

Introduction to Neural Networks

Introduction to Artificial Intelligence - 2023 Summer

Jul 27, 2023
Thu 4 PM

Kwangwoon University MI:RU
Artificial Intelligence Study



Agenda

In this course, you will learn

Part 1 – Quick Review of Partial Derivation and Chain Rule

Part 2 – Perceptron & Artificial Neuron Networks

Part 3 – Linear Regression & Backpropagation

Part 4 – Recurrent Neural Networks (RNN)



Quick review of Partial Derivative and Chain Rule

Partial derivative

$f'_x, f_x, \partial_x f, D_x f, D_1 f, \frac{\partial}{\partial x} f, \text{ or } \frac{\partial f}{\partial x}$.

Notation for Partial derivative of variable

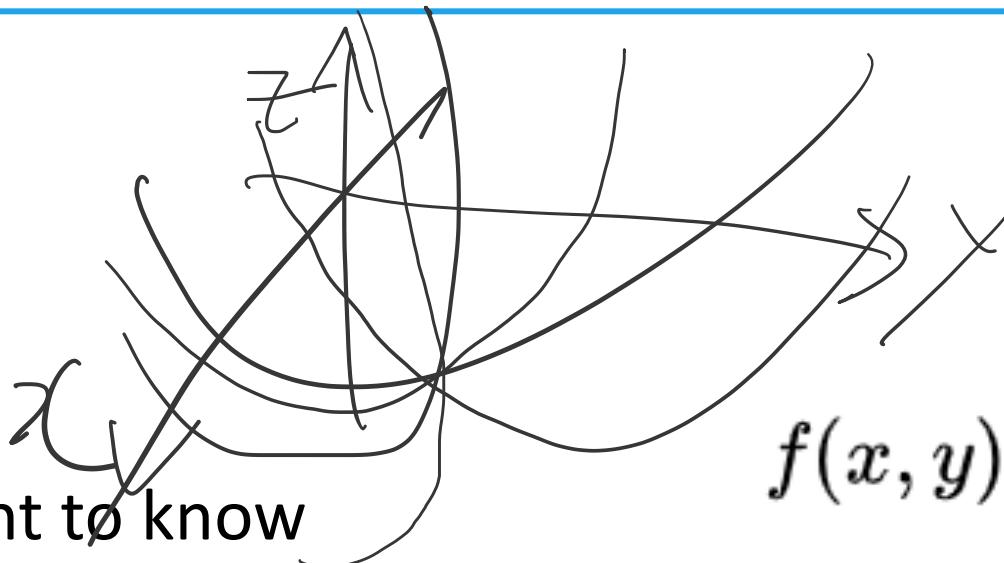
x

Quick review of Partial Derivative and Chain Rule

Partial derivative

What we want to know

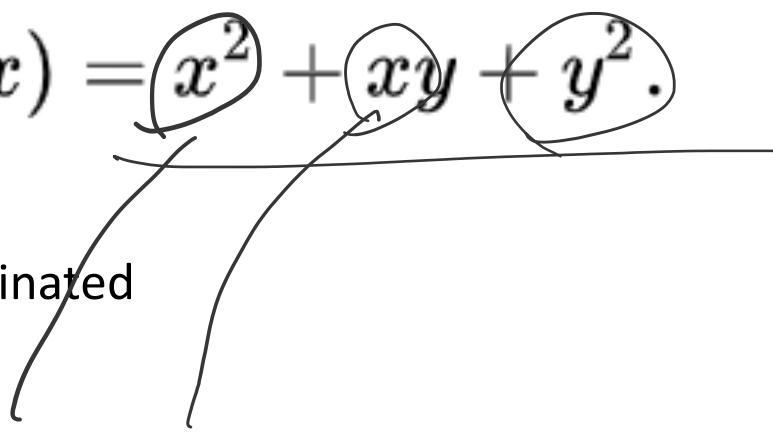
- Change rate of ONLY x
- Then, handle y as a constant
=> We can get the change rate of x ONLY!



$$f(x, y) = f_y(x) = x^2 + xy + y^2.$$

Handle y as a constant
=> Then, y will be eliminated

$$\frac{\partial f}{\partial x}(x, y) = 2x + y.$$





Quick review of Partial Derivative and Chain Rule

Chain Rule

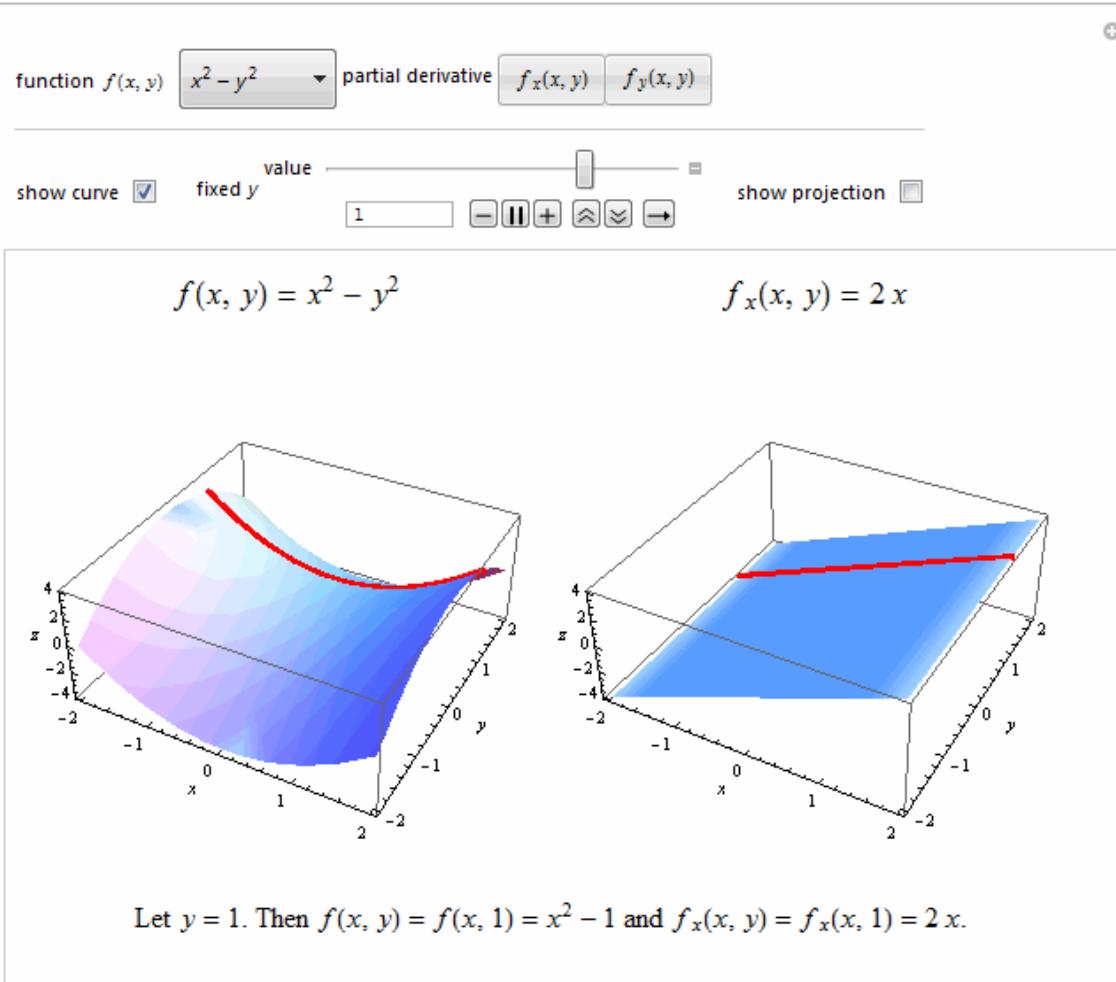
Chain rule

$$\frac{dz}{dx} = \frac{\frac{dz}{dt}}{\frac{dx}{dt}} \quad \text{OR} \quad \frac{dz}{dt} = \frac{dz}{dx} \cdot \frac{dx}{dt}$$

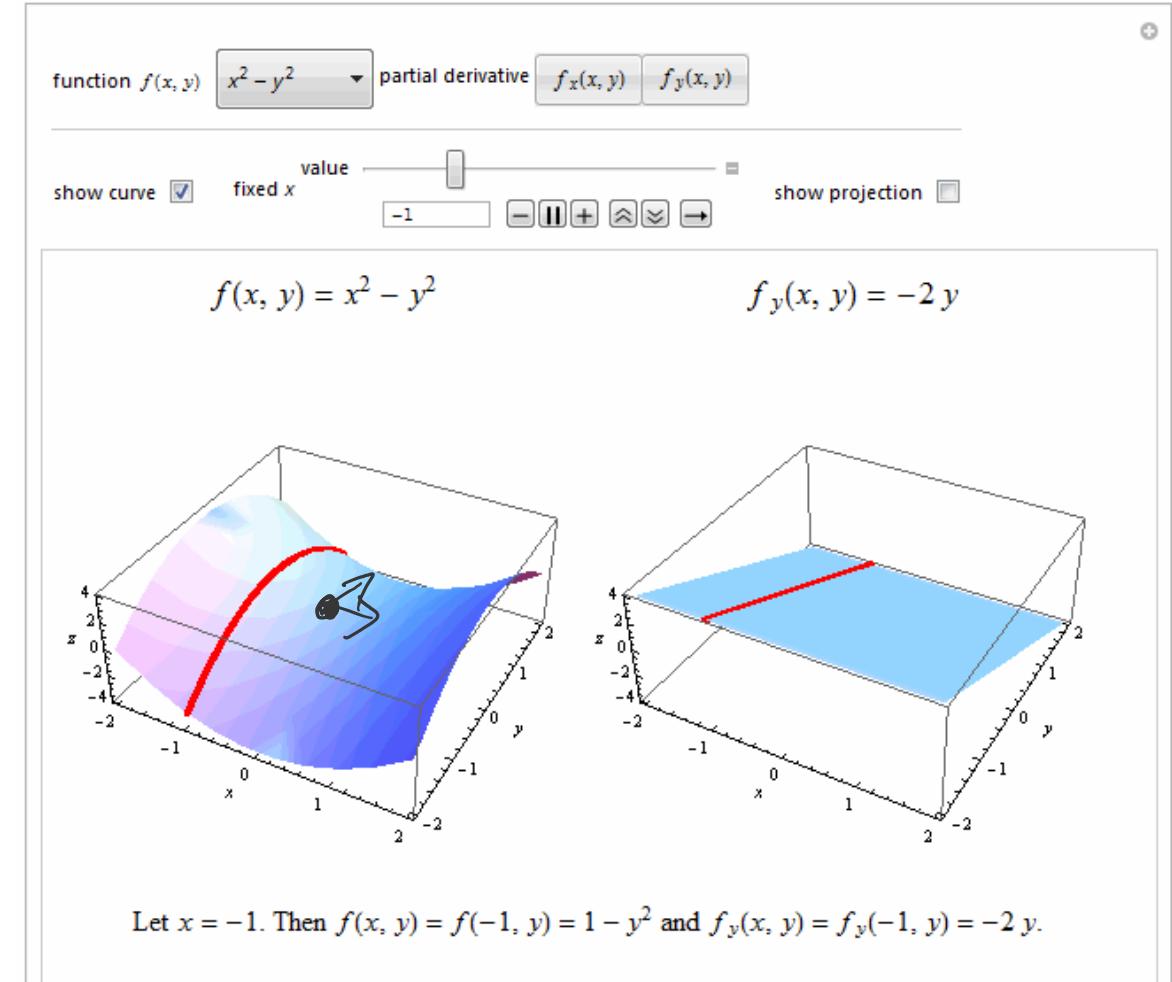
$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx} = 2 \cdot 4 = 8$$

Linear Regression & Backpropagation

Quick review of Partial derivation



$f_x(x, y)$: y is fixed

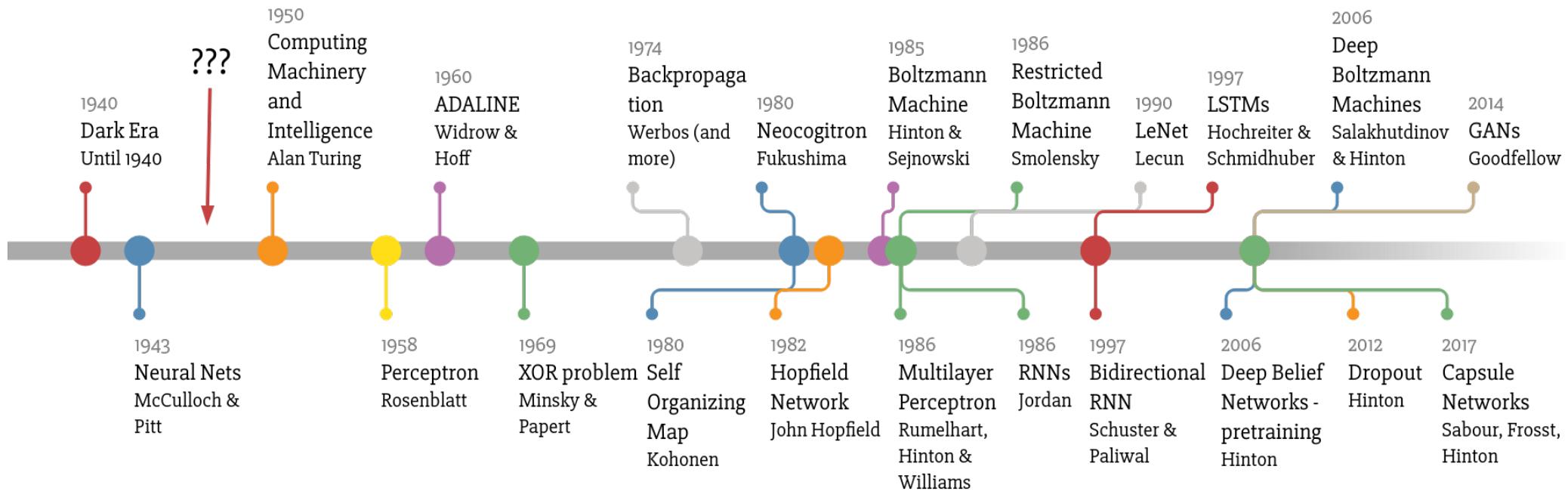


$f_y(x, y)$: x is fixed

Perceptron & Artificial Neural Networks

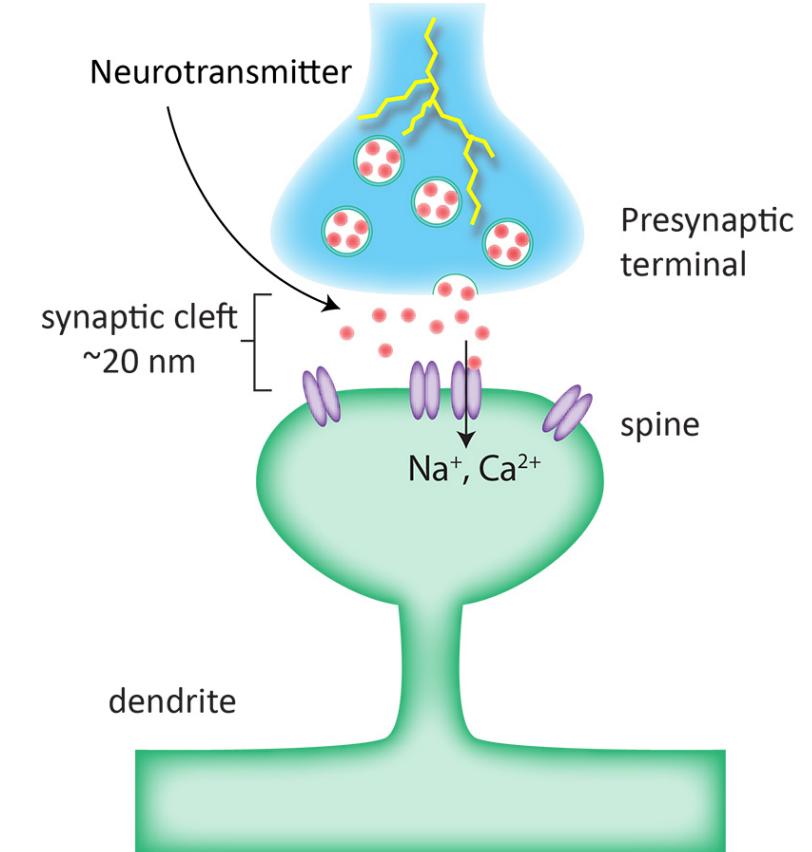
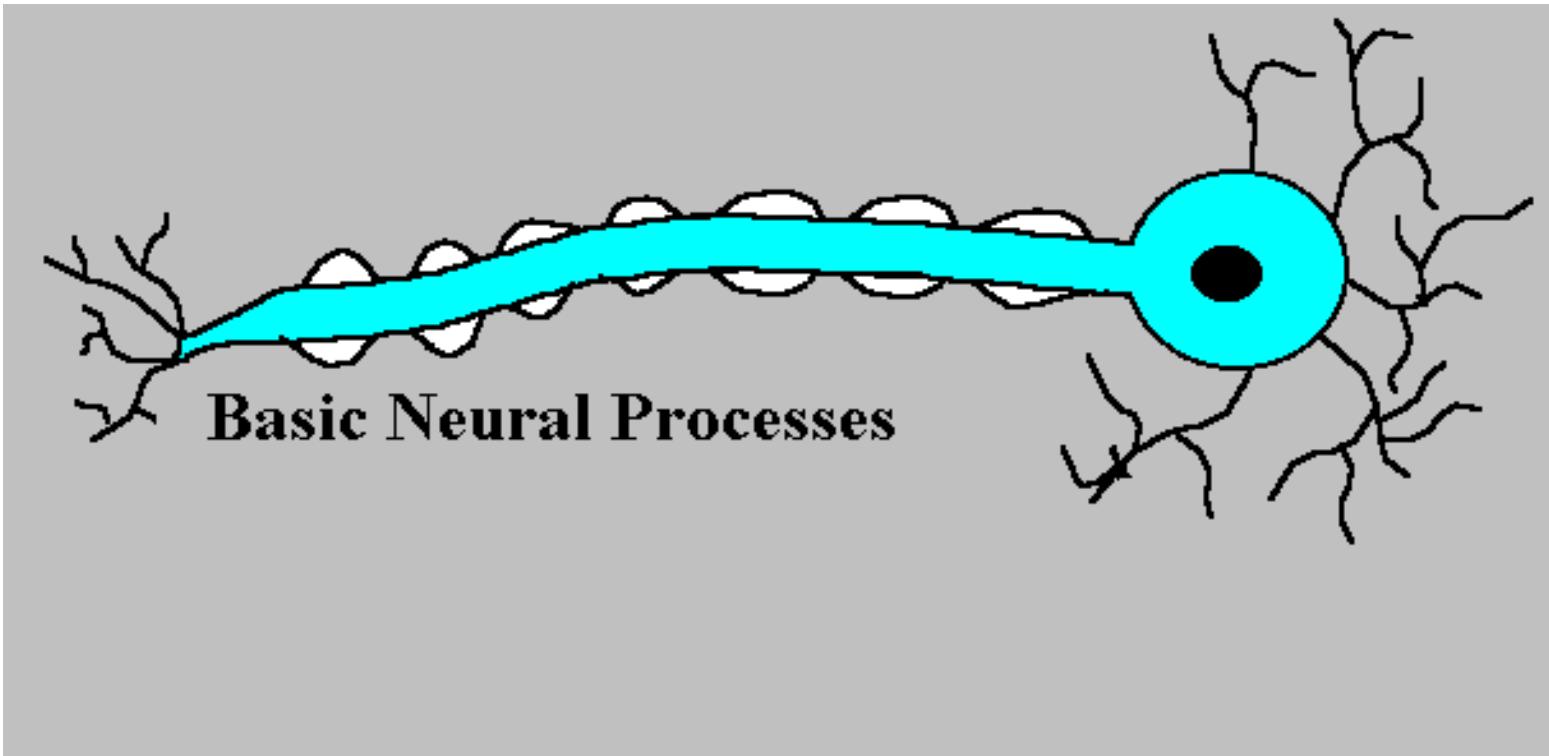
Neuron

Deep Learning Timeline



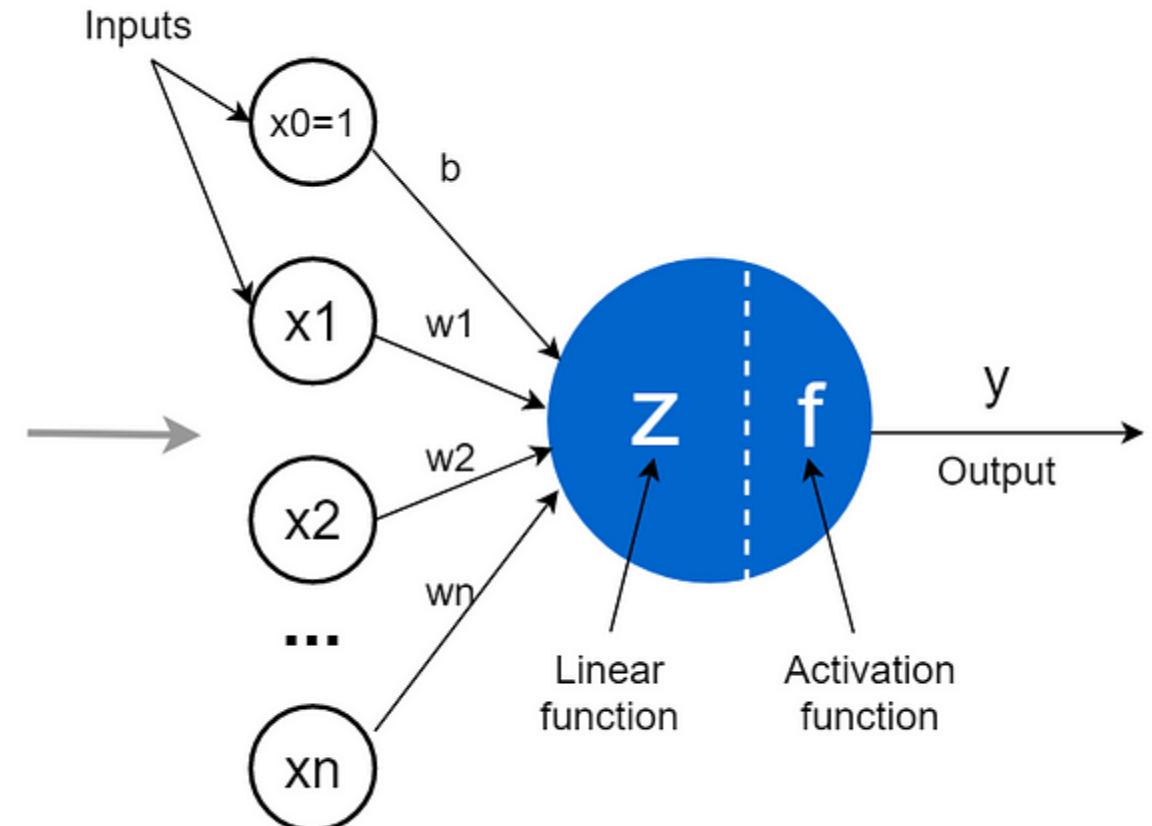
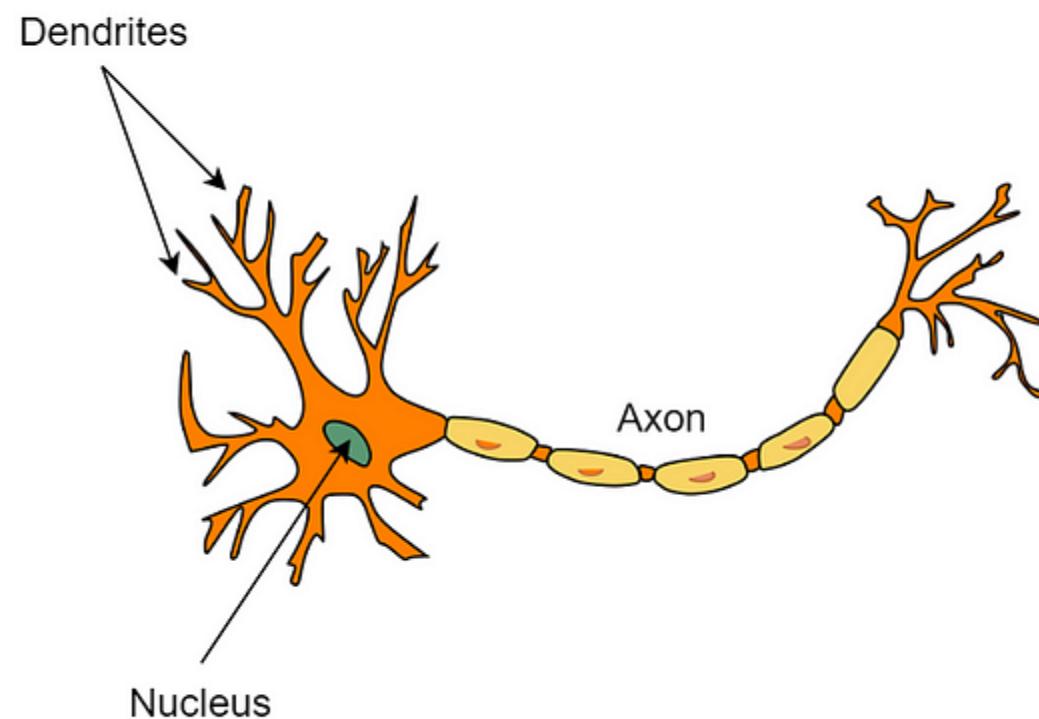
Perceptron & Artificial Neural Networks

Neuron



Perceptron & Artificial Neural Networks

Perceptron



Perceptron & Artificial Neural Networks

The first Perceptron machine

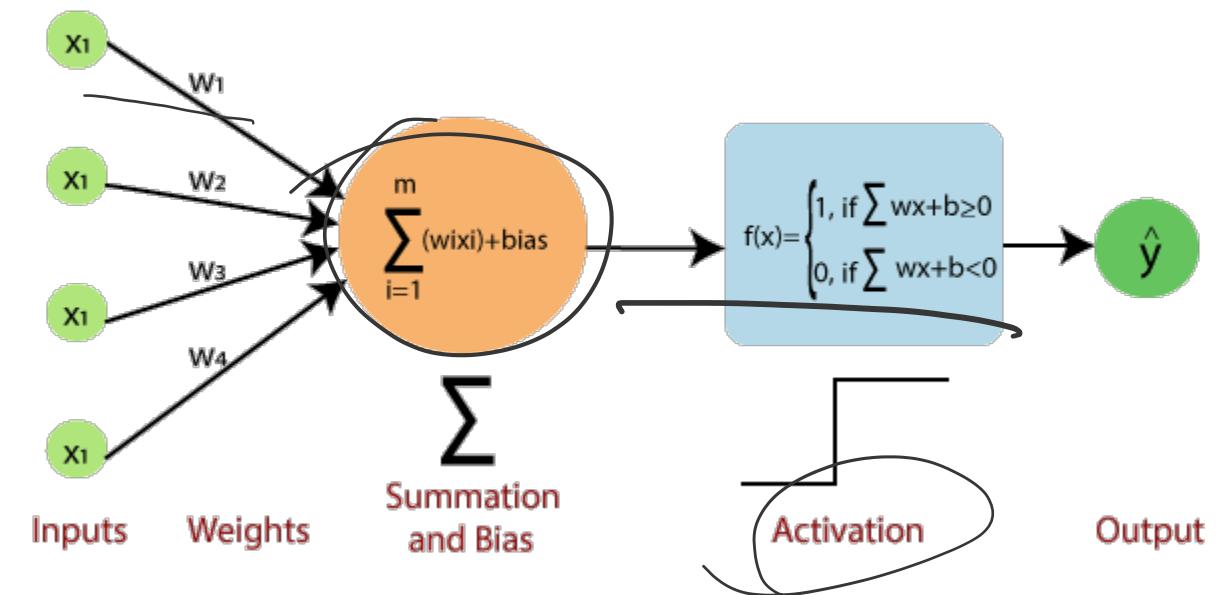
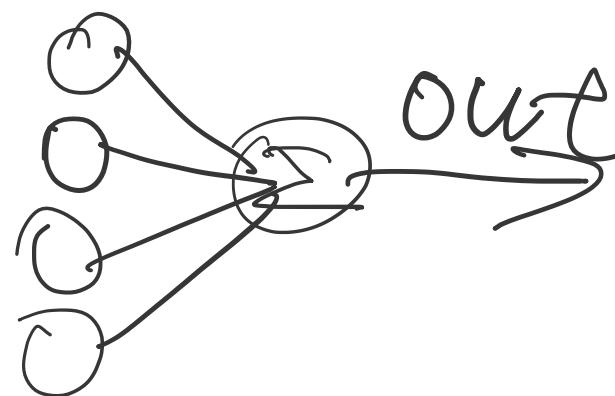
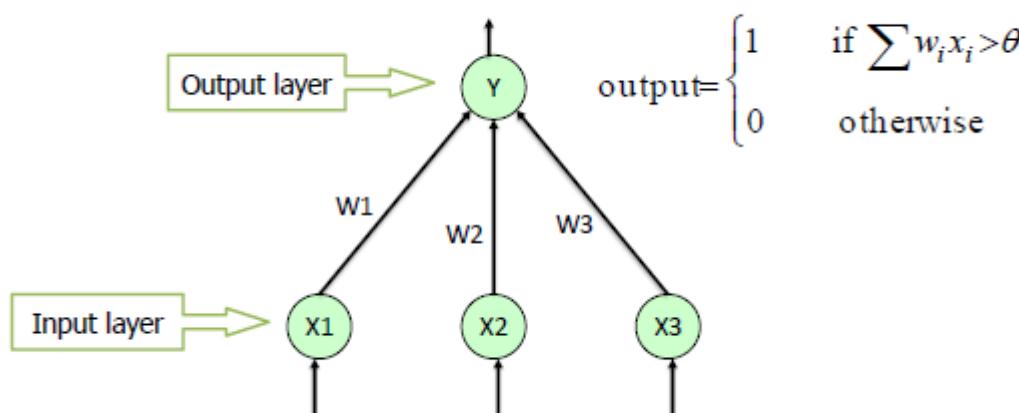


[One-Page Schoolhouse: Perceptron \(ronkowitz.blogspot.com\)](http://ronkowitz.blogspot.com)

Perceptron & Artificial Neural Networks

Single-Layer Perceptron

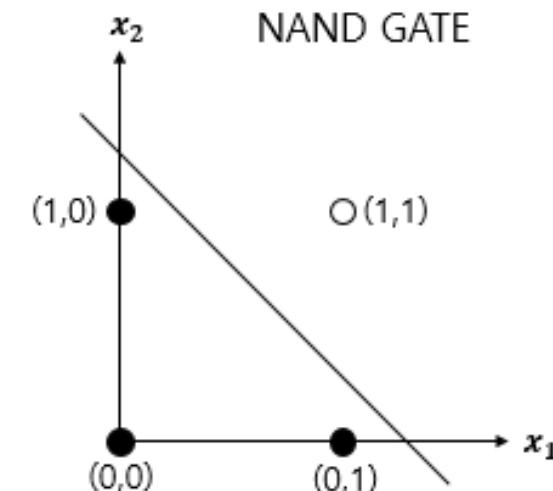
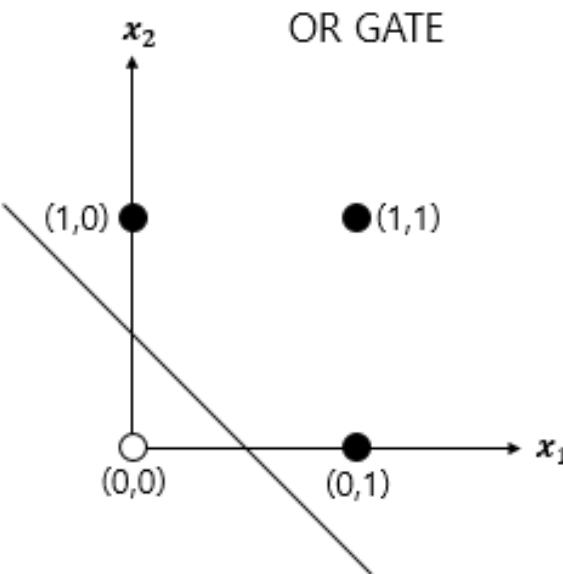
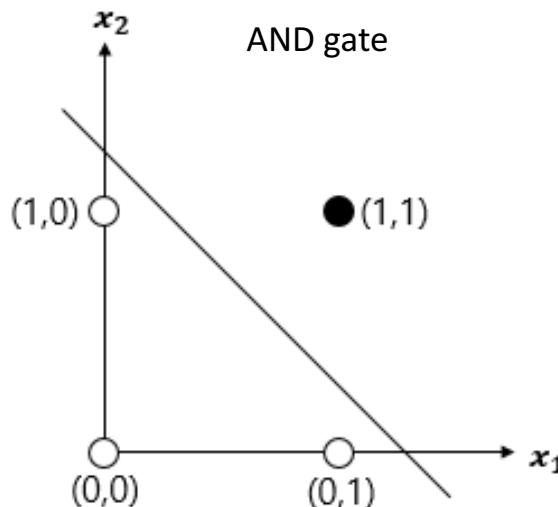
Single Layer Perceptron



Perceptron & Artificial Neural Networks

Classification using Single-Layer Perceptron

- A perspective of Logic gates

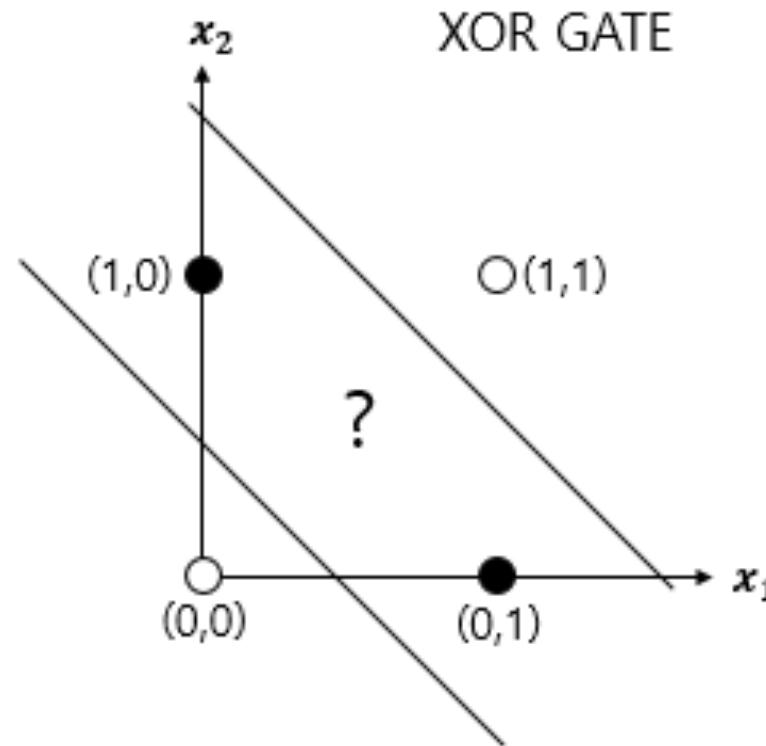


No problem occurred

Perceptron & Artificial Neural Networks

Classification using Single-Layer Perceptron

- A perspective of Logic gates



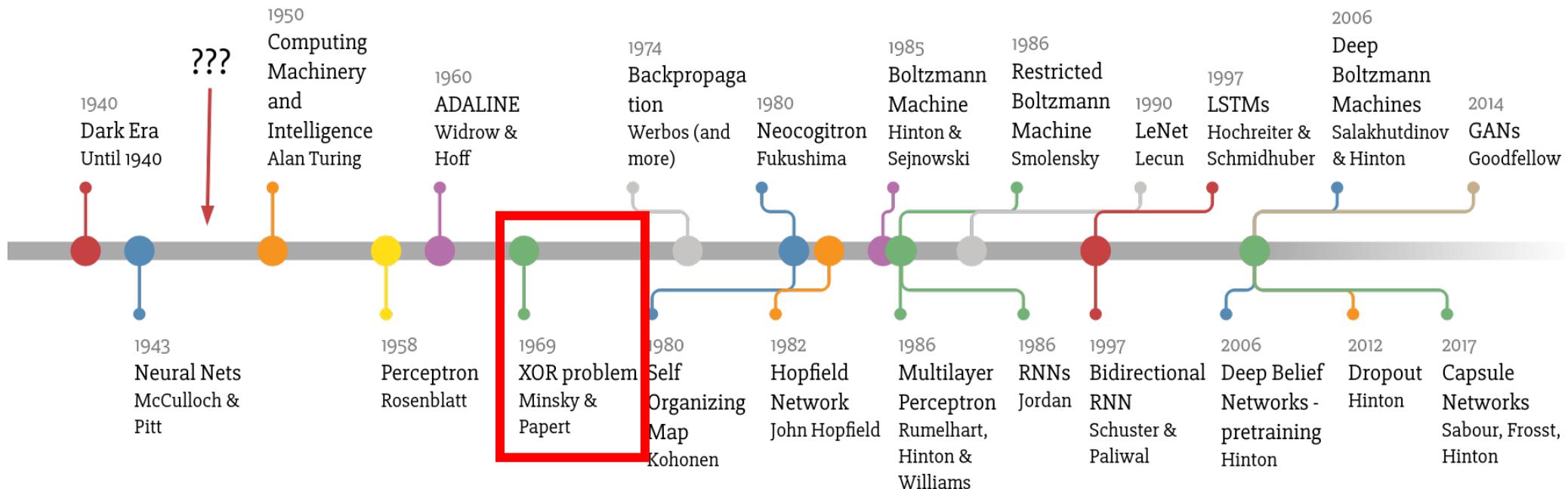
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

Problem occurred!!!

Perceptron & Artificial Neural Networks

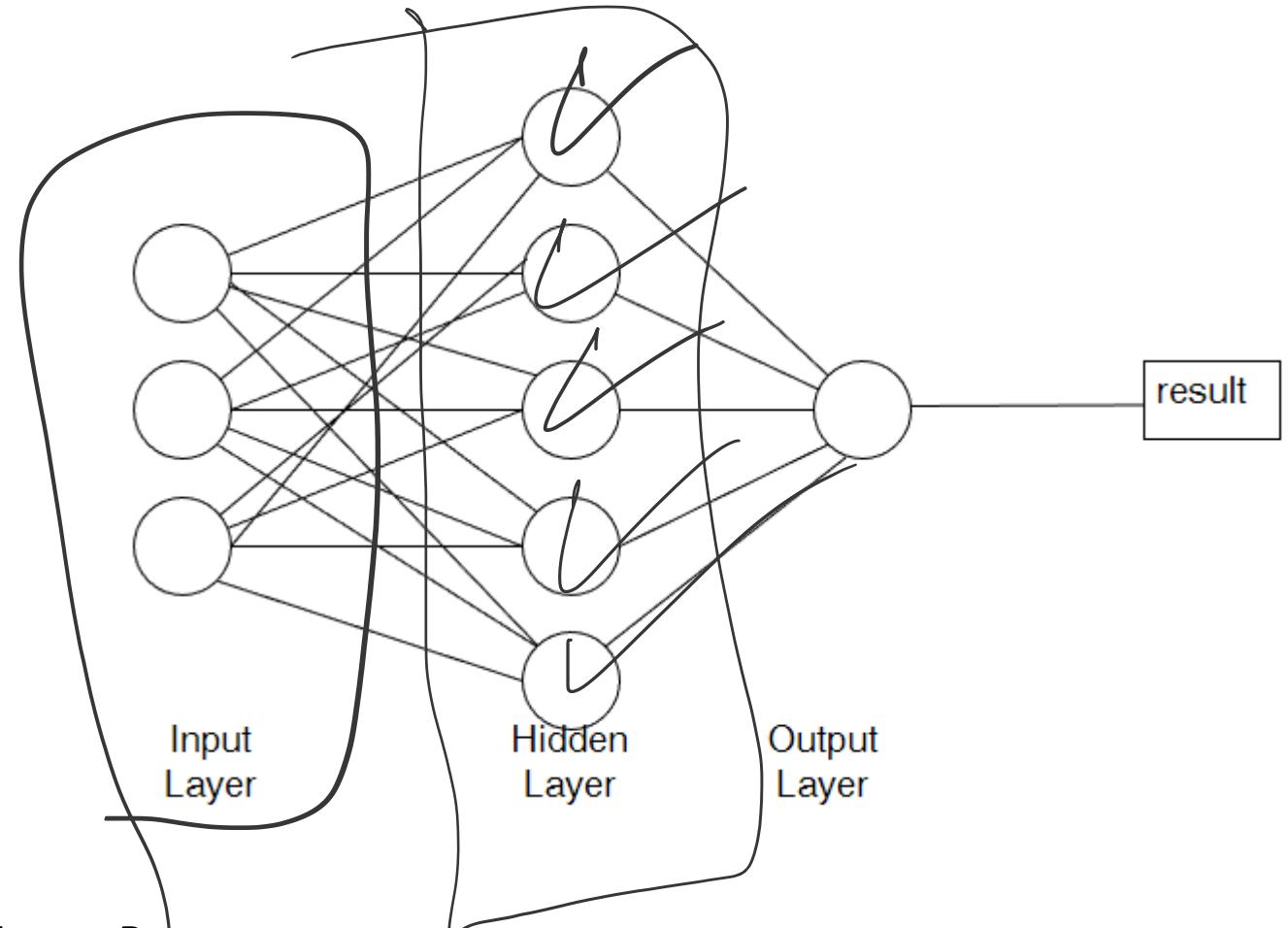
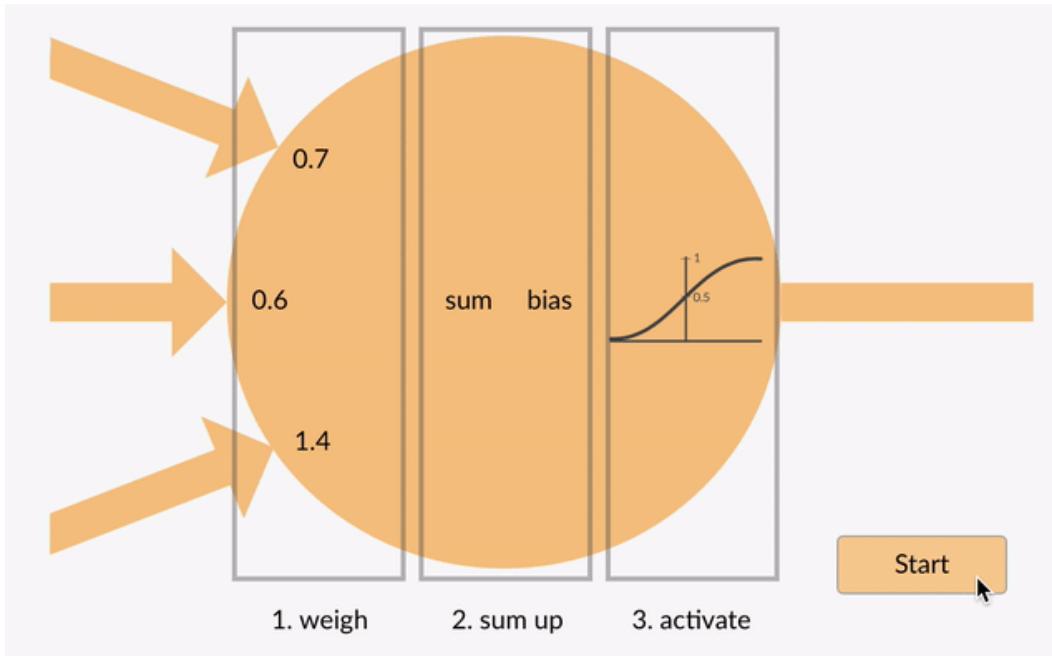
Neuron

Deep Learning Timeline



Perceptron & Artificial Neural Networks

Multi-Layer Perceptron (MLP)



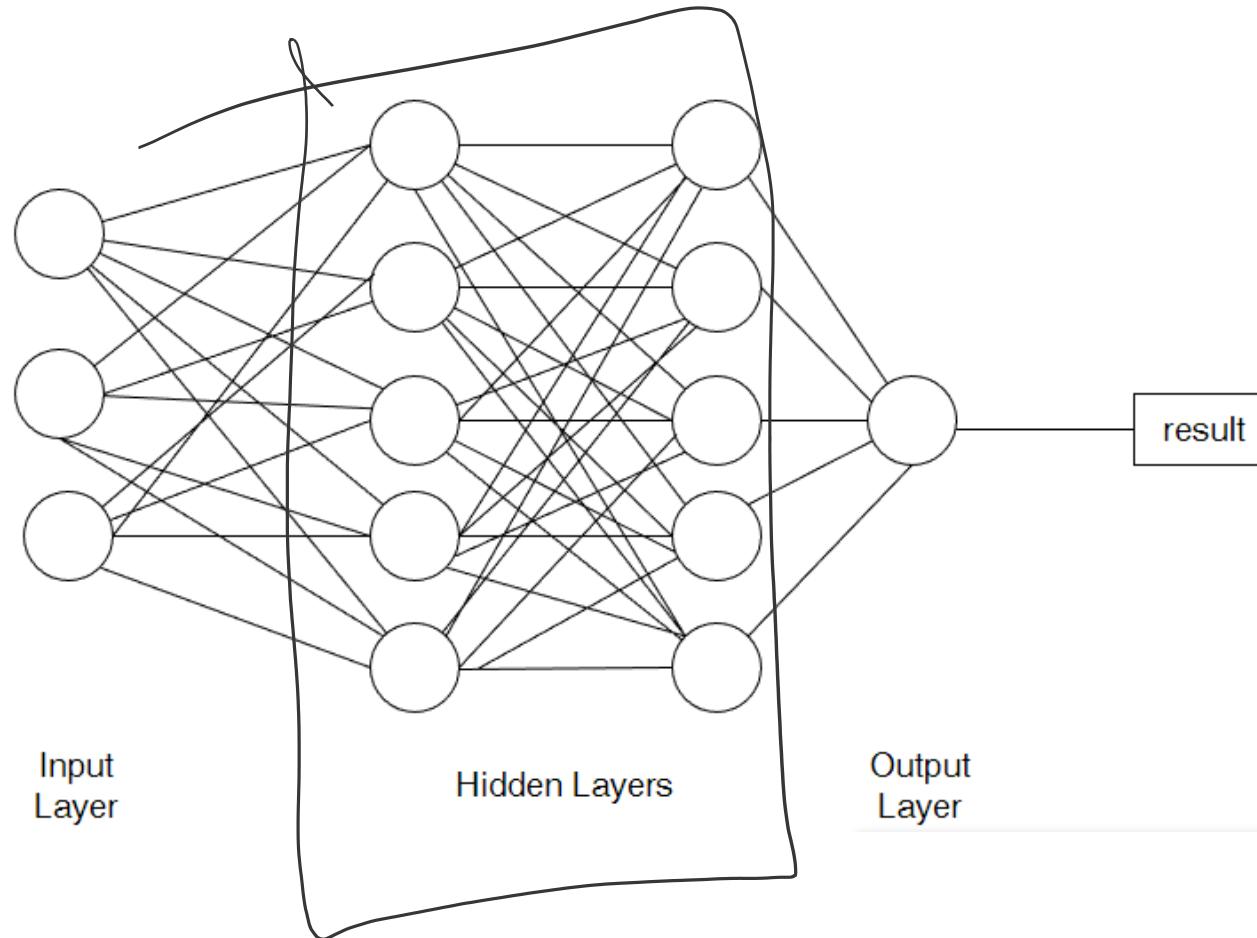
Multi-Layer Perceptron

==

Sequentially connected Single Layer Perceptron

Perceptron & Artificial Neural Networks

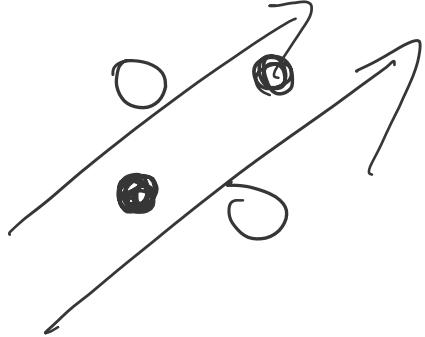
Deep Learning?



Deep Learning
==
Training a Deep Neural Network

Perceptron & Artificial Neural Networks

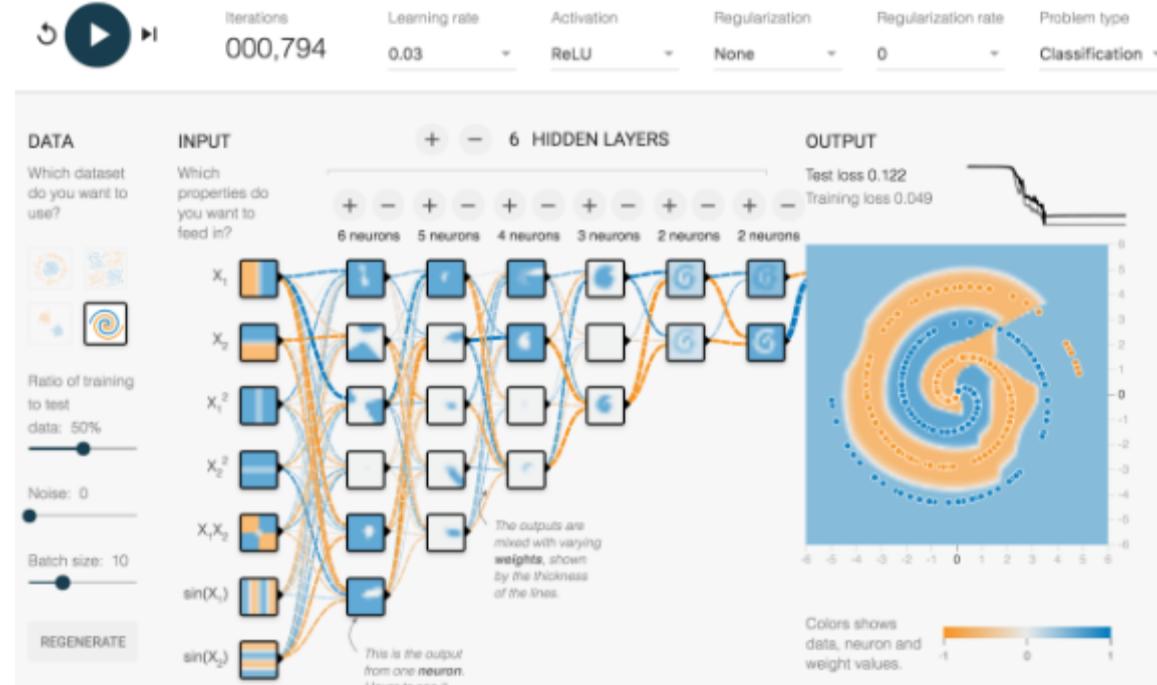
Multi-Layer Perceptron (MLP)



Structure	Regions	XOR	Meshed regions
single layer	Half plane bounded by hyperplane		
two layer	Convex open or closed regions		
three layer	Arbitrary (limited by # of nodes)		

Perceptron & Artificial Neural Networks

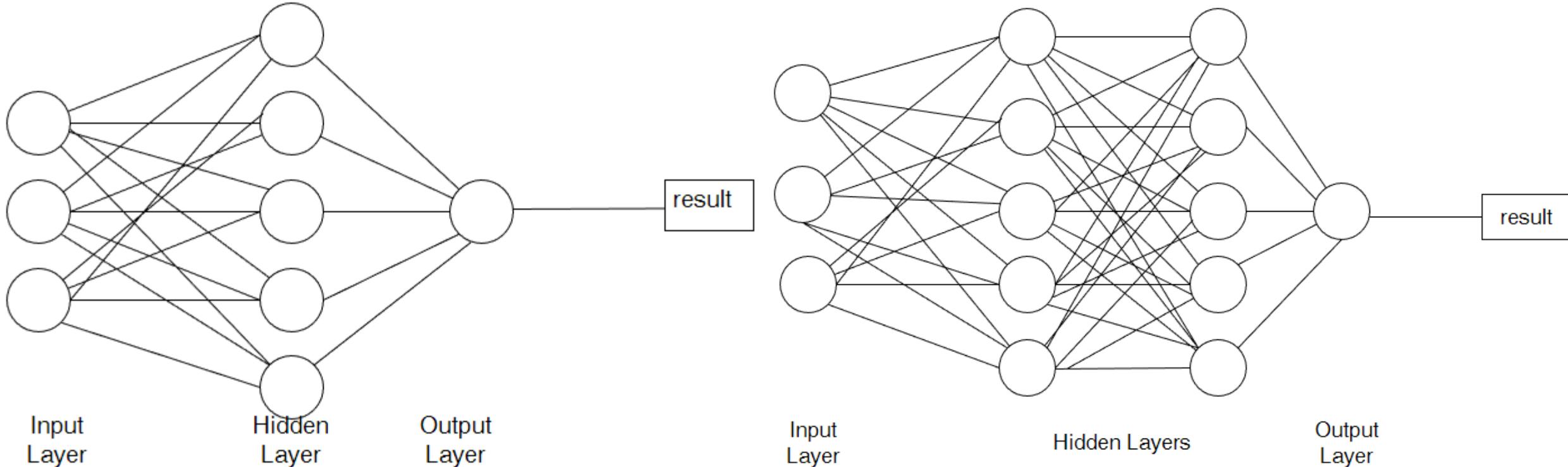
Deep Neural Network



[A Neural Network Playground \(tensorflow.org\)](https://tensorflow.org)

Perceptron & Artificial Neural Networks

Deep Neural Network (DNN)



Sequentially connected Multi-Layer Perceptron

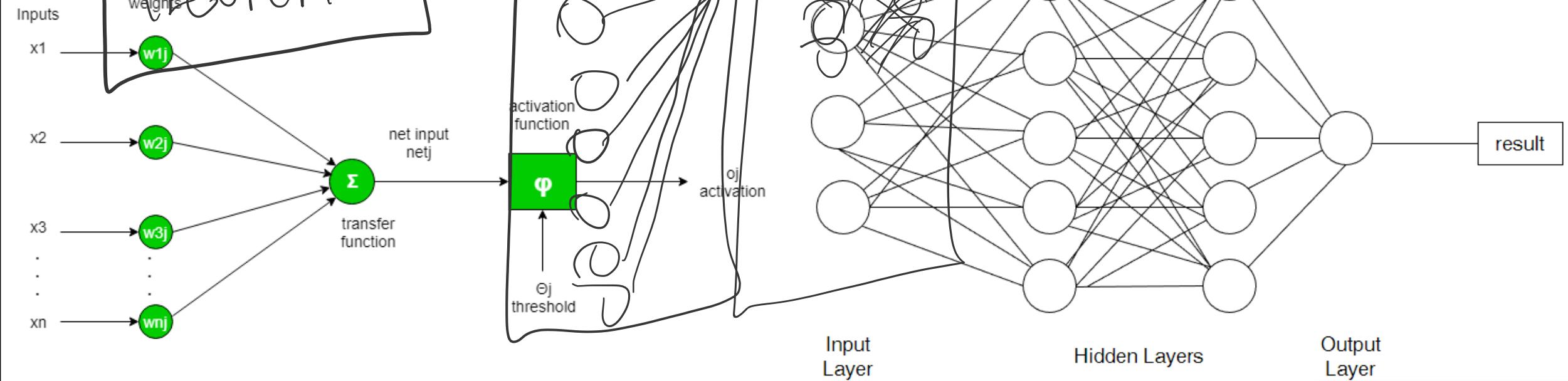
==

Deep Neural Network (DNN)

Perceptron & Artificial Neural Networks

Deep Neural Network (DNN)

Universal
approximation
theorem

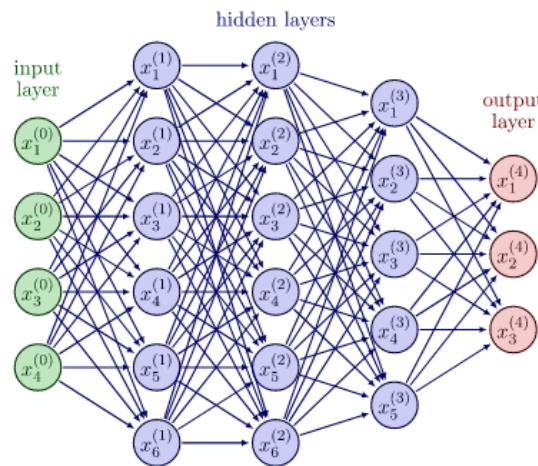


Structure of Neuron being used in Deep Neural Network

Activation function is included

Deep Neural Network (DNN)

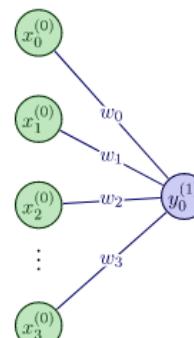
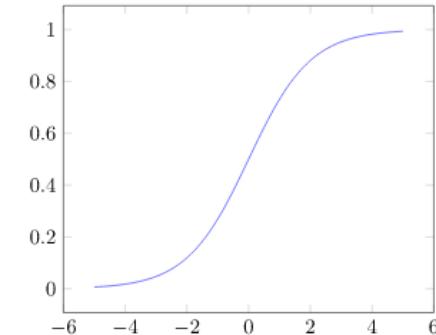
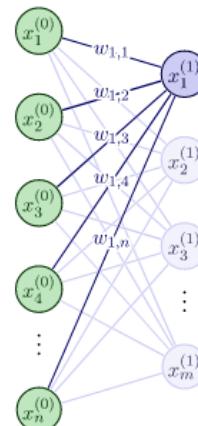
The Universal Approximation Theorem



$$W^{(l)} = \begin{pmatrix} w_{1,1}^{(l)} & w_{1,2}^{(l)} & \dots & w_{1,n}^{(l)} \\ w_{2,1}^{(l)} & w_{2,2}^{(l)} & \dots & w_{2,n}^{(l)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m,1}^{(l)} & w_{m,2}^{(l)} & \dots & w_{m,n}^{(l)} \end{pmatrix}$$

$$\psi(x) := \sigma \left(\sum_{i=1}^n x_i w_i - b \right) = \sigma(w^\top \cdot x - b),$$

$$\int_{x \in \mathbb{R}^n} \sigma(w^\top x - b) d\mu(x) = 0 \quad \forall w \in \mathbb{R}^n, b \in \mathbb{R}$$



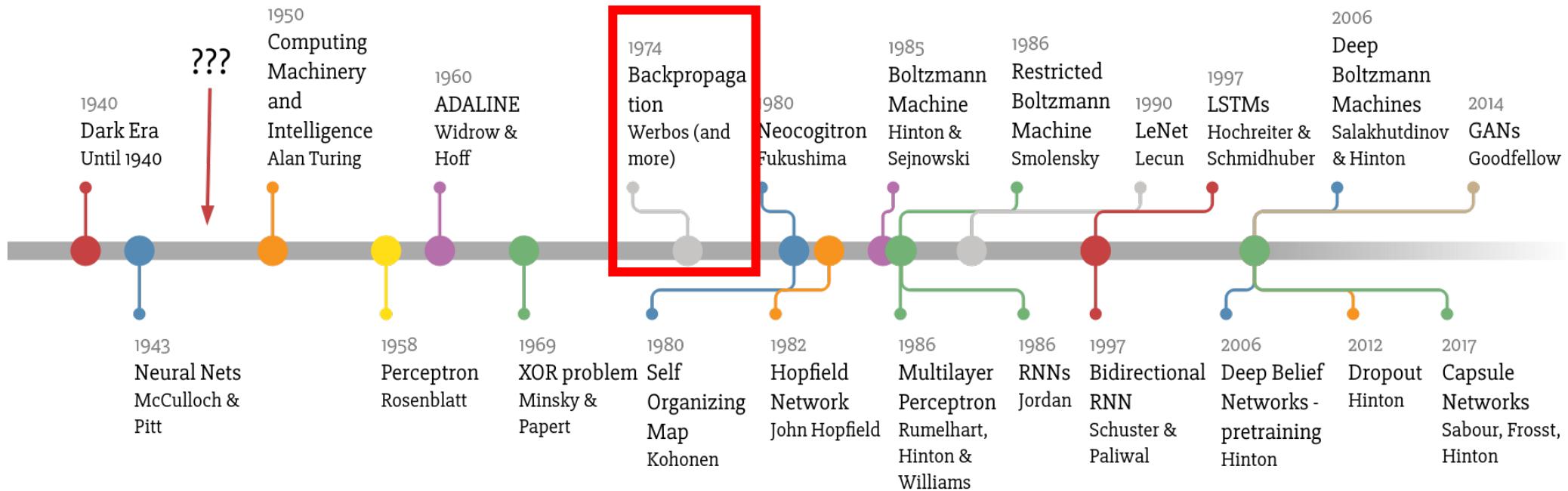
How does DNN works?

[Universal approximation theorem - Wikipedia](#)

Linear Regression & Backpropagation

Backpropagation

Deep Learning Timeline

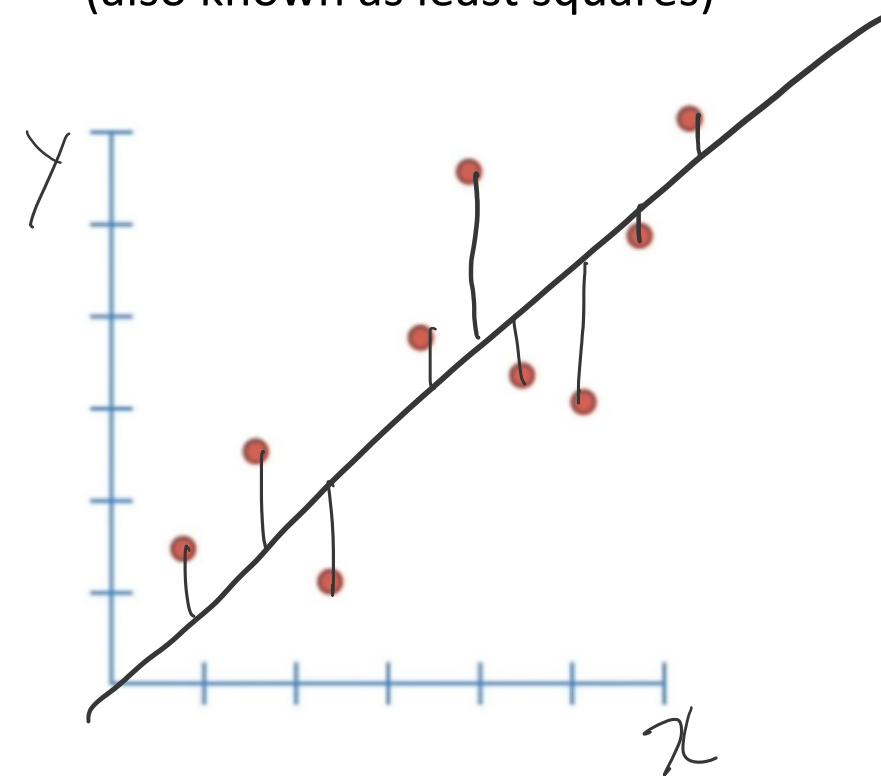




Linear Regression & Backpropagation

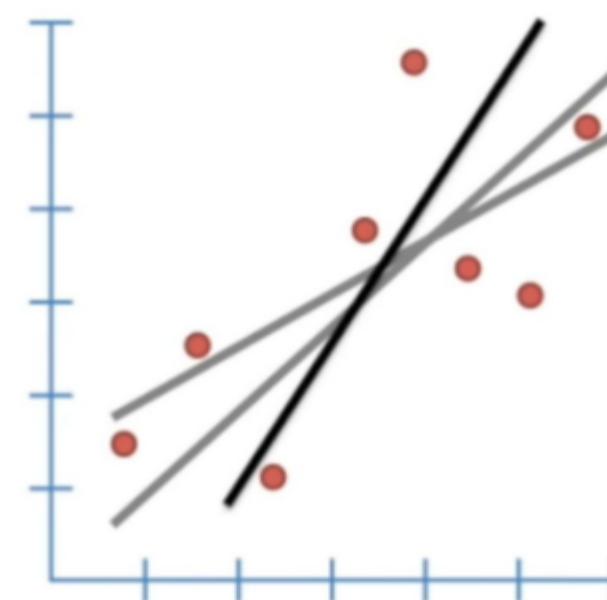
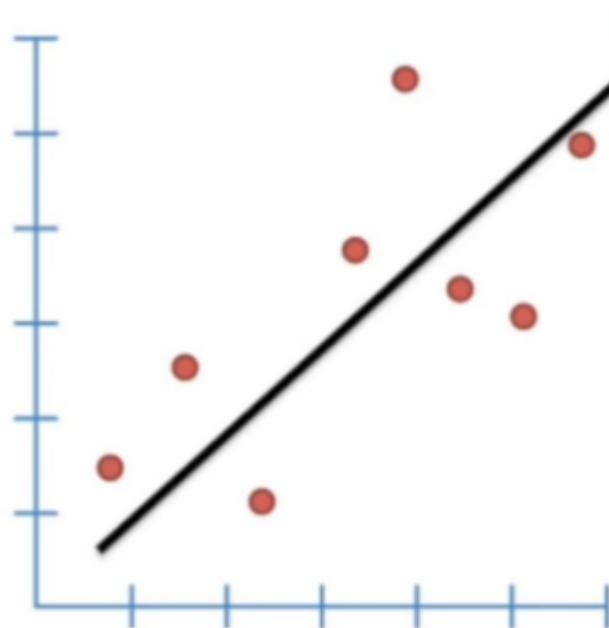
Linear Regression

Linear Regression
(also known as least squares)



Linear Regression & Backpropagation

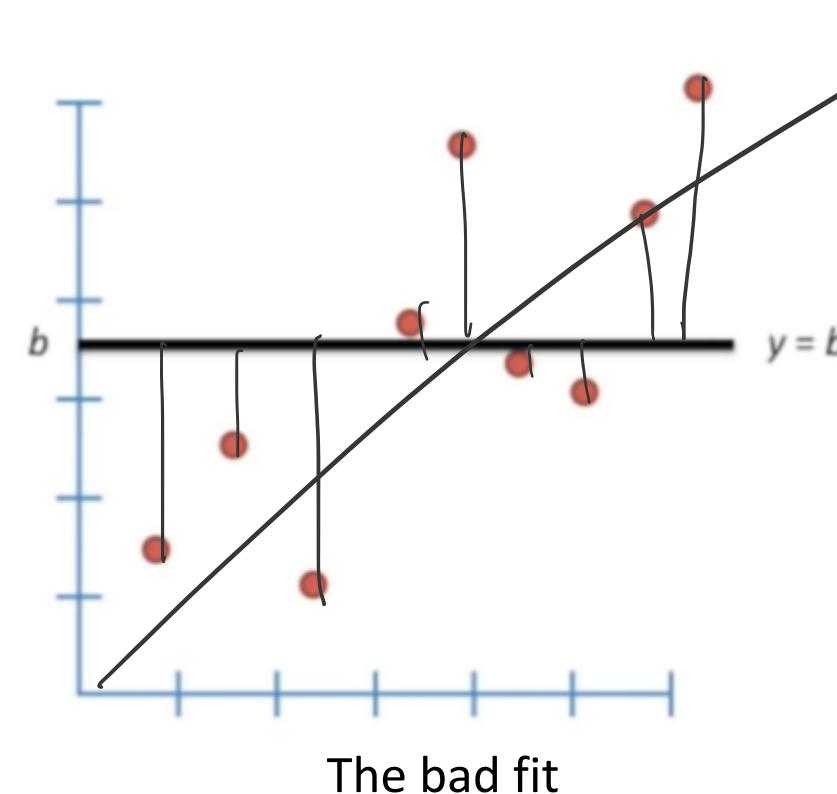
Linear Regression



Fitting a line so we can see what the trend is

Linear Regression & Backpropagation

Linear Regression



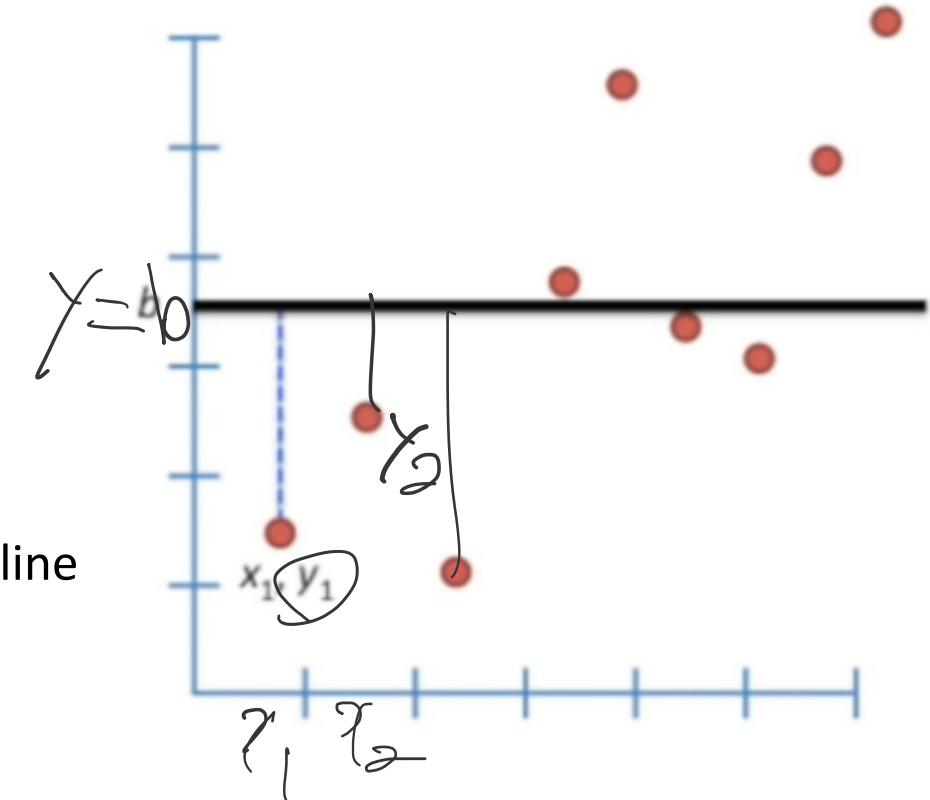
How do we estimate the fitting?

Linear Regression & Backpropagation

Linear Regression

We can measure how well this line fits the data by seeing how close it is to the data points

The distance between the line and 1st data point is $b - y_1$



Linear Regression & Backpropagation

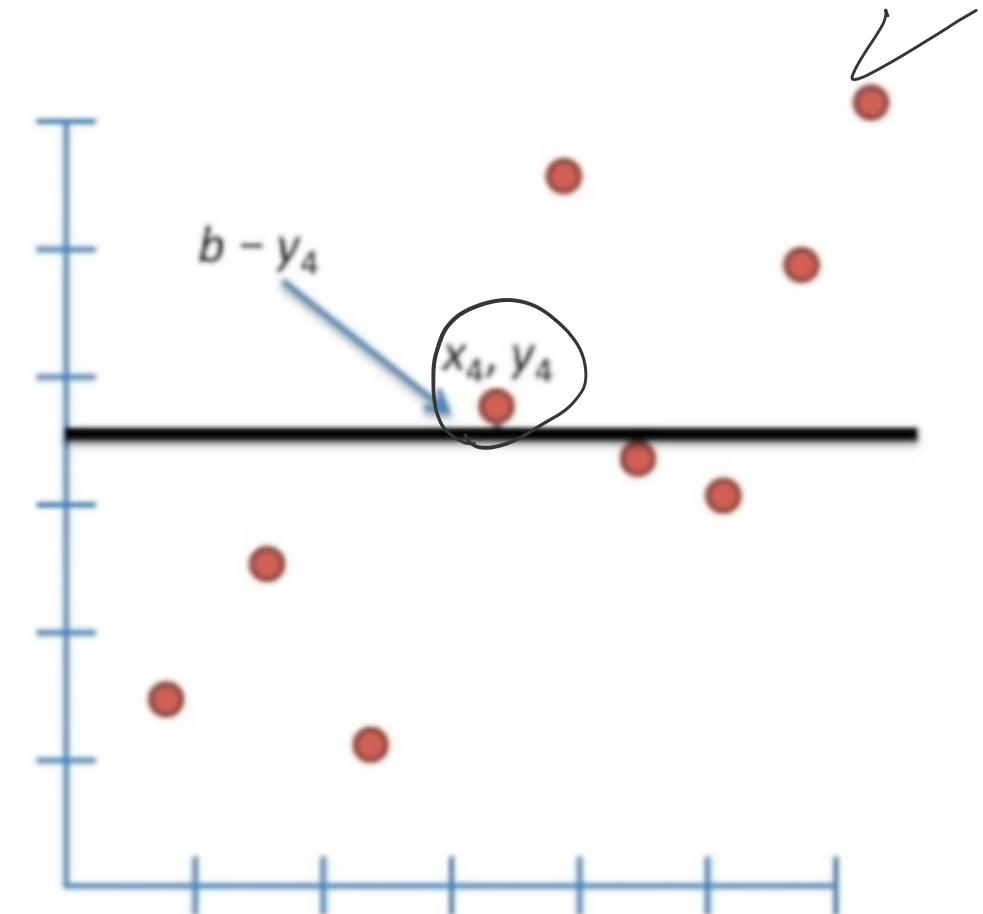
Linear Regression

$$\text{Score} = (b - y_1) + (b - y_2) + (b - y_3) + (b - y_4)$$

$$y_4 > b$$

$\Rightarrow b - y_4$ will be negative.

\Rightarrow Problem?

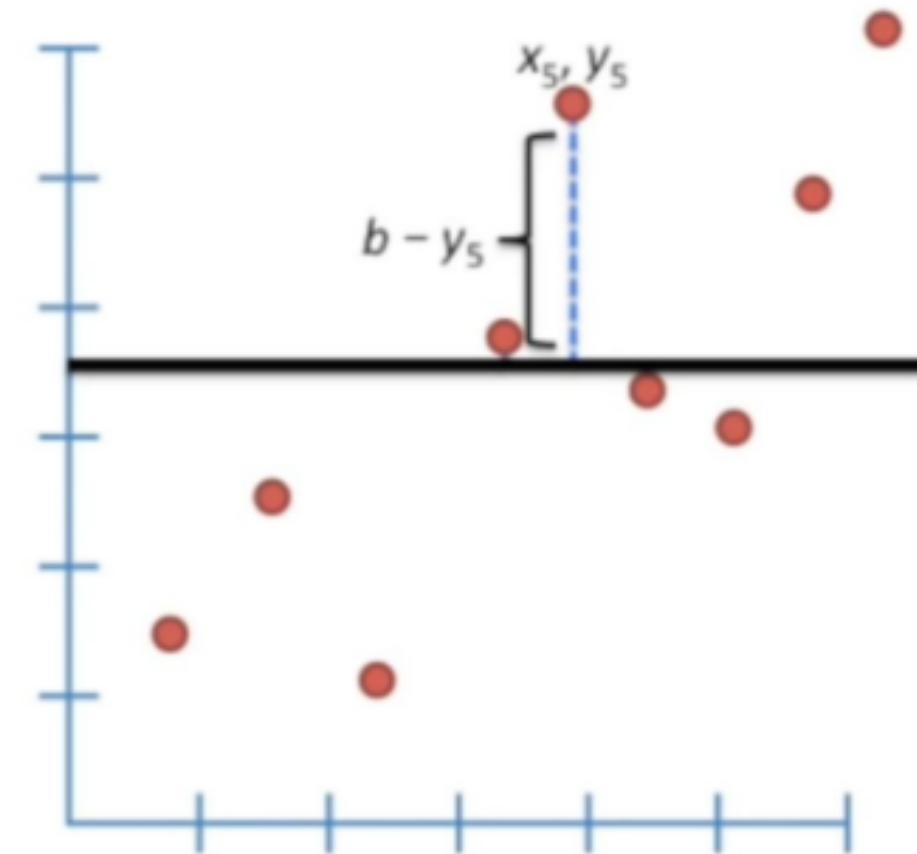


Linear Regression & Backpropagation

Linear Regression

$$\text{Score} = (b - y_1) + (b - y_2) + (b - y_3) + (b - y_4) + (b - y_5)$$

Solution?



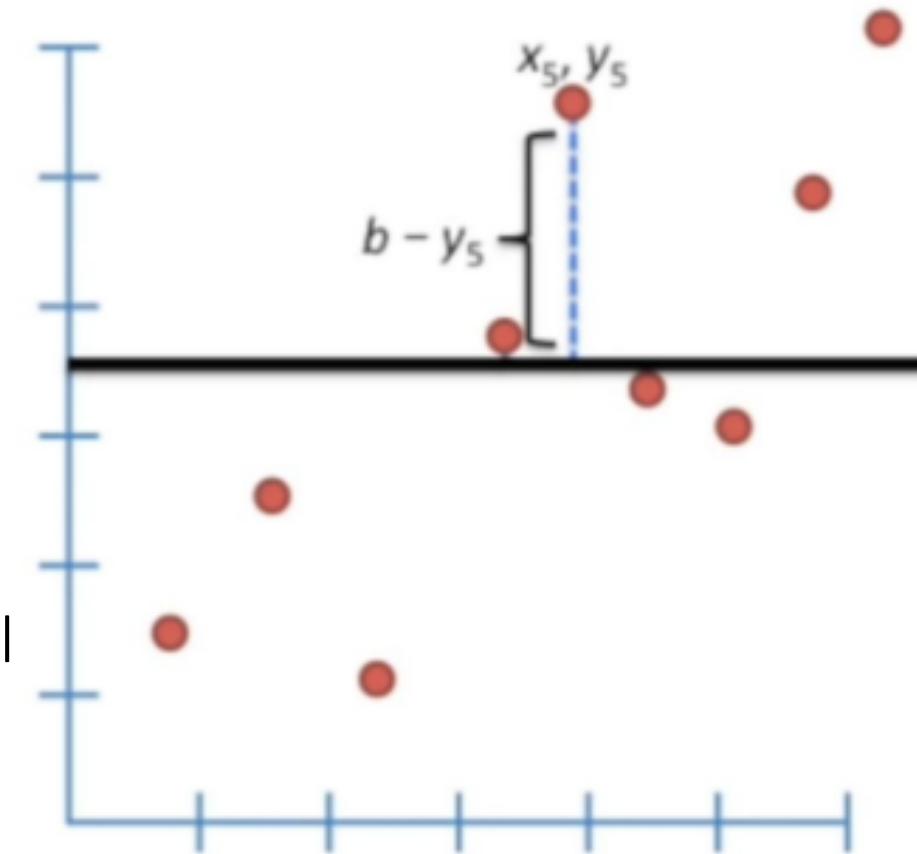
Linear Regression & Backpropagation

Linear Regression

$$\text{Score} = (b - y_1) + (b - y_2) + (b - y_3) + (b - y_4) + (b - y_5)$$

Solution?

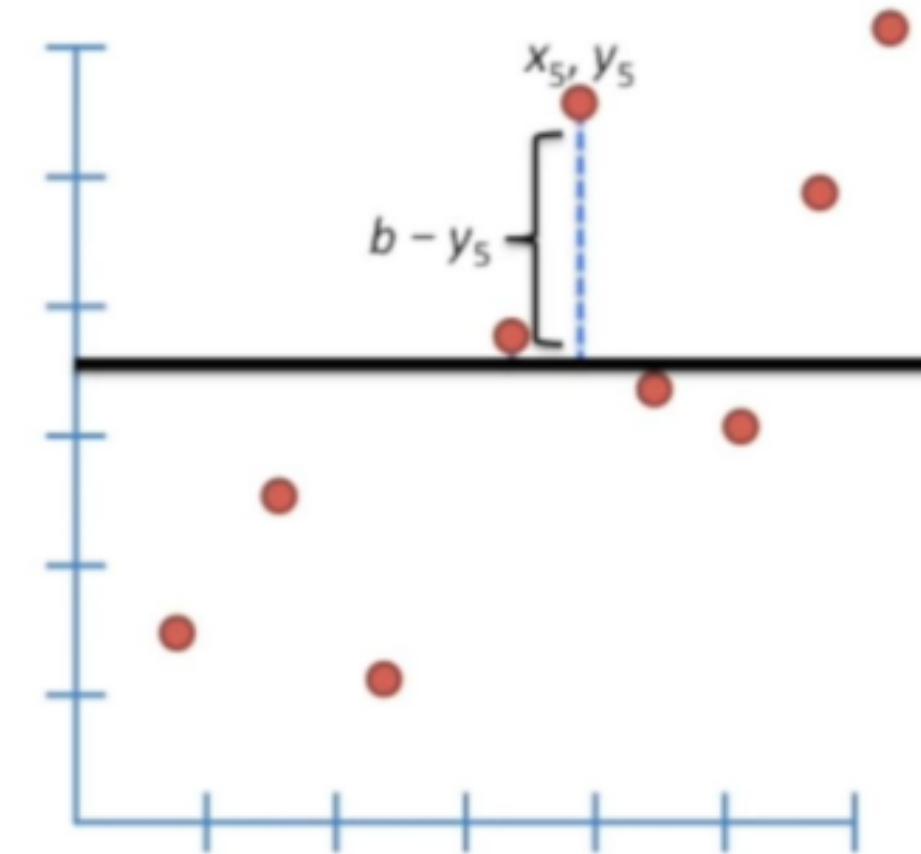
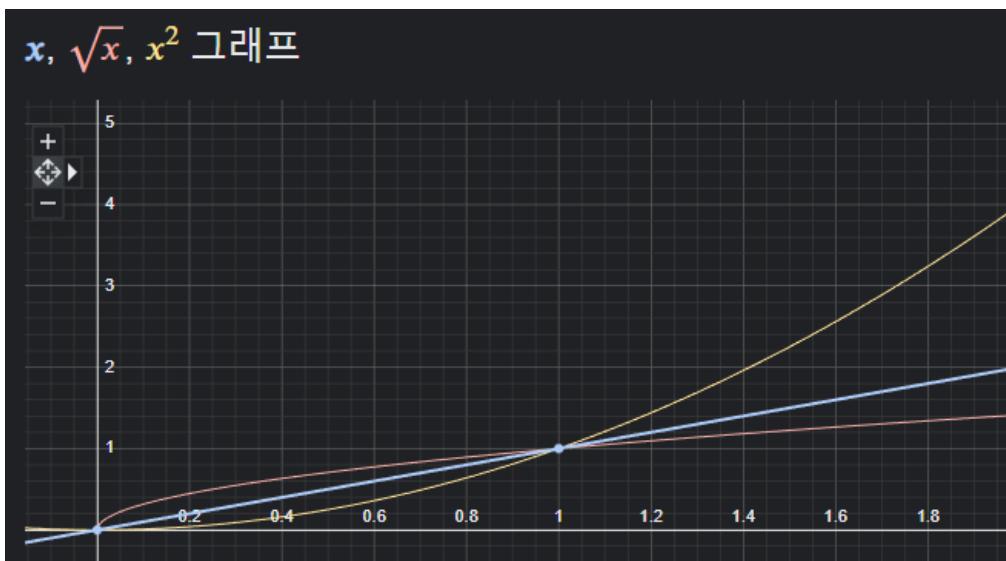
$$\text{Score} = |(b - y_1)| + |(b - y_2)| + |(b - y_3)| + |(b - y_4)| + |(b - y_5)|$$



Linear Regression & Backpropagation

Linear Regression

$$\text{Score} = (b - y_1) + (b - y_2) + (b - y_3) + (b - y_4) + (b - y_5)$$

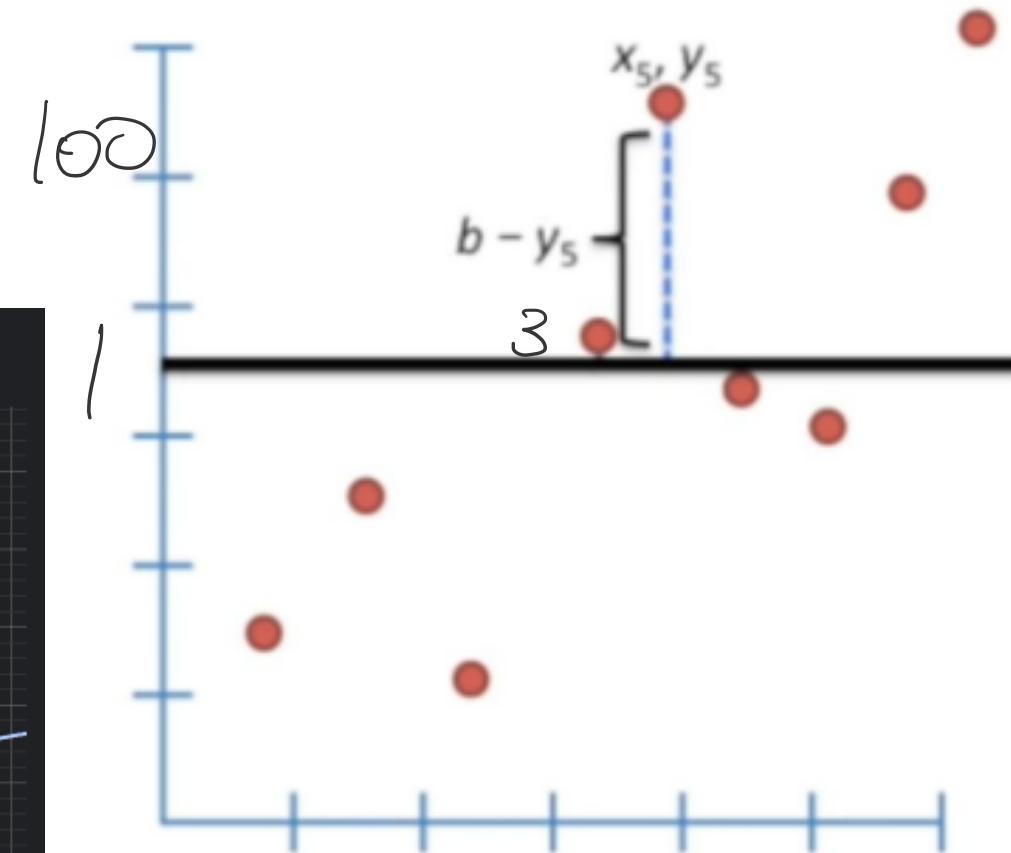
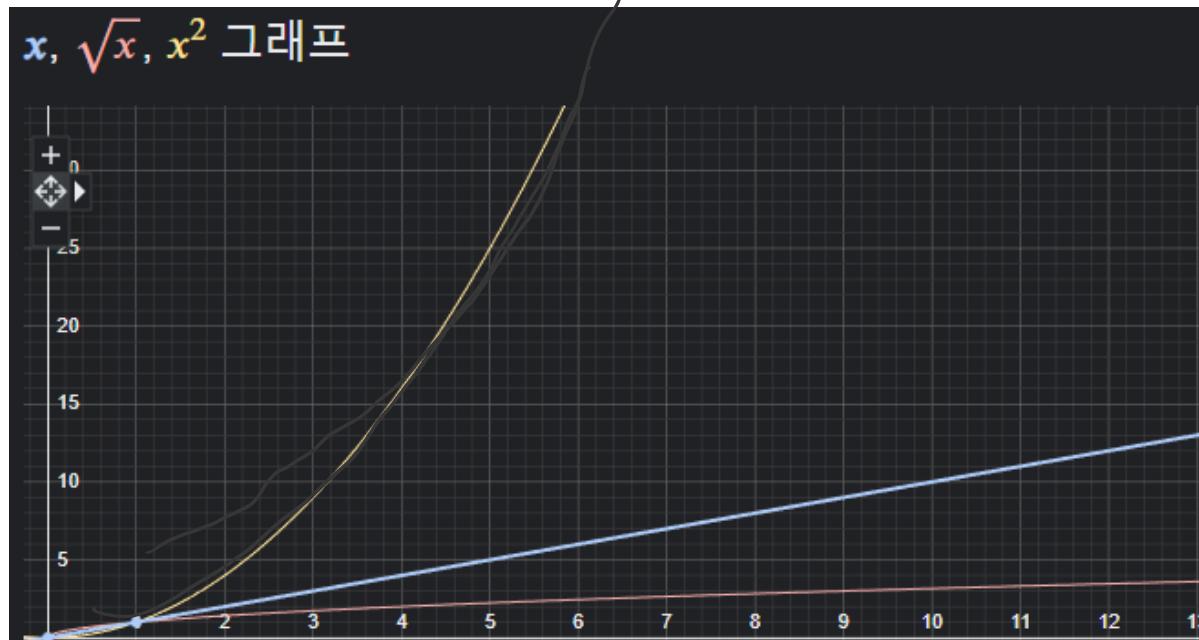


Linear Regression & Backpropagation

Linear Regression

$$\text{Score} = (b-y_1) + (b-y_2) + (b-y_3) + (b-y_4) + (b-y_5)$$

Solve this mathematical problem efficiently!



Linear Regression & Backpropagation

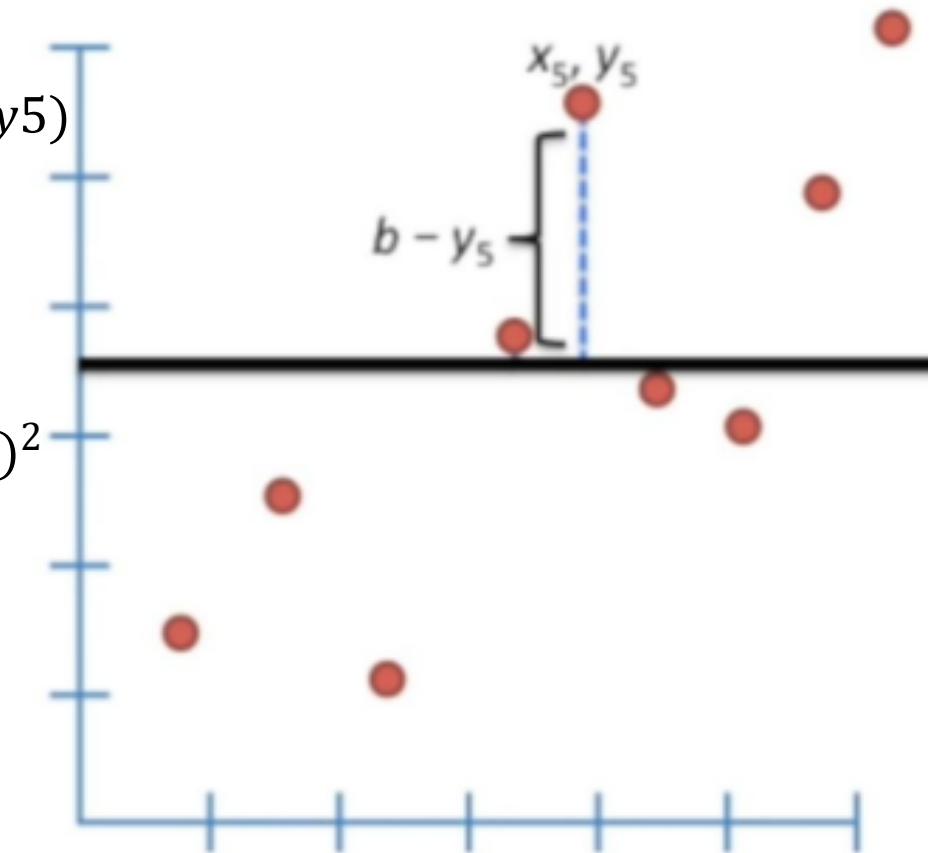
Linear Regression

$$Score = (b - y_1) + (b - y_2) + (b - y_3) + (b - y_4) + (b - y_5)$$

Solve this mathematical problem efficiently!

$$Score = (b - y_1)^2 + (b - y_2)^2 + (b - y_3)^2 + (b - y_4)^2 + (b - y_5)^2$$

We can square each term!

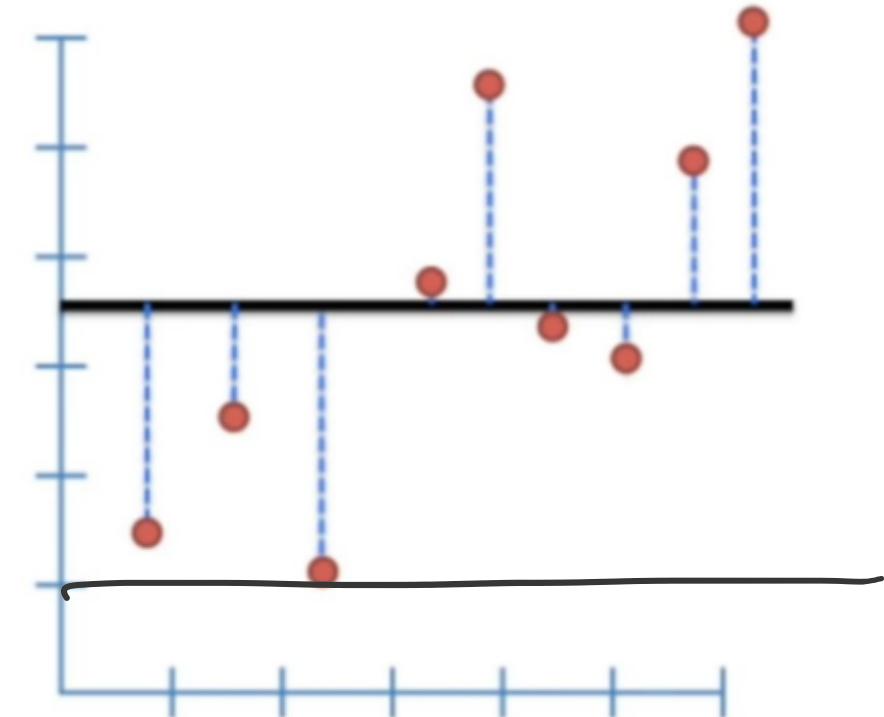


Linear Regression & Backpropagation

Linear Regression

Let's rotate the line a little bit

$$\begin{aligned} \text{Score} &= (b - y_1)^2 + (b - y_2)^2 + \dots + (b - y_4)^2 + (b - y_9)^2 \\ &= 23.45 \end{aligned}$$



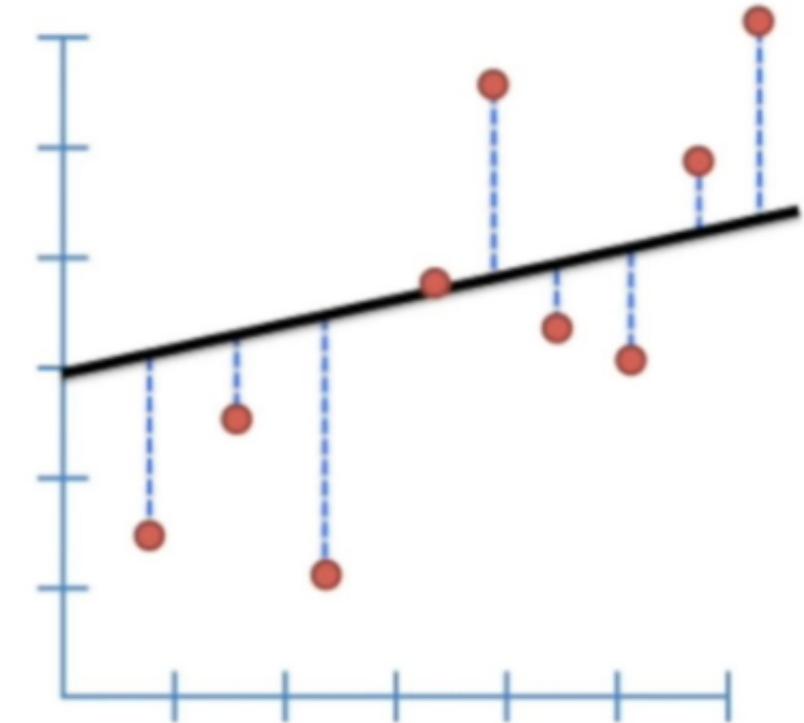
Linear Regression & Backpropagation

Linear Regression

Let's rotate the line a little bit

$$\begin{aligned} \text{Score} &= (b - y_1)^2 + (b - y_2)^2 + \dots + (b - y_4)^2 + (b - y_9)^2 \\ &= \cancel{23.45} \end{aligned}$$

The better result



Linear Regression & Backpropagation

Linear Regression

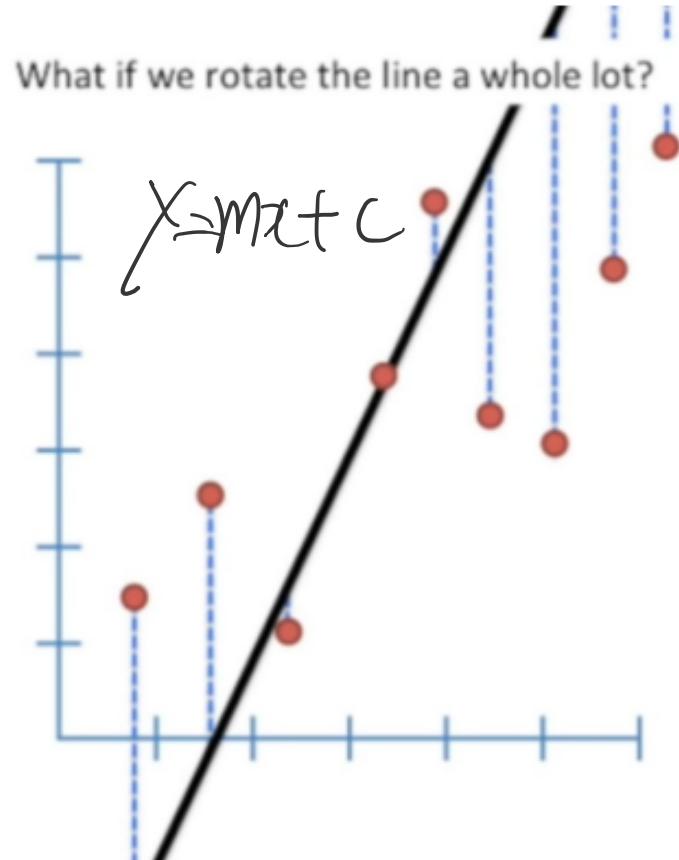
What if we rotate too much?

$$\text{Score} = (b - y_1)^2 + (b - y_2)^2 + \dots + (b - y_4)^2 + (b - y_9)^2$$

= 34.56

여기서 훈련하는 것인가
Learning rate * error

What if we rotate the line a whole lot?

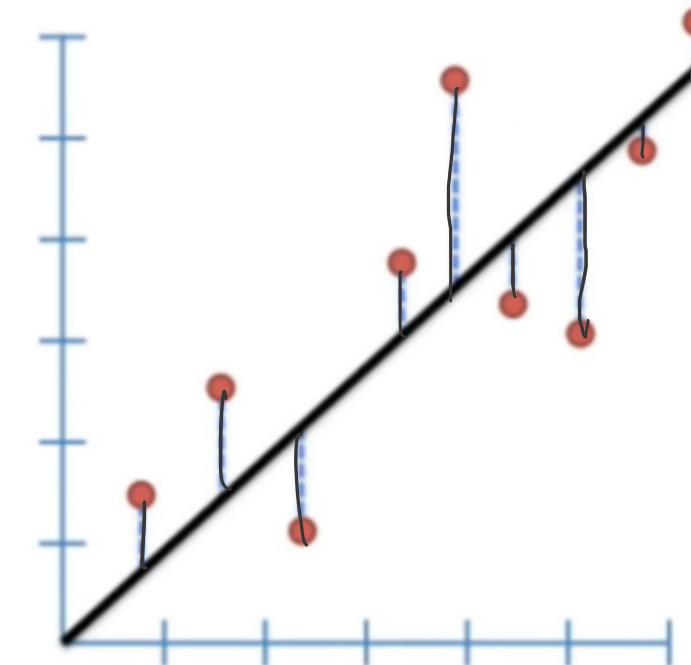


Linear Regression & Backpropagation

Linear Regression Equation expression

Generic line equation: $y = a * x + b$

Goal: Find optimal values for "a" and "b" so that we minimize the sum of squared residuals



Linear Regression & Backpropagation

Linear Regression

Equation expression

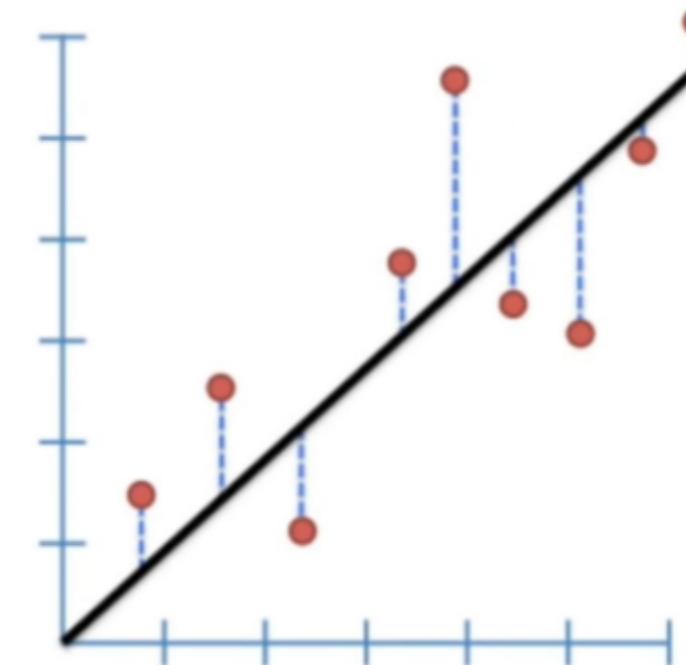
Generic line equation: $y = a * x + b$

Goal: Find optimal values for “a” and “b” so that we minimize the sum of squared residuals

Mean Squared Error (MSE)

Sum of squared residuals

$$= ((a * x_1 + b) - y_1)^2 + ((a * x_2 + b) - y_2)^2 + \dots + ((a * x_9 + b) - y_9)^2$$



Linear Regression & Backpropagation

Linear Regression
Equation expression

Minimize

$$r_1^2 + r_2^2 + \dots + r_q^2$$

Generic line equation: $y = a * x + b$

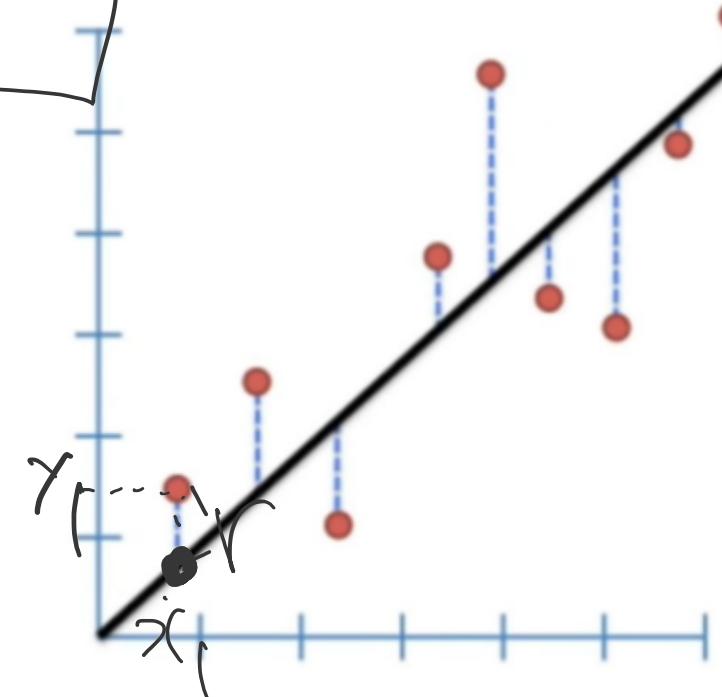
Goal: Find optimal values for "a" and "b" so that we minimize the sum of squared residuals

Sum of squared residuals

$$= ((a * x_1 + b) - y_1)^2 + ((a * x_2 + b) - y_2)^2 + \dots + ((a * x_9 + b) - y_9)^2$$

$(a * x_1 + b)$: The value of the line at x_1

y_1 : observed value at x_1

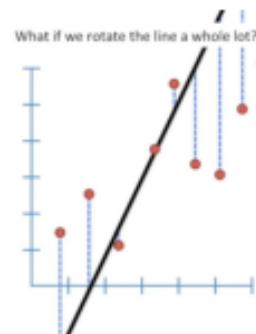
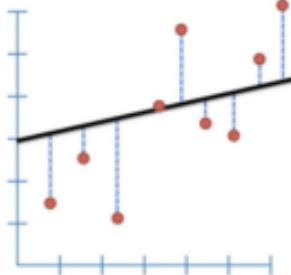
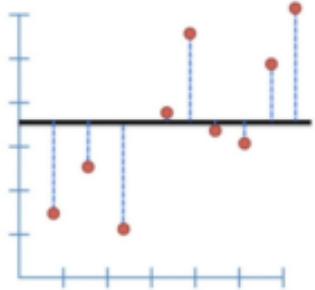


Since we want the line that give us the smallest sum of squares, this methods for finding the best value of "a" and "b" is called "least squares".

Linear Regression & Backpropagation

Linear Regression

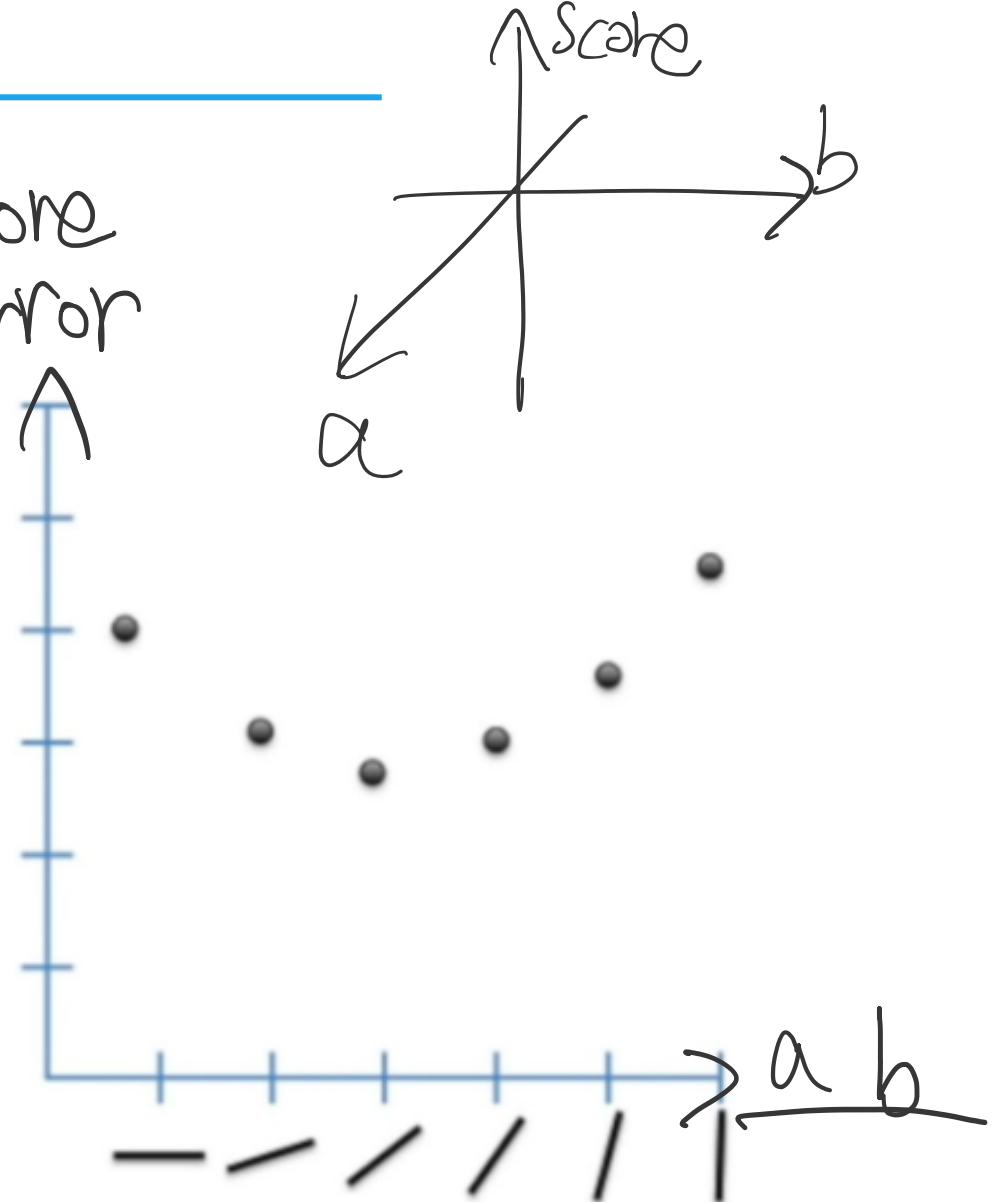
Visualization of each rotation



Score
error

A hand-drawn diagram of a coordinate system. The vertical axis is labeled "Score" and the horizontal axis is labeled "a".

Sum of
squared
residuals



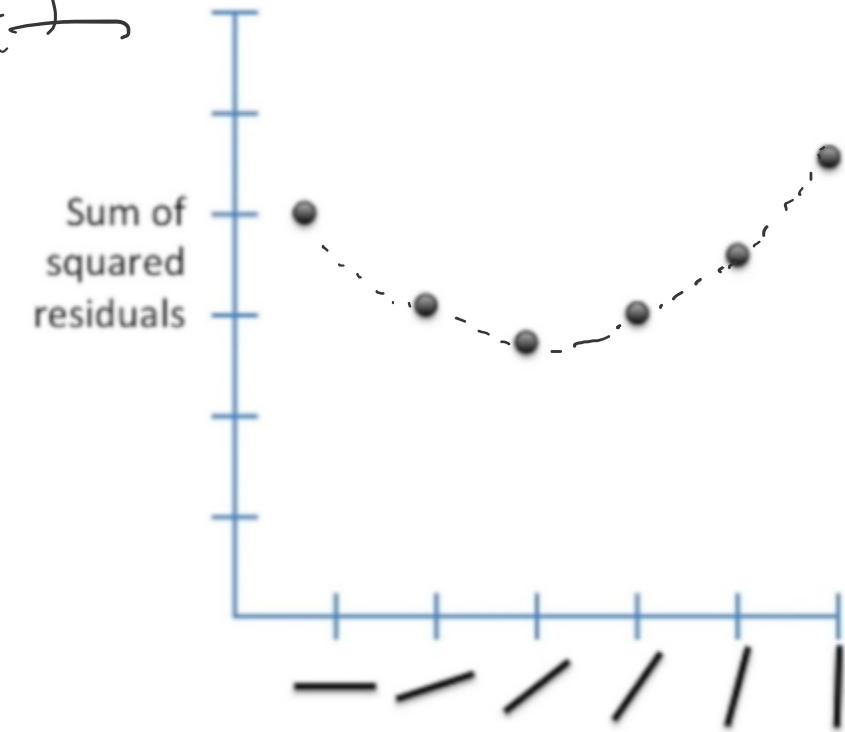
Linear Regression & Backpropagation

Linear Regression

Visualization of each rotation

3 eu 3 l
3 eu 3 l

How do we find the optimal rotation for the line?



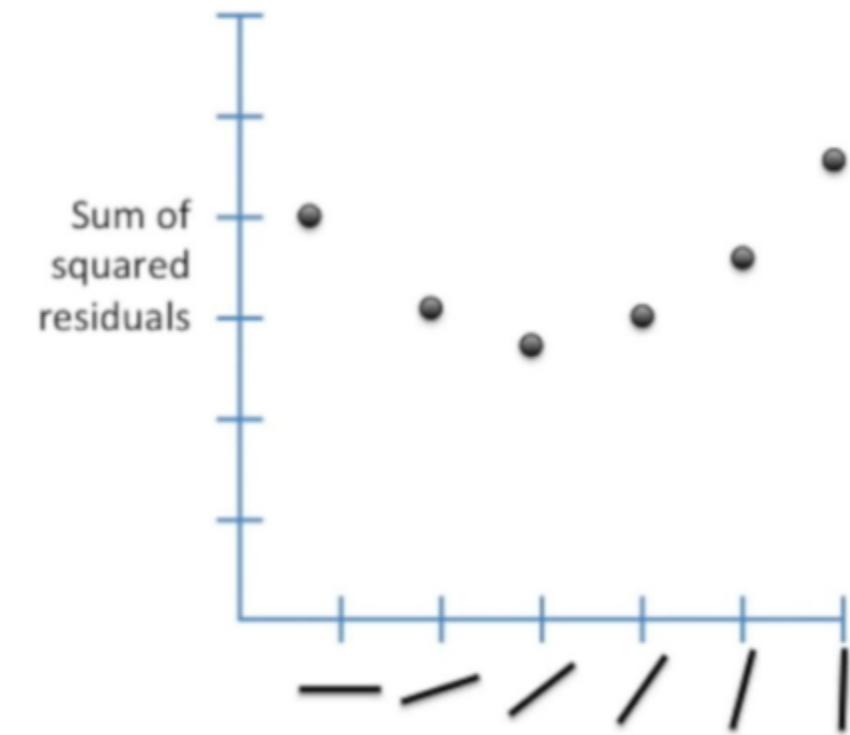
Linear Regression & Backpropagation

Linear Regression

Visualization of each rotation

How do we find the optimal rotation for the line?

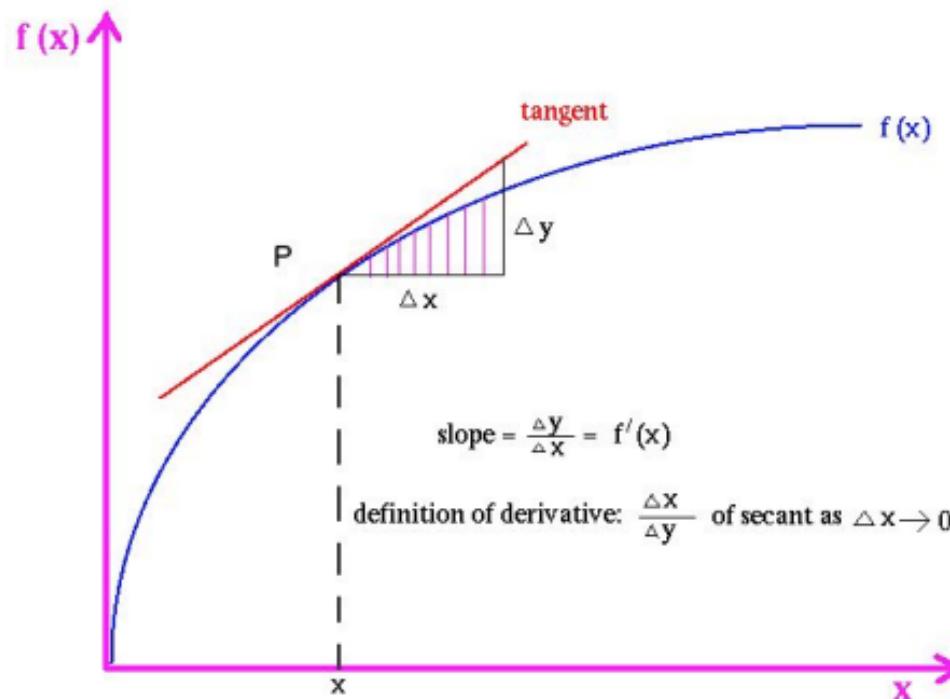
Derivatives!



Linear Regression & Backpropagation

Linear Regression

Derivative



Linear Regression & Backpropagation

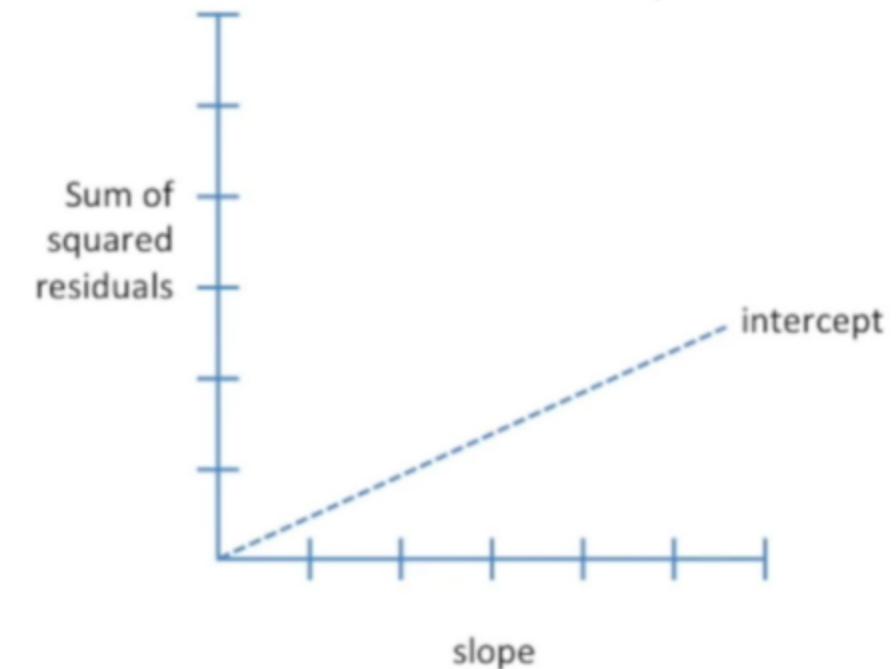
Linear Regression

Two independent variables

$$y = a * x + b$$

We have two independent variables that needs to be optimized

- 1) a : the slope
- 2) b : the intercept



Linear Regression & Backpropagation

Linear Regression

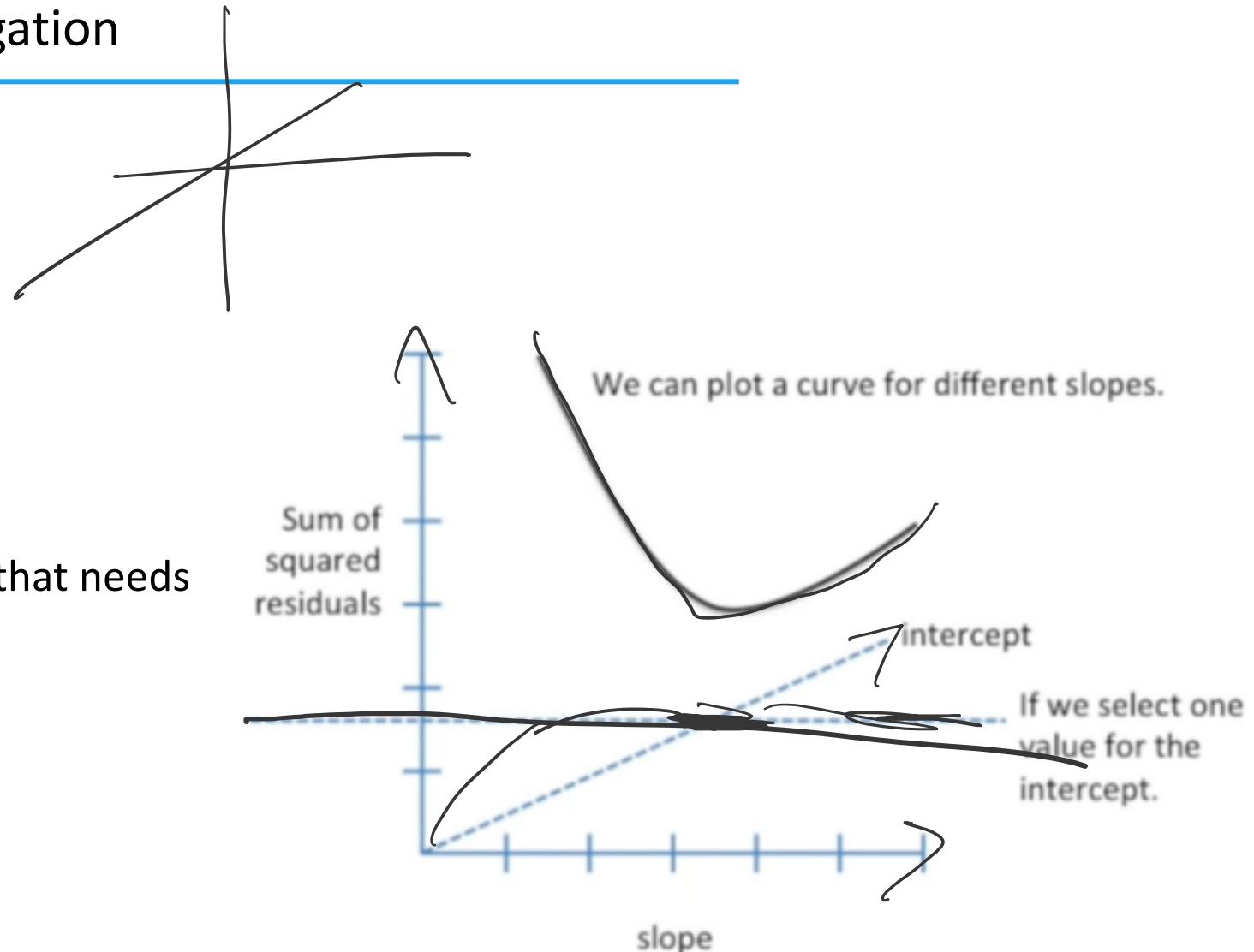
If the intercept has fixed?

=> Partial derivative

$$y = a * x + b$$

We have two independent variables that needs to be optimized

- 1) a : the slope
- 2) b : the intercept



Linear Regression & Backpropagation

Linear Regression

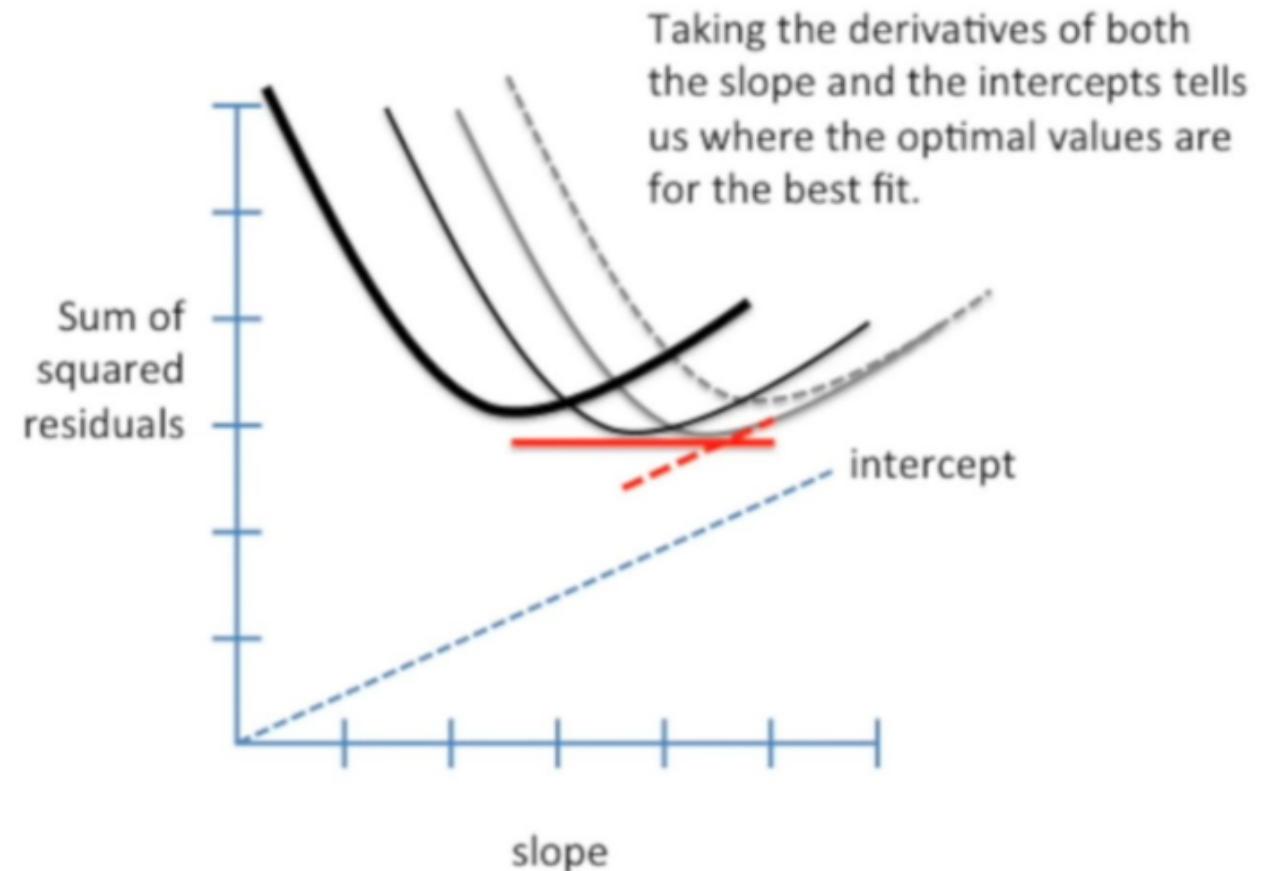
Take the derivative for both independent variables

=> Total derivative

$$y = a * x + b$$

We have two independent variables that needs to be optimized

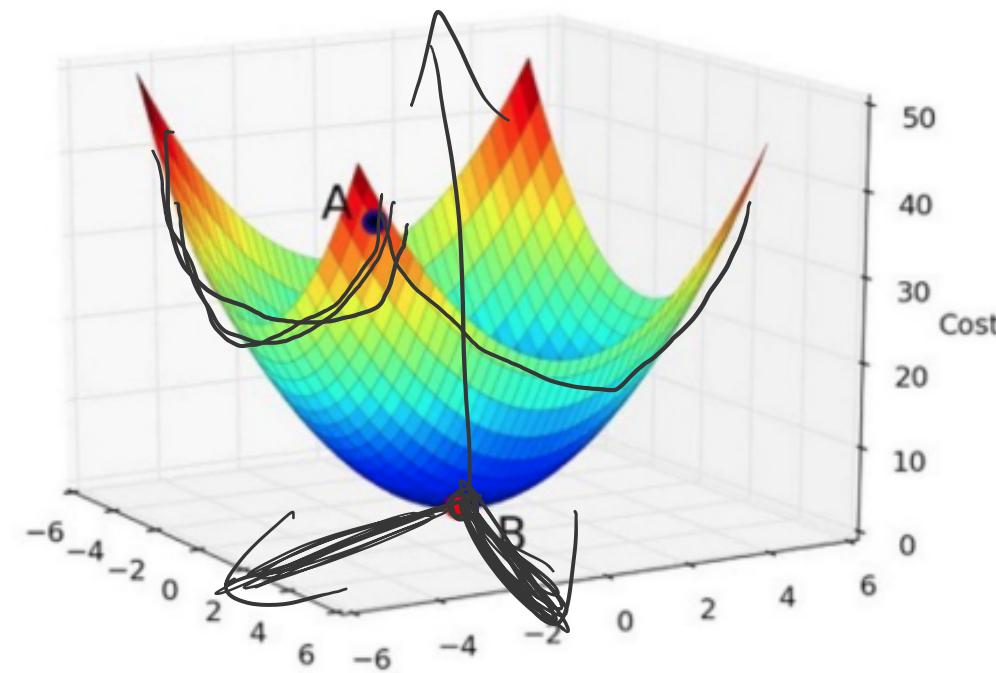
- 1) a : the slope
- 2) b : the intercept



Linear Regression & Backpropagation

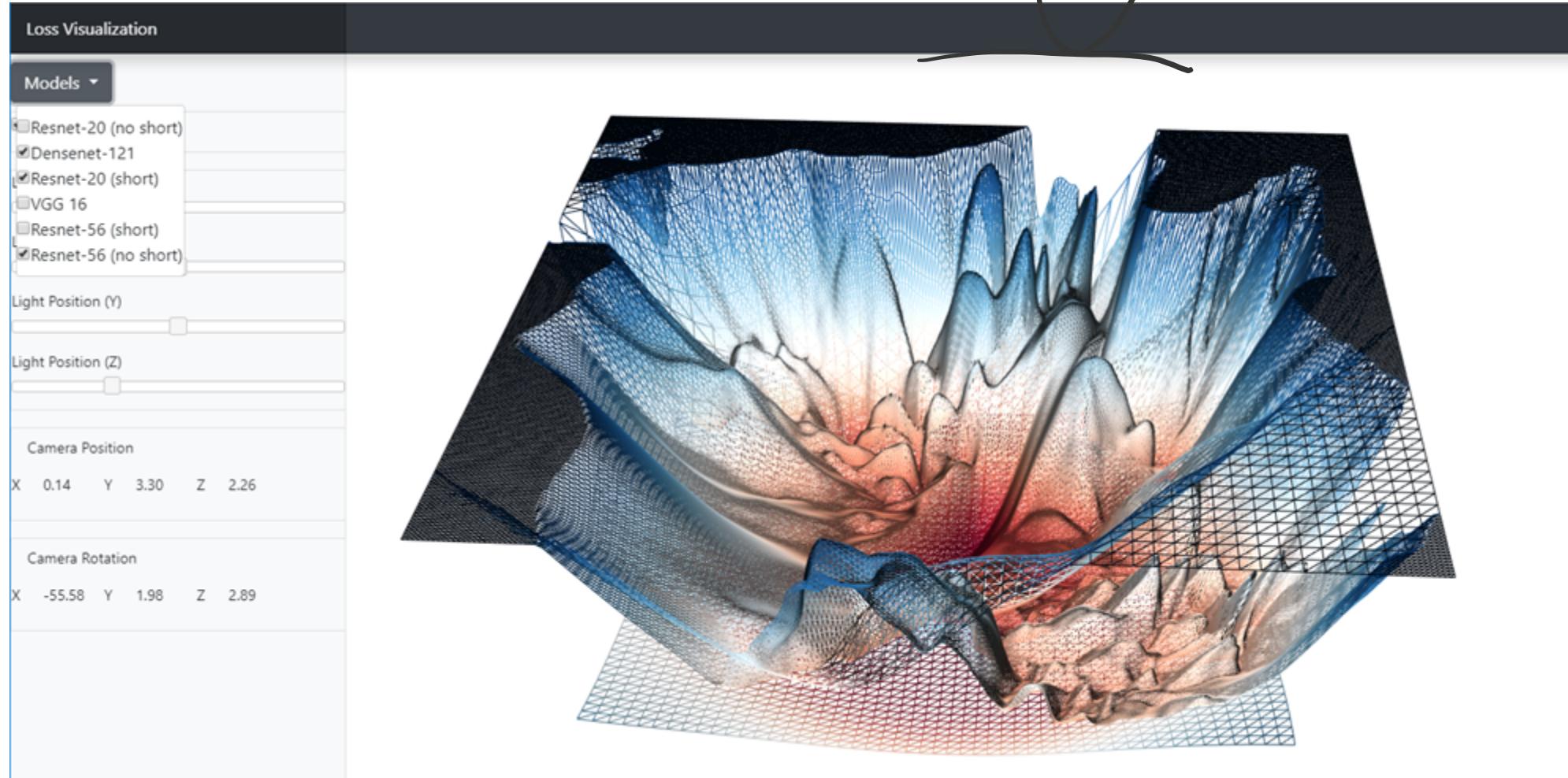
Linear Regression

Advanced topic: Gradient Descent



Back to ANN

Deep Neural Network (DNN)

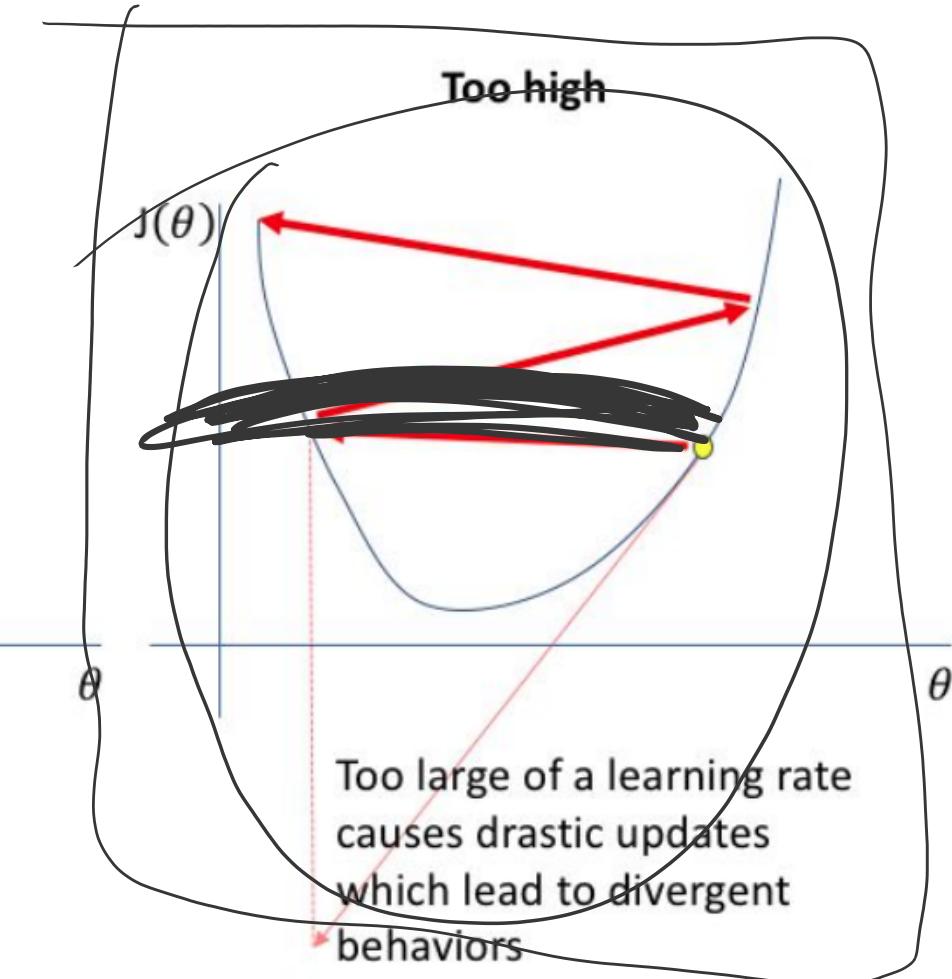
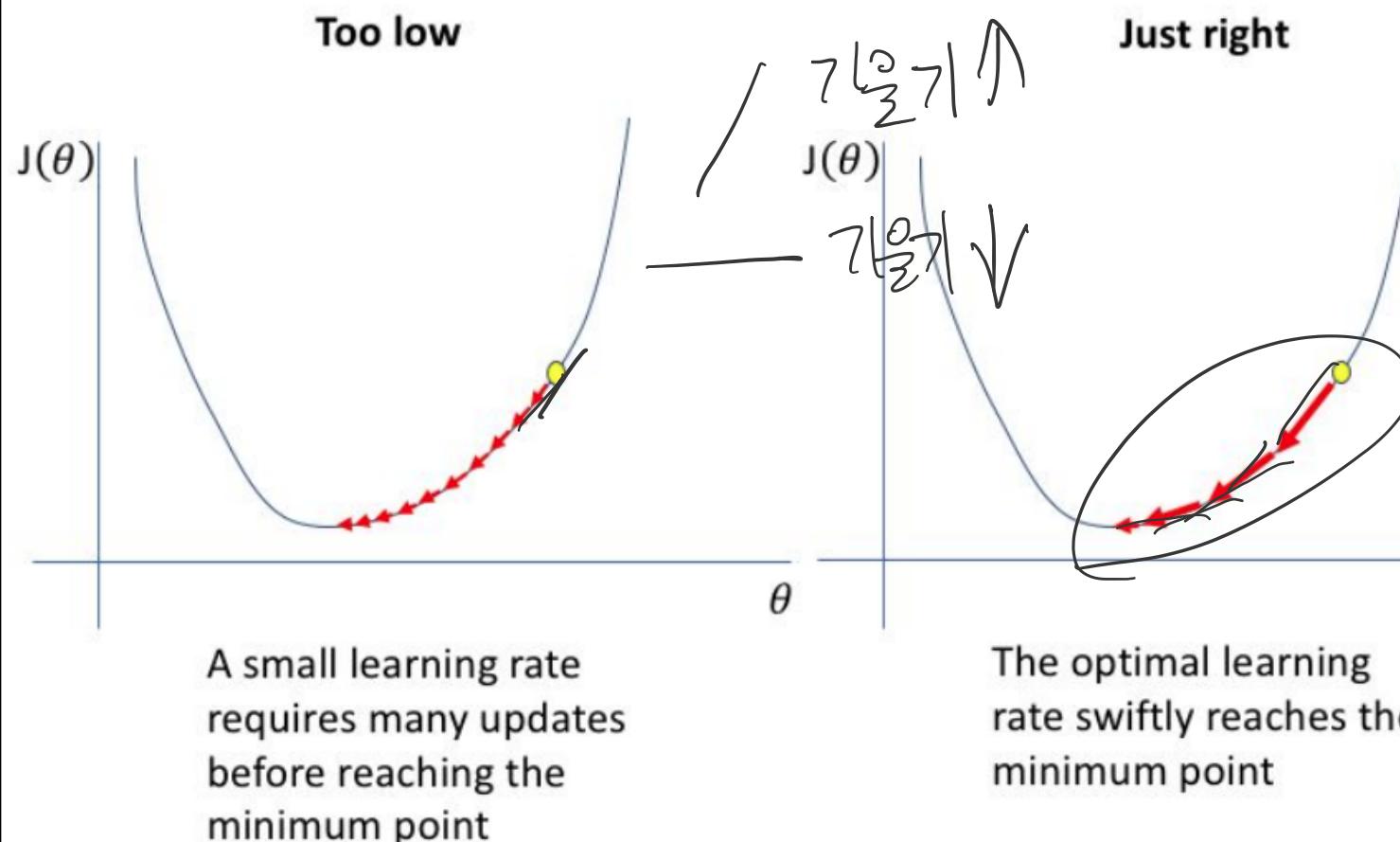


Visualized the loss map of Deep
Neural Network

Linear Regression & Backpropagation

Linear Regression

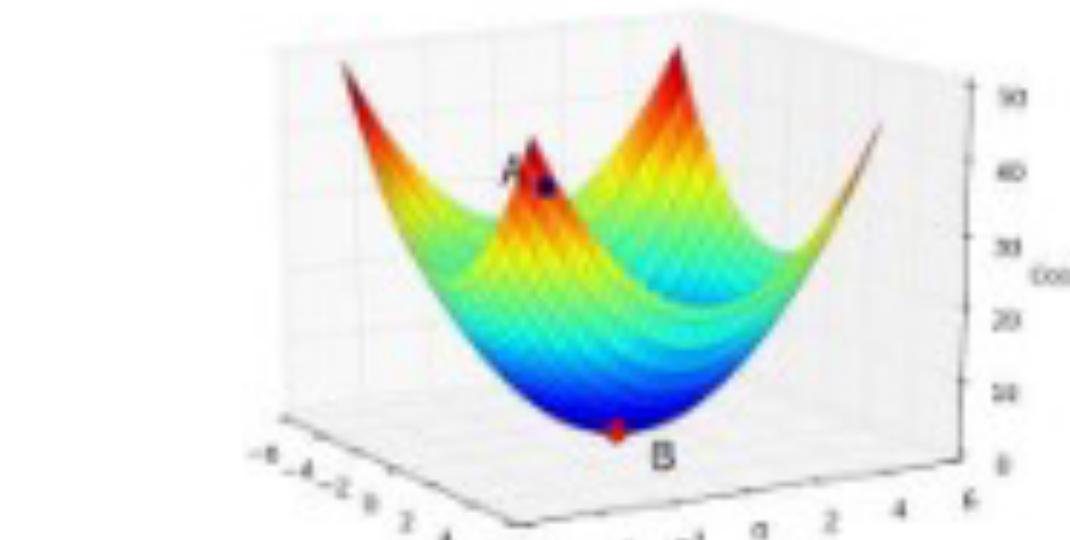
Advanced topic: Gradient Descent



Linear Regression & Backpropagation

Linear Regression

Advanced topic: Gradient Descent

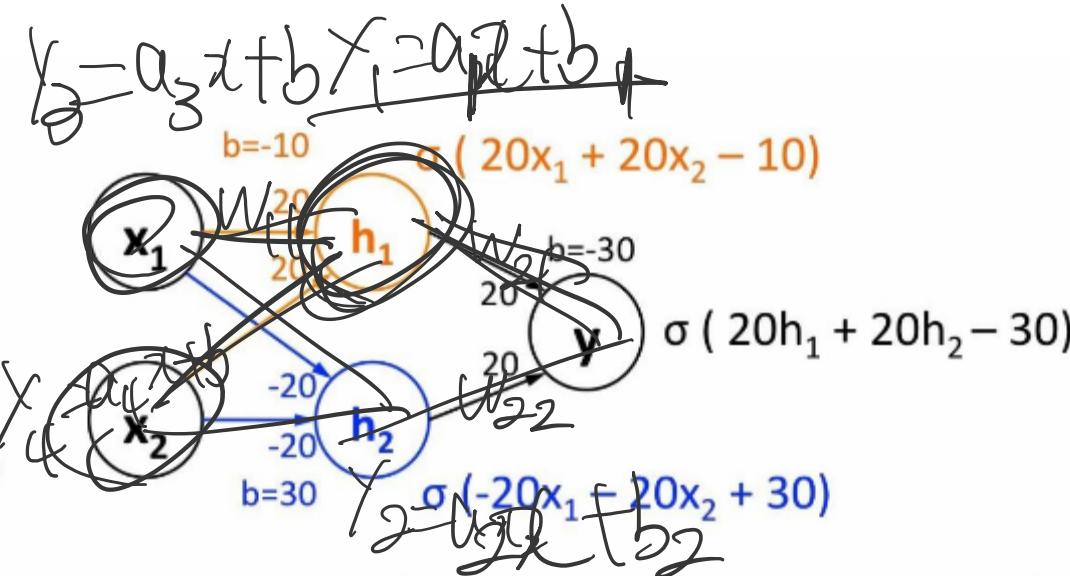
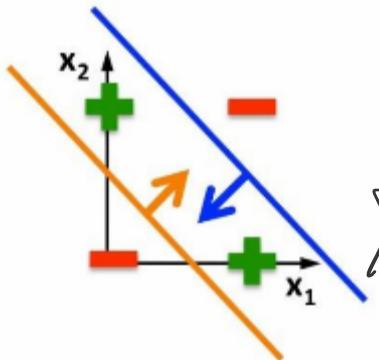


Linear Regression & Backpropagation

Linear Regression

Advanced topic: Backpropagation

Linear classifiers
cannot solve this



$$\sigma(20*0 + 20*0 - 10) \approx 0$$

$$\sigma(20*1 + 20*1 - 10) \approx 1$$

$$\sigma(20*0 + 20*1 - 10) \approx 1$$

$$\sigma(20*1 + 20*0 - 10) \approx 1$$

$$\sigma(-20*0 - 20*0 + 30) \approx 1$$

$$\sigma(-20*1 - 20*1 + 30) \approx 0$$

$$\sigma(-20*0 - 20*1 + 30) \approx 1$$

$$\sigma(-20*1 - 20*0 + 30) \approx 1$$

$$\sigma(20*0 + 20*1 - 30) \approx 0$$

$$\sigma(20*1 + 20*0 - 30) \approx 0$$

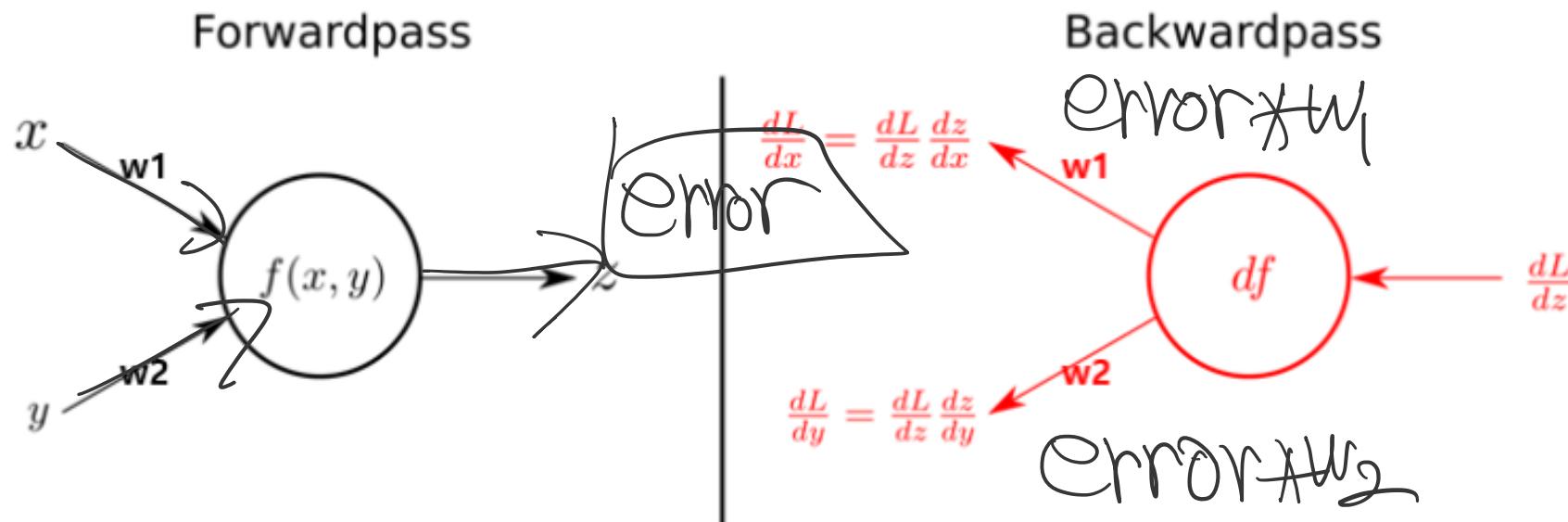
$$\sigma(20*1 + 20*1 - 30) \approx 1$$

$$\sigma(20*1 + 20*1 - 30) \approx 1$$

Linear Regression & Backpropagation

Linear Regression

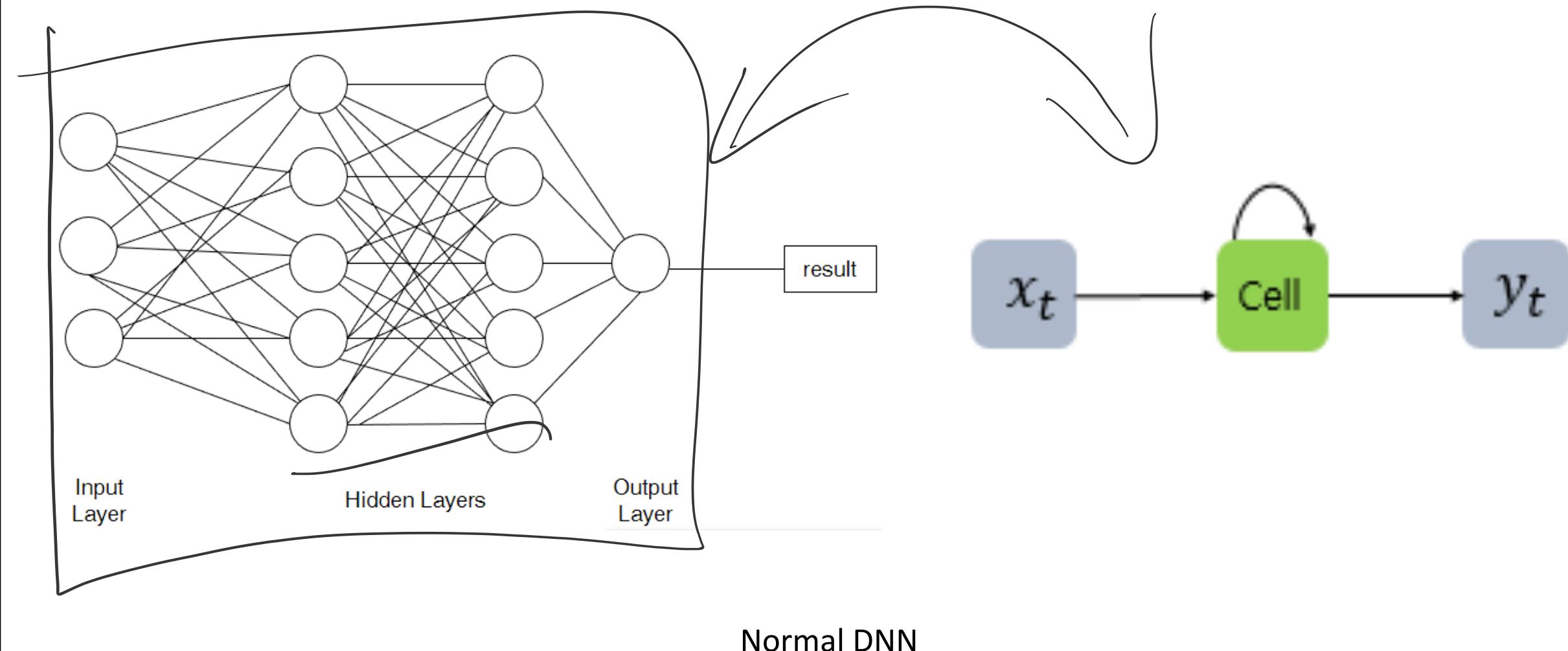
Advanced topic: Backpropagation



Calculate through the Chain rule!

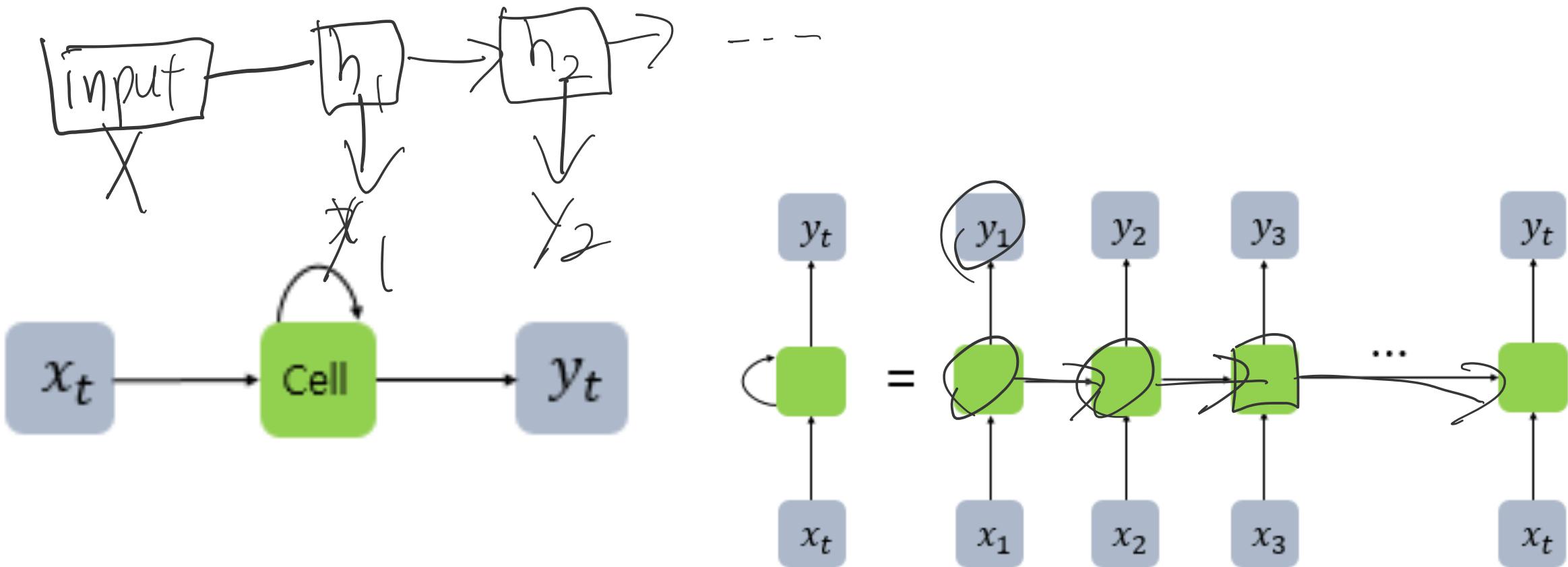
Recurrent Neural Networks (RNN)

What is RNN?



Recurrent Neural Networks (RNN)

What is RNN?

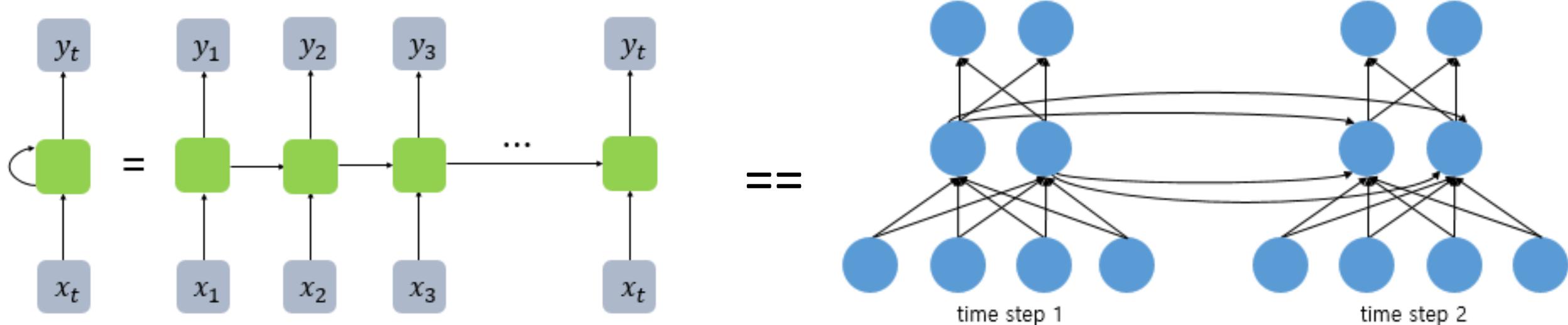


Output is calculated sequentially!
(or Recursively)

Recurrent Neural Networks (RNN)

Visualization of RNN

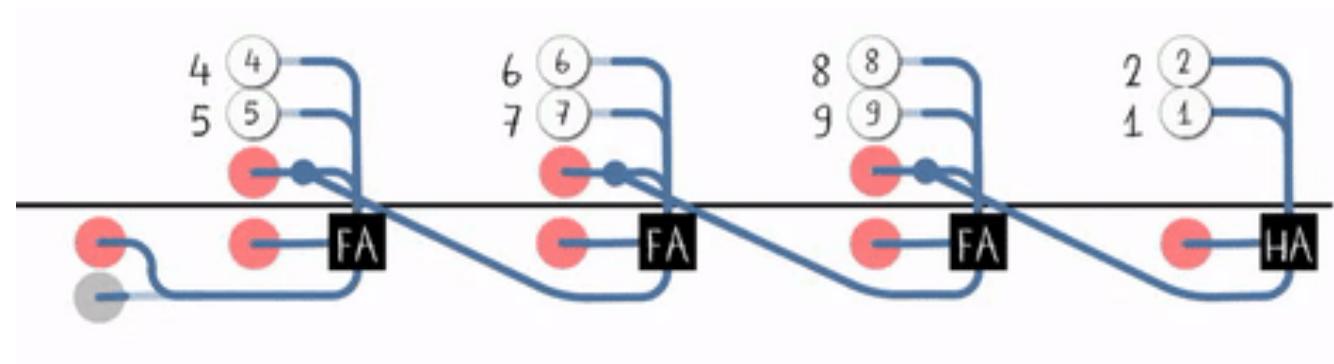
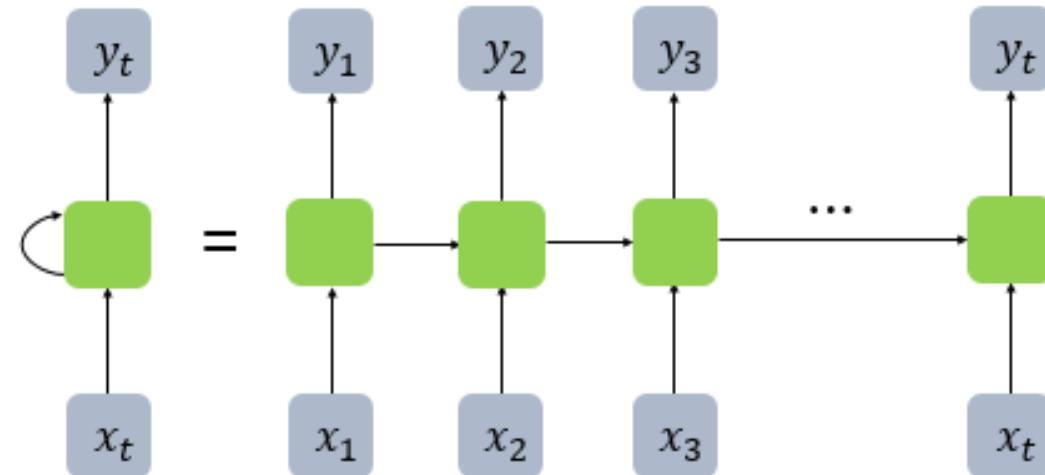
$$y_1 \propto x_1 \quad y_2 \propto x_1 * f + x_2$$



Recurrent Neural Networks (RNN)

ETC

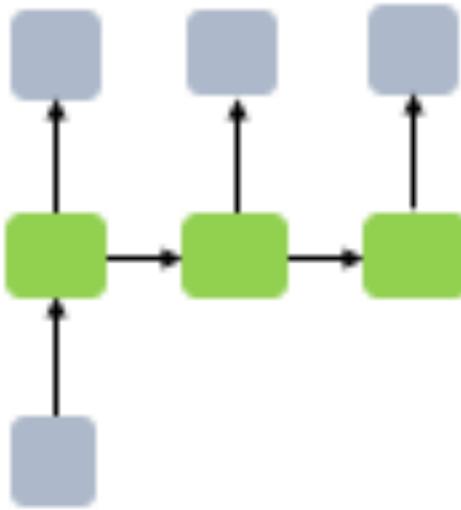
LSTM (1997)
GRU



Ripple Carry Adder

Recurrent Neural Networks (RNN)

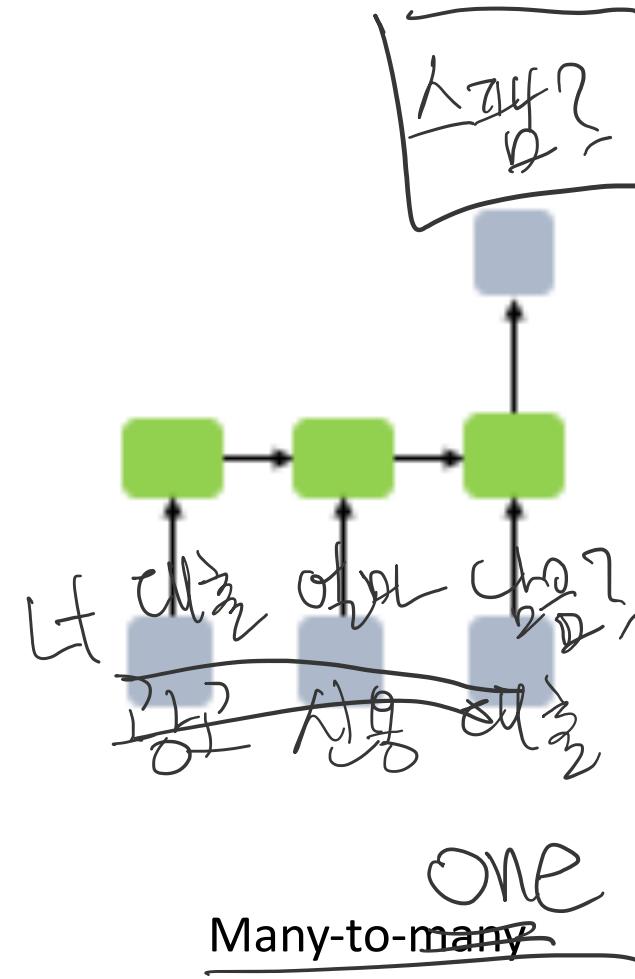
RNN: Examples



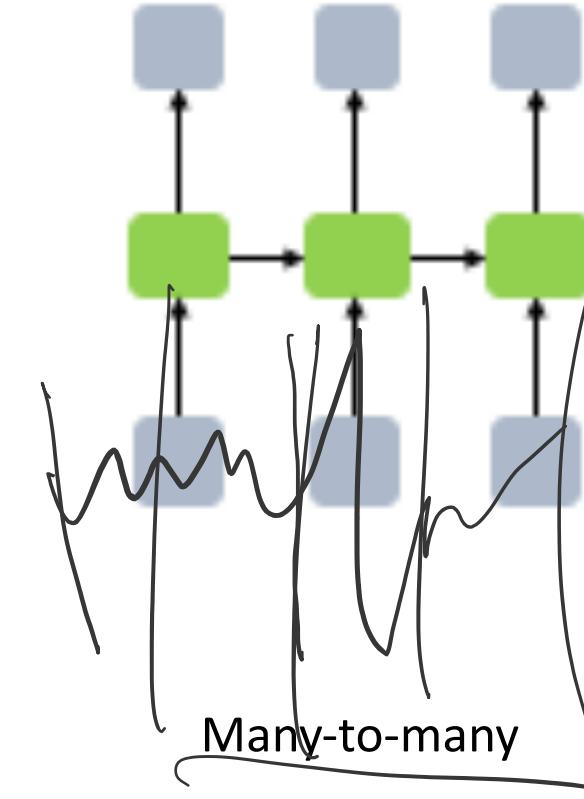
One-to-many

Ex)

~~Input: Word~~
~~Output: Emotions~~



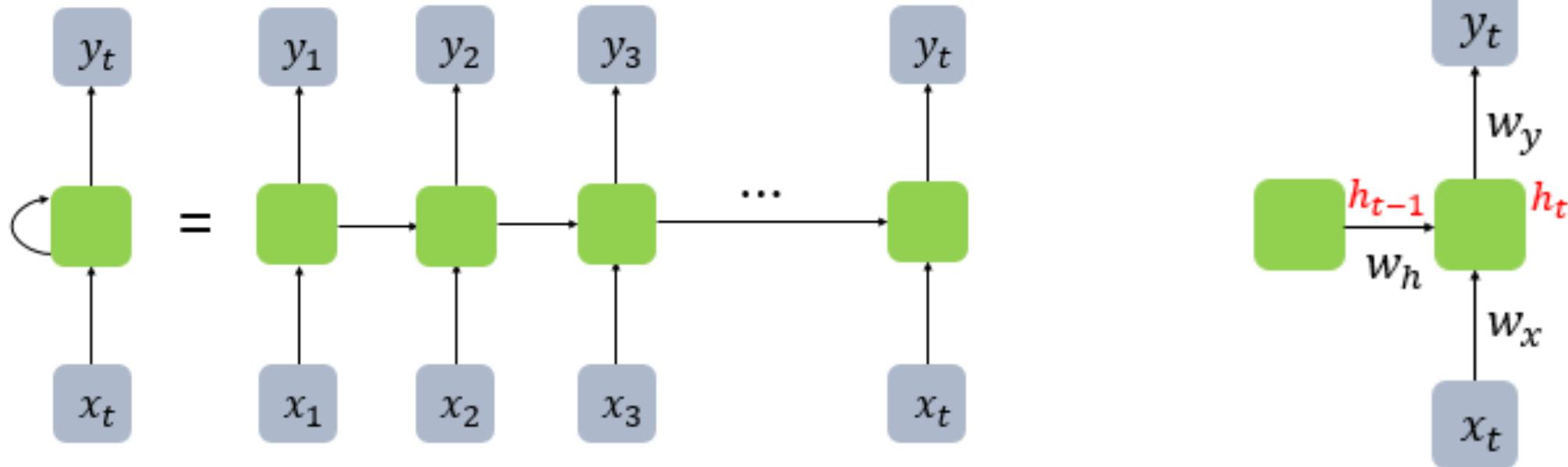
Many-to-many
Input: Words
Output: Classes



Many-to-many
Input: Stock price
Output: Stock price

Recurrent Neural Networks (RNN)

Equations for RNN



$$\text{Hidden } h_t = \tanh(W_x x_t + W_h h_{t-1} + b)$$

$$\text{Output } y_t = f(W_y h_t + b)$$

f is a nonlinear activation function.

d : Dimension of the input (vector)

D_h : Size of the hidden state

x_t : $(d \times 1)$

W : $(D_h \times d)$

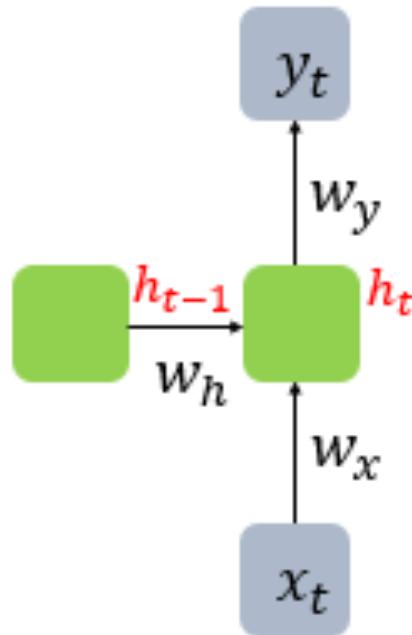
W_h : $(D_h \times D_h)$

h_{t-1} : $(D_h \times 1)$

b : $(D_h \times 1)$

Recurrent Neural Networks (RNN)

Equations for RNN



Batch size: 1
 D_h and d : 4

$$\tanh \left(W_h \begin{matrix} h_{t-1} \\ D_h \times 1 \end{matrix} + W_x \begin{matrix} x_t \\ d \times 1 \end{matrix} + b \begin{matrix} \\ D_h \times 1 \end{matrix} \right) = \begin{matrix} h_t \\ D_h \times 1 \end{matrix}$$

Visualization of RNN calculation

Hidden $h_t = \tanh(W_x x_t + W_h h_{t-1} + b)$
Output $y_t = f(W_y h_t + b)$
 f is a nonlinear activation function.

d : Dimension of the input (vector)
 D_h : Size of the hidden state
 x_t : $(d \times 1)$
 W : $(D_h \times d)$
 W_h : $(D_h \times D_h)$
 h_{t-1} : $(D_h \times 1)$
 b : $(D_h \times 1)$

Practice: Partial Derivatives for Linear Equations

Equations for RNN

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

$$= C + 1$$

MSE = mean squared error

n = number of data points

Y_i = observed values

\hat{Y}_i = predicted values

$$\frac{1}{n} \sum_{i=1}^n (Y_i - (mx_i + c))^2$$

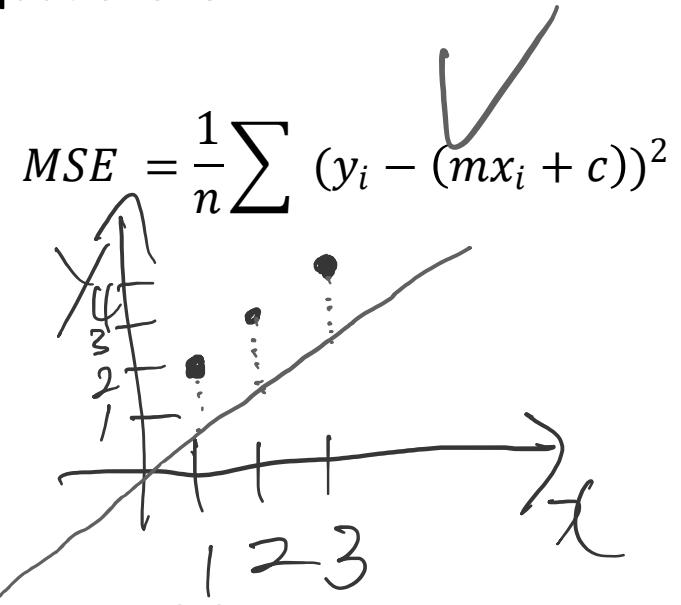
$$y = mx + c$$

$$\frac{\partial MSE}{\partial m} = \frac{2}{n} \sum -x_i(y_i - (mx_i + c))$$

$$\frac{\partial MSE}{\partial c} = \frac{2}{n} \sum - (y_i - (mx_i + c))$$

Practice: Partial Derivatives for Linear Equations

Equations for RNN



Problem

$$(x_1, y_1) = (1, 2)$$

$$(x_2, y_2) = (2, 3)$$

$$(x_3, y_3) = (3, 4)$$

$$y = mx + c$$
$$\frac{\partial MSE}{\partial m} = \frac{2}{n} \sum -x_i(y_i - (mx_i + c))$$
$$\frac{\partial MSE}{\partial c} = \frac{2}{n} \sum -(y_i - (mx_i + c))$$

Initial condition
m:1, c:0

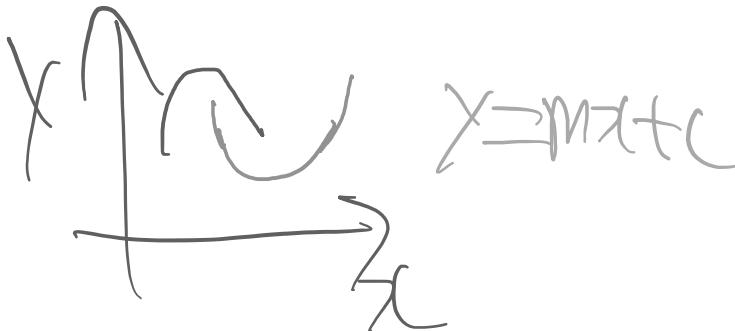
$y = x$

Learning rate = 0.01

Practice: Partial Derivatives for Linear Equations

Equations for RNN

$$MSE = \frac{1}{n} \sum (y_i - (mx_i + c))^2$$



$$y = mx + c$$

$$\frac{\partial MSE}{\partial m} = \frac{2}{n} \sum -x_i(y_i - (mx_i + c))$$

$$\frac{\partial MSE}{\partial c} = \frac{2}{n} \sum -(y_i - (mx_i + c))$$

$$y = 1.04x + 0.02$$

Initial condition
 $\underbrace{m:1, c:0}_{y = x}$

Problem

$$(x_1, y_1) = (1, 2)$$

$$(x_2, y_2) = (2, 3)$$

$$(x_3, y_3) = (3, 4)$$

$$MSE = \frac{1}{3} [(2-1)^2 + (3-2)^2 + (4-3)^2] = 1$$

$$\frac{\partial MSE}{\partial m} = \frac{2}{3} [(-1)(2-1) - 2(3-2) - 3(4-3)] = -2 - 4 \cancel{- 0.02}$$

$$\frac{\partial MSE}{\partial c} = \frac{2}{3} [-(2-1) - (3-2) - (4-3)] = -2$$

-0.04
-2 × 0.01
-0.02



References

[https://youtube.com/playlist?list=PLblh5JKOoLUIzaEkCLIUxQFjP
llapw8nU](https://youtube.com/playlist?list=PLblh5JKOoLUIzaEkCLIUxQFjPllapw8nU)

[Section 01 \(massey.ac.nz\)](#)