

2장 : Entropy, Relative, Mutual

* Entropy

$$H(X) = - \sum_{x \in X} p(x) \log p(x) = E[-\log p(x)]$$

• 특성 함수 g 에 대해 $g(X)$ 도 R.V.

1) Non-negativity : $H(X) \geq 0$

PF) $-\log p(x) > 0$

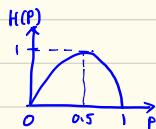
2) $H_b(X) = \log_b a H_a(X)$

PF) $\log_b p(x) = \log_a p(x) \cdot \log_b a$

- Example i) Bernoulli RV's entropy

- $X = \begin{cases} 0 & \text{확률 } p \\ 1 & \dots \\ 1-p \end{cases}$

- $H(X) = -p \log p - (1-p) \log(1-p) = H(p)$



binary entropy function

* Joint Entropy

- R.V. X, Y 에 대한 Joint Entropy

- $H(X, Y) = -E[\log p(X, Y)] = -\sum_x \sum_y p(x, y) \log p(x, y)$

• 의미: 마치 2개의 R.V 가 모여, 하나의 RV가 되는 것.

- Marginal Distribution

- $\sum_{y \in Y} p(x, y) = p(x)$

- ex) $\begin{array}{c|cc|c} & 0 & 1 & p(x) \\ \hline 0 & \frac{1}{2} & \frac{1}{2} & \frac{3}{4} \\ 1 & \frac{1}{8} & \frac{1}{8} & \frac{1}{4} \end{array} \quad \text{marginal sum.}$

* Conditional Entropy

- $(X, Y) \sim p(x, y)$ 일 때 $-H(Y|X=x)$

conditional probability
 $\cdot P(Y|X) = \frac{P(X, Y)}{P(X)}$

- $H(Y|X) = -\sum_x p(x) \sum_y p(y|x) \log p(y|x)$

$$= -\sum_x \sum_y p(x, y) \log p(y|x) = -E_{x,y} [\log p(y|x)]$$

ex)

X	1	2
1	$\frac{1}{2}$	$\frac{1}{4}$
2	$\frac{1}{8}$	$\frac{1}{8}$

$$H(Y|X) = ?$$

$$H(Y|X) = - \sum_{x,y} p(x,y) \log p(y|x)$$

$$\textcircled{1} \quad H(Y|X) = - \sum_{x,y} p(x,y) \log p(y|x)$$

$$= p(y=1|x=1) = \frac{\frac{1}{2}}{\frac{1}{2} + \frac{1}{8}} = \frac{4}{5}, \quad p(y=1|x=2) = \frac{\frac{1}{4}}{\frac{1}{4} + \frac{1}{8}} = \frac{2}{3}$$

$$p(y=2|x=1) = \frac{\frac{1}{8}}{\frac{1}{2} + \frac{1}{8}} = \frac{1}{5}, \quad p(y=2|x=2) = \frac{1}{3}$$

$$H(Y|X) = - \left(\frac{1}{2} \log \frac{4}{5} + \frac{1}{5} \log \frac{1}{5} + \frac{1}{4} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3} \right)$$

$$\textcircled{2} \quad H(Y|X) = \boxed{\sum_x p(x) H(Y|x)}$$

$$= \boxed{(p(x=1) H(Y|x=1) + p(x=2) H(Y|x=2))}$$

$$= \boxed{\left(\frac{2}{3} H\left(\frac{4}{5}\right) + \frac{1}{3} H\left(\frac{1}{3}\right) \right)}$$

$$H(Y|X) = \sum_x p(x) H(Y|X=x)$$

* Chain Rule

$$H(X,Y) = H(X) + H(Y|X)$$

$$\text{pf)} \quad H(X,Y) = - \sum_{x,y} p(x,y) \log p(x,y)$$

$$= \log p(x) \cdot p(y|x)$$

$$= \log p(x) + \log p(y|x) = H(X) + H(Y|X)$$

$$H(X,Y) = H(X) + H(Y|X)$$

$$= H(Y) + H(X|Y) > H(Y|X) \neq H(X|Y)$$

* Relative Entropy \Rightarrow KL divergence

- 두 probability distribution의 distance 측정.

- RV의 분포를 P 라고 가정했을 때, length = $H(P)$

그러면 실제 분포를 우리가 모른다고, 우리가 가정해버리면, 필요 bit = $H(Q) = H(P) + D(P||Q)$

- KL Divergence

$$\cdot D(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} = E_P \left[\log \frac{P(x)}{Q(x)} \right]$$

• If $P=Q$, $D(P||Q)=0$

• $D(P||Q) \neq D(Q||P)$: Asymmetric

Penalty

* Mutual Information $\Rightarrow I(X;Y)$

- 2개의 RV X, Y 에 대해서 joint prob $P(x,y)$ 와 marginal prob $P(x), P(y)$ 를 가질 때,

$$\cdot I(X;Y) = D(P(x,y) || P(x)P(y))$$

$$= \sum_{x,y} P(x,y) \log \frac{P(x,y)}{P(x)P(y)}$$

$$\cdot I(X;Y) = \sum_{x,y} P(x,y) \log \frac{P(x,y)}{P(x)P(y)}$$

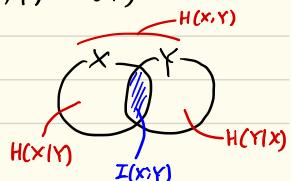
$$~~~~~ = \log \frac{P(x,y)P(y)}{P(x)P(y)} = \log \frac{P(x|y)}{P(x)} = \log P(x|y) - \log P(x)$$

$$I(X;Y) = - \sum_{x,y} P(x,y) \log P(x) - \left(- \sum_y P(x,y) \log P(x|y) \right)$$

$$= H(X) - H(X|Y)$$

$$\cdot I(X;Y) = H(X) - H(X|Y) = H(X) - (H(X,Y) - H(Y))$$

$$= H(X) + H(Y) - H(X,Y) \Rightarrow$$



* Chain rule for entropy

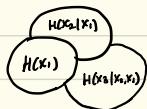
- 어떤 개의 RV $X_1, \dots, X_n \sim p(x_1, \dots, x_n)$ 일 때

$$\therefore H(X) = H(X|Y) + I(X;Y)$$

$$\cdot H(X_1, \dots, X_n) = \sum_{i=1}^n H(X_i | X_1, \dots, X_{i-1}) \leq H(X_1) + H(X_2) + \dots + H(X_n)$$

\nwarrow 확률은 품임. 엔트로피는 토심.

\hookrightarrow 등호는 independent.



* Conditional Mutual Information

- Z 가 주어졌을 때, X 와 Y 의 conditional MI.

$$\cdot I(X;Y|Z) = H(X|Z) - H(X|Y,Z)$$

$\hookrightarrow I(X;Y)$ 보다 클수도, 작을 수도 있음.

- Chain rule of MI

$$\cdot I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_1, \dots, X_{i-1})$$

$\underbrace{\quad}_{\text{MI}}$

* Conditional Relative Entropy

- $p(y|x)$ 와 $q(y|x)$ 의 relative entropy.

$$\approx \sum_x p(x) \frac{1}{q(x)} \log \frac{p(y|x)}{q(y|x)}$$

$$\cdot D(p(y|x) || q(y|x)) = \sum_x E_x [D(p(y|x=x) || q(y|x=x))]$$

$$= \sum_x \sum_y p(x,y) \log \frac{p(y|x)}{q(y|x)} = E_{x,y} \left(\log \frac{p(y|x)}{q(y|x)} \right)$$

\hookrightarrow 확률임.

$$\cdot D(p(x,y) || q(x,y)) = D(p(x) || q(x)) + D(p(y|x) || q(y|x))$$

$$\text{PF) } \frac{p(x,y)}{q(x,y)} = \frac{p(y|x)p(x)}{q(y|x)q(x)} = \frac{p(y|x)}{q(y|x)} \cdot \frac{p(x)}{q(x)}$$

↑
 $\log \rightarrow$ 확률임.

아래 볼록

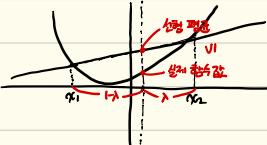
* Convex functions

- 모든 $x_1, x_2 \in (a, b)$ 와 $0 \leq \lambda \leq 1$ 에 대해, 아래 그림에서 $f(x)$ 는 (a, b) 에서 convex function이다.

$$\bullet f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2) \quad \text{--- Theorem 2.6.1.}$$

한국어

선형평균



• $\lambda=0$ or 1 일 때 equal.

위로 볼록

- $f(x)$ 가 convex면 $-f(x)$ 은 concave.

- $f(x)$ 의 second derivative가 positive면, $f(x)$ 는 convex.

증명) Taylor series : $f(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x^*)}{2}(x-x_0)^2$

→ equality를 만족하는 x^* 가 존재.

$$x_0 = \lambda x_1 + (1-\lambda)x_2 \text{ 라 하면}$$

$$1) x=x_1 \text{ 이면 } f(x_1) \geq f(x_0) + f'(x_0) [(1-\lambda)(x_1-x_2)] \quad \text{--- ①}$$

$$2) x=x_2 \text{ 이면 } f(x_2) \geq f(x_0) + f'(x_0) [\lambda(x_2-x_1)] \quad \text{--- ②}$$

$$\text{①} \times \lambda + \text{②} \times (1-\lambda) \therefore \lambda f(x_1) + (1-\lambda)f(x_2) \geq f(x_0) = f(\lambda x_1 + (1-\lambda)x_2) \quad \square$$

* Jensen's Inequality

- $f(x)$ 가 convex function, X 가 RV면

$$\bullet E[f(X)] \geq f(E[X]) \quad \text{위의 식에서 } \lambda, 1-\lambda \text{가 RV의 probability인 것 됨.}$$

PF) X 의 값: x_1, x_2, \dots , 확률: p_1, p_2, \dots 일 때.

$$1. P_1 f(x_1) + P_2 f(x_2) \geq f(P_1 x_1 + P_2 x_2), P_1 + P_2 = 1 \text{인가? } P_1 = \lambda \text{면 2.6.1이 의해 } 0.$$

$$2. \sum_{i=1}^{k-1} P_i f(x_i) \geq f\left(\sum_{i=1}^{k-1} P_i x_i\right), \sum_{i=1}^{k-1} P_i = 1 \text{ 일 때 증명.}$$

$$\begin{aligned} \sum_{i=1}^k P_i f(x_i) &= P_k f(x_k) + (1-P_k) \sum_{i=1}^{k-1} \frac{P_i}{1-P_k} f(x_i) \quad (\because \sum_{i=1}^{k-1} P_i = 1 - P_k) \\ &= P_k f(x_k) + (1-P_k) \sum_{i=1}^{k-1} P'_i f(x_i) \quad (P'_i = \frac{P_i}{1-P_k}, \sum_{i=1}^{k-1} P'_i = 1) \\ &\geq P_k f(x_k) + (1-P_k) \sum_{i=1}^{k-1} P'_i x_i \\ &\geq f(P_k x_k + (1-P_k) \sum_{i=1}^{k-1} P'_i x_i) = f\left(\sum_{i=1}^k P_i x_i\right) \quad \square \end{aligned}$$

2.6.1

* Information Inequality

- 두 Probability distribution $p(x), q(x)$ 이면 다음

$$\cdot D(p \parallel q) \geq 0. \text{ 등로는 } p(x) = q(x) \text{ 일 때 } 0.$$

PF) Using Jensen's Inequality

- \log ; concave function

$$- D(p \parallel q) = - \sum_x p(x) \log \frac{p(x)}{q(x)}$$

$$= \sum_x p(x) \log \frac{q(x)}{p(x)}$$

$$\leq \log \sum_x p(x) \cdot \frac{q(x)}{p(x)} = \log \sum q(x) = \log 1 = 0.$$

$\hookrightarrow (\because \log \text{ is concave function. } \sum p(x) = 1)$

- Non-negativity of MI Corollary (특별 case) 증명

$$\cdot I(X; Y) = D(p(x,y) \parallel p(x)p(y)) \geq 0, \text{ 등로는 } p(x), p(y) \text{가 independent 일 때.}$$

$$\text{PF) } D(p \parallel q) = 0 \Leftrightarrow p = q \text{ 이므로}$$

$$I(X; Y) = 0 \Leftrightarrow p(x,y) = p(x)p(y) : X, Y \text{가 서로 independent.}$$

$$\cdot I(X; Y | Z) \geq 0 \text{ 등로 : } X, Y \text{ independent given } Z.$$

$$\Rightarrow H(X) = \log |X| - D(p \parallel u)$$

* Uniform distribution \rightarrow RU X 에 의해 생성되는 symbol 종류를 갖는 집합.

- $|X|$: X 의 원소의 개수라 하면

$$\cdot H(X) \leq \log |X|. \text{ 등로 : } X \text{가 } X \text{에 대해 uniform distribution.}$$

$$\text{PF) } u(x) = \frac{1}{|X|} : \text{uniform 확률 분포.}$$

$p(x) : X$ 에 대한 확률 분포 라 하면

$$D(p \parallel u) = \sum p(x) \log \frac{p(x)}{u(x)} = \sum p(x) \log p(x) - \sum p(x) \log u(x)$$

$$= -H(X) - \sum p(x) \log \frac{1}{|X|}$$

$$= -H(X) + \log |X| \underbrace{\sum p(x)}_{=1} =$$

$$= -H(X) + \log |X| \geq 0$$

$\therefore H(X) \leq \log |X|. \text{ 등로는 } p=u \text{ (uniform).}$

* Conditional Entropy reduce.

- $H(X|Y) \leq H(X)$. 등호 : X, Y independent.

Pf) $I(X;Y) = H(X) - H(X|Y) \geq 0$.

- 특히 $RV Y$ known 이면, entropy가 줄어듬.

* Independence Bound and Entropy

- $X_1, \dots, X_n \sim p(x_1, \dots, x_n)$ 이면

- $H(X_1, \dots, X_n) \leq \sum_{i=1}^n H(X_i)$. 등호 : X_i all independent

Pf) $H(X_1, \dots, X_n) = \sum_{i=1}^n H(X_i | X_1, \dots, X_{i-1}) \leq \sum_{i=1}^n H(X_i)$,

* Log sum inequality

- Non-negative number a_1, \dots, a_n or b_1, \dots, b_n all $\neq 0$

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left(\sum_{i=1}^n a_i \right) \log \frac{\sum a_i}{\sum b_i}$$

Pf) $f(t) = t \log t \rightarrow f''(t) = \frac{1}{t} \log e > 0 \text{ for all } t > 0$.

$$\Rightarrow \sum a_i f(t_i) = f(\sum a_i t_i), \quad a_i = \frac{b_i}{\sum b_i}, \quad t_i = \frac{a_i}{b_i}, \quad \sum b_i = \lambda \text{ or } t$$

$$\Rightarrow \sum_i \frac{b_i}{\sum b_i} \cdot \frac{a_i}{b_i} \log \left(\frac{a_i}{b_i} \right) = \left(\sum_i \frac{b_i}{\sum b_i} \cdot \frac{a_i}{b_i} \right) \log \left(\sum_i \frac{b_i}{\sum b_i} \cdot \frac{a_i}{b_i} \right)$$

$$\Rightarrow \sum_i \frac{a_i}{\lambda} \log \frac{a_i}{b_i} = \sum_i \frac{a_i}{\lambda} \log \frac{\sum a_i}{\sum b_i}$$

$$\Rightarrow \therefore \sum_i a_i \log \frac{a_i}{b_i} = \left(\sum a_i \right) \log \frac{\sum a_i}{\sum b_i}$$

* Convexity of Relative Entropy

- $D(p \parallel q)$ is convex in pair (p, q) : $(p_1, q_1), (p_2, q_2)$ 에 대해

$$= D(\lambda p_1 + (1-\lambda)p_2 \parallel \lambda q_1 + (1-\lambda)q_2) = \lambda D(p_1 \parallel q_1) + (1-\lambda) D(p_2 \parallel q_2)$$

PF) log sum equality ; $\sum a_i \log \frac{a_i}{b_i} \geq \sum a_i \log \frac{\sum a_i}{\sum b_i}$ --- ①

$$a_1 = \lambda p_1, a_2 = (1-\lambda)p_2, b_1 = \lambda q_1, b_2 = (1-\lambda)q_2 \text{ 증명}$$

$$\begin{aligned} \text{①: } (\lambda p_1 + (1-\lambda)p_2) \log \frac{\lambda p_1 + (1-\lambda)p_2}{\lambda q_1 + (1-\lambda)q_2} &\leq \lambda p_1 \log \frac{\lambda p_1}{\lambda q_1} + (1-\lambda)p_2 \log \frac{(1-\lambda)p_2}{(1-\lambda)q_2} \\ &= \lambda p_1 \log \frac{p_1}{q_1} + (1-\lambda)p_2 \log \frac{p_2}{q_2} \end{aligned}$$

$$\text{위의 식을 모든 } x \text{에 대해 더함} : \sum \rightarrow \sum p(x) \log \frac{p(x)}{q(x)} = D(p \parallel q)$$

$$\Rightarrow D(\lambda p_1 + (1-\lambda)p_2 \parallel \lambda q_1 + (1-\lambda)q_2) \leq \lambda D(p_1 \parallel q_1) + (1-\lambda) D(p_2 \parallel q_2)$$

* Concavity of Entropy

- $H(p)$ is concave function of p .

$$\cdot H(\lambda p_1 + (1-\lambda)p_2) \geq \lambda H(p_1) + (1-\lambda)H(p_2)$$

$$\text{PF) } H(p) = \underbrace{\log |x|}_{\text{concave}} - \underbrace{D(p \parallel u)}_{\text{convex}} \Rightarrow \text{concave function.}$$

* Mutual Information and Input distribution

$$-(X, Y) \sim p(x, y) = p(x)p(y|x)$$

① $p(y|x)$ 가 주어졌을 때, $I(X; Y)$ 는 $p(x)$ 에 대한 concave function이다.
 (channel transaction probability) (변환 가능) \Rightarrow 최대 값이 있다.

$$\text{PF) } I(X; Y) = H(Y) - H(Y|X) = H(Y) - \sum_x p(x) H(Y|x)$$

\downarrow
 $p(y|x)$ fix 이고

$H(Y)$ 는 $p(Y)$ 에 대한 concave function이다

$H(Y)$ 는 $p(x)$ 에 대한 concave.

$$\therefore p(y) = \sum_x p(x) p(y|x)$$

) 따라서
 $I(X; Y)$ 는 concave

② $p(x)$ 가 주어졌을 때, $I(X; Y)$ 는 $p(y|x)$ 에 대한 convex function이다.

(채널의 변화)

RU: $Y_1, Y_2 \rightarrow Y$ 는 Y_1, Y_2 의 mixed.
 PF) 2개의 channel $p_1(y|x), p_2(y|x)$ 이 있다고 하면,

RU 자체 의해 한 channel이 선택된다 흥미. ($\lambda \in [0, 1] \ni p_2$)

$$\begin{aligned} I(X; Y, Z) &= I(X; Y|Z) + I(X; Z) \\ &= I(X; Z|Y) + I(X; Y) \end{aligned}$$

Z 는 X 와 상관 없음. $I(X; Z) = 0$

$$\therefore I(X; Y) \leq I(X; Y|Z) = \lambda I(X; Y_1) + (1-\lambda) I(X; Y_2) \quad \therefore \text{convex.}$$

* Markov Chain

- 3개의 RV X, Y, Z 에 대해 Z 는 Y 의 영향만 받고, Y 는 X 의 영향만 받을 때, 즉

$$\cdot p(x, y, z) = p(x)p(y|x)p(z|y, x) \rightarrow p(x)p(y|x)p(z|y)$$

이면 X, Y, Z 는 그 순으로 Markov Chain이며, $X \rightarrow Y \rightarrow Z$ 로 표기.

- $p(x, z|y) = p(x|y)p(z|y)$ Y 가 주어지면 x, z 는 independent.

- $X \rightarrow Y \rightarrow Z \Leftrightarrow Z \rightarrow Y \rightarrow X$

- $Z = f(Y)$ 면 $X \rightarrow Y \rightarrow Z$.

* Data Processing Inequality

- $X \rightarrow Y \rightarrow Z$ 면, $I(X;Y) \geq I(X;Z)$

$$\begin{aligned} \text{pf)} \quad I(X;Y,Z) &= I(X;Y|Z) + I(X;Z) \\ &= \underline{I(X;Z|Y)} + I(X;Y) \end{aligned}$$

여기서 $X \rightarrow Y \rightarrow Z$ 이므로 $I(X;Z|Y) = 0$: Y 가 주어지면 X, Z 는 independent

$$\text{그러므로 } I(X;Y) = I(X;Z) + I(X;Y|Z) \geq I(X;Z) \quad \square$$

$$\text{또한 } I(X;Y) \geq I(X;Y|Z)$$

같은 MI 가 줄어든다.

* Fano's Inequality

- X 를 Y 를 estimate 한 후에 error boundary 를 제공.

- X (sent) $\rightarrow Y$ (observed)

. $\hat{X} = g(Y)$: decoding (복호)

. $P_e = \Pr(\hat{X} \neq X)$

- Error

$$E = \begin{cases} 1 & \text{if } \hat{X} \neq X \\ 0 & \text{if } \hat{X} = X \end{cases}$$

- $\hat{X} = g(Y)$ 이므로 $X \rightarrow Y \rightarrow \hat{X}$

• Fano's Inequality : $H(P_e) + P_e \log |X| \geq H(X|\hat{X}) \geq H(X|Y)$

pf) $X \rightarrow Y \rightarrow \hat{X}$. E 또한 RV.

$$\begin{aligned} H(E, X|\hat{X}) &= H(E|X, \hat{X}) + H(X|\hat{X}) \quad H(X|\hat{X}) \geq H(X|Y) \\ &= \underline{H(X|E, \hat{X})} + \underline{H(E|\hat{X})} \end{aligned}$$

$$\cdot H(E|\hat{X}) \leq H(E) = H(P_e)$$

$$\cdot H(E|X, \hat{X}) : X, \hat{X} 면 E 를 알 수 있음. = 0.$$

$$\cdot H(X|E, \hat{X}) = \Pr(E=0) H(X|\hat{X}, E=0) + \Pr(E=1) H(X|\hat{X}, E=1) \quad \leq \log |X|$$

$$\leq (1-P_e) \cdot 0 + P_e \cdot \log |X|$$

$$\therefore H(X|\hat{X}) \leq H(P_e) + P_e \log |X|, \quad H(X|Y) \leq H(P_e) + P_e \log |X|.$$

예외로, $H(X|\hat{X}) \leq H(X)$

* Fano 정리

- RV X, Y 에 대해 $p = \Pr(X \neq Y)$ 라 하면

$$H(X|Y) \leq H(p) + p \log |X|$$

- Corollary : More strict

• $\hat{X} : Y \rightarrow X$, $p_e = \Pr(X \neq \hat{X})$ 라 하면

$$H(X|Y) \leq H(p_e) + p_e \log (\underbrace{|X|-1}_{\text{error}})$$

그자신이 빠졌다는 소리.