

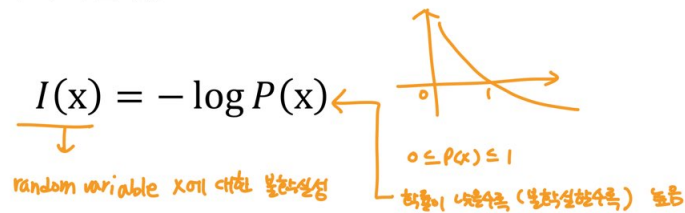
# 7. Information and Entropy

## Information

불확실성(Uncertainty)을 나타내는 값으로 random variable  $X$ 에 대한 불확실성이다.

확률이 낮을수록 불확실성이 높다. 왜냐하면  $P(x)$ 는  $[0, 1]$ 인데  $-\log$ 는 0에 가까울수록 무한, 1에 가까울수록 0이기 때문이다.

을 나타내는 값



## Entropy

정보량의 기대값(평균)으로 분포의 평균적인 uncertainty를 나타내는 값이다.

이 Entropy를 통해 분포의 형태를 예측해 볼 수 있는데 entropy가 작을수록 평평하고, 클수록 한 군데에 집중된 sharp한 형태이다.

$$H(P) = -\mathbb{E}_{x \sim P(x)}[\log P(x)]$$

## Cross Entropy

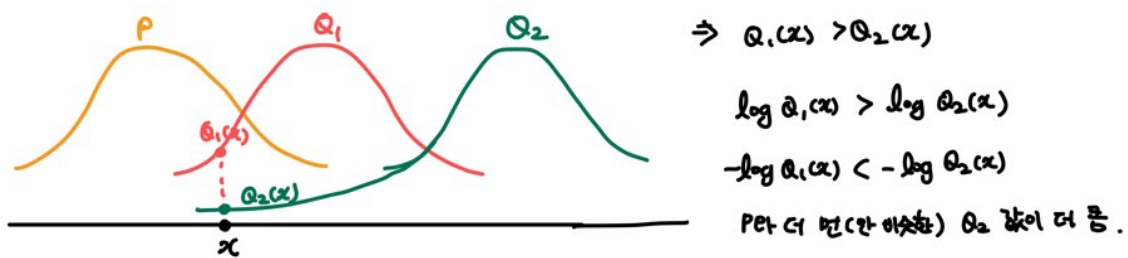
분포  $P$ 의 관점에서 본 분포  $Q$ 의 정보량의 평균이다.

⇒아래 식은  $P(X)$ 에서 sampling 한  $x$ 를  $Q(x)$  분포에 넣어주고 평균을 구하는 것이다.

여기서  $P$ 를 참인 확률분포,  $Q$ 를 신경망 확률분포라 생각하면 적용하기 쉽다.

$\downarrow$  항  $\downarrow P_\theta$  : neural network  
 $H(P, Q) = -\mathbb{E}_{x \sim P(x)} [\log Q(x)] = -\int p(x) \log q(x) dx$   
 $\underbrace{\hspace{10em}}$   
 $P(x)$ 에서 sampling 한  $x$ 를  $Q(x)$ 함수에 넣어줌  
 monte carlo  
 $\approx -\frac{1}{n} \sum_{i=1}^n \log p(x_i)$

두 분포가 비슷할수록 작은 값을 가진다.



## DNN Optimization using Cross Entropy

$$\begin{aligned}
 \mathcal{L}(\theta) &= -\mathbb{E}_{x \sim P(x)} \left[ \mathbb{E}_{y \sim P(y|x)} [\log P(y|x; \theta)] \right] \\
 &\stackrel{\text{monte carlo}}{\approx} -\frac{1}{N \cdot k} \sum_{i=1}^N \sum_{j=1}^k \log P(y_{i,j} | x_i; \theta) \\
 &\stackrel{k=1}{\approx} -\frac{1}{N} \sum_{i=1}^N \log P(y_i | x_i; \theta), \text{ if } k = 1. \\
 &\stackrel{\text{one-hot encoding}}{\approx} -\frac{1}{N} \sum_{i=1}^N y_i^T \cdot \log \hat{y}_i \\
 &\text{Softmax layer output}
 \end{aligned}$$

## KL-Divergence and Cross Entropy

KL-Divergence와 Cross Entropy를 seta로 미분하면 같다.

$$\begin{aligned}
 \text{KL}(p||p_\theta) &= -\mathbb{E}_{x \sim p(x)} \left[ \log \frac{p_\theta(x)}{p(x)} \right] \\
 &= - \int p(x) \log \frac{p_\theta(x)}{p(x)} dx \\
 &= - \int p(x) \log p_\theta(x) dx + \int p(x) \log p(x) dx \\
 &= \underbrace{H(p, p_\theta)}_{\text{Cross Entropy}} - \underbrace{H(p)}_{\text{Entropy}}
 \end{aligned}$$

$$\nabla_\theta \text{KL}(p||p_\theta) = \nabla_\theta H(p, p_\theta) - \cancel{\nabla_\theta H(p)}$$