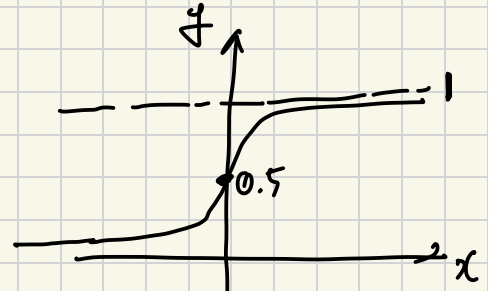


softmax → model 예측이 해석 가능성 제공 (0~1 사이로 mapping)

→ 이진분류에서 class의 점수 = s_1 인 경우
" 2의 점수 = s_2

answer = class1인 경우 $s_1 - s_2$ 가 ↑ (1에 가까워)
" class2인 경우 $s_1 - s_2$ 가 ↓ (0에 가까워)



$$\sigma(s_1 - s_2) = \frac{1}{1 + e^{-(s_1 - s_2)}}$$

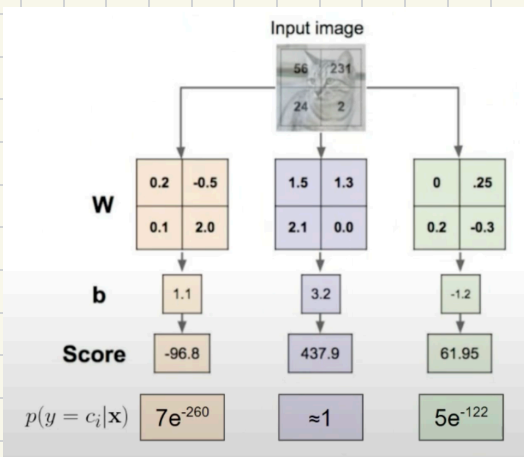
→ 확률 계산 사용 (① 모든 p의 sum = 1
② 모든 p는 $0 \leq p \leq 1$ 이다.)

$$\begin{aligned} \approx, p(y=1|x) &= \sigma(s_1 - s_2) = \frac{1}{1 + e^{-(s_1 - s_2)}} = \frac{e^{s_1}}{e^{s_1} + e^{s_2}} \\ p(y=0|x) &= \sigma(s_2 - s_1) = \frac{1}{1 + e^{-(s_2 - s_1)}} = \frac{e^{s_2}}{e^{s_1} + e^{s_2}} \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{sum} = 1.$$

→ Generalize.

$$p(y=k|x) = \frac{e^{s_k}}{\sum_{i=1}^n e^{s_i}} \rightarrow \text{해당 index.}$$

$p(y=3|x) = 0.8$ 이라고 해서 class=3 예측의 80%가 정답이란게 아니다. Just confidence 역할.
But threshold 지정이 필요하다.



다음과 같이 last layer에 softmax classifier를 달아서 Model의 예측을 해석가능하도록 만들 수 있다.

< parameter 학습 >

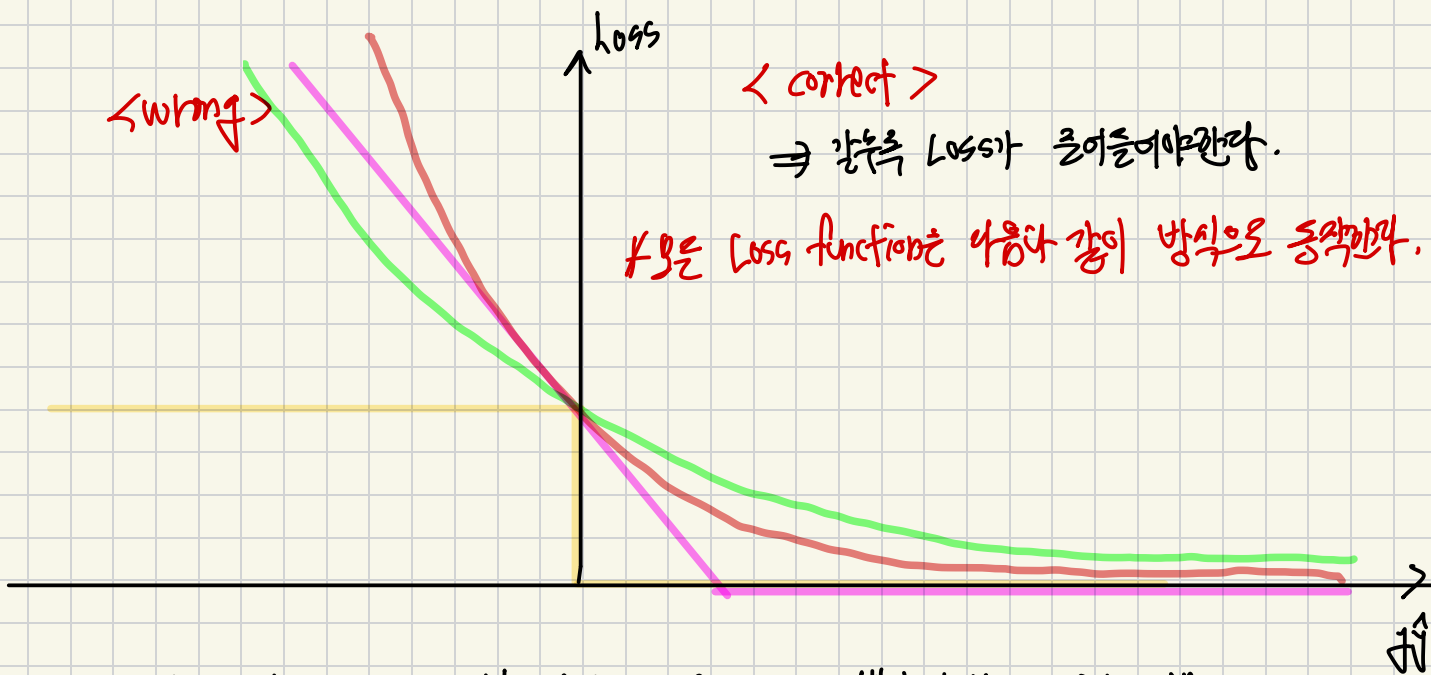
Human이 Model의 틀만 설정. 나머지는 Model이 Data로 부터 학습한다.
 answer를 보고 손실의 의미를 W를 조정한다.

Loss function : Model이 잘 학습하는지 수치화 해준다.

$L(y, \hat{y})$ 를 계산한다.

$y = \{+1, -1\}$ 인 경우. (\hat{y} 의 부호를 가지고 예측한다)

즉, $y\hat{y} > 0$ 인 경우 = model의 예측이 맞은 경우
 $y\hat{y} < 0$ 인 경우 = model의 예측이 틀린 경우.



— : 가장 이상적인 Loss : 맞으면 $loss = 0$ 틀리면 $loss = k$ \Rightarrow 미분이 안되는 point가 너무 critical해서 사용X.

— : $\log loss = \log(1 + e^{-y\hat{y}})$ \approx Logistic Regression Loss

— : $e^{-y\hat{y}}$
 \hookrightarrow Noise에 취약

— : Hinge Loss : $\max(0, 1 - y\hat{y})$ \Rightarrow 미분X point가 많은 critical X여서 사용X. \rightarrow 잘못 사용

$\hat{=} SVM$

$y = \{0, 1\}$ 인 정수

$p(y=0|x) = \hat{y}$, $p(y=1|x) = 1 - \hat{y}$ 이다.

$y = [0, 0, 0, 1, 0, 0]$
 $\hat{y} = [0.2, 0, 0, 0.7, 0.1, 0.5]$ 다음과 같은 형태를 갖는다.

* y 와 \hat{y} 의 값을 본다면 비슷하면 model이 잘 예측했다고 생각 가능.

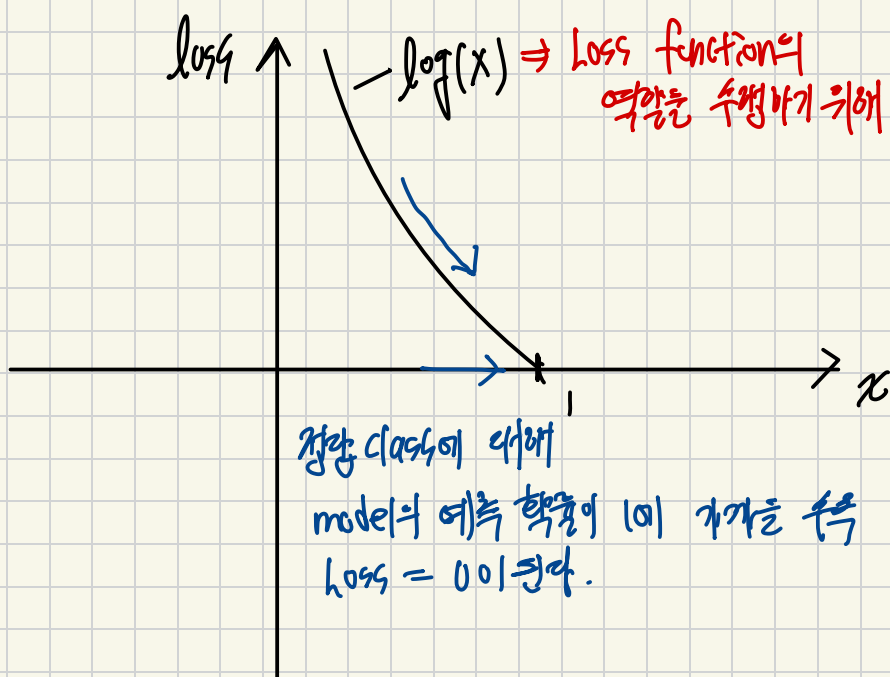
"Cross Entropy"

$$L = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \underline{y_{ik}} \log(\underline{\hat{y}_{ik}})$$

\therefore 정답에 대해 얼마나 확신할 수 있는지를 나타낸다.

$\underline{y_{ik}}$ 정답값, $\underline{\hat{y}_{ik}}$ model 예측값.

\hookrightarrow 정답 = 1
 else = 0

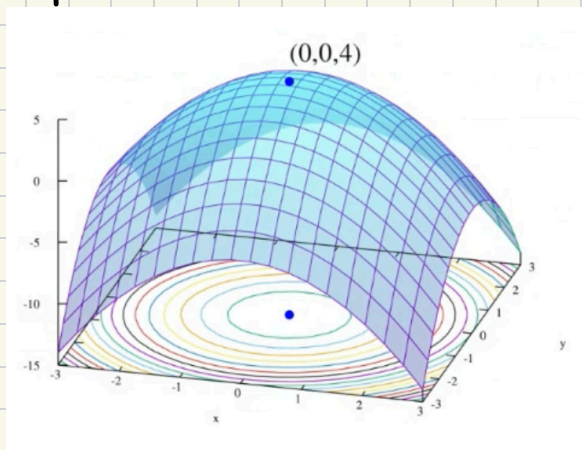


"KL-Divergence"

$$D_{KL}(P||Q) = \sum_i p(i) \log \frac{p(i)}{Q(i)}$$

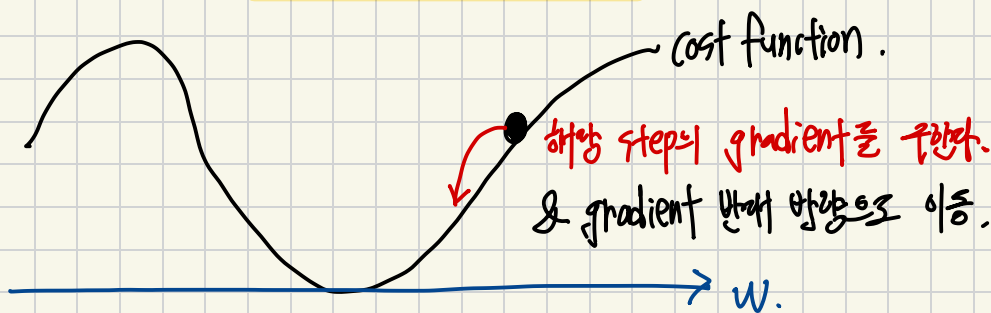
P 과 Q 의 확률분포 사이의 차이를 수치화 해준다.

Optimization



Loss function = Cost function $\hat{=}$ minimize.

* Cost function에서 x 와 y 를 model의 가중치이고, z 는 Loss이다.



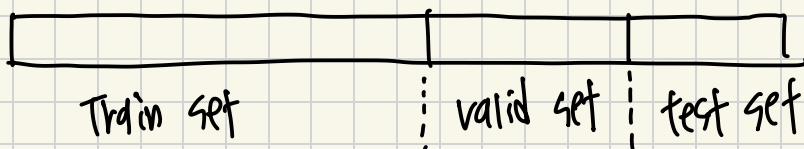
$$\theta^{\text{new}} = \theta^{\text{old}} - \eta \cdot \frac{\nabla \text{cost-function}}{\nabla \theta^{\text{old}}}$$

mini-batch \Rightarrow N 개의 sample만 사용해서 이동하자. (N 이 작을수록 정확. But Cost \uparrow)

Cross Validation

Train 성능 \uparrow 이 목표가 아냐. 일반적인 data를 잘 예측하는 model이 목표야.

hyperparameter 설정을 위한 validation set을 만들어야함.



valid set도 마지막에 Train 하듯이 이동할수록 성능도 있음. (validation overfitting 가능성 있음)

