

<KL-Divergence>

정의: 두 확률 분포의 차이를 계산할 때 사용되는 함수.

어떤 이상적인 분포에 대해, 그 분포를 근사하는 다른 분포를 사용해 샘플량을 할 때 발생할 수 있는 정보 엔트로피 차이를 계산한다.

* 정보 엔트로피 차이는 Cross-Entropy를 이용해 구할 수 있다.

먼저 Cross-Entropy를 수식으로 표현하면 아래와 같다.

$$H(p, q) = \sum_i p_i \log_2 \frac{1}{q_i} = - \sum_i p_i \log_2 q_i$$

여기서 p_i 는 target 값 (확률에 대한 참값), q_i 는 우리가 학습한 모델에서 도출된 결과값이다. 즉 우리가 어떤 q_i 를 학습하고 있고 이 방향이 올바른 방향으로 진행된다면, Cross-entropy는 작아지는 방향으로 진행될 것이다.

위 Cross-Entropy를 이용해 KL-Divergence를 구하면 아래와 같다.

$$\begin{aligned} H(p, q) &= - \sum_i p_i \log(q_i) \\ &= - \sum_i p_i \log(q_i) - \underbrace{\sum_i p_i \log(p_i)}_{= H(p)} + \sum_i p_i \log(p_i) \end{aligned}$$

$$= H(p) + \sum_i p_i \log(p_i) - \sum_i p_i \log(q_i)$$

$$= H(p) + \sum_i p_i \log\left(\frac{p_i}{q_i}\right)$$

p 의 엔트로피에 이만큼 더해진 것이
Cross-entropy

이만큼 더해진 것이 바로
 p, q 의 정보량 차이, 즉 KL-divergence가 된다.

KL-Divergence는 $KL(p||q)$ 로 쓰여, 위 식에서 유도된 $KL(p||q)$ 는 아래와 같다

$$KL(p||q) = H(p, q) - H(p)$$

또한 KL-divergence의 정확한 식은 아래와 같다.

$$KL(p||q) = \begin{cases} \sum_i p_i \log \frac{p_i}{q_i} & \text{or} & \sum_i p_i \log \frac{q_i}{p_i} & (\text{discrete}) \\ \int p(x) \log \frac{p(x)}{q(x)} dx & \text{or} & - \int p(x) \log \frac{q(x)}{p(x)} dx & (\text{continuous}) \end{cases}$$

$KL(p||q)$ 의 특성을 가장 대표적인 2개는 다음과 같다.

① $KL(p||q) > 0$

② $KL(p||q) \neq KL(q||p)$.

examples)

1) P 와 Q 2개의 확률 분포가 다음과 같이 나와있을 때,
 $KL(P||Q)$ 와 $KL(Q||P)$ 를 구하는 과정은 아래와 같다.

x	0	1	2
Distribution $P(x)$	$9/25$	$12/25$	$4/25$
Distribution $Q(x)$	$1/3$	$1/3$	$1/3$

$$1) KL(P||Q) = \sum P(x) \cdot \ln\left(\frac{P(x)}{Q(x)}\right) = 9/25 \cdot \ln \frac{9/25}{1/3} + 12/25 \cdot \ln \frac{12/25}{1/3} + 4/25 \cdot \ln \frac{4/25}{1/3}$$

$$2) KL(Q||P) = \sum Q(x) \cdot \ln\left(\frac{Q(x)}{P(x)}\right) = 1/3 \cdot \ln \frac{1/3}{9/25} + 1/3 \cdot \ln \frac{1/3}{12/25} + 1/3 \cdot \ln \frac{1/3}{4/25}$$

2) 자주 사용되는 gaussian distribution에 대한 KL-divergence는
아래와 같다.

$$KL(P||Q) = \int P(x) \cdot \ln \frac{P(x)}{Q(x)} dx$$

$$P(x) = N(\mu_P, \Sigma_P) = (2\pi)^{-k/2} \cdot |\Sigma_P|^{-1/2} \cdot e^{(-\frac{1}{2}(x-\mu_P)^T \Sigma_P^{-1} (x-\mu_P))}$$

$$Q(x) = N(\mu_Q, \Sigma_Q) = (2\pi)^{-k/2} \cdot |\Sigma_Q|^{-1/2} \cdot e^{(-\frac{1}{2}(x-\mu_Q)^T \Sigma_Q^{-1} (x-\mu_Q))}$$

P 와 Q 모두 k -dimension 이라 가정하면

$$\begin{aligned} KL(P||Q) &= E_P[\log(P) - \log(Q)] = E_P\left[\frac{1}{2} \log \frac{|\Sigma_Q|}{|\Sigma_P|} - \frac{1}{2} (x-\mu_P)^T \Sigma_P^{-1} (x-\mu_P) + \frac{1}{2} (x-\mu_Q)^T \Sigma_Q^{-1} (x-\mu_Q)\right] \\ &= \frac{1}{2} \log \frac{|\Sigma_Q|}{|\Sigma_P|} - \frac{1}{2} E[(x-\mu_P)^T \Sigma_P^{-1} (x-\mu_P)] + \frac{1}{2} E[(x-\mu_Q)^T \Sigma_Q^{-1} (x-\mu_Q)] \end{aligned}$$

먼저, $(x-\mu_P)^T \Sigma_P^{-1} (x-\mu_P)$ 는 trace 연산을 통해 구할 수 있다. 즉

$$(x-\mu_P)^T \Sigma_P^{-1} (x-\mu_P) = \text{tr}\{(x-\mu_P)(x-\mu_P)^T \Sigma_P^{-1}\} \text{ 이다.}$$

$$\begin{aligned} \therefore E[(x-\mu_P)^T \Sigma_P^{-1} (x-\mu_P)] &= E[\text{tr}\{(x-\mu_P)(x-\mu_P)^T \Sigma_P^{-1}\}] \\ &= \text{tr}\{E[(x-\mu_P)(x-\mu_P)^T] \cdot \Sigma_P^{-1}\} \end{aligned}$$

$$\text{이 때 } E[(x-\mu_p)(x-\mu_p)^T] = \Sigma_p \quad \text{이므로}$$

$$\text{tr}\{E[(x-\mu_p)(x-\mu_p)^T] \cdot \Sigma_p^{-1}\} = \text{tr}\{\Sigma_p \Sigma_p^{-1}\} = \text{tr}\{I_k\} = k \text{ 가 된다.}$$

$$E[(x-\mu_g)^T \Sigma_g^{-1} (x-\mu_g)] = (\mu_p - \mu_g)^T \Sigma_g^{-1} (\mu_p - \mu_g) + \text{Tr}\{\Sigma_g^{-1} \Sigma_p\} \quad \text{이다.}^{[1]}$$

$$\text{위 식을 정리하면 } KL(p||g) = \frac{1}{2} \left(\log \frac{|\Sigma_g|}{|\Sigma_p|} - k + (\mu_p - \mu_g)^T \Sigma_g^{-1} (\mu_p - \mu_g) + \text{tr}\{\Sigma_g^{-1} \Sigma_p\} \right)$$

가 나오,

만약 $g \sim N(0, I)$ 라면

$$KL(p||g) = \frac{1}{2} [\mu_p^T \mu_p + \text{tr}\{\Sigma_p\} - k - \log|\Sigma_p|] \text{ 가 된다.}$$

[1] Matrix Cookbook section B.2, eq 380.

8.2.2 Mean and variance of square forms

Mean and variance of square forms: Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$

$$E(\mathbf{x}\mathbf{x}^T) = \Sigma + \mathbf{m}\mathbf{m}^T \quad (377)$$

$$E[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \text{Tr}(\mathbf{A}\Sigma) + \mathbf{m}^T \mathbf{A} \mathbf{m} \quad (378)$$

$$\begin{aligned} \text{Var}(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \text{Tr}[\mathbf{A}\Sigma(\mathbf{A} + \mathbf{A}^T)\Sigma] + \dots \\ &\quad + \mathbf{m}^T (\mathbf{A} + \mathbf{A}^T) \Sigma (\mathbf{A} + \mathbf{A}^T) \mathbf{m} \end{aligned} \quad (379)$$

$$E[(\mathbf{x} - \mathbf{m}')^T \mathbf{A} (\mathbf{x} - \mathbf{m}')] = (\mathbf{m} - \mathbf{m}')^T \mathbf{A} (\mathbf{m} - \mathbf{m}') + \text{Tr}(\mathbf{A}\Sigma) \quad (380)$$