



Lecture 2. Linear Classifier

HI



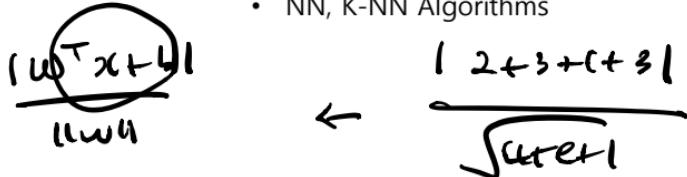
Review

1. What is Machine Learning?

- Definition
- Fields of ML
- Narrow down to Image Classification

2. Making a Model I

- Narrow down to Iris Classification
- Data-Driven Approach
 - NN, K-NN Algorithms



Handwritten note: $(w^T x + b) / \|w\|$

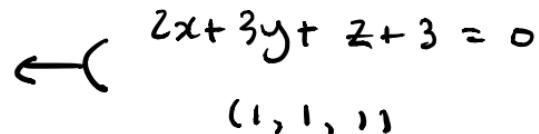
Diagram showing a fraction where the numerator is $w^T x + b$ and the denominator is $\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$. An arrow points from the left towards the fraction.

3. How to Test a Model

- Hyperparameter
- Cross Validation

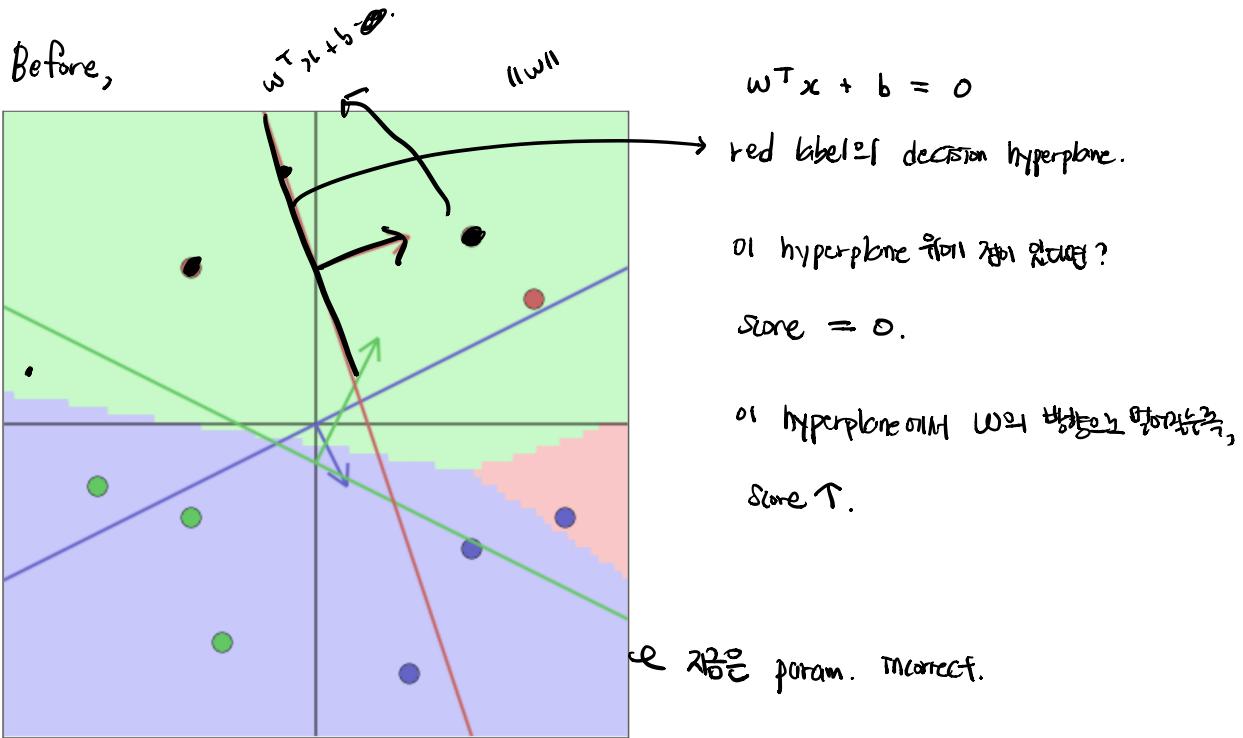
4. Making a Model II

- Limitations of Data-Driven Approach
- Parametric Approach
- Linear Classifier : Algebraic & Geometric

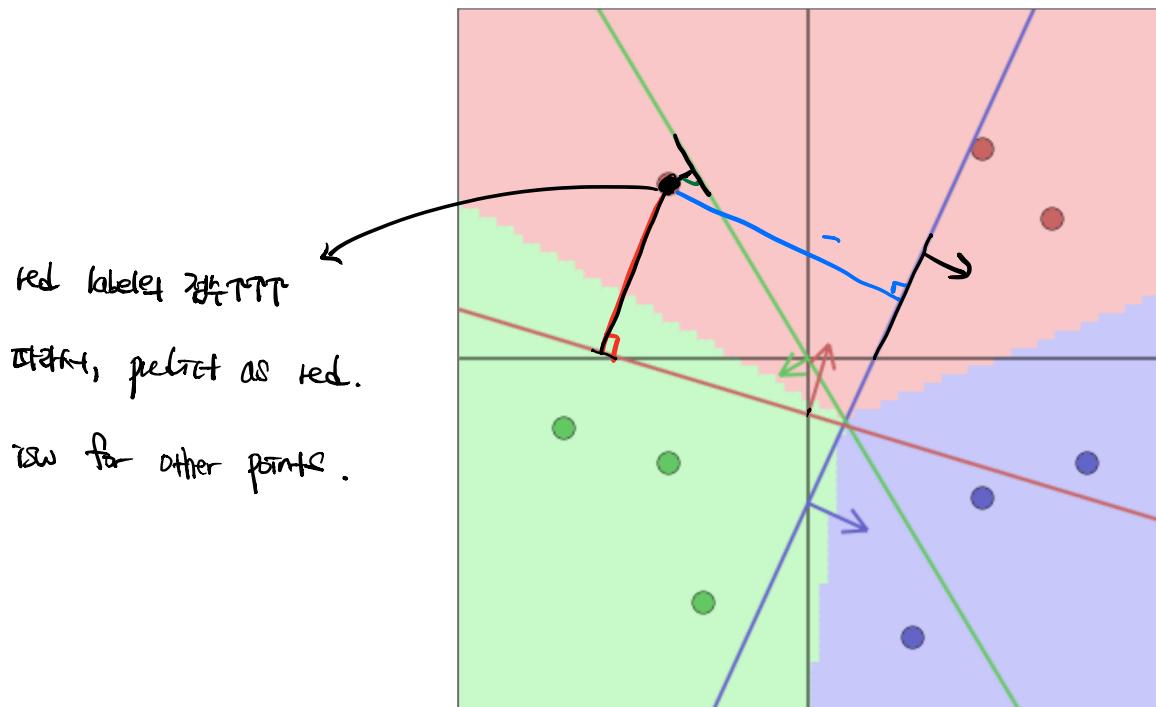


Handwritten note:
 $2x + 3y + 2 + 3 = 0$
 $(1, 1, 1)$

Diagram showing a handwritten linear equation $2x + 3y + 2 + 3 = 0$ and a point $(1, 1, 1)$.



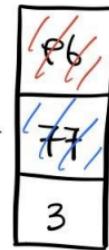
After,



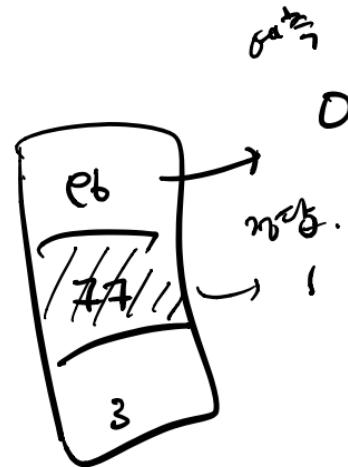
Review

$x = [x_1, x_2, \dots, x_n]$ $w = [w_1, w_2, \dots, w_n]$

+ feature \rightarrow predict ans: 1 -----
but, actually ans: 2 -----



\neq



Review

Questions

- Train이 잘 되었는지 판단할 수치적 척도 필요

: define a **Loss Function** that quantifies our unhappiness with the scores across the training data

- Parameter를 update하는 algorithm 필요

: come up with a way of efficiently finding the parameters that minimize the **Loss Function**

Today's Contents

- 1. Loss**
- 2. Loss Function**
- 3. Regularization**

Loss

Loss

130.

→ 130 개 예측 결과.



130개 정답.

How is a Model Optimized / Updated?

1. Training Set의 Data들을 Linear Classifier에 통과시켜서, 그 결과들을 정답과 비교
2. 1에서의 결과를 바탕으로 Parameter의 값을 Update
3. 다시 1로

Loss

170

5.

2. 1에서의 결과를 바탕으로 Parameter의 값을 Update

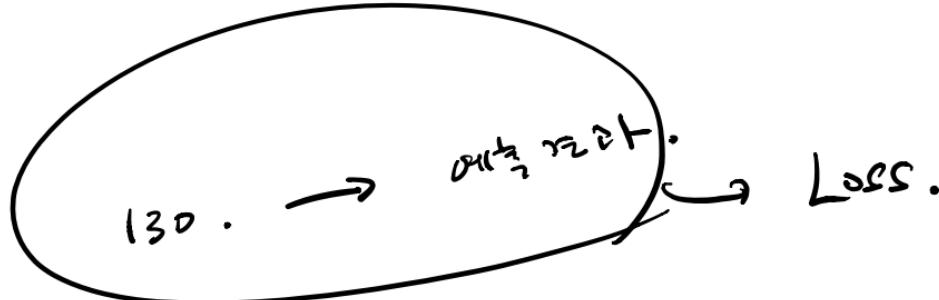
정답과의 차이가 '많다'면?

정답과의 차이가 "작아" 서 거의 비슷하다면?

x feature → predict ans : 1 -----
but, actually ans = 2 -----

111
777
3

Loss



Loss Over Dataset : 현재 Model의 Parameter가 주는 결과가 얼마나 부정확한지를 숫자로 표현

We want to, (MINIMIZE) / MAXIMIZE) the LOSS

Loss

130.

How do we calculate the Loss Over (Training) Dataset?

Training dataset $\{(x_i, y_i)\}_{i=1}^N$ on data,

feature
label

Loss over the dataset은,

$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, w, b), y_i)$$

loss function
ith data's prediction 결과
ith dataset 정답

$\{(x_i, y_i)\}_{i=1}^N$

↓
feature

label

(...,)
(0 , 1 , 2)

$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, w, b), y_i)$$

L_i 예측값.
 x_i w b (내려온다)

Loss Function

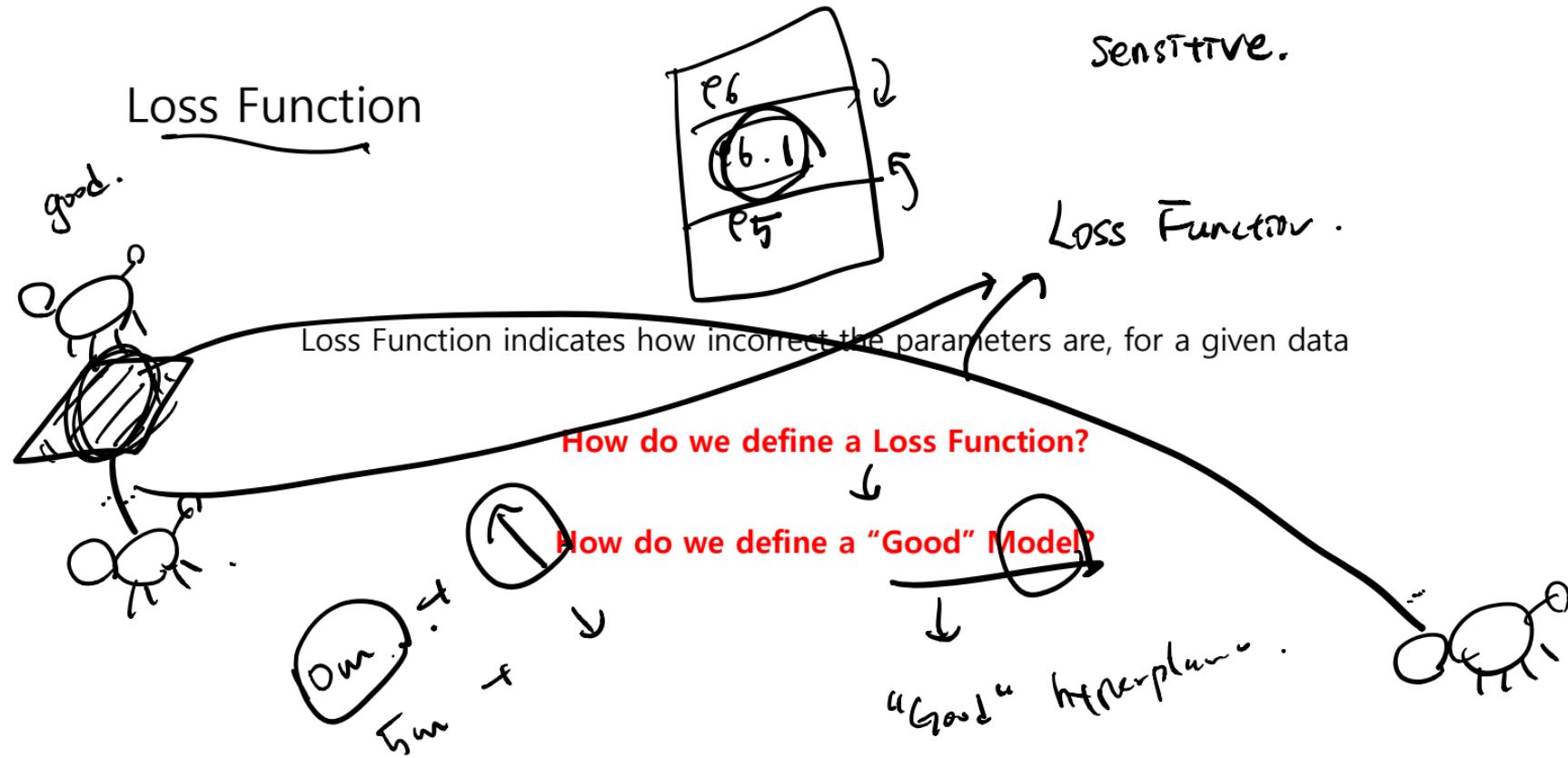
Training dataset $\{(x_i, y_i)\}_{i=1}^N$ or others,
feature \xrightarrow{f} label \xrightarrow{g}

Loss over the dataset \sum_i

$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, w, b), y_i)$$

loss function \xrightarrow{h} i th dataset prediction 결과 \xrightarrow{j} i th dataset 정답 \xrightarrow{k}

We calculate the Loss via,
Loss Function



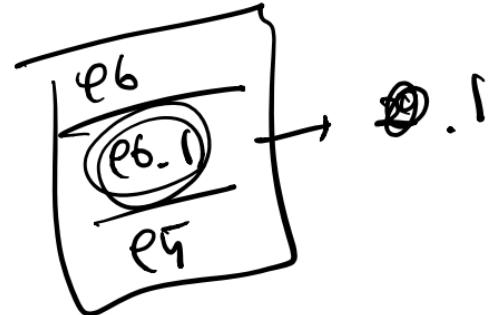
Loss Function

Two Popular Loss Functions,

1. Multiclass SVM Loss
2. Cross-Entropy Loss

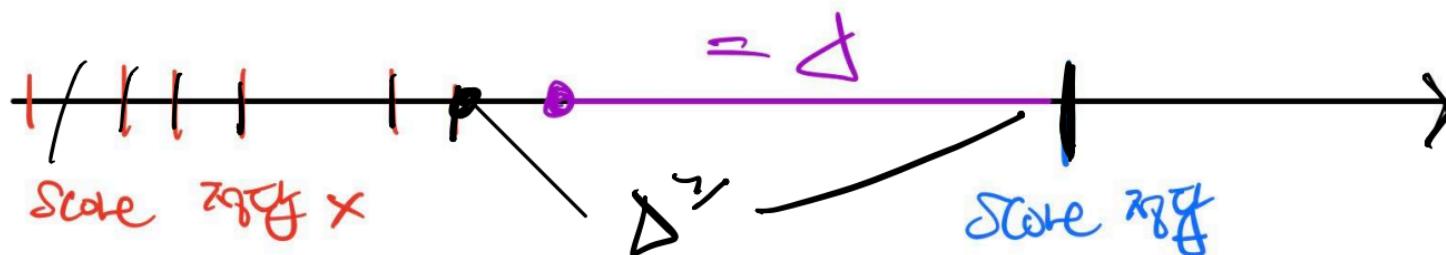
Loss Function : Multiclass SVM Loss

$$\Delta = 10.$$



"The SVM Loss is set up so that
the SVM wants the correct class for each input
to have a score higher than the incorrect classes by some fixed margin Δ ."

정답 label의 score가 나머지 incorrect label의 score보다 Δ 만큼 크기를 바람.



Loss Function : Multiclass SVM Loss

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + \Delta)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

250g Labeled 25kg

ex. $L_{\bar{i}} \leftarrow$ 예측 결과 $s_{\bar{i}}$
정답 $y_{\bar{i}}$.

$$L = (2 \cdot e_0 + 12 \cdot e) / 3$$

Loss Function : Multiclass SVM Loss



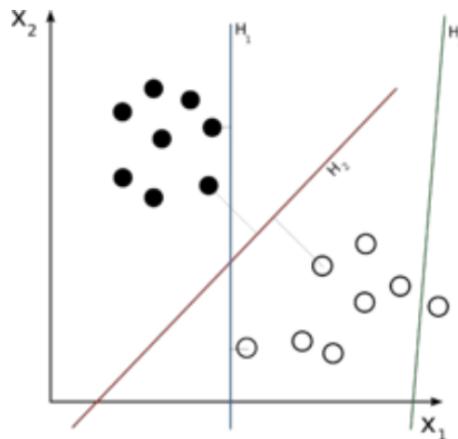
Loss

$$\max(0, 1.3 - 4 \cdot e + 1) + \max(0, 2.0 - 4 \cdot e + 1) = 0$$

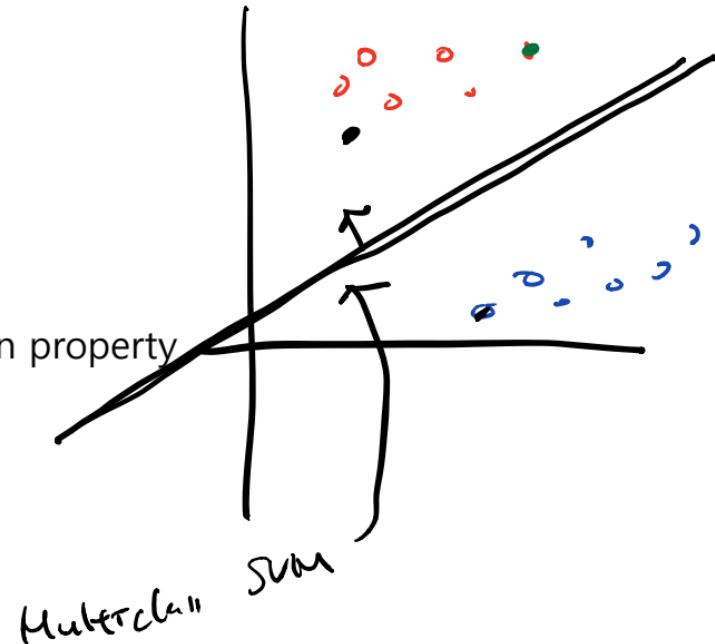
More on SVM (Optional)

Multiclass SVM Loss의 objective는,

Try to find the decision hyperplane with max-margin property



** For more detail,
https://ko.wikipedia.org/wiki/%EC%84%9C%ED%8F%AC%ED%8A%B8_%EB%B2%A1%ED%84%B0_%EB%A8%B8%EC%8B%A0
(cs229.stanford.edu/notes/cs229-notes3.pdf)

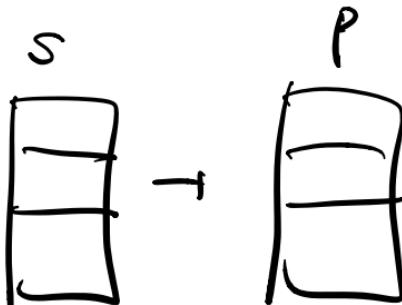


②

Loss Function : Cross-Entropy Loss

An approach based on probability

Softmax function maps the scores to probabilities



softmax function :

$$\frac{e^{f_j}}{\sum_i e^{f_i}}$$

정답 score
e^{score} 의 sum

$e \in [0, 1]$

$$\frac{e^1}{\text{sum}} + \frac{e^2}{\text{sum}} + \frac{e^3}{\text{sum}} = 1$$

...
...
...

1
2
3

$$\frac{e^3}{e^1+e^2+e^3} = 1$$

Loss Function : Cross-Entropy Loss

softmax function :

$$\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \text{ 정답 score} \approx e^{\text{score}} \text{ 의 sum}$$

↓

interpret as, $P(y_i | x_i; w, b) = \frac{e^{f_{y_i}}}{\sum_j e^{f_j}}$

(parameter는 w, b . x_i 가 input일 때, y_i 가 나올 확률)

Input.

$$P(y_i | x_i; w, b) = \dots$$

y_i w, b .

Loss Function : Cross-Entropy Loss

$$p_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

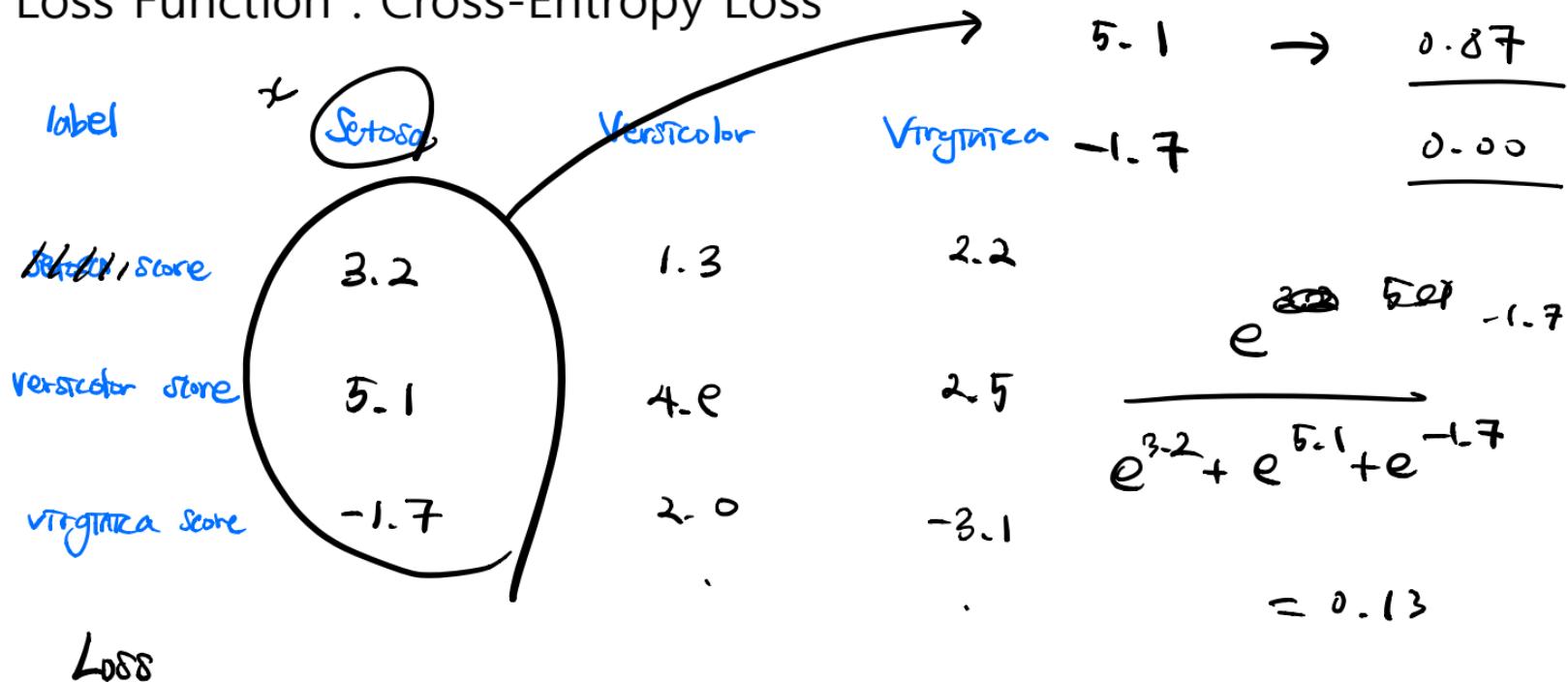
제일 높은 차원에 대한 확률을 계산하는 식입니다.

$$\vec{x}_n \rightarrow \hat{y} = \text{argmax } p_i$$

$$L_i = -\log(p_i)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

Loss Function : Cross-Entropy Loss



Loss Function : Cross-Entropy Loss

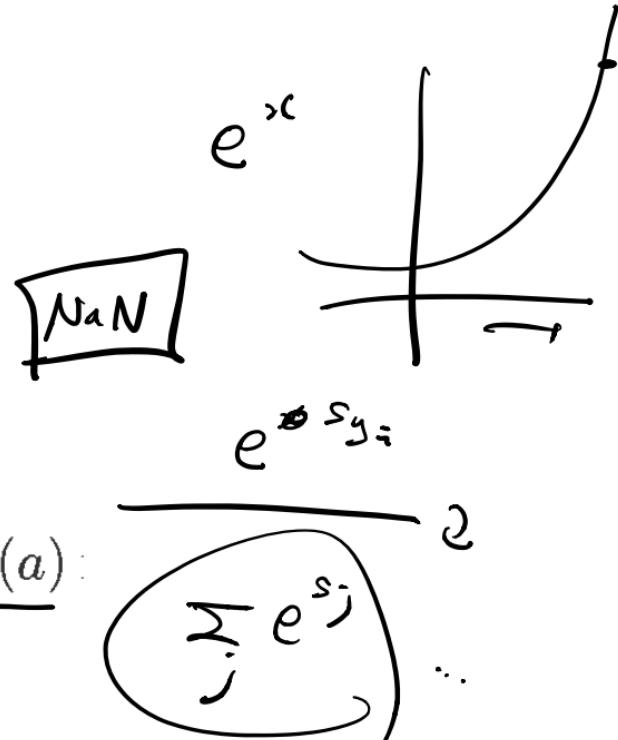
In reality,

$$\begin{aligned}
 p_i &= \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}} \times C \\
 &= \frac{Ce^{a_i}}{C \sum_{k=1}^N e^{a_k}} \\
 &= \frac{e^{a_i + \log(C)}}{\sum_{k=1}^N e^{a_k + \log(C)}}
 \end{aligned}$$

"Stable" Softmax.

$$\underline{\log(C) = -\max(a)}$$

$$-\max(s)$$



Loss Function : Cross-Entropy Loss

$$\sum \text{CE}$$

이때, Unknown Parameter W와 b를 어떻게 estimate?

via **MLE(Maximum Likelihood Estimation)**

More on MLE (Optional)

$$x_1 \rightarrow y_1$$

$$x_2 \rightarrow y_2$$

$$S \rightarrow P$$

conditional probability를 계산할 때,

(unknown parameter : W, b) 를.

We want to find W, b that maximizes the (Log) Likelihood Function,

so define the loss as,

$$L_i = -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

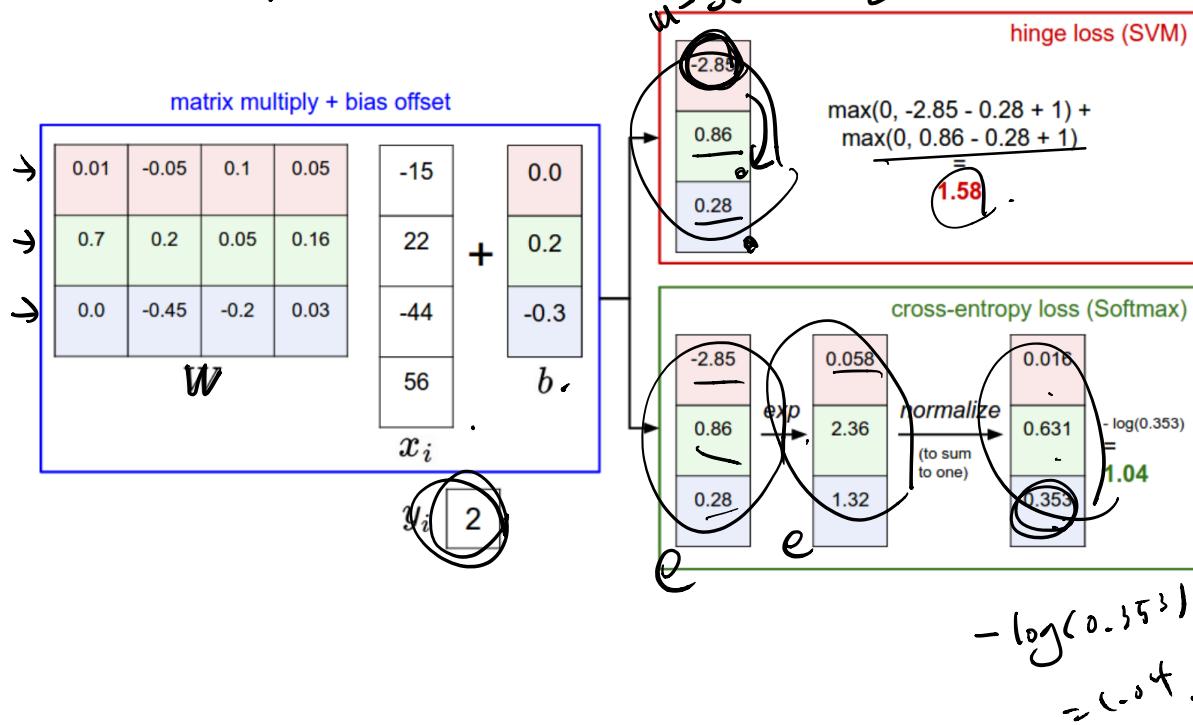
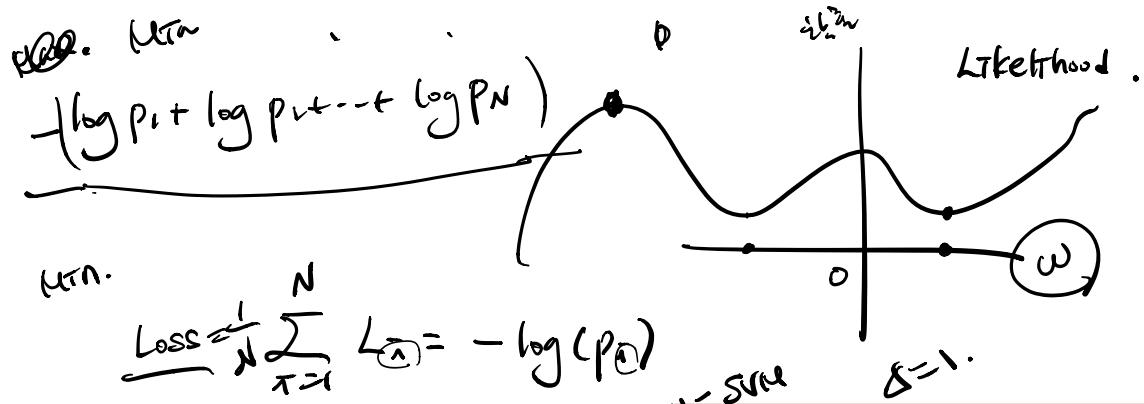
$$S(x_i, y_i) ? \underset{i=1}{\overset{n}{\dots}}$$

$$p_1 \times p_2 \times \dots \times p_n$$

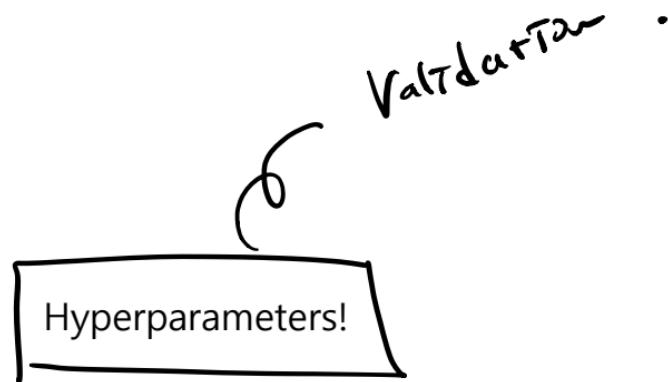
ω' , b'

ω'' , b''

** For more detail, https://en.wikipedia.org/wiki/Multinomial_logistic_regression/
<http://jaejunyoo.blogspot.com/2018/02/minimizing-negative-log-likelihood-in-kor-3.html>
<https://ratsgo.github.io/deep%20learning/2017/09/24/loss>



Actually, these are....



Regularization

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

But is this enough...?

loss L 을 최소화하는 weight W 를 찾는 것이 목적.

이때, $L = 0$ 이 되게 하는 W' 가 있다면,

임의의 실수 $k > 1$ 에 대해, kW' 도 $L = 0$ 을 만족.

따라서, W is “NOT” uniquely determined!

Mult - SVM.

$$L = 0.$$

$$\omega$$

$$2\omega'$$

$$3\omega' \quad 100\omega'$$



Regularization

$k\omega'$

among all $k\omega'$, small ones are preferred (reasons discussed later)

so, add regularization term to discourage large weight

$$L = \underbrace{\frac{1}{N} \sum_{n=1}^N L_n(f(x_n, w, b), y_n)}_{\text{Loss Function}} + \lambda R(w)$$

Regularization Term

λ : regularization strength

$$L_1 = \sum_k \sum_e |w_{ke}|$$

$$R(w) = \sum_k \sum_e |w_{ke}| \quad \text{-- } L_1$$

$$L_2 = \sum_k \sum_e w_{ke}^2$$

$$= \sum_k \sum_e w_{ke}^2 \quad \text{-- } L_2$$

Regularization Term.

$$\textcircled{1} \quad \underline{R(w)}$$

L_1

w

L_2

\overline{w}

Regularization

$$\chi = (1, 1, 1, 1)$$

①

②

더 작은 w 가 갖는 이점?

$$\text{ex) } \chi = [1.1.1.1]^T$$

$$w_1 = (1.0.0.0)$$

$$1+b \quad \frac{1}{\sqrt{5}}$$

$$w_2 = (.25, \dots)$$

$$1+b$$

$$w_1 = [1.0.0.0] . \quad w_2 = [\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}] \quad \text{같은 } \frac{1}{4} \times 4$$

\checkmark^2

$$\text{면, } w_1\chi = w_2\chi = 1$$

$$\text{but, } R(w_1) = 1 > R(w_2) = \frac{1}{4}$$

$$R(w_1) = 1$$

$$R(w_2) = \frac{1}{4}$$

Regularization

$$L(Q, \dots)$$

$$\omega_1 = \underline{(1.0.0.0)} \dots \omega_r = \underline{c \cdot 2^k}, \dots 1$$

Regularization을 통해,

No input dimension can have a very large influence on the scores all by itself

Regularization

train data의
w

If an input dimension has a very large influence on the scores all by itself,

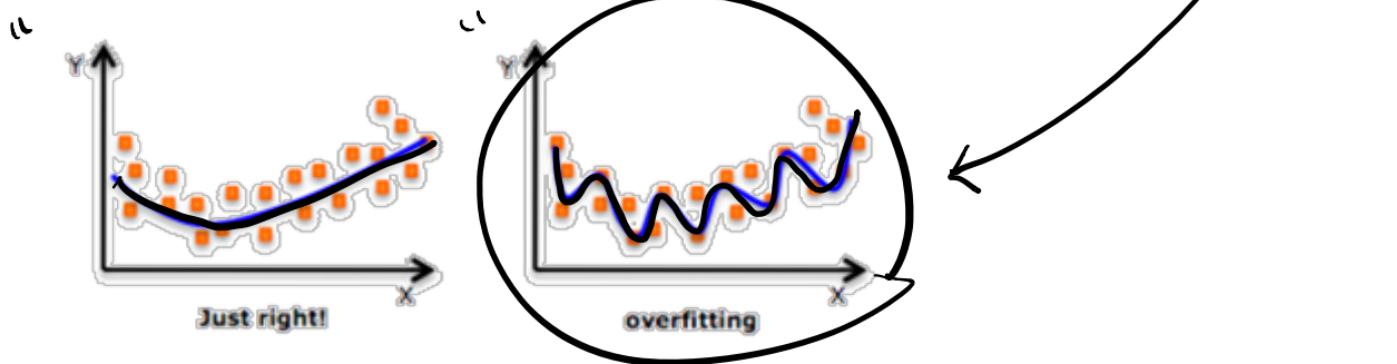
then, model이 training data의 noise까지 학습해버릴 수 있다!

Model performing TOO WELL on training data

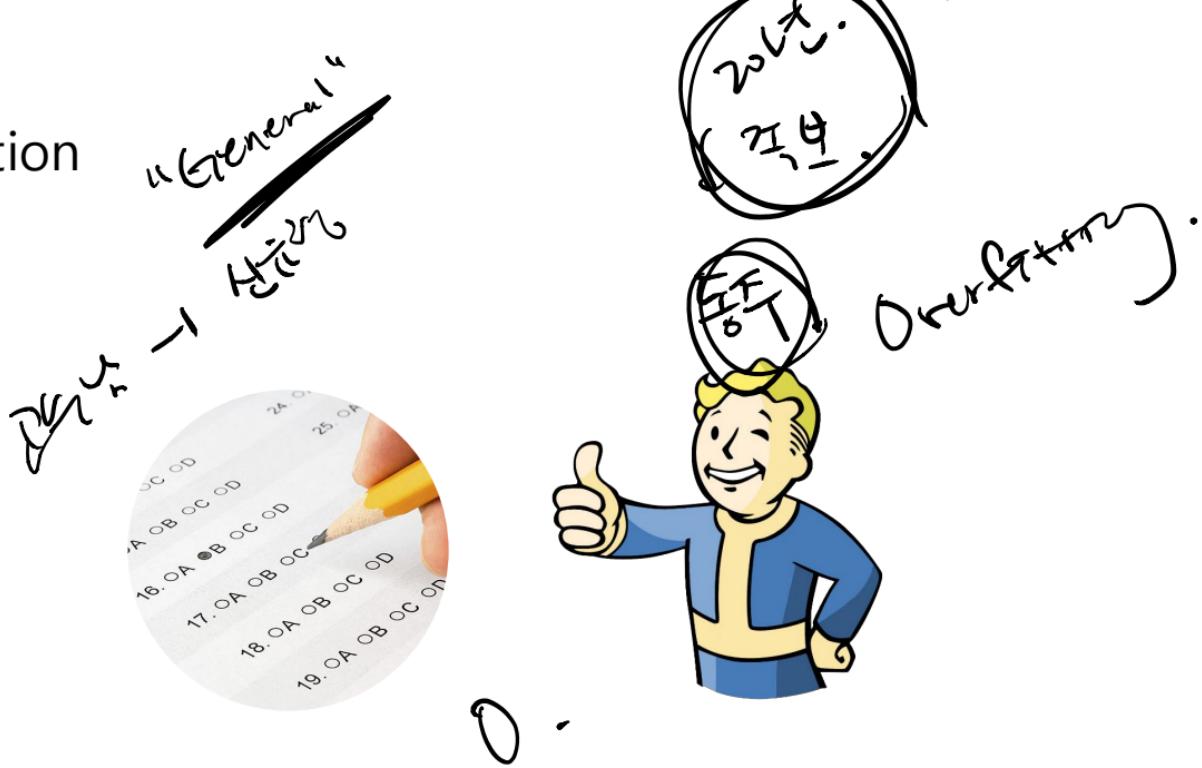
Regularization

We are building a **GENERAL** model for classification

Thus, doing **TOO WELL** on the training data is not desired



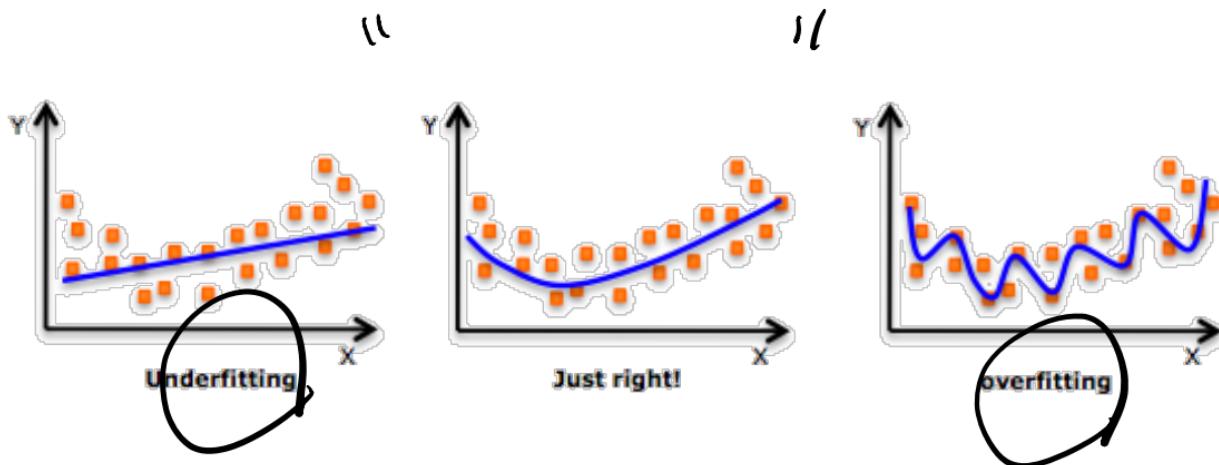
Regularization



Regularization : Overfitting



Regularization



Regularization

$$L^2 \quad R(W) = \underbrace{\sum_k \sum_l W_{kl}^2}_{\text{L2}} \times \lambda \rightarrow L.$$

L2 Regularizer favors W that is spread out

$$\begin{aligned} w_1 &= (1.0, -1.0) \\ w_2 &= (1.1, 1.1, 1.1, 1.1) \end{aligned}$$

Regularization : Overfitting

by $R(w)$, we discourage large w = prevent some input dim. from having too much influence on the output

$\backslash w$

= prevent the model from learning the noises

= prevent overfitting

(to some extent)

Review

w, b .

1. Loss

$y \neq w^T x + b$

2. Loss Function

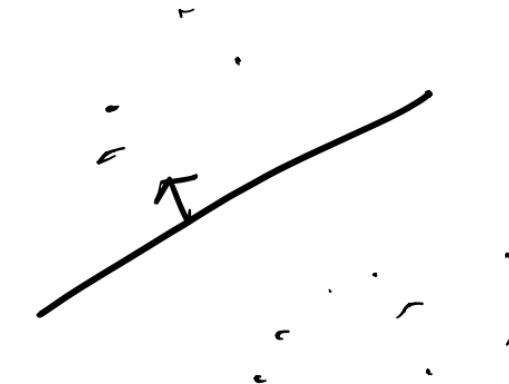
out of many not.

- ① Multiclass SVM Loss
- ② Softmax

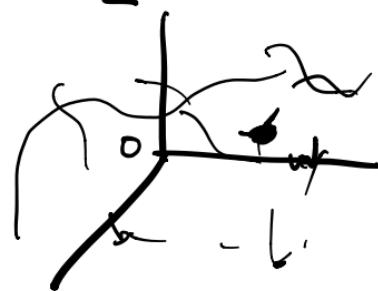
MLE

3. Regularization

- Why it is needed : to discourage large weight matrix
- Overfitting



$$L \leq 0$$



Preview on Next Class

Questions

- Train이 잘 되었는지 판단할 수치적 척도 필요

: define a **Loss Function** that quantifies our unhappiness with the scores across the training data

- Parameter를 update하는 algorithm 필요

: come up with a way of efficiently finding the parameters that minimize the **Loss Function**

“Loss”

Loss Function scalar.

$w, b \rightarrow \text{scalar}$

Optimization and Backpropagation

Loss Function은 Scalar Function

Thus, Loss Function의 input W, b에 대해

Gradient를 구해 주어 W, b에서 빼면,

Loss가 줄어드는 방향으로 학습하지 않을까?