$$P_\theta(x) = p(x;\theta) = p(x|\theta)$$

확률분포 $\theta$에서 $x$라는 점이 가지는 확률

$$P_\theta(y|x) = p(y|x;\theta) = p(y|x,\theta)$$

확률분포 $\theta$가 있을 때 $x$가 주어졌을 때 $y$의 확률값

MLE

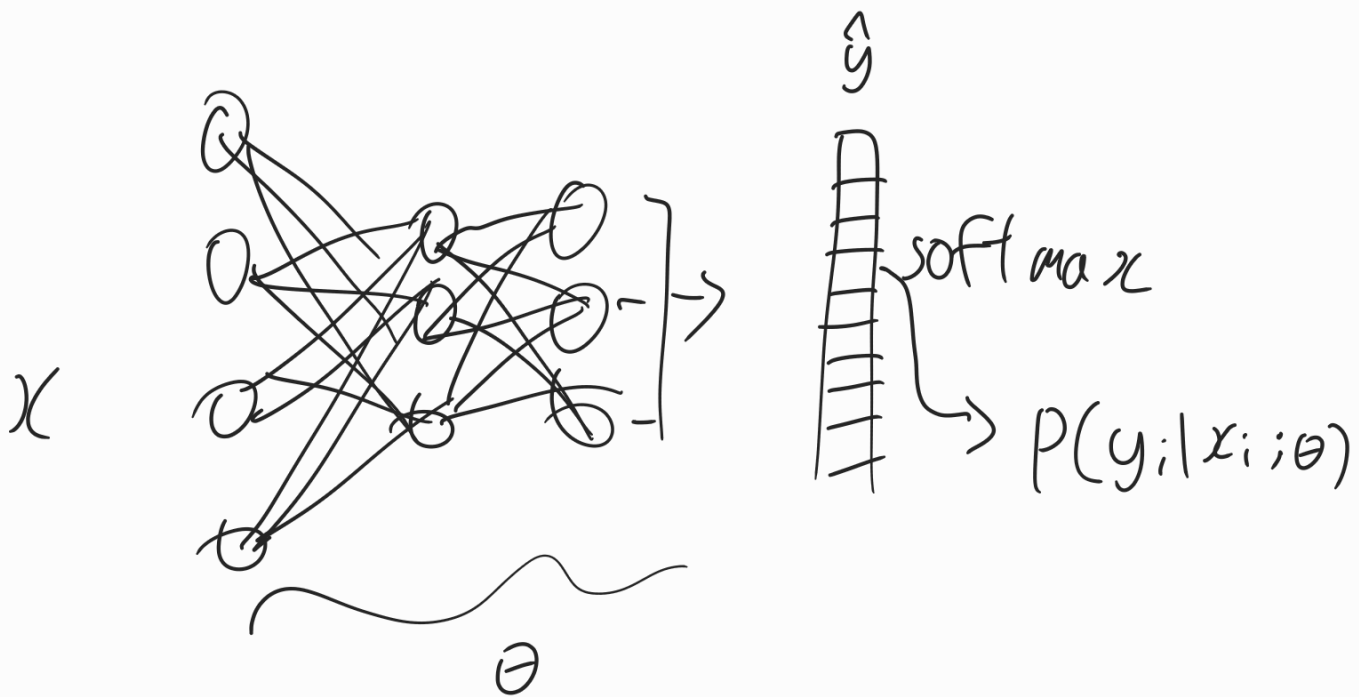$$\theta = \{ w_1, b_1, w_2, b_2 \cdots \}$$

Gradient Ascent를 통해 도출

하지만 대부분 딥러닝 프레임워크는 Gradient descent만 지원

No problem    -1 곱 하면 됨

따라서 Maximization 에서 minimization 문제로짐?

Negative log Likelihood (NLL)

$X$

$\hat{y}$

softmax

$\rightarrow P(y_i | x_i ; \theta)$

$\theta$

$$-\sum_{i=1}^{N} \log p(y_i | x_i ; \theta) = -\sum_{i=1}^{n} y_i^T \cdot \log \hat{y}_i$$

← cross entropy Loss

자 MLE

$$P(y_i | x_i ; \theta)$$

확률 분포 $\theta$ ( DL에서는 weight )에서 $X$과 주어졌을 때의 $Y$의 확률

한번의 시행에서 사건 A가 일어날 확률은 $P$라 할 때, 이 시행을 N회 반복한 독립시행에서 사건 A가 $r$번 일어날 확률은 $nC_r P^r (1-P)^{n-r}$

MLE $L(\theta) =$ $\left( \prod_{i=1}^{N} P(y_i | x_i ; \theta) \right)$

해당 관찰 결과가 발생할 확률

여기에 $X$을 $t$로 반복기 위해 Data underflow 를 방지하기위해

LogL($\theta$) 로그 Maximum likelyhood 로 변환

$$\sum_{i=1}^{N} \log P(Y_i|x_i;\theta)$$

(이 식을 최대로 만드는 $\theta$가 가능도함수 $L(\theta)$도 최대로 만듦)

MLE의 목적 $L(\theta)$가 가장 큰 곳의 $\theta$ 즉 $L(\theta)$ 를 미분했을 때 값이 0 이 되는 $\theta$를 찾는것

$$\hat{\theta} = \underset{\theta \in \Theta}{\arg\max} \sum_{i=1}^{N} \log P(Y_i|x_i;\theta)$$

그런데  이 식은 Gradient Ascent 를 통해 가중치 조절

$\Downarrow$

$$\hat{\theta} = \underset{\theta \in \Theta}{\arg\min} - \sum_{i=1}^{N} \log P(Y_i|x_i;\theta) \text{ 로변환}$$

$\longrightarrow$ Gradient descent

Model 의 결과물 $\hat{y}$

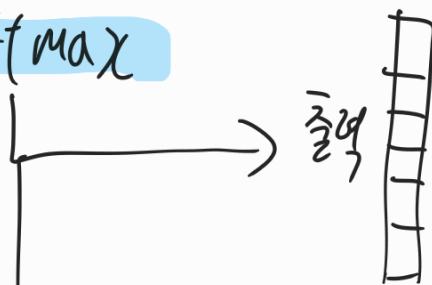X 가 파라미터 $\theta$를 거쳐서 나온 출력 $\hat{y}$
$\hookrightarrow$ 확률

$P(y_i|x_i;\theta)$
$\hookleftarrow$ 확률 분포 $\theta$ ( DL에서는 weight ) 에서 X가 주어졌을 때 Y의 확률
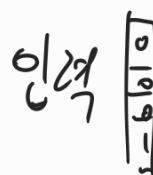
$\hat{y}_i = p(y_i|x_i;\theta)$

$$-\sum_{i=1}^{N} \log p(y_i|x_i;\theta) = -\sum_{i=1}^{N} \log \hat{y}_i$$

자 여기서 softmax

$\longrightarrow$ 출력 │ 각 클래스 별로 확률 제공

$\longrightarrow$ 인덱 │ one-hot vector

$\hat{y}_i = y_i^T \times \hat{y}_i$
$\hookrightarrow$ i번째 1 나머지 0 해당 index만 생존

$$-\sum_{i=1}^{N} \log P(y_i | x_i; \theta) = -\sum_{i=1}^{N} y^T_i \cdot \log \hat{y}_i$$

cross entrophy loss

평균내는 건 상수