



with tensorflow

구현을 위한 딥러닝

■ 고려대학교 물리학과 한승희

■ mod96@naver.com

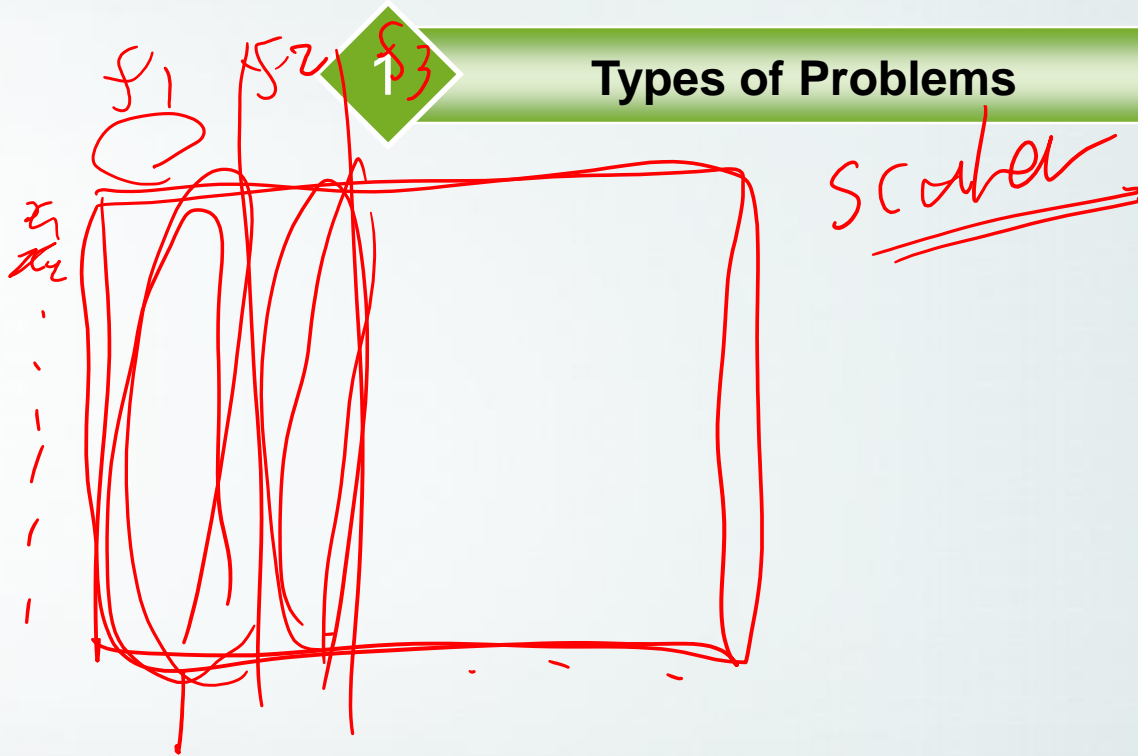
Contents

- 1 Types of Problems
- 2 Theoretical Background
- 3 Bias – Variance Tradeoff
- 4 About data
- 5 Implementation

gradient vanishing/exploding

Contents

Types of Problems



Types of Problems - 1

- ❖ Supervised Learning
- ❖ Unsupervised Learning
- ❖ Semi-Supervised Learning

Supervised Learning

- ❖ Majority of algorithms. Machine is trained using well-labeled data. (inputs and outputs are matched)
- ❖ Ex> Classification, Regression



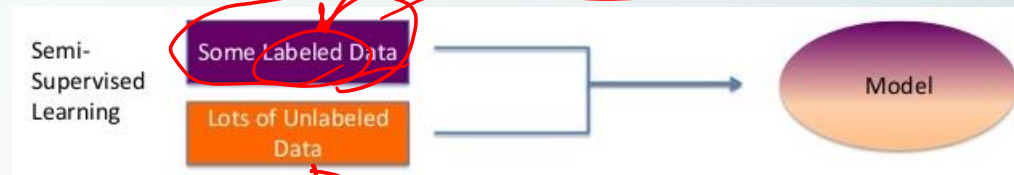
Unsupervised Learning

- ❖ Learning happens without supervision. Only inputs are used to create a model.
- ❖ Ex> Clustering



Semi - Supervised Learning

- ❖ Some data is labeled, some not. Since clean, perfectly labeled datasets aren't easy to come by, good for real world data.

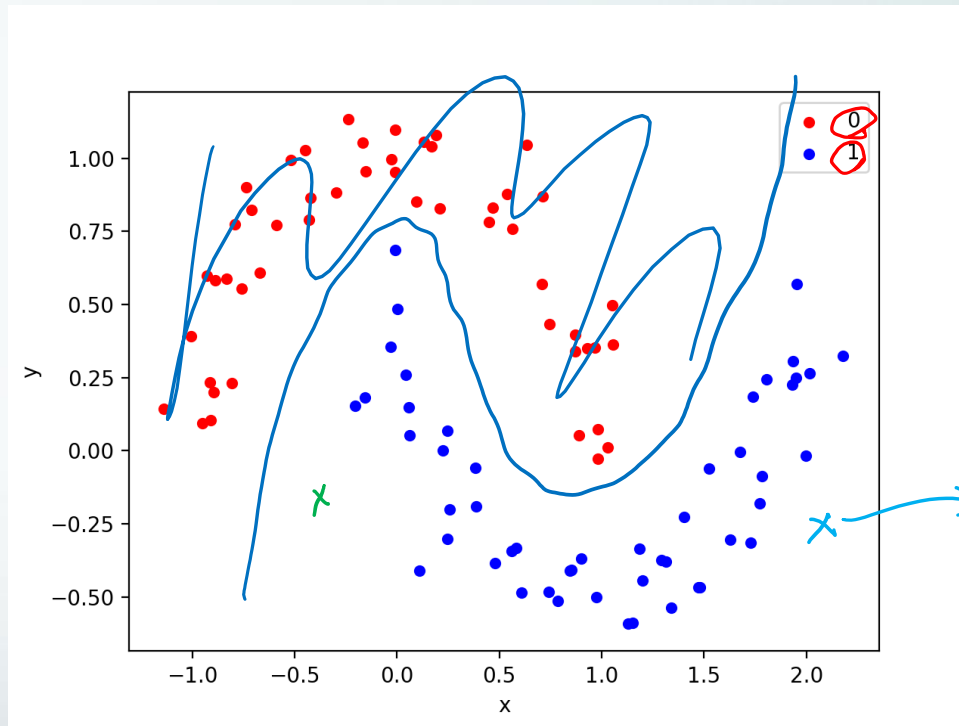


Types of Problems - 2

- ❖ Classification
 - ❖ Regression
 - ❖ Clustering
 - ❖ Etc.
- Sup.
- Un-sup.

Classification

- ❖ Predicts discrete number of values. The data is categorized under different labels



Classification

0 ~ 1 Standardization

가장

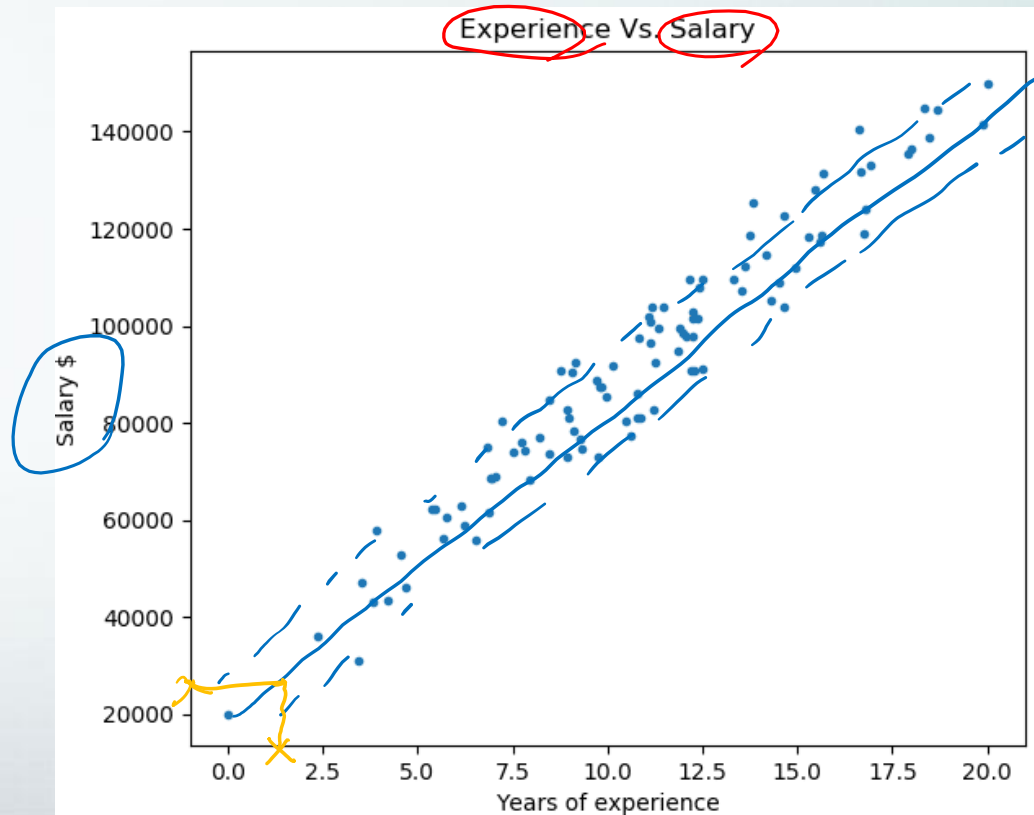
Sepal length →	Sepal width →	Petal length →	Petal width →	Species →
5.2	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.3	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.6	3.4	1.4	0.3	<i>I. setosa</i>
7.0	3.2	4.7	1.4	<i>I. versicolor</i>
6.4	3.2	4.5	1.5	<i>I. versicolor</i>
6.9	3.1	4.9	1.5	<i>I. versicolor</i>
5.5	2.3	4.0	1.3	<i>I. versicolor</i>
6.5	2.8	4.6	1.5	<i>I. versicolor</i>
5.7	2.8	4.5	1.3	<i>I. versicolor</i>
6.3	3.3	4.7	1.6	<i>I. versicolor</i>
4.9	2.4	3.3	1.0	<i>I. versicolor</i>
6.6	2.9	4.6	1.3	<i>I. versicolor</i>
6.3	3.3	6.0	2.5	<i>I. virginica</i>
5.8	2.7	5.1	1.9	<i>I. virginica</i>
7.1	3.0	5.9	2.1	<i>I. virginica</i>
6.3	2.9	5.6	1.8	<i>I. virginica</i>
6.5	3.0	5.8	2.2	<i>I. virginica</i>
7.6	3.0	6.6	2.1	<i>I. virginica</i>
4.9	2.5	4.5	1.7	<i>I. virginica</i>

W

?? ~

Regression

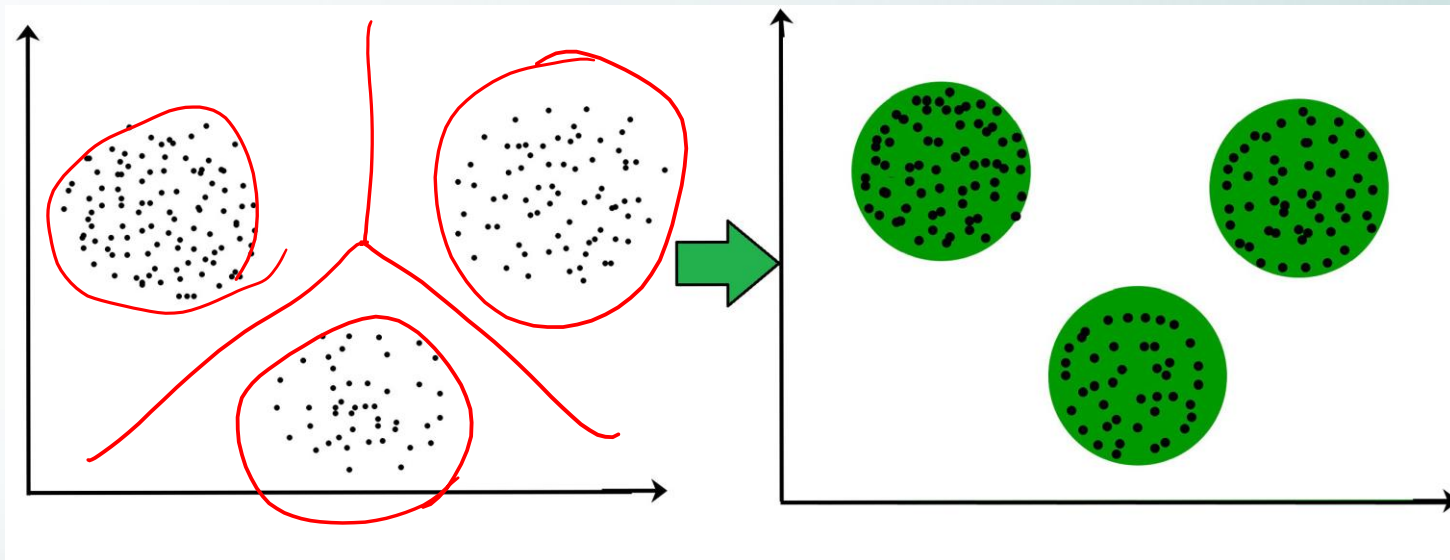
- ❖ Predicts continuous values output. Mostly analysis using statistical model which is used to predict the numeric data instead of labels.



이론
↓
실험
↓
데이터

Clustering

- ❖ Task of partitioning the dataset into groups, called clusters. This splits the data in such a way that points within single cluster are very similar and points in different are different.

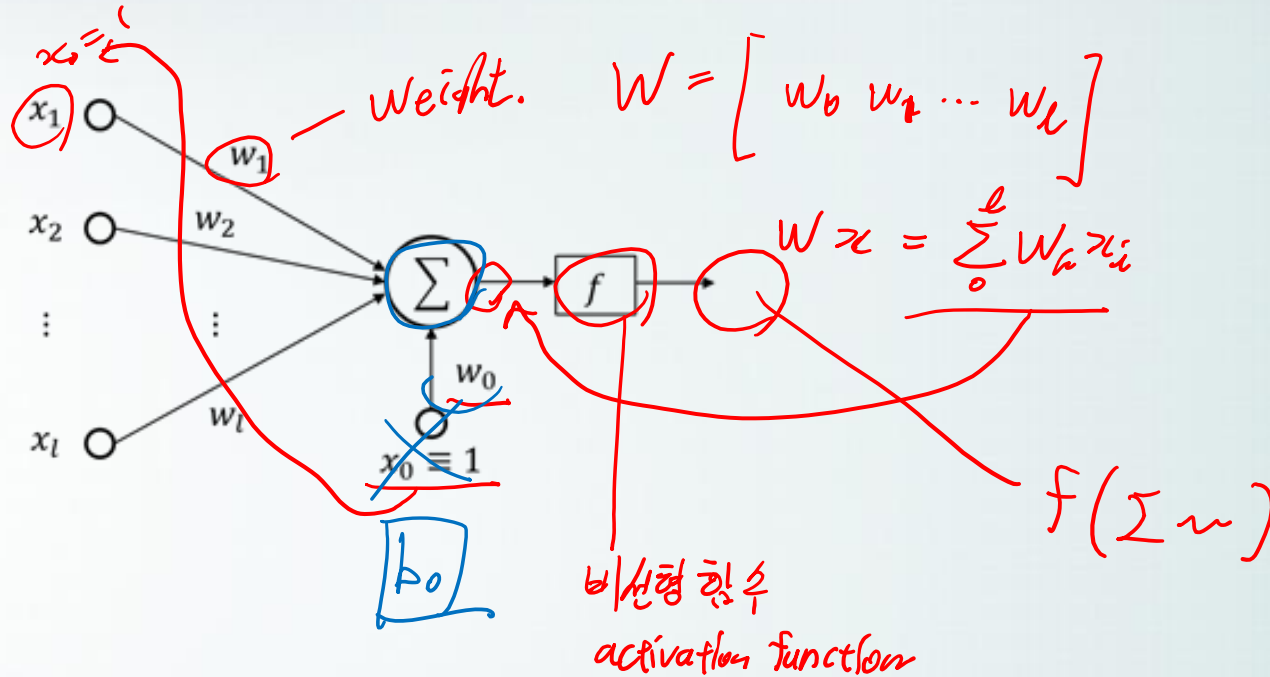
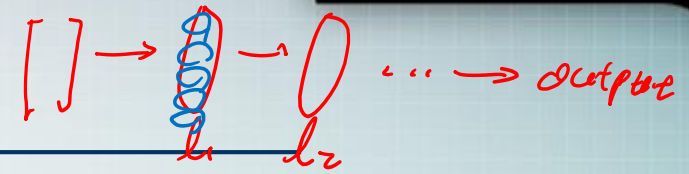


Contents

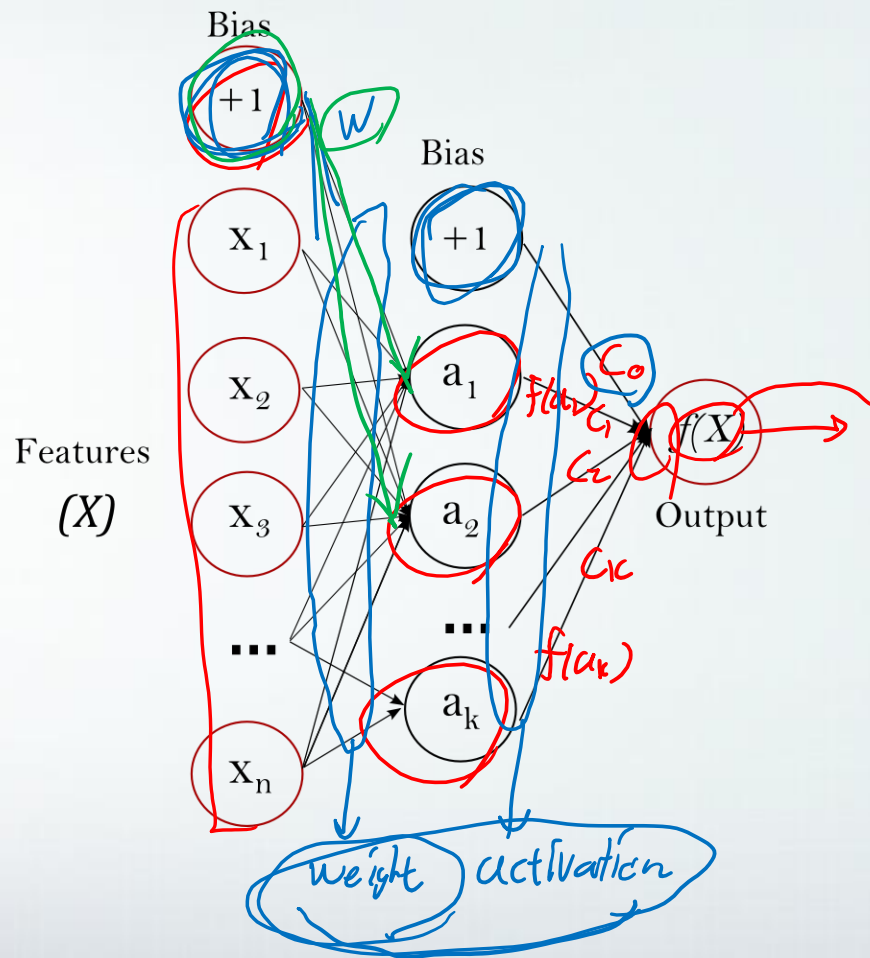
2

Theoretical Background

Perceptron (Neuron)

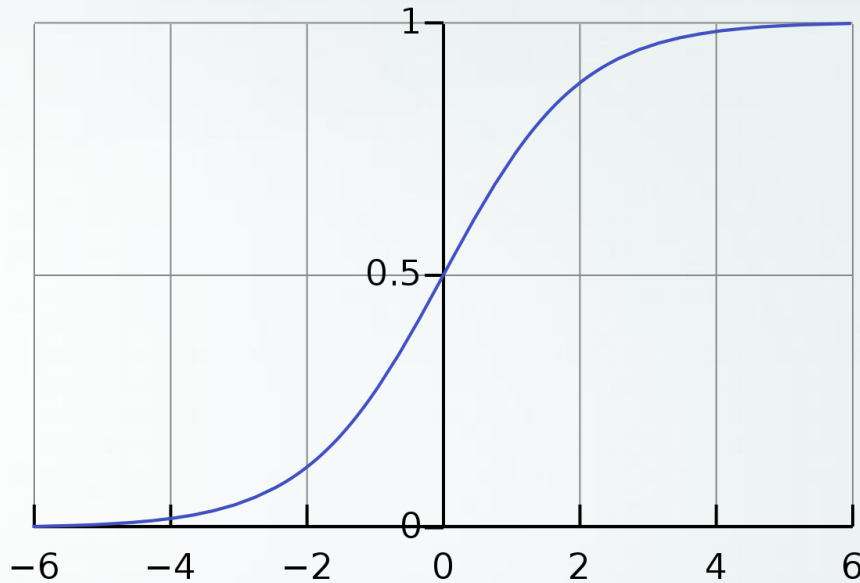


Multi Layer Perceptron



layer (tf. keras. layers. Dense)

Activation ftn (Sigmoid)



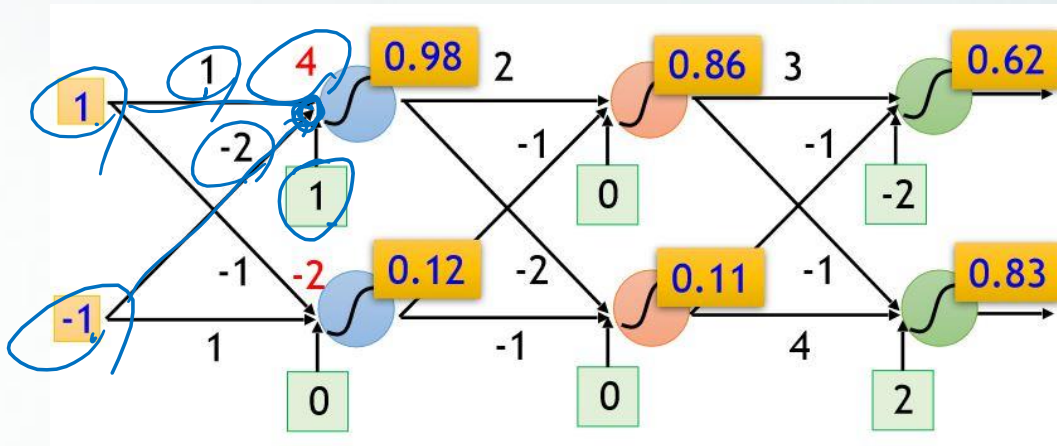
$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad : 0 \sim 1$$

$$\sigma'(x) = - \frac{-e^{-x}}{(1 + e^{-x})^2}$$

$$= \frac{1}{1 + e^{-x}} \times \frac{e^{-x}}{1 + e^{-x}}$$

$$= \sigma(x) \cdot (1 - \sigma(x)) \quad : 0 \sim 1$$

Example



But...Why?

Universal Approximation thm

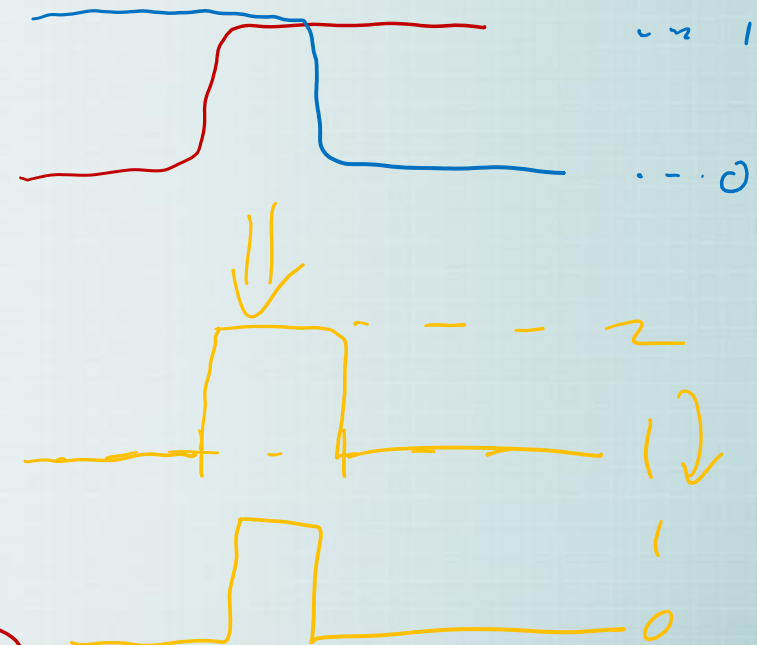
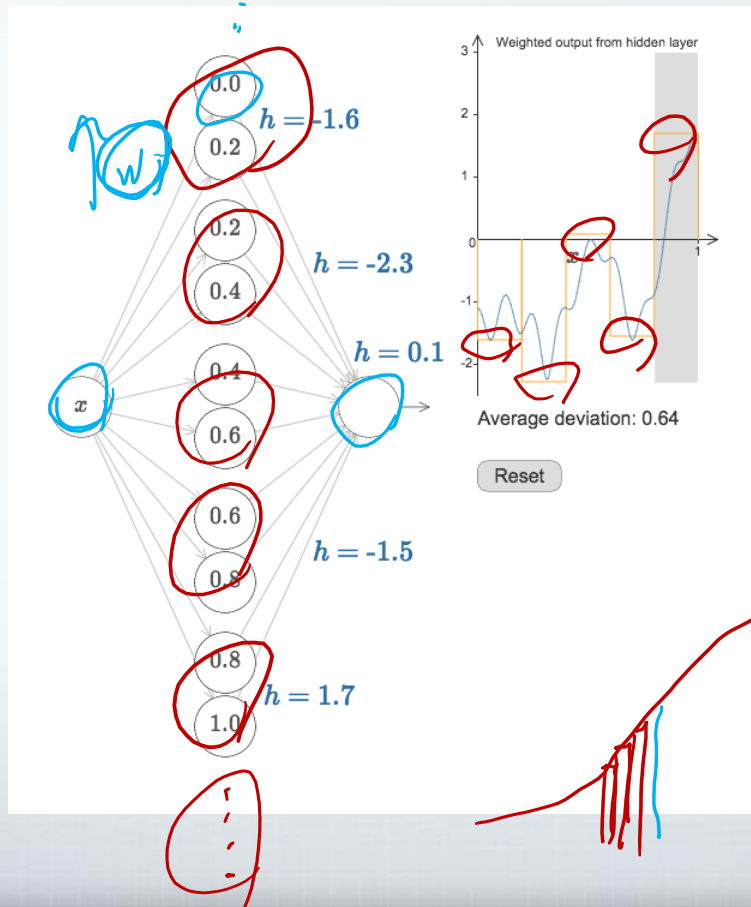
- ❖ The sum of the following form can approximate any continuous function F on $[0, 1]^n$ to any degree of accuracy:

- ❖
$$F(x) \approx \sum_k c_k \sigma(\sum_{i=1}^n w_{ki} x_i + b_k)$$

Diagram illustrating the universal approximation theorem. It shows two vertical rectangles representing hidden units. The left rectangle has an input x_i and the right one has an input x_l . Lines connect these inputs to a central point where a handwritten formula is written: $\sigma(\sum_{i=1}^n w_{ki} x_i + b)$. A blue arrow points from this formula to a blue curve representing the target function $F(x)$. A red curve below it represents the approximation, with a vertical double-headed arrow indicating the error between the two curves.

Universal Approximation thm 10x10x10

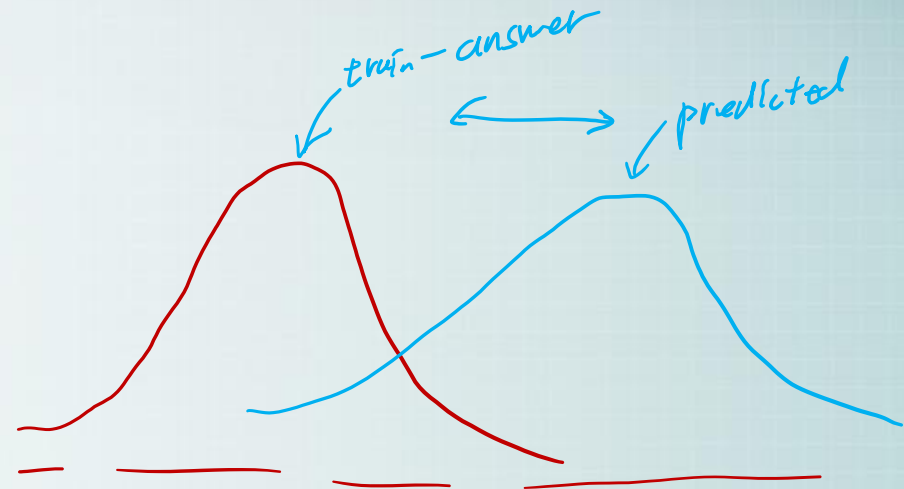
❖ See what happens with perceptron



Training – Loss ftn?

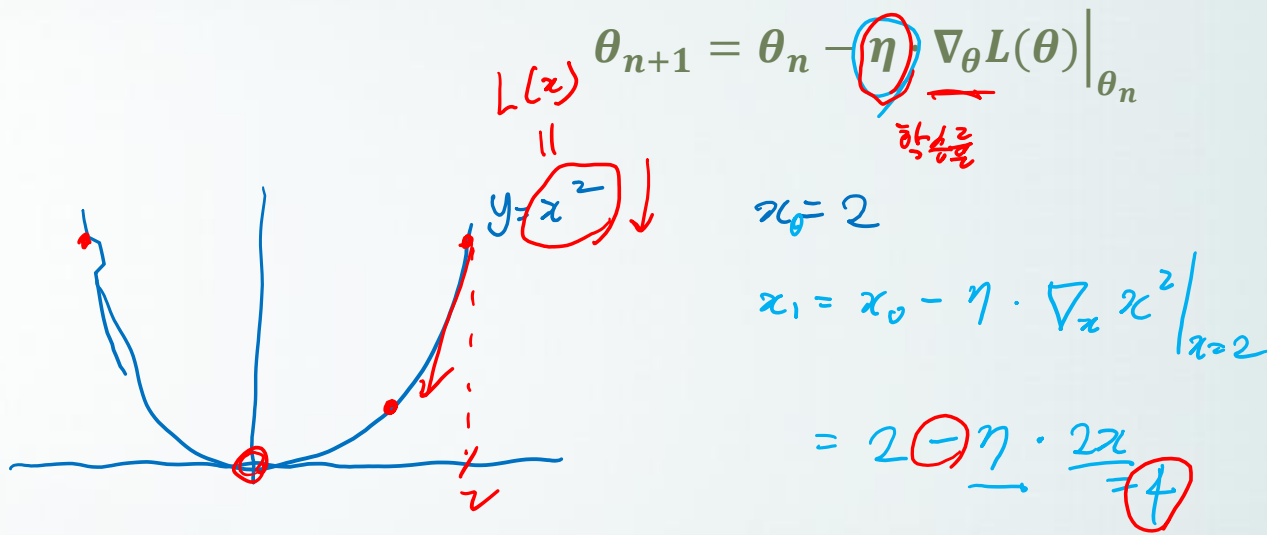
❖ Loss ftns

- Mean Squared Error
- Mean Absolute Error
- KL-Divergence
- Categorical Cross Entropy
- etc.



Training - Optimizer

❖ Gradient Descent & Backpropagation



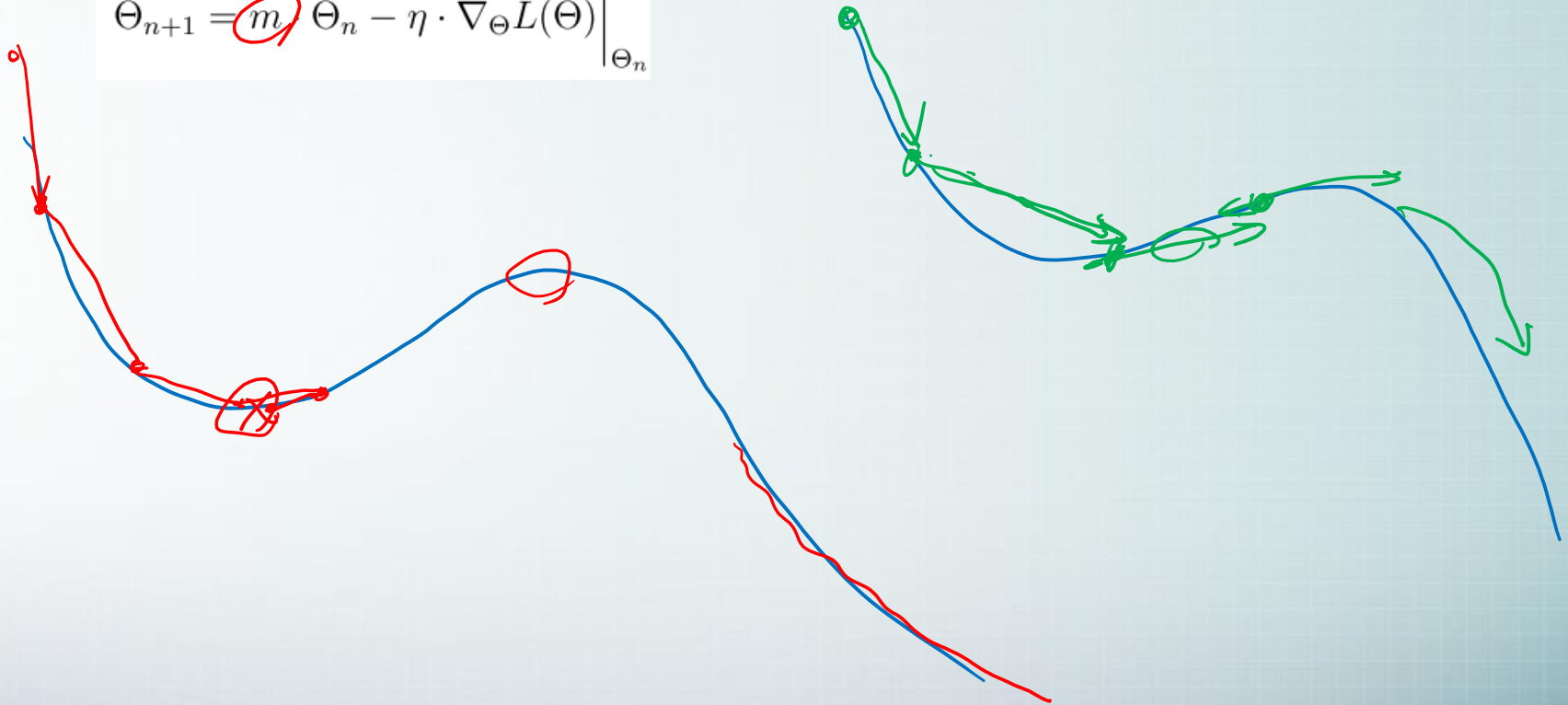
$$\eta = 0.5 \Rightarrow \underline{0}$$

$$\eta = \frac{1}{x} \Rightarrow x_1 = 1$$

Training - Optimizers

❖ Momentum

$$\Theta_{n+1} = m \Theta_n - \eta \cdot \nabla_{\Theta} L(\Theta) \Big|_{\Theta_n}$$



Training - Optimizers

❖ Adagrad (adaptive gradient)

$$\Theta_{n+1} = \Theta_n - \frac{\eta}{\sqrt{G_n + \epsilon}} \nabla_{\Theta} L(\Theta) \Big|_{\Theta_n}$$
$$G_n = \sum_{i=1}^n \left[\nabla_{\Theta} L(\Theta) \Big|_{\Theta_i} \right]^2$$

Training - Optimizers

❖ RMSProp (root mean square propagation)

$$\Theta_{n+1} = \Theta_n - \frac{\eta}{\sqrt{G_n + \epsilon}} \nabla_{\Theta} L(\Theta) \Big|_{\Theta_n}$$
$$G_n = \underbrace{\gamma \cdot G_{n-1}} + \underbrace{(1 - \gamma)} \left[\underbrace{\nabla_{\Theta} L(\Theta) \Big|_{\Theta_n}} \right]^2$$

Training - Optimizers

❖ **Adam** (adaptive moment estimation) *바탕지수, overfit ↓. / AdamW / SWE*

$$\begin{aligned}\hat{m}_{\Theta} &= \frac{m_{\Theta}^{n+1}}{1 - (\beta_1)^{n+1}} \quad \text{where} \quad m_{\Theta}^{n+1} = \beta_1 m_{\Theta}^n + (1 - \beta) \nabla_{\Theta} L(\Theta) \Big|_{\Theta_n} \\ \hat{G}_{\Theta} &= \frac{G_{\Theta}^{n+1}}{1 - (\beta_2)^{n+1}} \quad \text{where} \quad G_{\Theta}^{n+1} = \beta_2 G_{\Theta}^n + (1 - \beta_2) \left[\nabla_{\Theta} L(\Theta) \Big|_{\Theta_i} \right]^2 \\ \Theta_{n+1} &= \Theta_n - \eta \frac{\hat{m}_{\Theta}}{\sqrt{\hat{G}_{\Theta} + \epsilon}}\end{aligned}$$

❖ **About learning rate...**

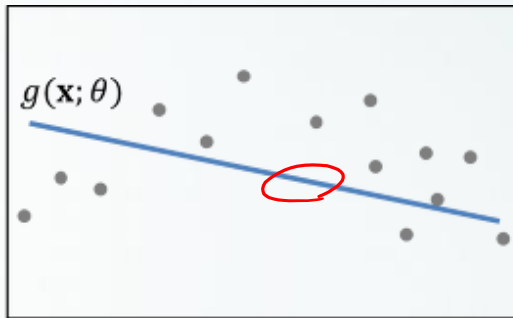
- [Google exercise](#)

Contents

3

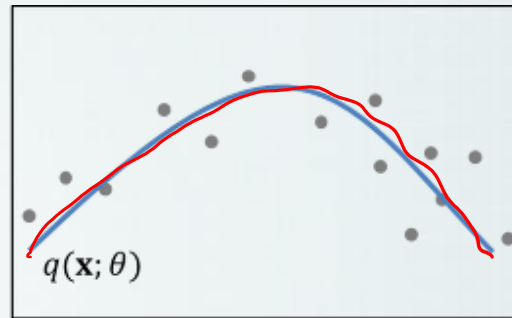
Bias – Variance Tradeoff

Bias vs Variance

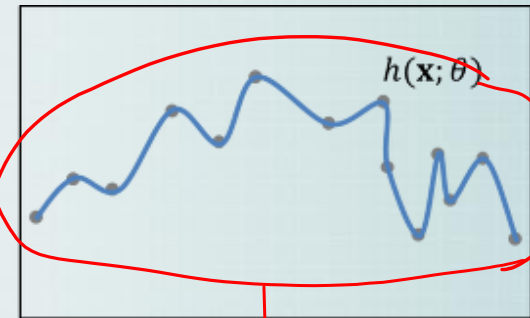


(a) Underfit

linear



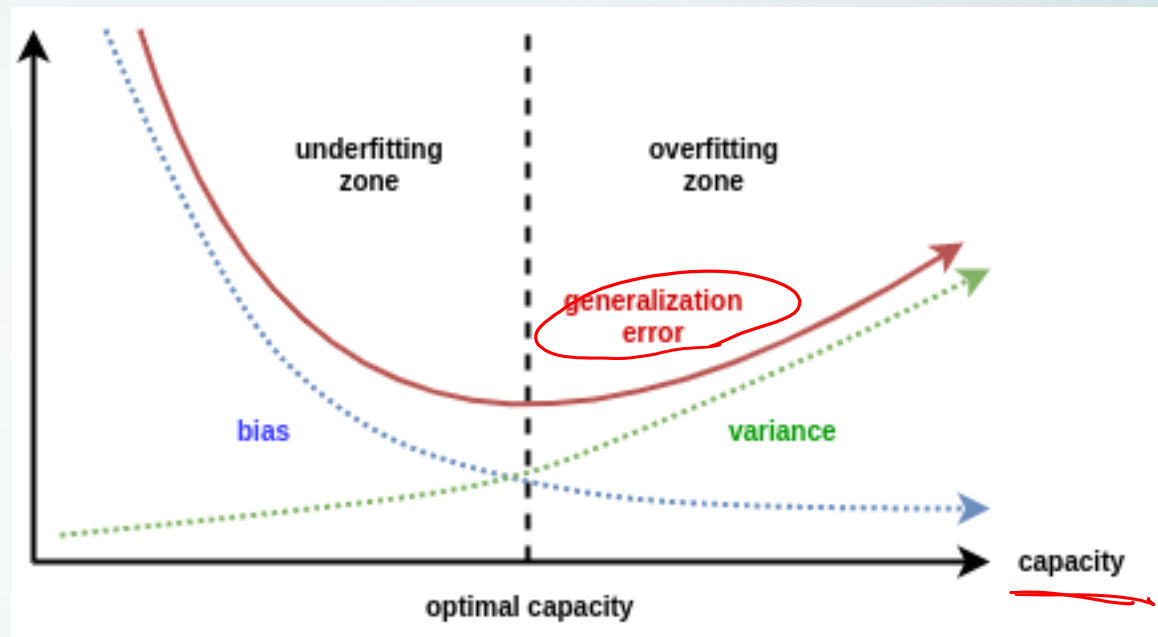
(b) Ideal fit



(c) Overfit

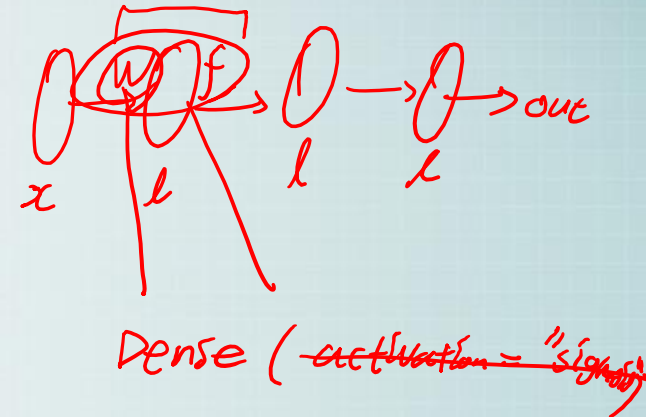
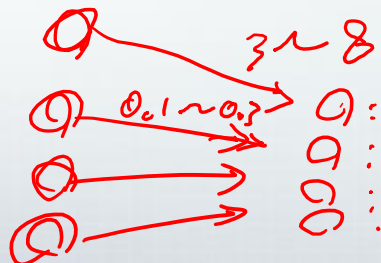
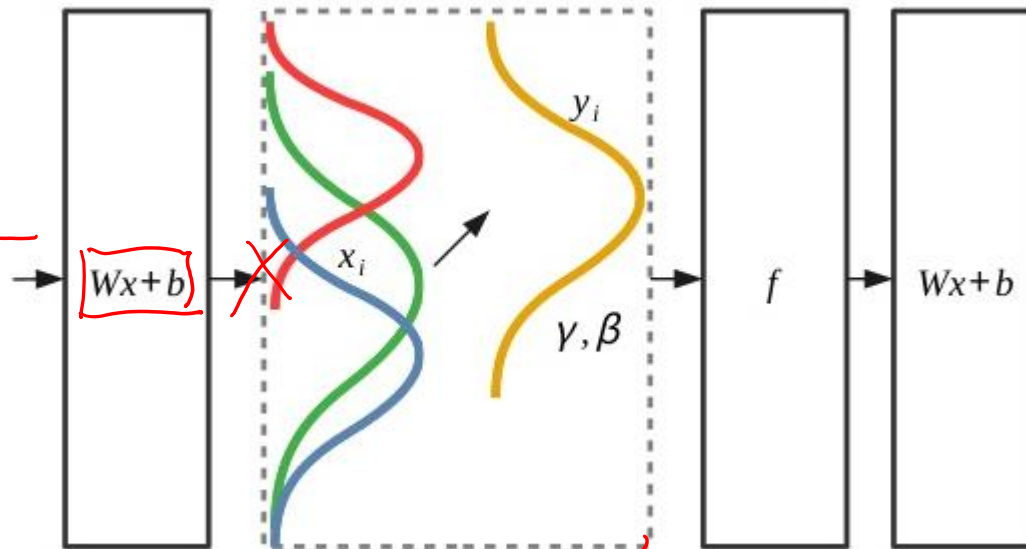
high dim = data > H_f

Bias vs Variance



Bias vs Variance : BN

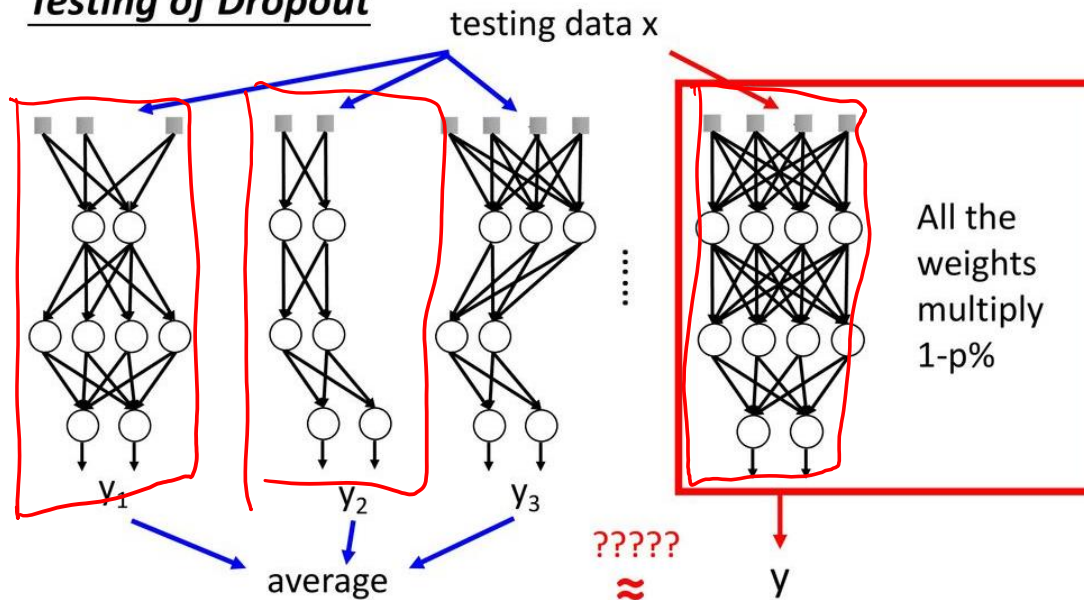
Ensure the output statistics of a layer are fixed.



Bias vs Variance : dropout

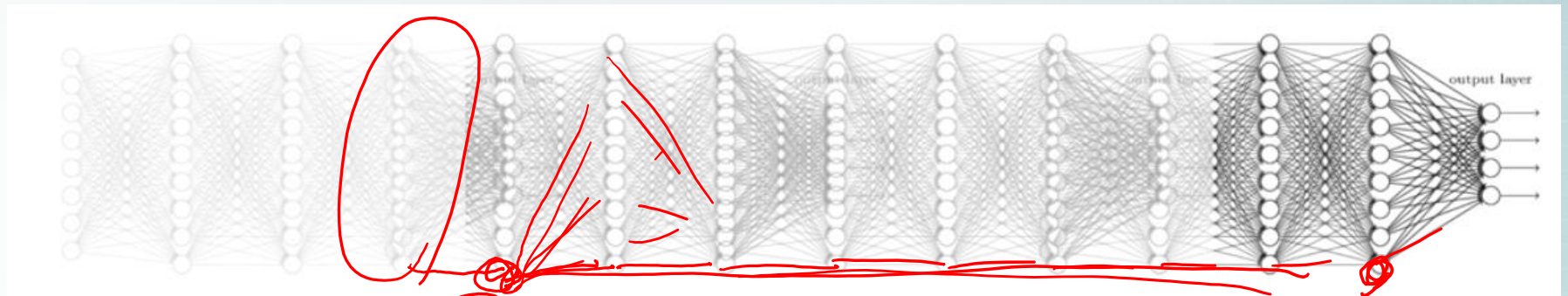
Dropout is a kind of ensemble.

Testing of Dropout



dense (acc -)
dropout (ratio = 0.2)

Bias vs Variance : vanishing gradient



$$\boxed{\nabla_w L(w)}$$

(f)

$$a: 0 \sim 1$$

$$\underline{a': 0 \sim 1} \downarrow$$

$$\nabla_w L(w) \big|_{w=w_h}$$

$$w_{n+1} = w_h + \nabla_w L(w)$$

Vanishing Gradient : relu

|||

Activation function	Equation	Example	1D Graph
Unit step (Heaviside)	$\phi(z) = \begin{cases} 0, & z < 0, \\ 0.5, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant	
Sign (Signum)	$\phi(z) = \begin{cases} -1, & z < 0, \\ 0, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant	
Linear	$\phi(z) = z$	Adaline, linear regression	
Piece-wise linear	$\phi(z) = \begin{cases} 1, & z \geq \frac{1}{2}, \\ z + \frac{1}{2}, & -\frac{1}{2} < z < \frac{1}{2}, \\ 0, & z \leq -\frac{1}{2}, \end{cases}$	Support vector machine	
Logistic (sigmoid)	$\phi(z) = \frac{1}{1 + e^{-z}}$	Logistic regression, Multi-layer NN	
Hyperbolic tangent	$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	Multi-layer Neural Networks	
Rectifier, ReLU (Rectified Linear Unit)	$\phi(z) = \max(0, z)$	Multi-layer Neural Networks	
Rectifier, softplus	$\phi(z) = \ln(1 + e^z)$	Multi-layer Neural Networks	

selu
 elu
 leaky-relu

$y = x$
 $x > 0$
 $\Rightarrow 1$

Contents

4

About data

Normalization vs Standardization

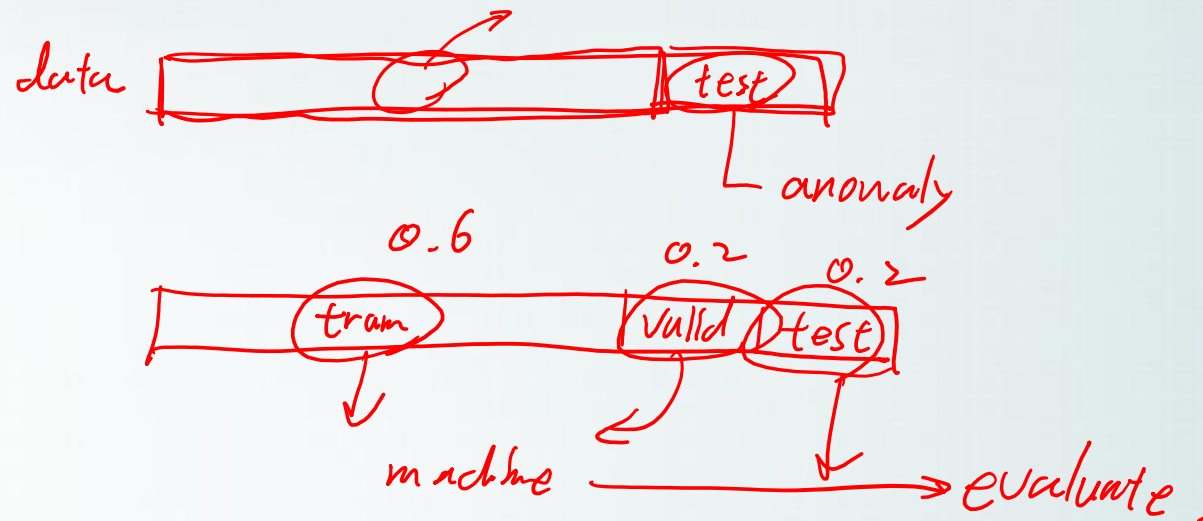
❖ **Normalization :**

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

❖ **Standardization :**

$$X' = \frac{X - \mu}{\sigma}$$

Train-Valid-Test split



Contents

5

Implementation

Pima-Indian Diabetes dataset

- ❖ <https://www.kaggle.com/kumargh/pimaindiandabetescsv>

Credit Card Fraud Detection

❖ <https://www.kaggle.com/mlg-ulb/creditcardfraud/>