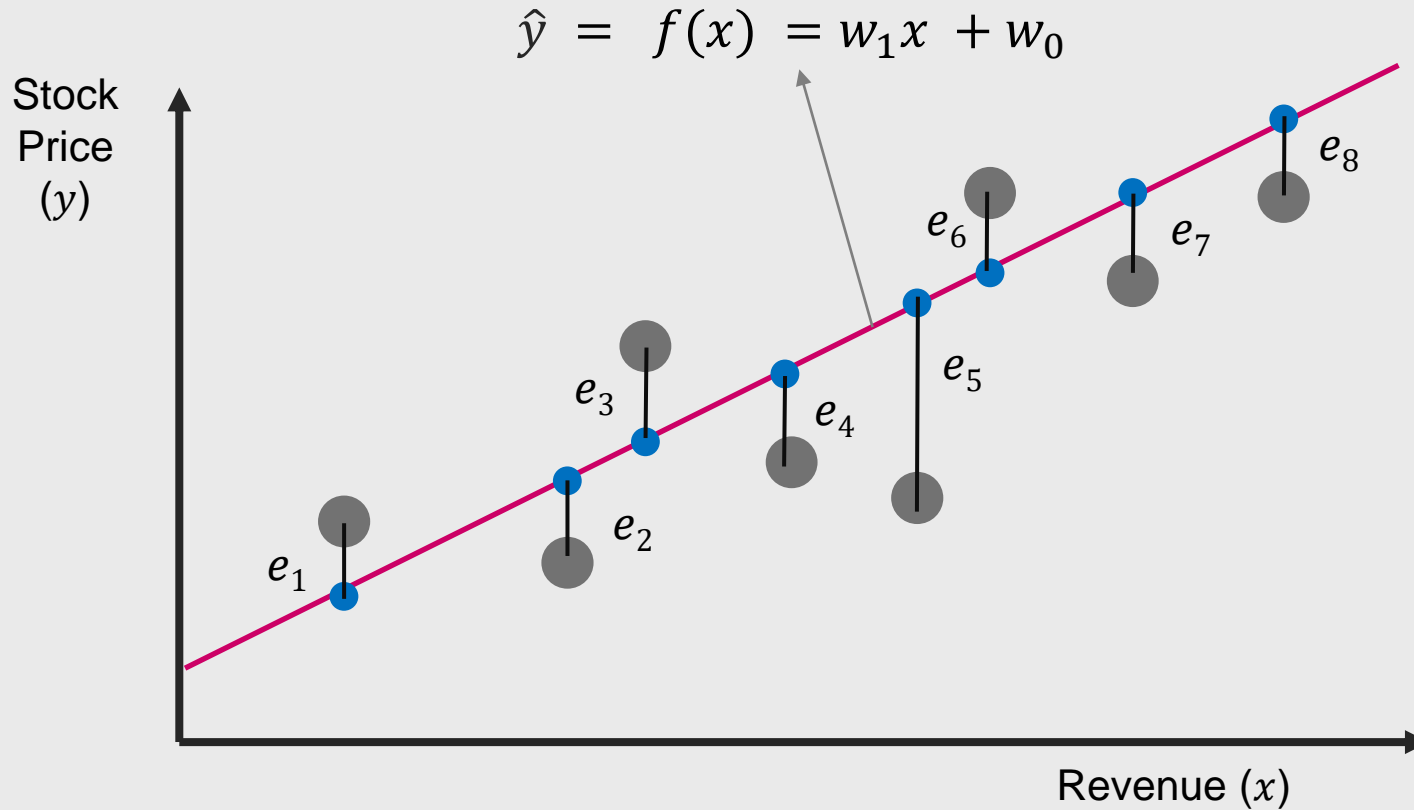
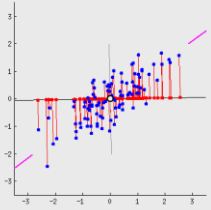
#	Inputs (X)		Target (Y)	
	Area (m ²)	Distance(mile)	Price (£Bn)	
Ex. 0	76.85	17.27	0.15	
Ex. 1	76.97	19.54	0.5	
Ex. 2	77.10	18.51	0.76	
Ex. 3	85.28	46.09	0.23	
Ex. 4	85.42	35.83	0.6	
Ex. 5	88.02	2.59	0.67	
Ex. 6	77.25	6.34	0.89	
Ex. 7	77.49	6.98	0.2	
Ex. 8	85.81	12.18	0.55	
Ex. 9	98.81	2.18	9.45	

Discrete labeled data

#	Inputs (X)		Class (Y)	
	Length (cm)	Weight (kg)	Sales	
Ex. 0	23.2	3.2	Good	Red
Ex. 1	70.9	19.5	Bad	Green
Ex. 2	60.5	18.51	Bad	Green
Ex. 3	24.5	4.6	Good	Red
Ex. 4	110.0	35.83	Bad	Green
Ex. 5	23.8	3.7	Good	Red
Ex. 6	25.8	4.5	Good	Red
Ex. 7	24.7	4.9	Good	Red
Ex. 8	85.8	25.6	Bad	Green
Ex. 9	78.8	20.33	Bad	Green



Regression



- ✓ Best Fit
- ✓ Find the line (parameters of a line equation) that minimize the norm of the y errors
- ✓ (sum of the squares)

Error
 $e_i = \hat{y}_i - y_i$

$$e = \sum_{i=1}^8 (\hat{y}_i - y_i)^2$$

Loss function: Mean Squared Error, E

$$E = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

\hat{y}_i - predicted output

y_i - target output

n - number of examples in training/test set

Loss function: Mean Absolute Error, E

$$E = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

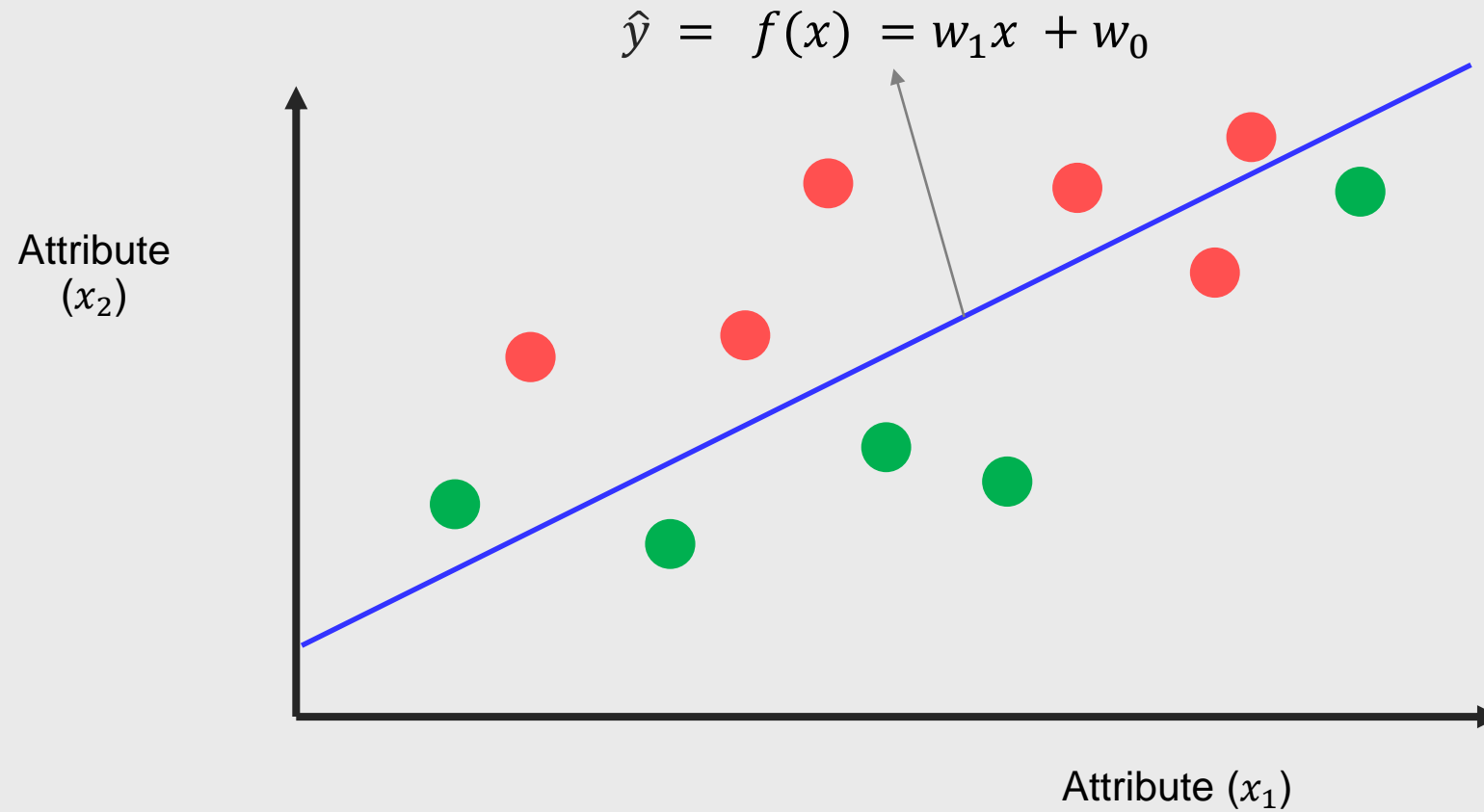
\hat{y}_i - predicted output

y_i - target output

n - number of examples in training/test set





Classification



- ✓ Best Fit
- ✓ Find the line (parameters of a line equation) that minimize the error (misclassification) rate

$$e = \frac{1}{n} \sum_{i=1}^n \hat{y}_i \neq y_i$$

Loss function: Misclassification rate, E

$$E = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i \neq y_i)$$


\hat{y}_i - predicted output

y_i - target output

n - number of examples in training/test set

Loss function: Log loss

It effectively works by updating network weights on correct classification and penalizing models for a misclassification

This part will be zero if $y_i = 0$

This part will be zero if $y_i = 1$

$$-\log P(y_i | \hat{y}_i) = -\left((y_i) \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right)$$

\hat{y}_i - predicted output

y_i - target output

n - number of examples in training/test set

Cross Entropy, E

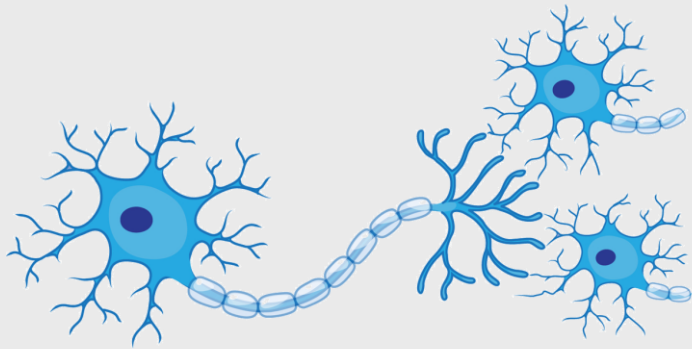
$$E = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

\hat{y}_i - predicted output distribution (SoftMax output)

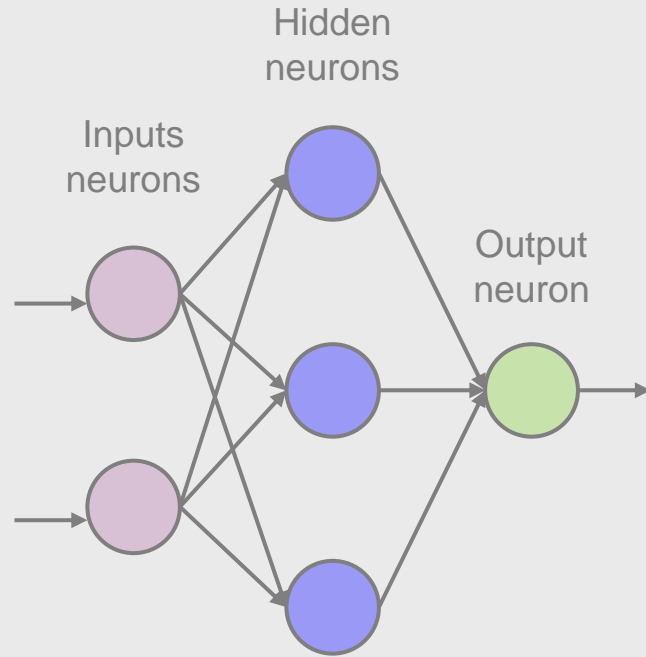
y_i - target output (target out distribution , one-hot encoding)

C - number of classes

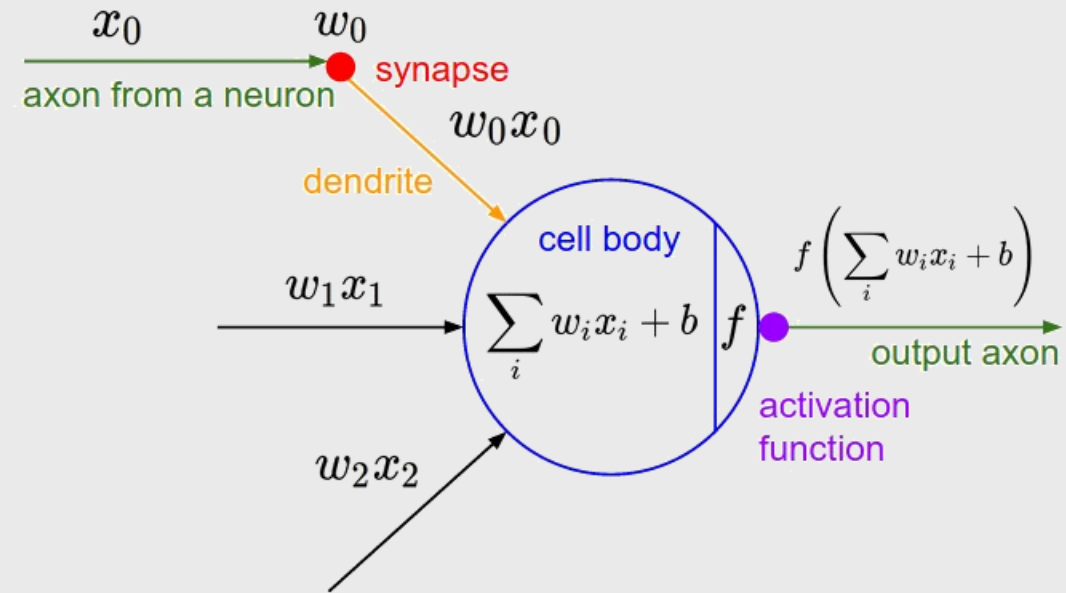
Neural Networks



Biological networks of neurons in human brains



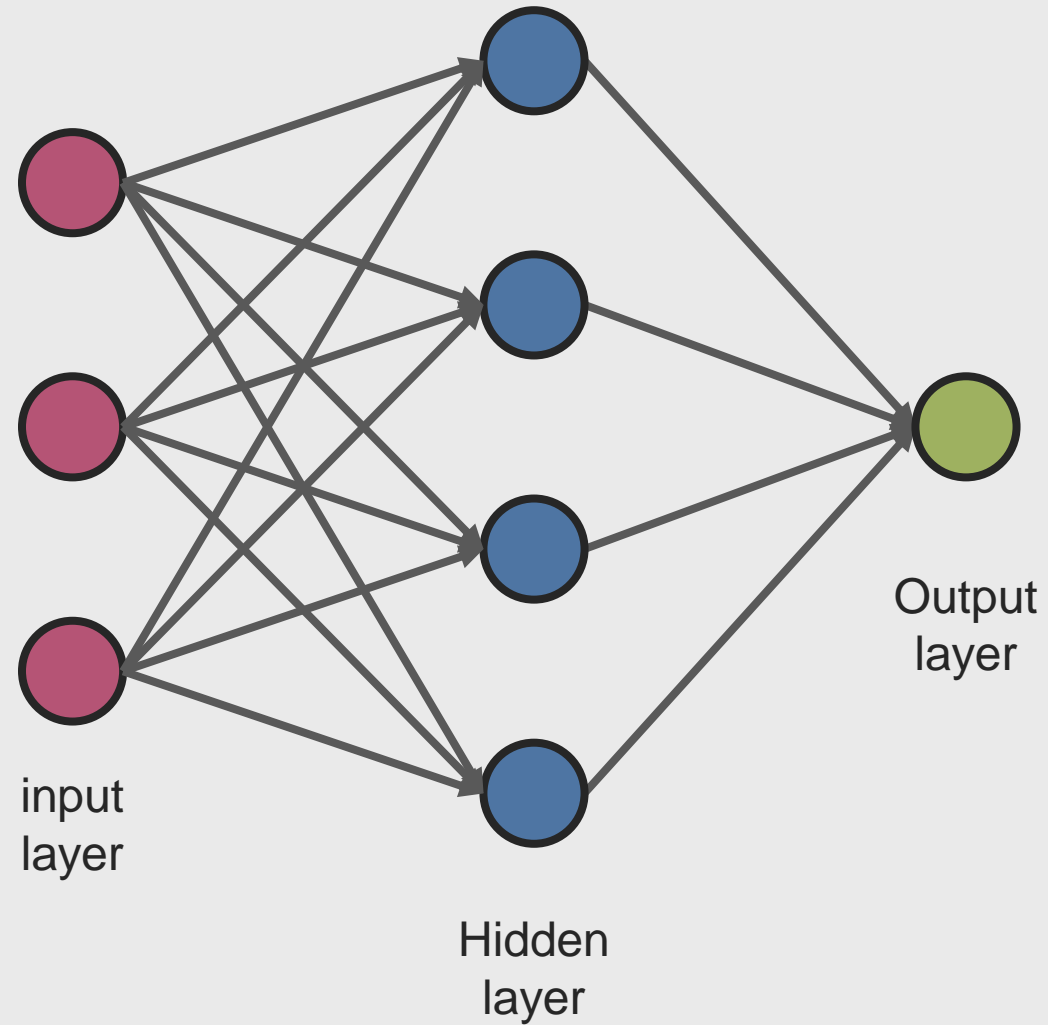
AI representation of biological neural networks



NEURAL NETWORK

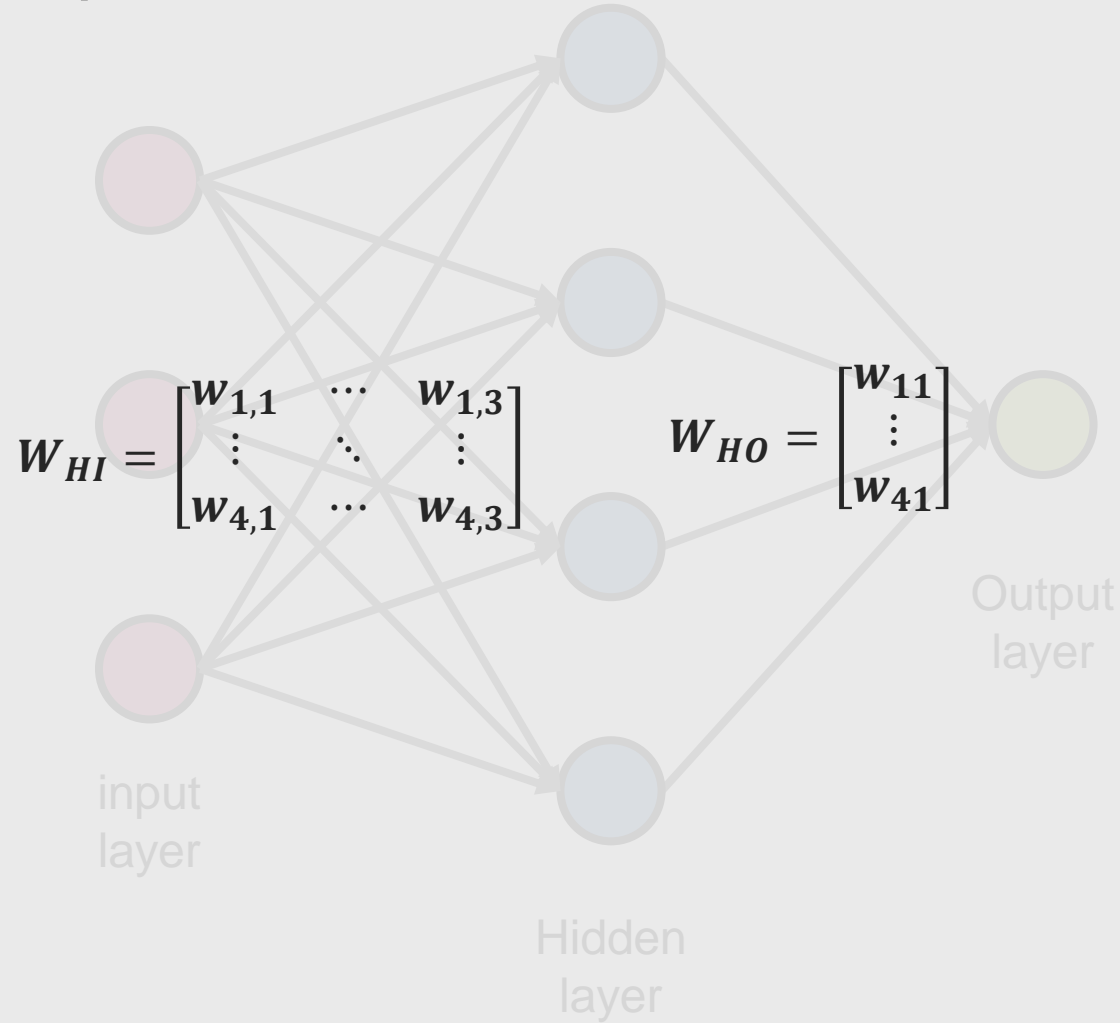
Architecture

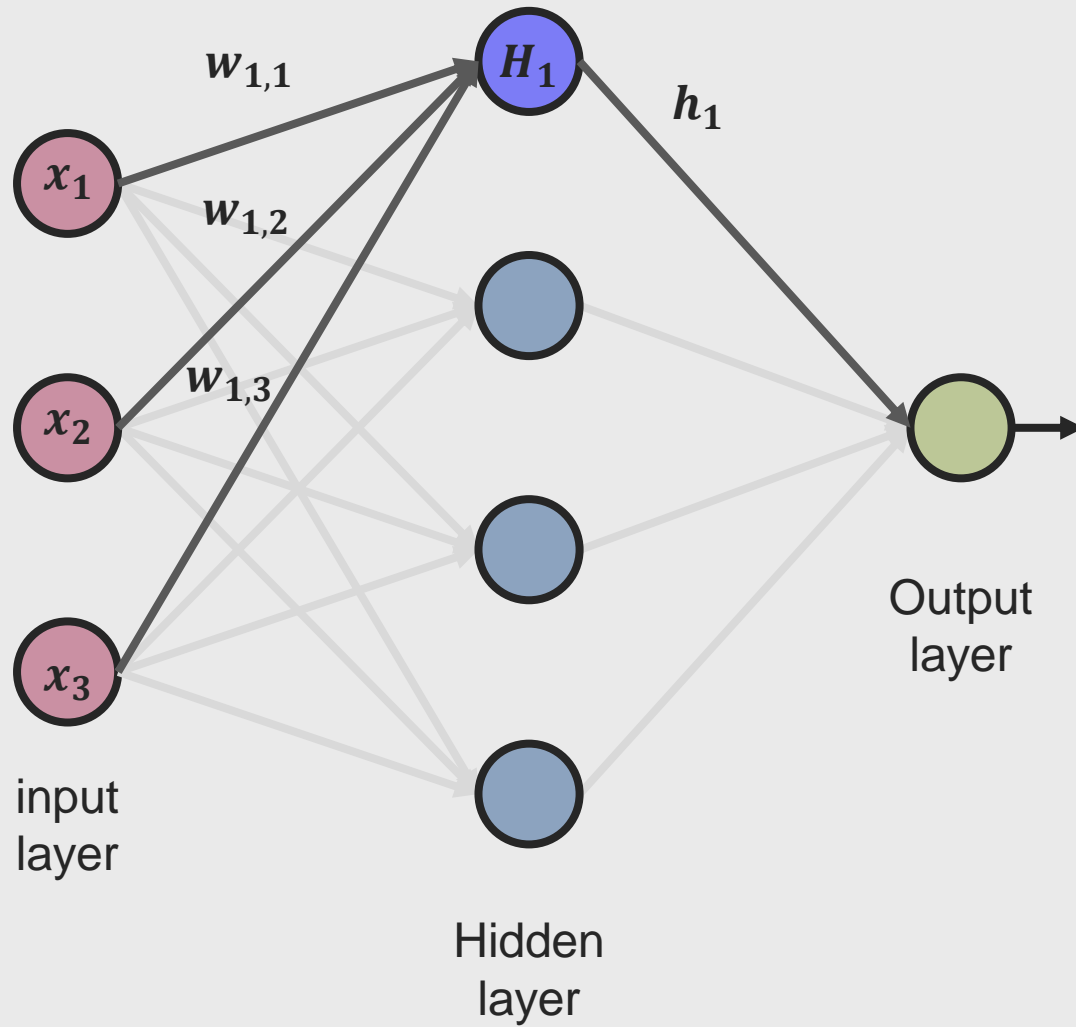
9:37 AM



NEURAL NETWORK

Weights (parameters)

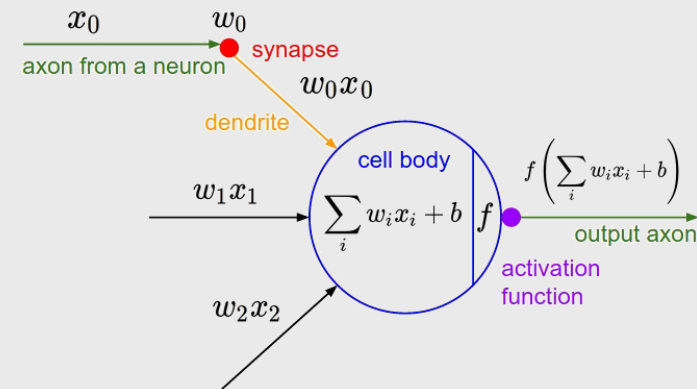




For n inputs, a hidden layer node's h_j output is expressed as:

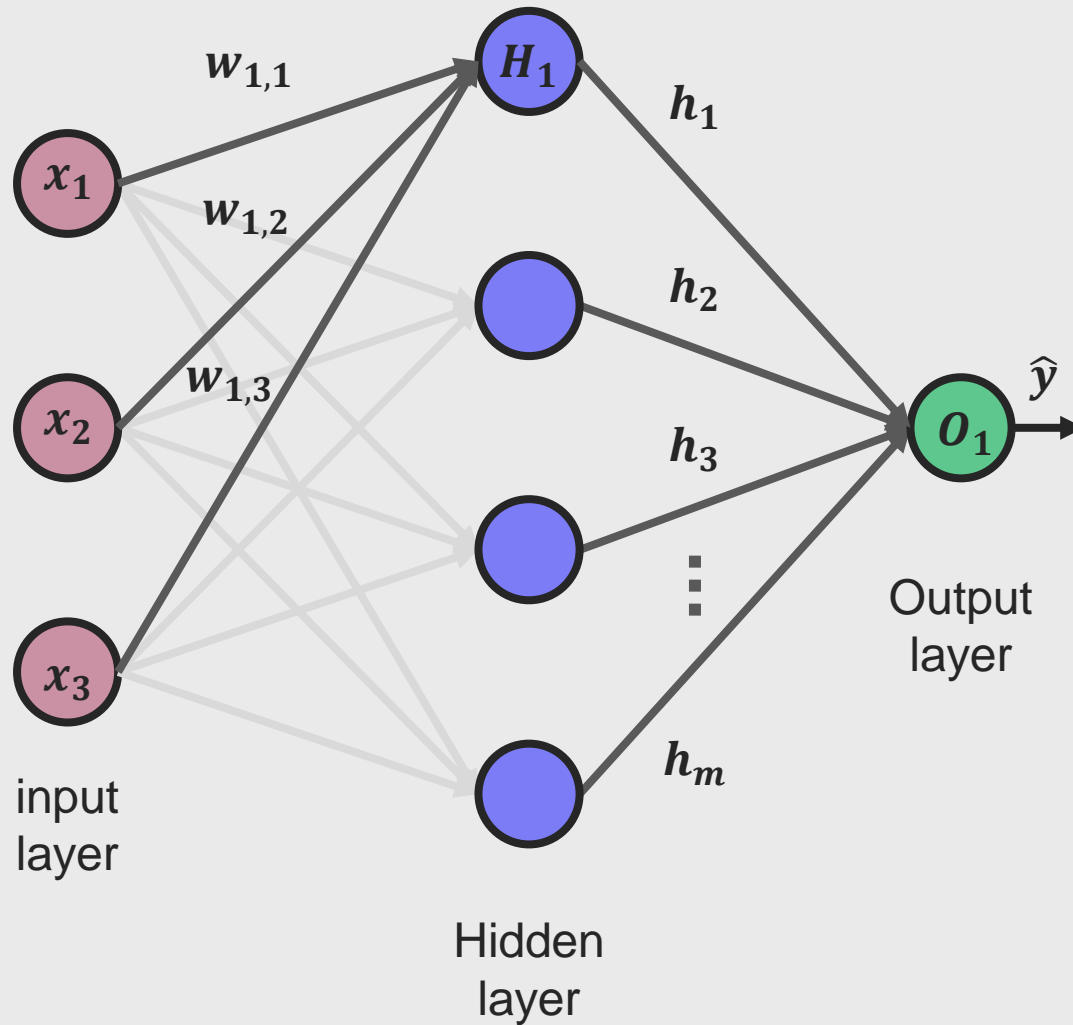
$$h_j = \varphi_h \left(\sum_{i=1}^n w_{ji} \cdot x_i \right)$$

Where φ_h is an activation function:



NEURAL NETWORK

Computation: Hidden layer



For n inputs, a hidden layer node's h_j output is expressed as:

$$h_j = \varphi_h \left(\sum_{i=1}^n w_{ji} \cdot x_i \right)$$

Where φ_h is an activation function:

For m hidden nodes and an output node, the output nodes output is expressed as:

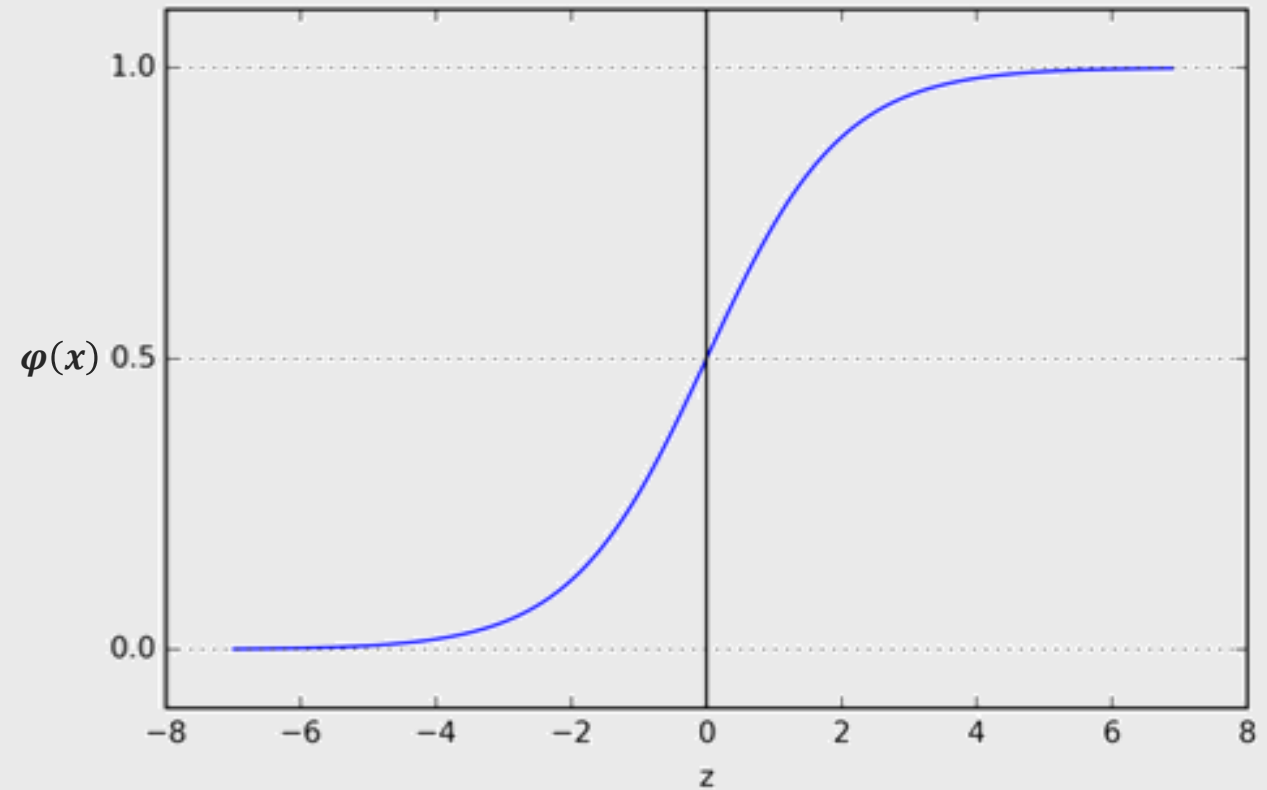
$$\hat{y} = \varphi_o \left(\sum_{j=1}^m w_{jk} \cdot h_j \right)$$

NEURAL NETWORK

Computation: Output layer

Sigmoid activation

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

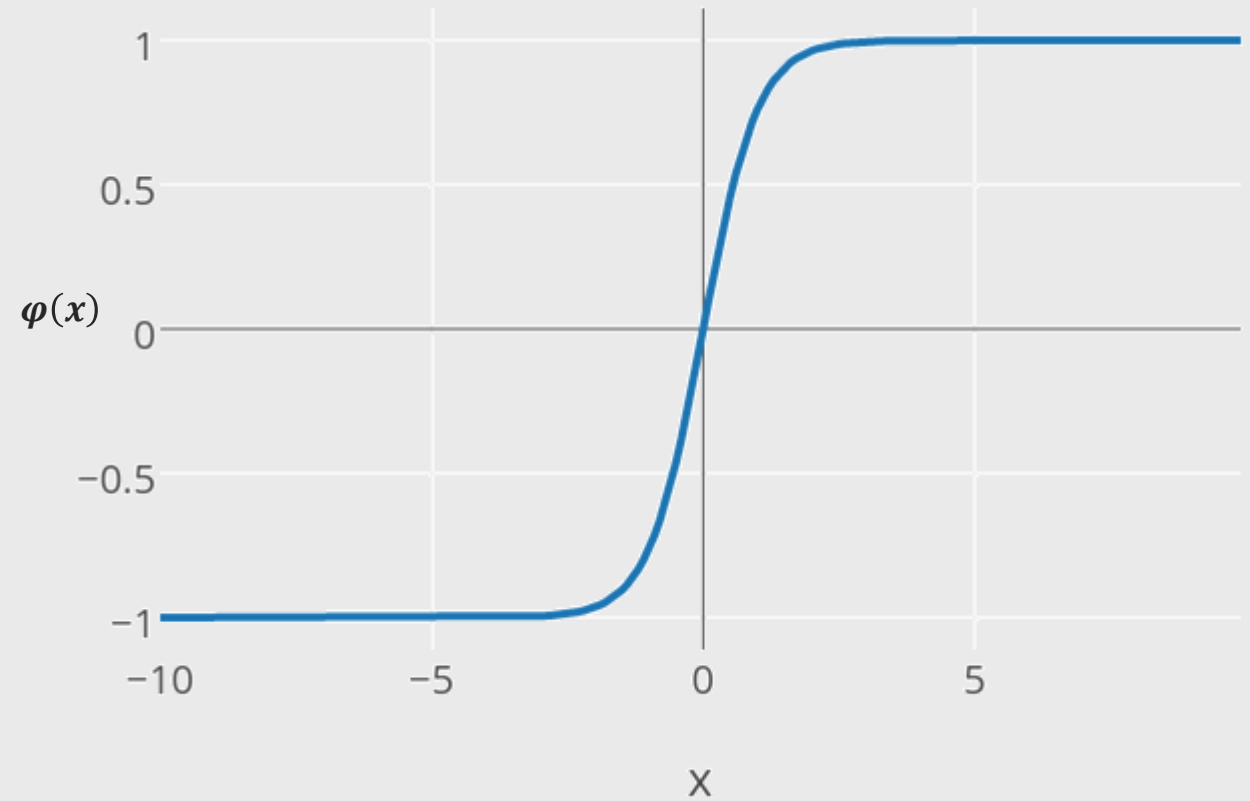


NEURAL NETWORK

Activation function

Tangent hyperbolic activation

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

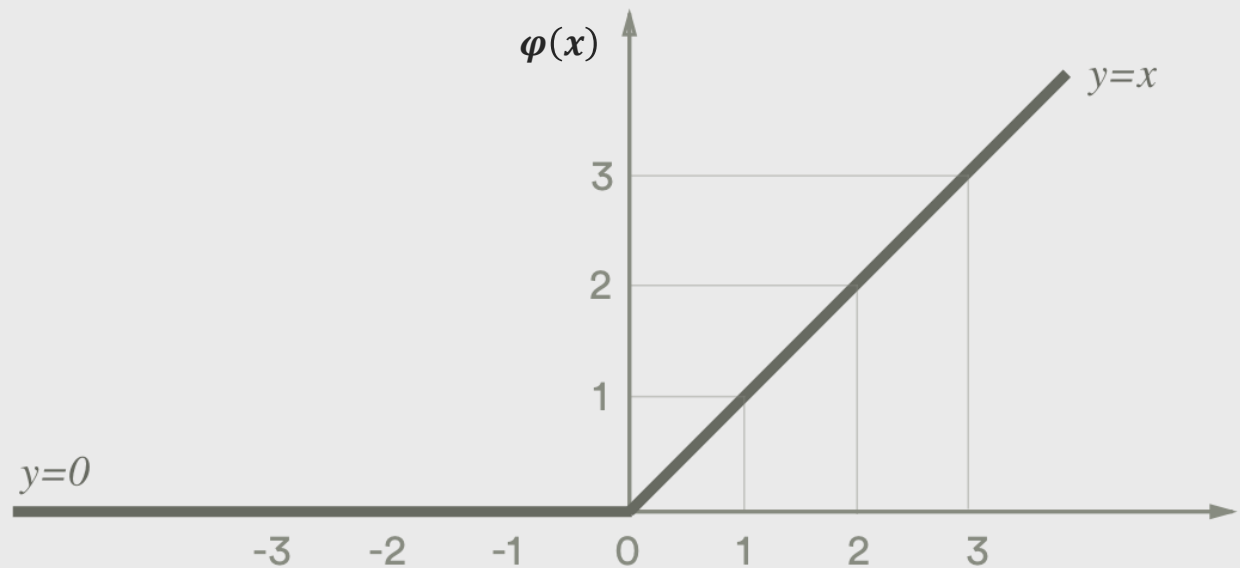


NEURAL NETWORK

Activation function

Rectified Linear Unit (ReLU)

$$\varphi(x) = \max(0, x)$$



Source: <https://medium.com/@danqing/a-practical-guide-to-relu-b83ca804f1f7>

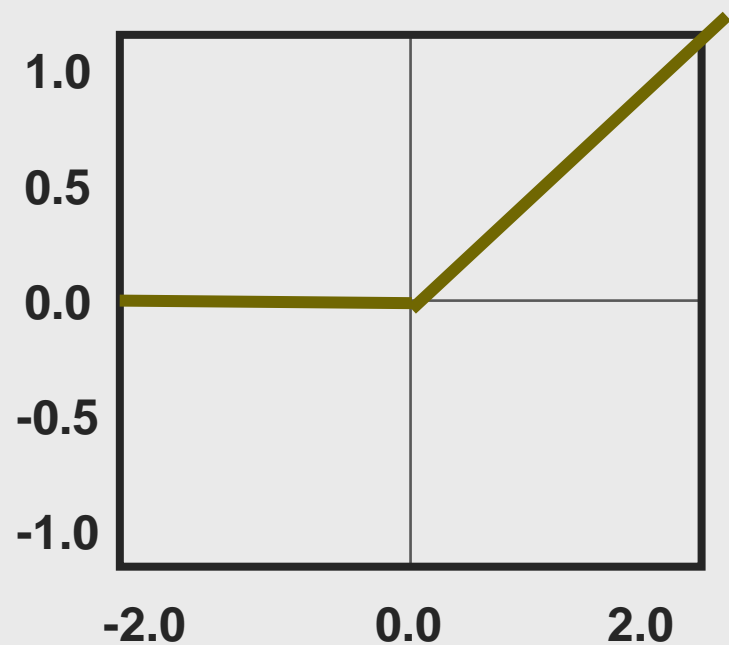
NEURAL NETWORK
Activation function

NEURAL NETWORK

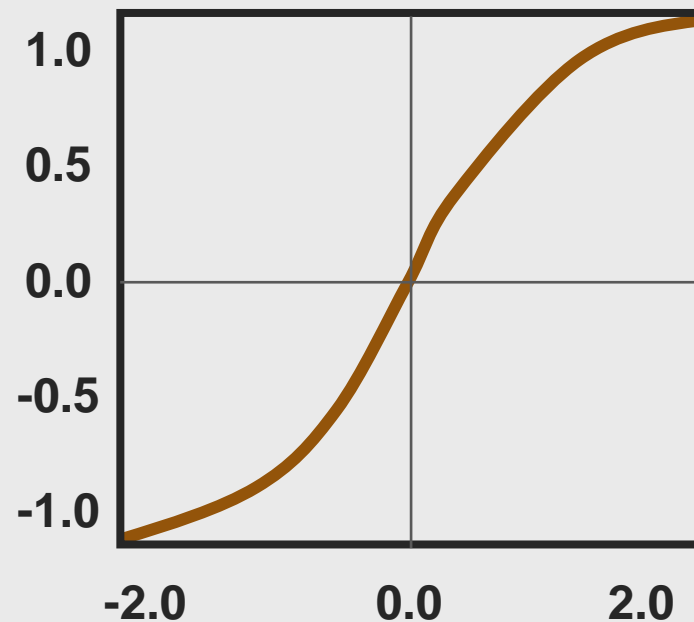
Activation functions

9:37 AM

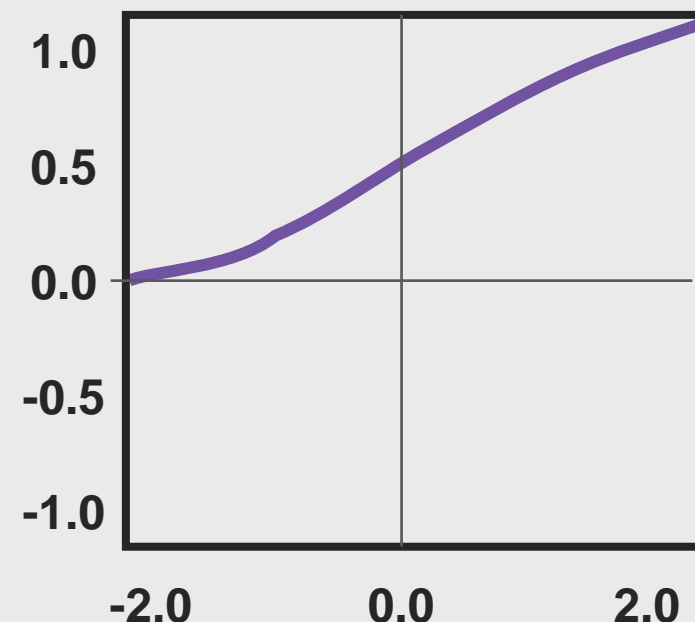
ReLU



Tanh

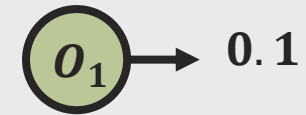


Sigmoid



SoftMax Activation

$$\varphi(x_i) = \frac{e^{x_i}}{\sum_j^k e^{x_j}} \text{ for } k \text{ units}$$

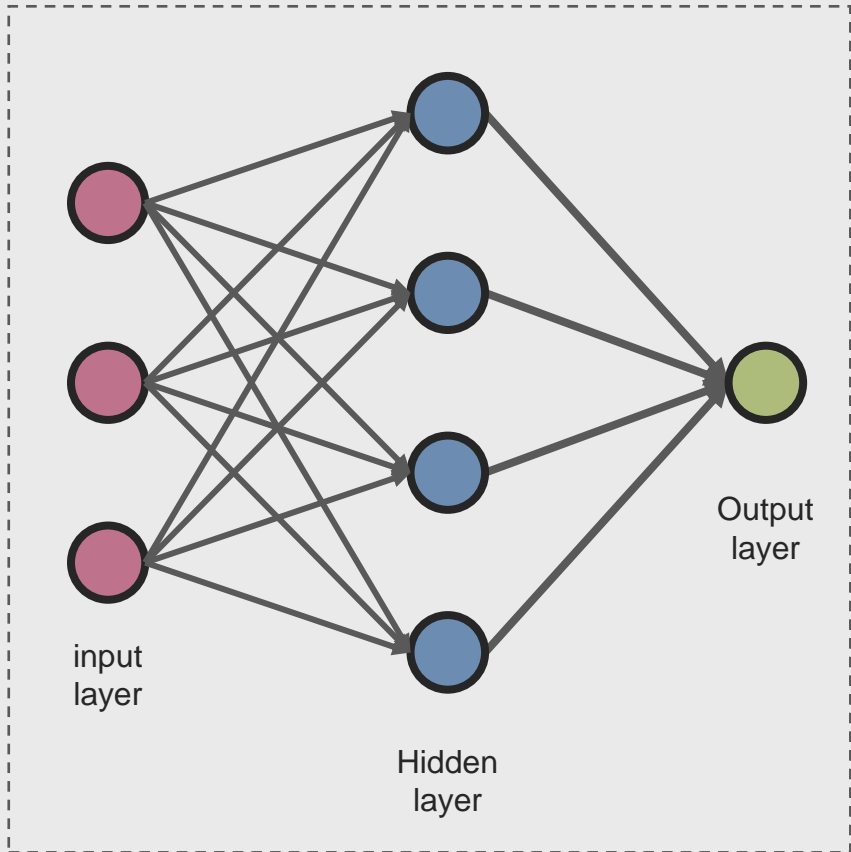


PROBABILITIES
DISTRIBUTION OF ALL
LABELS

NEURAL NETWORK
Activation function

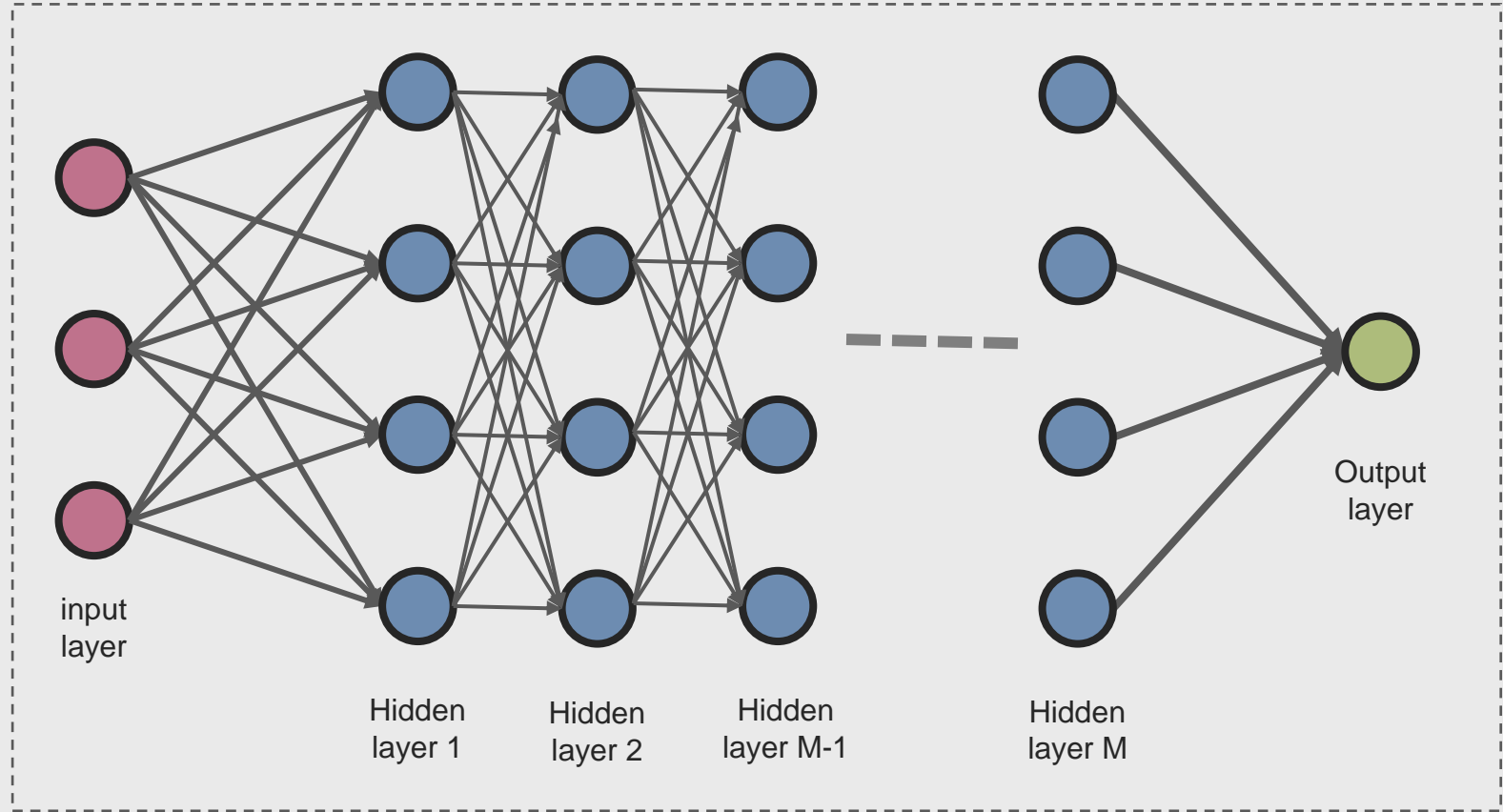
NEURAL NETWORK: Architecture

A regular neural network architecture



SHALLOW LEARNING

A deep neural network architecture



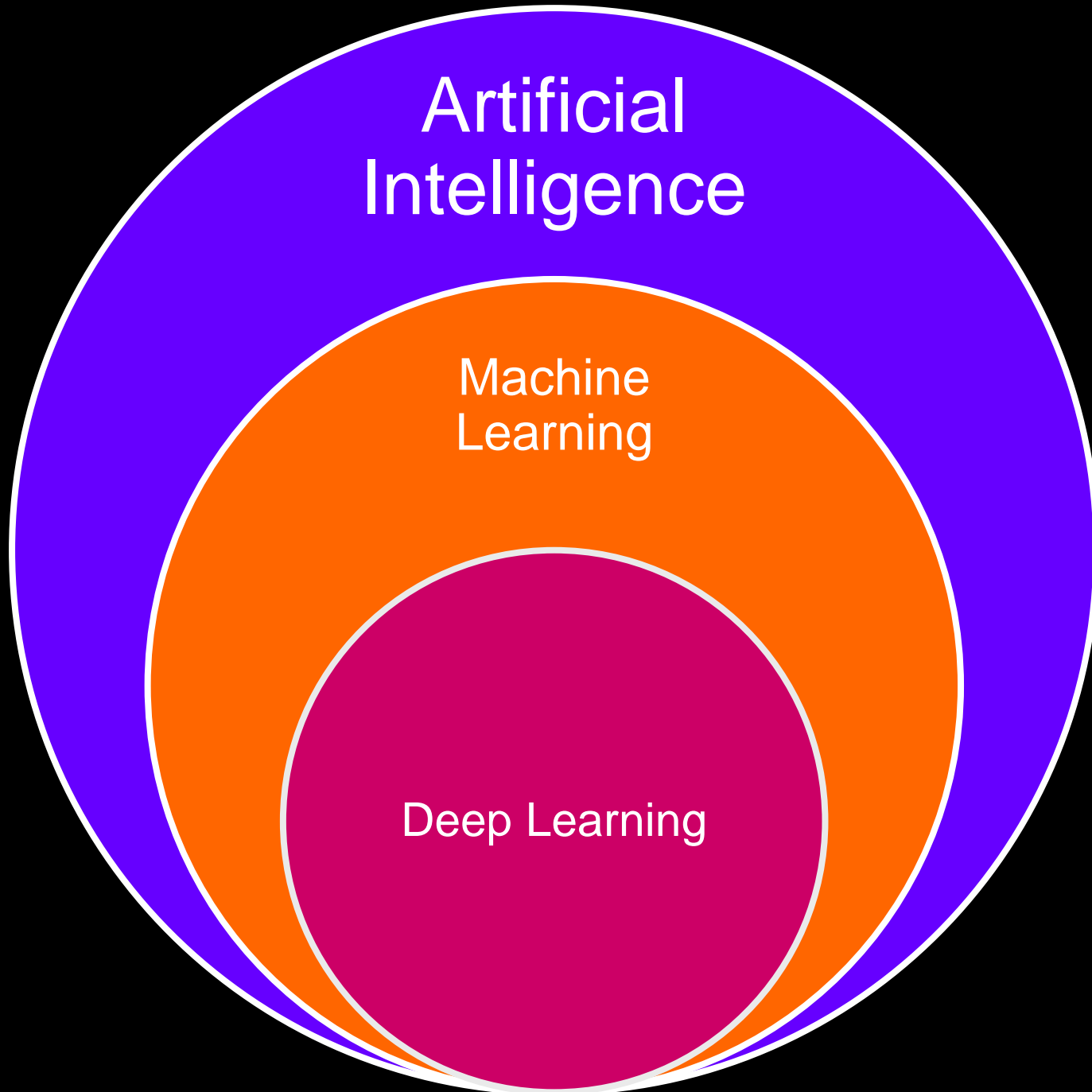
DEEP LEARNING

Deep Learning

Deep learning is an artificial intelligence function that imitates the workings of the human brain in processing data and creating patterns for use in decision making.

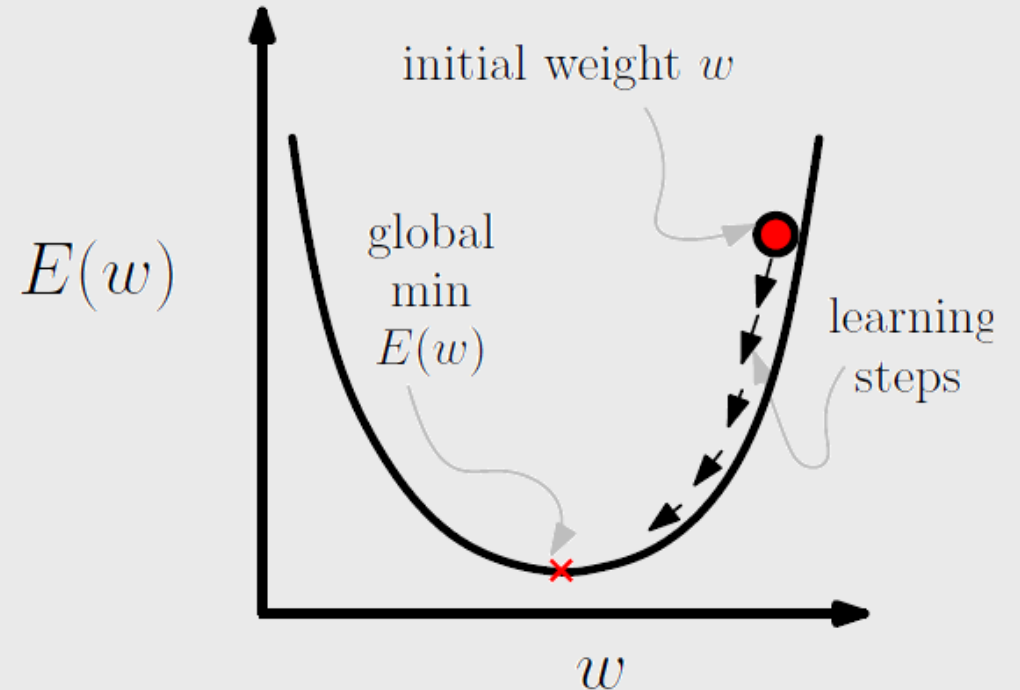
Deep learning is a **subset of machine learning in artificial intelligence** that has networks capable of learning supervised/unsupervised from data that is structured/unstructured or labelled/unlabelled.

Source: <https://www.investopedia.com/terms/d/deep-learning.asp>



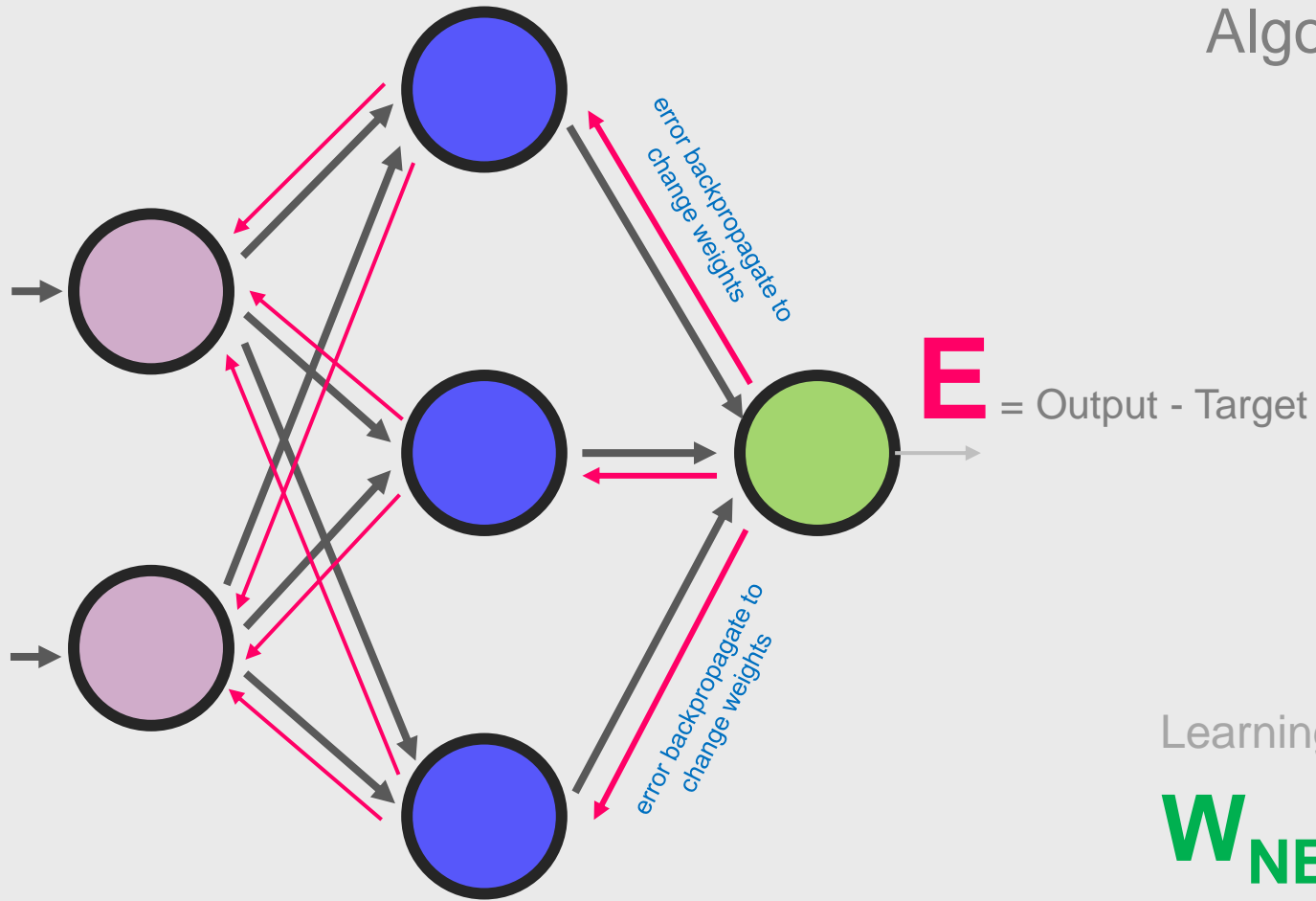
(DEEP) NEURAL NETWORK Optimisation

- Stochastic gradient descent
- Mini-batch gradient descent
- Batch gradient descent
- Backpropagation



BACKPROPAGATION

Algorithm for Updating Learning Systems

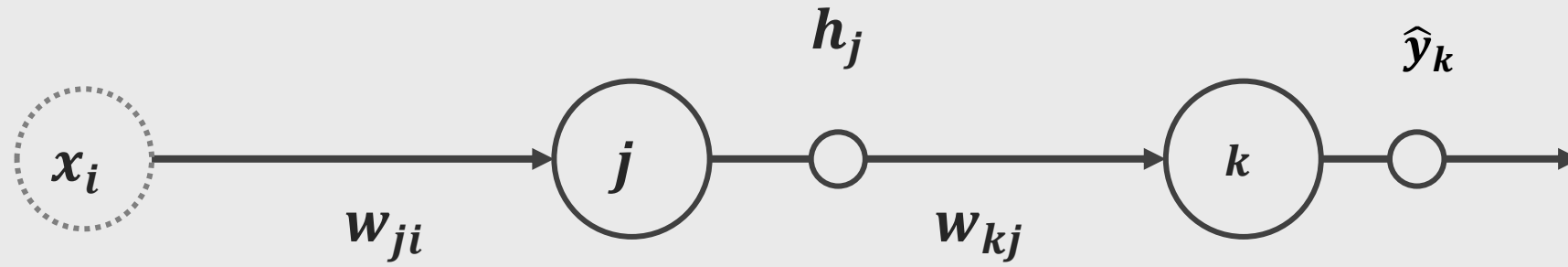


Learning Systems Update

$$W_{\text{NEW}} \leftarrow W_{\text{OLD}} + W_{\text{CHANGE}}$$

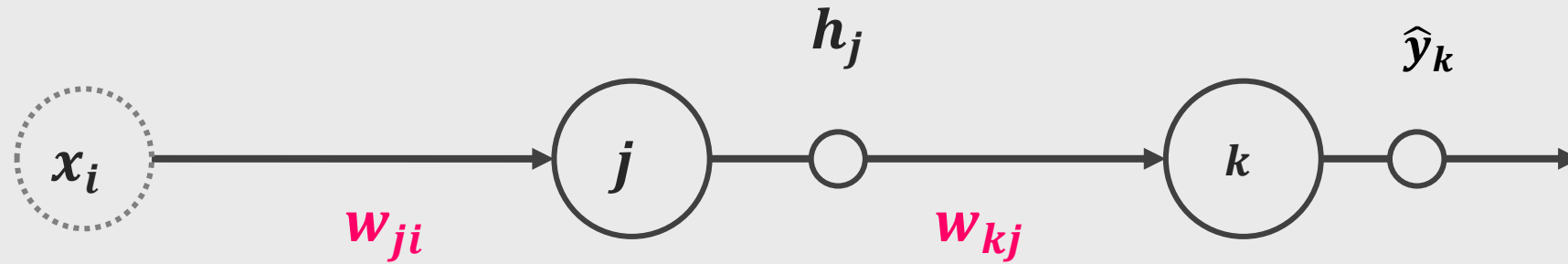
Backpropagation Algorithm

9:37 AM

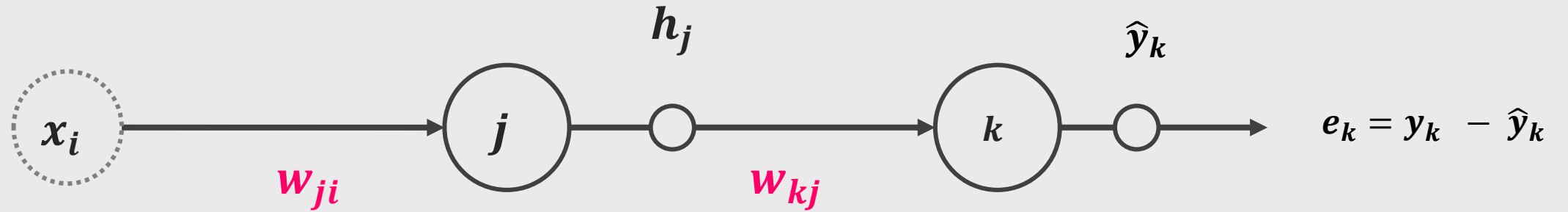


Backpropagation: Forward Pass

9:37 AM

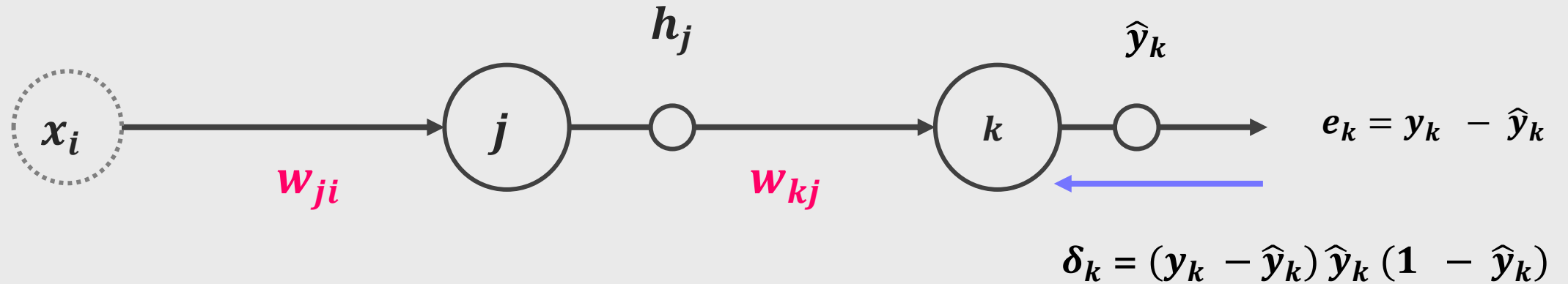


Backpropagation: Error at Output layer 9:37 AM



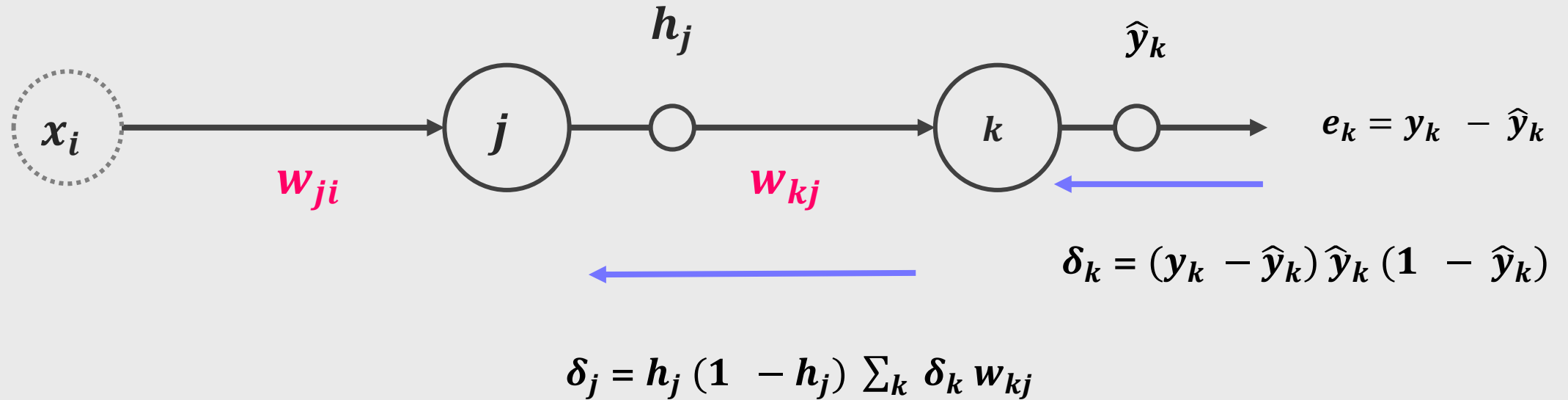
Backpropagation: Backward pass

Output layer delta (δ_k) considering sigmoidal output node(s)



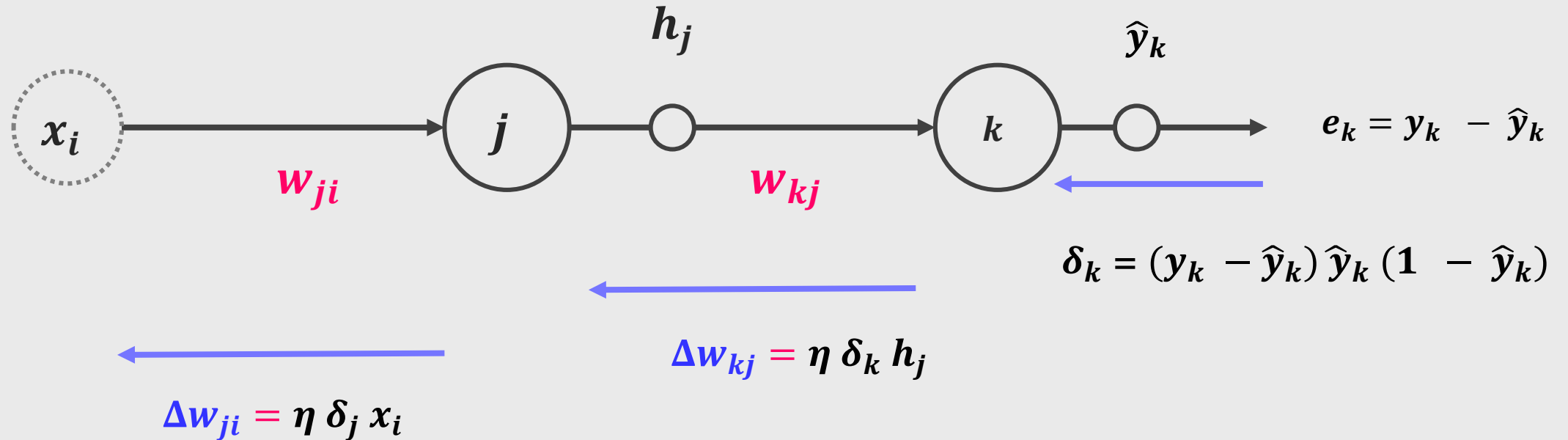
Backpropagation: Backward pass

Hidden layer delta (δ_j) considering sigmoidal hidden node(s)



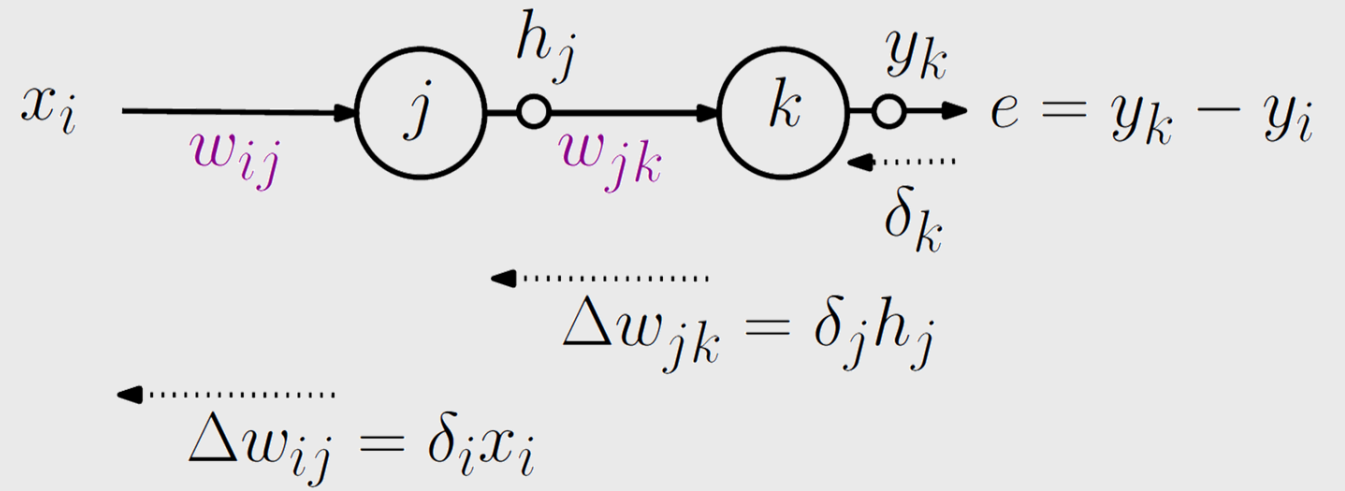
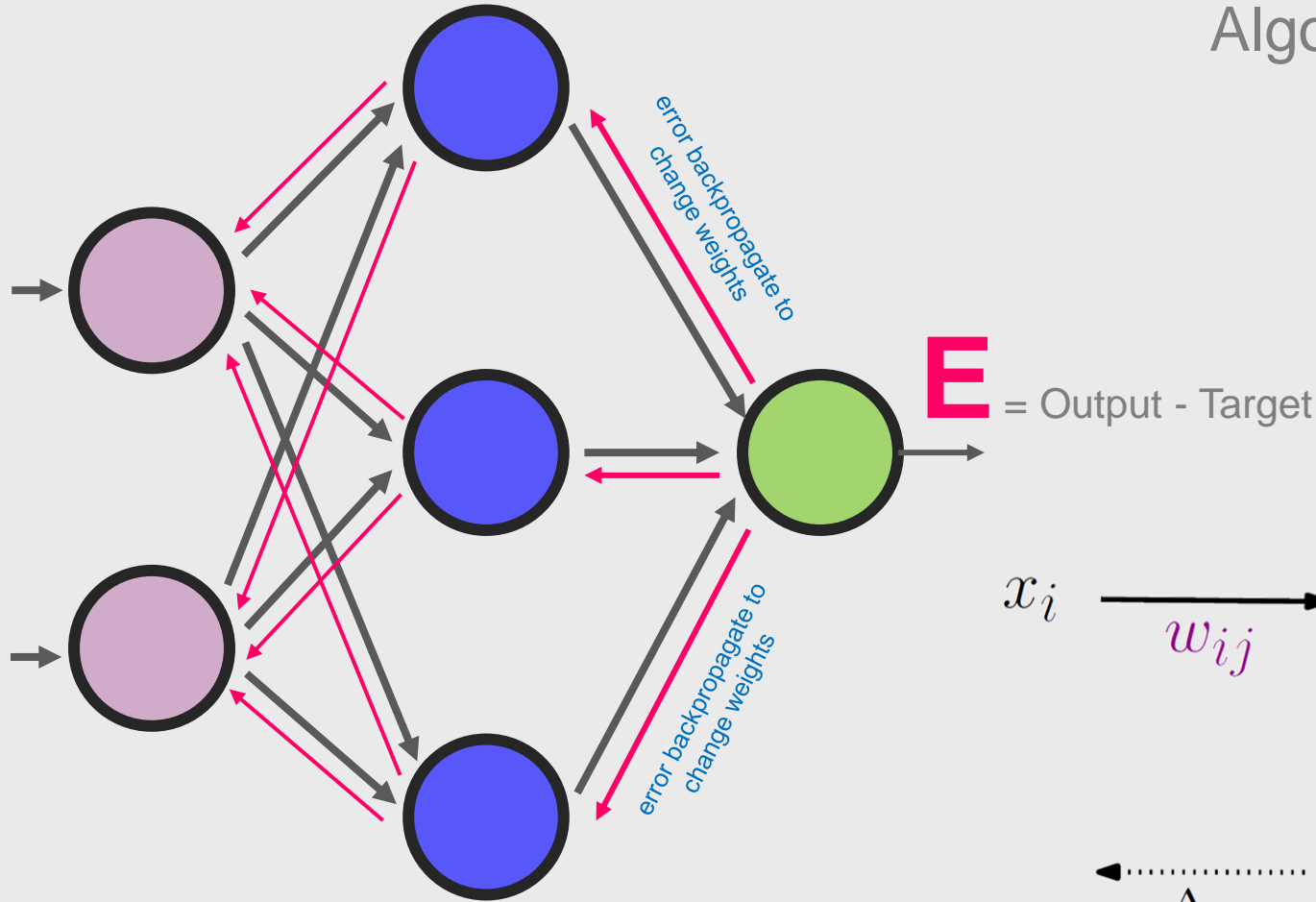
Backpropagation: Backward pass

Hidden layer delta (δ_j) considering sigmoidal hidden node(s)



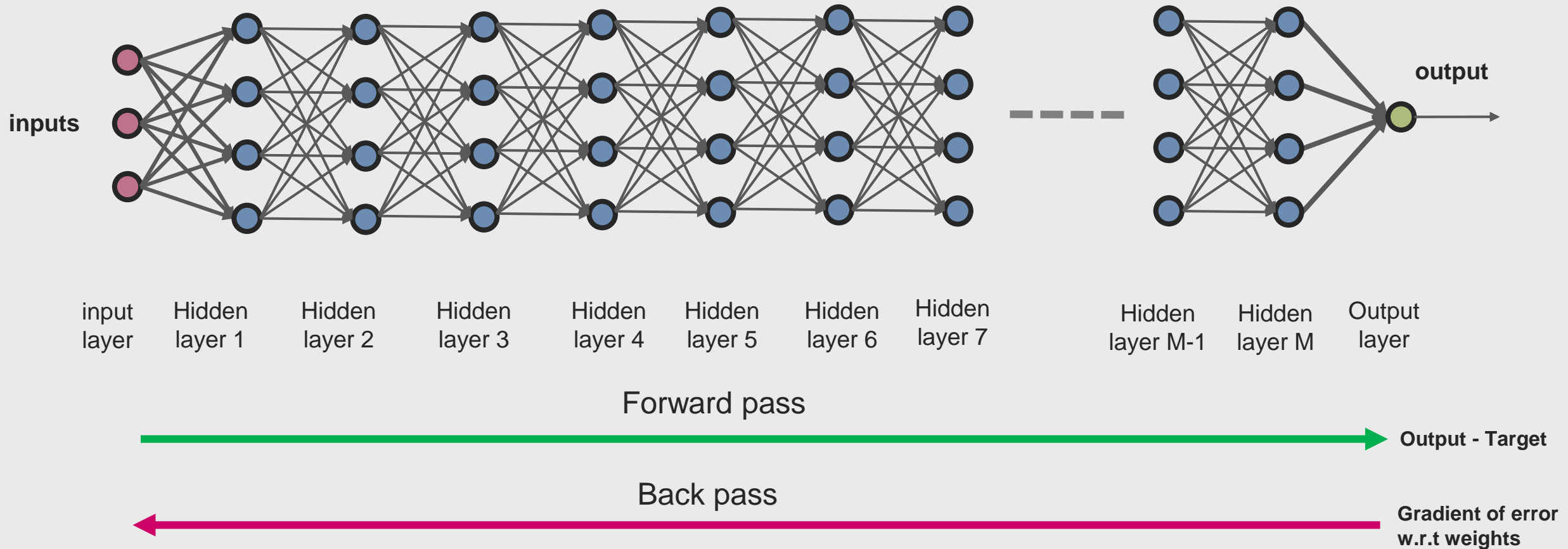
BACKPROPAGATION

Algorithm for Updating Learning Systems



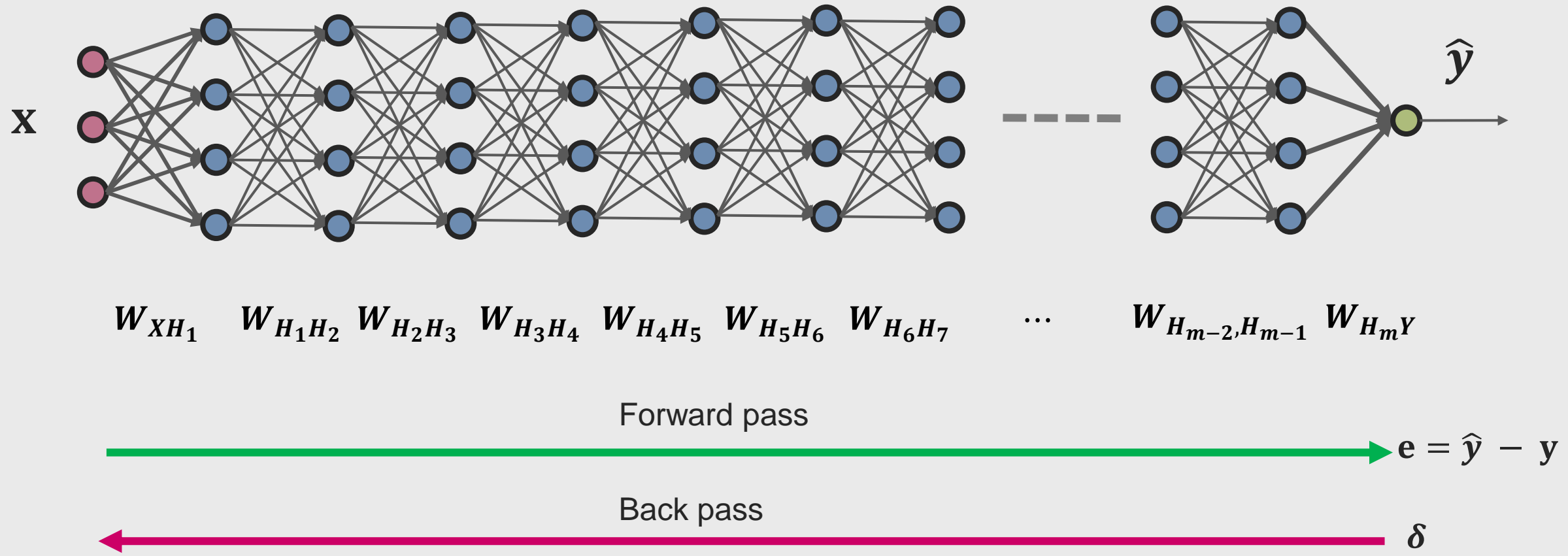
BACKPROPAGATION

Deep Learning

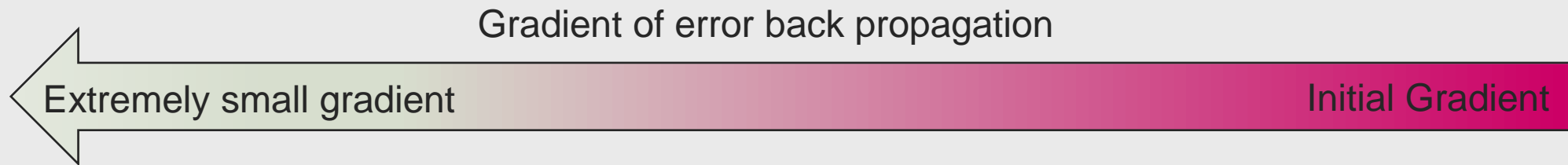
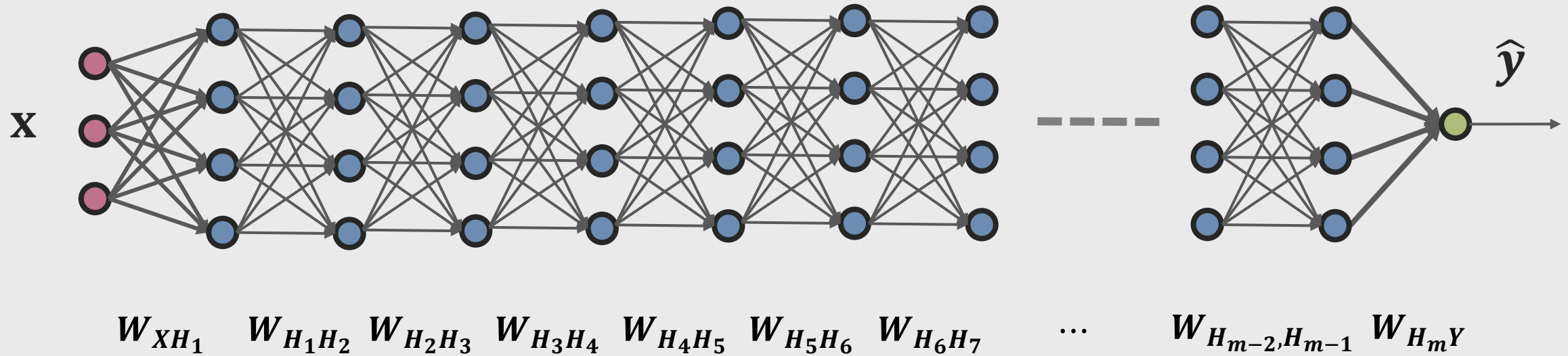


BACKPROPAGATION

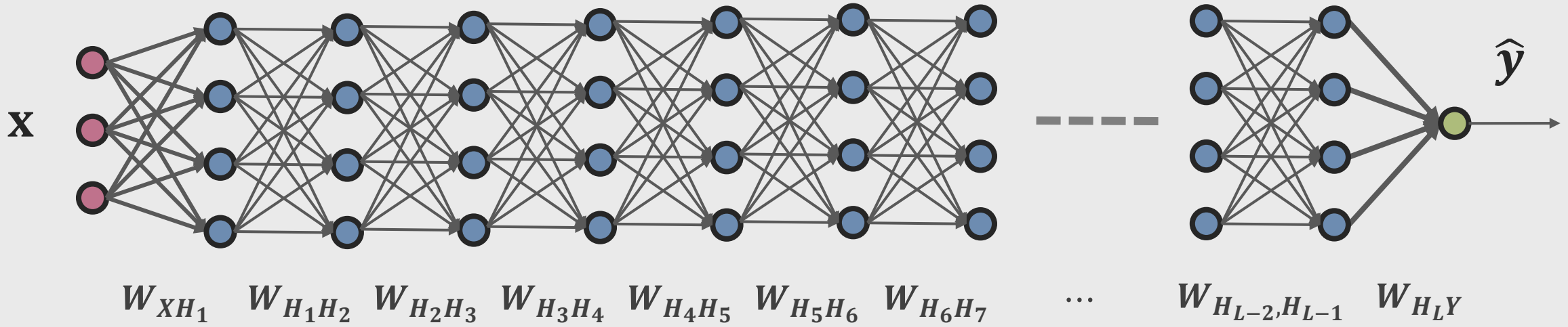
Deep Learning



Vanishing Gradient



Vanishing Gradient



Forward pass: $\hat{y} = \varphi_L(W_L \varphi_{L-1}(W_{L-1} \dots \varphi_3(W_3 \varphi_2(W_2 \varphi_1(W_1 x)))) \dots)$

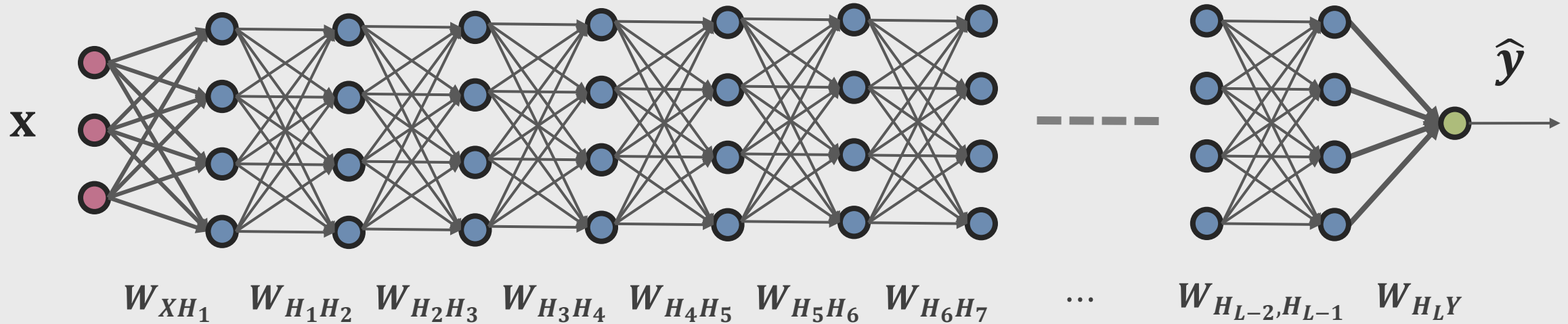
$\xrightarrow{\text{green arrow}} e = \hat{y} - y$

$$\hat{y} = W_L \begin{bmatrix} 0.5 & \dots & 0.0 \\ \vdots & \ddots & \vdots \\ 0.0 & \dots & 0.5 \end{bmatrix}^{L-1}$$

$w = 0.5^{L-1}$ for a large L this will be **extremely small**. That is, weight w is an exponentially **decreasing** function of L

This is caused by **sigmoid function** because its derivative lies between 0.0 and 0.25

Vanishing Gradient



Forward pass: $\hat{y} = \varphi_L(W_L \varphi_{L-1}(W_{L-1} \dots \varphi_3(W_3 \varphi_2(W_2 \varphi_1(W_1 \mathbf{x}))) \dots))$

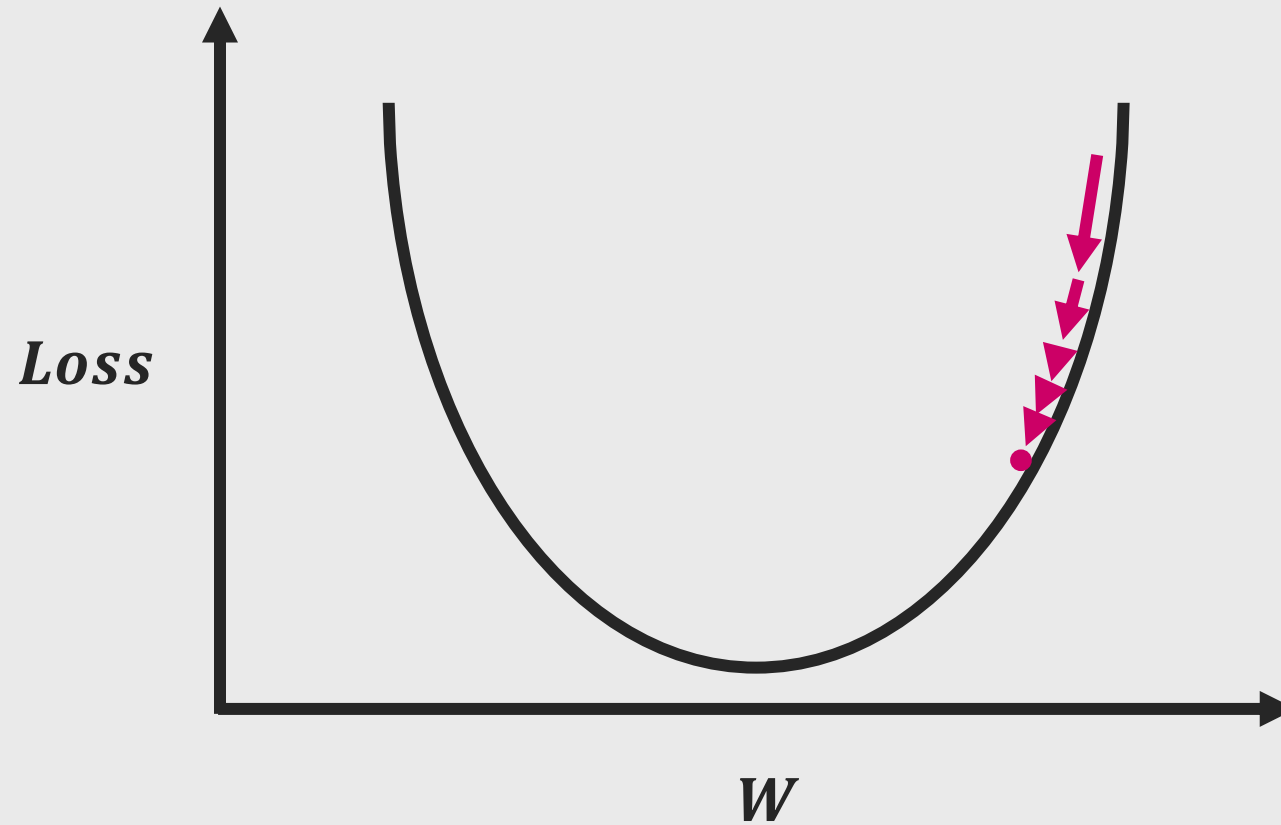
$\xrightarrow{\text{green arrow}} \mathbf{e} = \hat{y} - y$

$$\hat{y} = W_L \begin{bmatrix} 0.5 & \dots & 0.0 \\ \vdots & \ddots & \vdots \\ 0.0 & \dots & 0.5 \end{bmatrix}^{L-1}$$

$w = 0.5^{L-1}$ for a large L this will be **extremely small**. That is, weight w is an exponentially **decreasing** function of L

Solution:
Use of **ReLU** function
 $\varphi_1(x) = \max(0, x)$

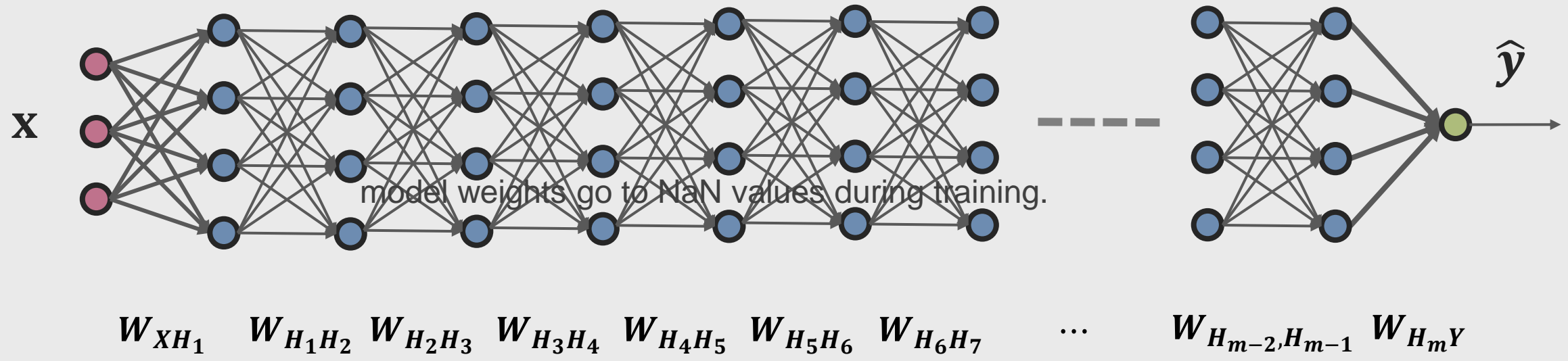
Vanishing Gradient



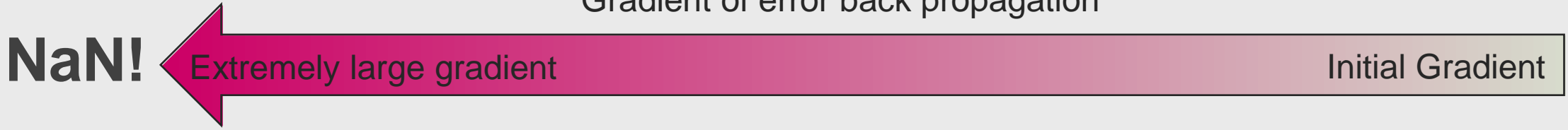
gradient become
very small

Convergence
virtually stops
because weights do
not change any more

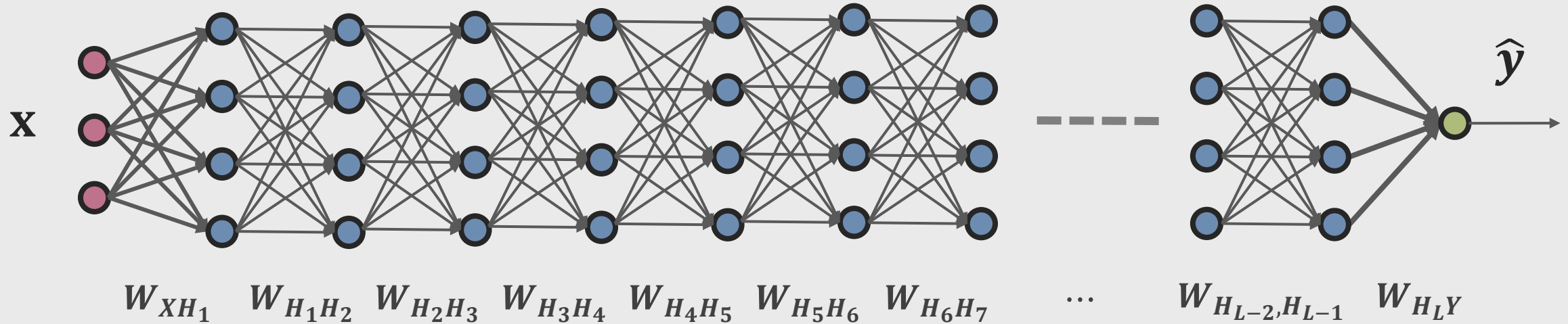
Exploding Gradient



Gradient of error back propagation



Exploding Gradient



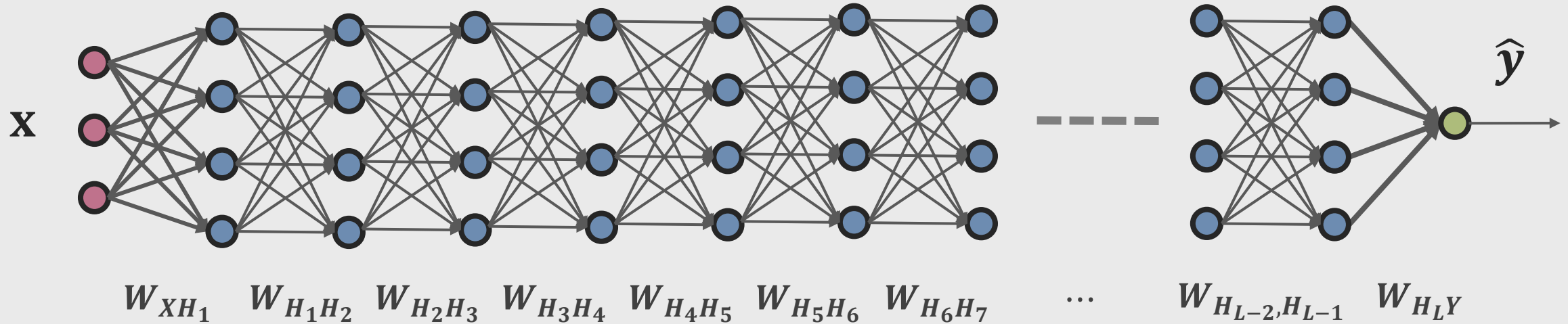
Forward pass: $\hat{y} = \varphi_L(W_L \varphi_{L-1}(W_{L-1} \dots \varphi_3(W_3 \varphi_2(W_2 \varphi_1(W_1 \mathbf{x}))) \dots))$

$$\hat{y} = W_L \begin{bmatrix} 1.5 & \dots & 0.0 \\ \vdots & \ddots & \vdots \\ 0.0 & \dots & 1.5 \end{bmatrix}^{L-1}$$

$w = 1.5^{L-1}$ for a large L this will be **extremely larger**. That is, weight w is an exponentially **increase** function of L

This is caused by **initialization of weights with large values.**

Exploding Gradient



Forward pass: $\hat{y} = \varphi_L(W_L \varphi_{L-1}(W_{L-1} \dots \varphi_3(W_3 \varphi_2(W_2 \varphi_1(W_1 \mathbf{x}))) \dots))$

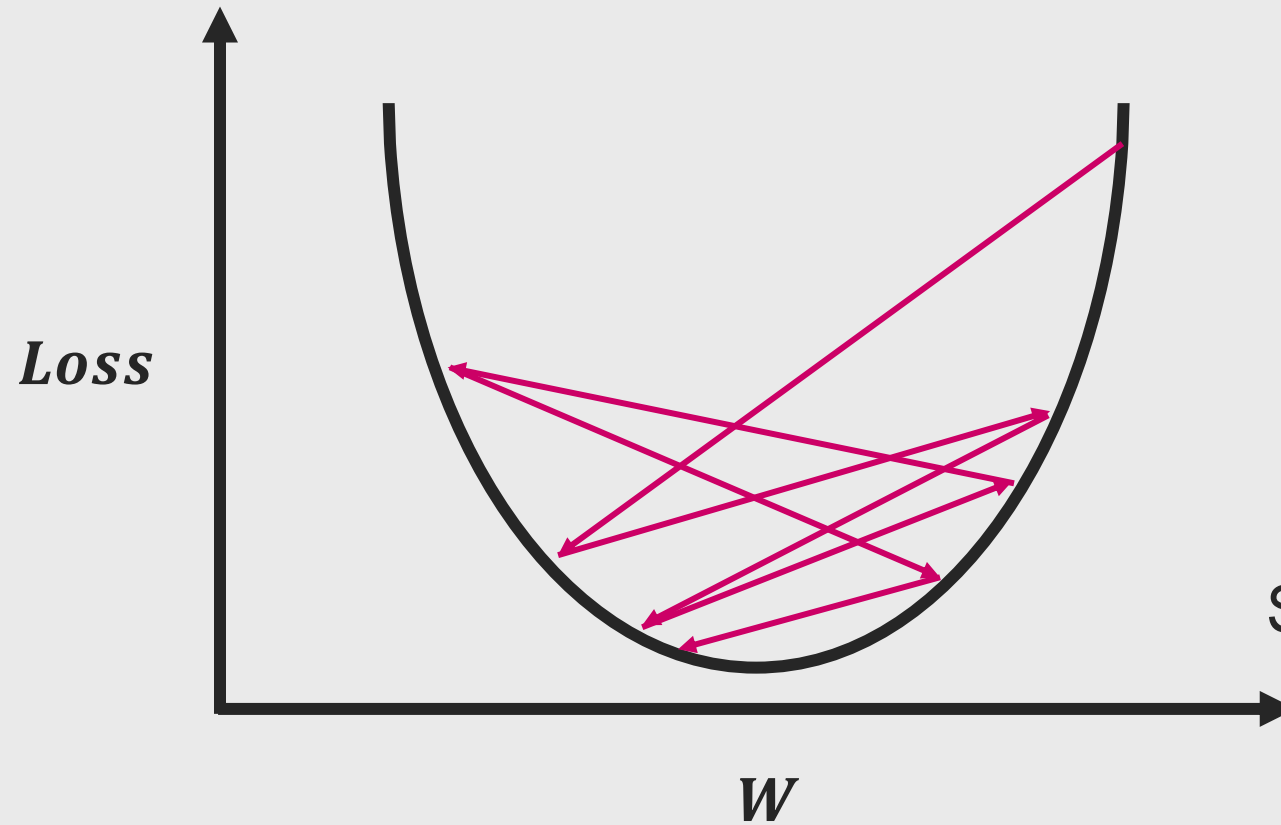
$\mathbf{e} = \hat{y} - y$

$$\hat{y} = W_L \begin{bmatrix} 1.5 & \dots & 0.0 \\ \vdots & \ddots & \vdots \\ 0.0 & \dots & 1.5 \end{bmatrix}^{L-1}$$

$w = 1.5^{L-1}$ for a large L this will be extremely larger. That is, weight w is an exponentially increase function of L

Solution:
Gradient clipping and/or better weight initialization

Exploding Gradient



Highly fluctuating gradient descent

Weight abruptly changes.

Skips Global minima

Part 4

Practical

Session

See attached document in Canvas

