

TODAY: BASICS OF IT. (contd.)

- Conditional Entropy, KL Divergence
- (In)equalities

Review of Entropy

- Suppose transmitting n i.i.d copies $X_1 \dots X_n$, $X_i \sim P_x$
 $\& P_x = (P_1, \dots, P_m)$ $\left[\sum_{i=1}^m P_i = 1, P_i \geq 0 \right]$
- "Budget" $l_i^* = \log \frac{1}{P_i}$ bits to transmit i
- Total Transmission cost will be $\sum_{i=1}^m P_i l_i^* = H(X)!$
- Could try to use any other $\{l_i\}_{i=1 \dots m}$. Why $l_1^* \dots l_m^*$?
- ① Not all l_i achievable. E.g. $l_1 = l_2 = l_3 = \dots = l_m = 1!$

Psst 1, Problem 4: l_i achievable in prefix-free way $\Rightarrow \sum 2^{-l_i} \leq 1$.

So if we let $q_i = 2^{-l_i}$ then

we can try to pretend $X_j \sim Q = (q_1, \dots, q_m)$
 $\&$ transmit according to Q "sub-distribution".

But presumably $\sum P_i l_i = \sum P_i \log \frac{1}{z_i} \geq \sum P_i \log \frac{1}{P_i}$

\uparrow sub-optimal \uparrow Optimal

Is the above really the case? Will see

Back to Basics

- Notation : $P_{xy}(\alpha, \beta) \triangleq \Pr [X = \alpha, Y = \beta]$

"Distribution": $\sum_{\alpha, \beta} P_{xy}(\alpha, \beta) = 1$, $P_{xy}(\alpha, \beta) \geq 0$.

- P_x, P_y = marginals $P_x(\alpha) = \sum P_{xy}(\alpha, \beta)$.

- $P_{y|x=\alpha}$ = conditional dist. $P_{y|x=\alpha}(\beta) = \frac{P_{xy}(\alpha, \beta)}{P_x(\alpha)}$

- Conditional Entropy = Expected Entropy after conditioning

$$\begin{aligned}
 H(Y|X) &= \mathbb{E}_{\alpha} [H(Y|_{X=\alpha})] \\
 &= \sum_{\alpha} P_x(\alpha) \cdot H(Y|_{X=\alpha}) \\
 &= \sum_{\alpha} P_x(\alpha) \cdot \sum_{\beta} P_{y|x=\alpha}(\beta) \cdot \log \frac{1}{P_{y|x=\alpha}(\beta)} \\
 &= \sum_{\alpha, \beta} P_{xy}(\alpha, \beta) \cdot \log \frac{P_x(\alpha)}{P_{xy}(\alpha, \beta)}
 \end{aligned}$$

AXIOMS OF ENTROPY

$$\textcircled{1} \quad H(X) \leq \log |\Omega| \quad \text{if } X \in \Omega$$

Equality iff $X \sim \text{Unif}(\Omega)$

$$\textcircled{2} \quad H(X, Y) = H(X) + H(Y|X)$$

$$\textcircled{3} \quad H(Y|X) \leq H(Y) \quad \text{"Conditioning reduces Entropy"}$$

Will prove above today

$$\textcircled{2}: H(X, Y) = \sum_{\alpha, \beta} P_{X,Y}(\alpha, \beta) \log \frac{1}{P_{X,Y}(\alpha, \beta)}$$

$$H(X) = \sum_{\alpha} P_X(\alpha) \log \frac{1}{P_X(\alpha)} = \sum_{\alpha, \beta} P_{X,Y}(\alpha, \beta) \log \frac{1}{P_X(\alpha)}$$

$$H(Y|X) = \sum_{\alpha, \beta} P_{X,Y}(\alpha, \beta) \log \frac{P_X(\alpha)}{P_{X,Y}(\alpha, \beta)}$$

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

$$\Leftarrow \log \frac{1}{P_{X,Y}(\alpha, \beta)} = \log \frac{1}{P_X(\alpha)} + \log \frac{P_X(\alpha)}{P_{X,Y}(\alpha, \beta)}$$

Corollary: $H(X_1, \dots, X_n) \quad X_i \text{ i.i.d. } \sim X$
 $= n \cdot H(X).$

The inequalities:

- ① $H(x) \leq \log |\Omega|$
- ② $H(x) < \log |\Omega|$ if $P_x \neq \text{Unif}$
- ③ $H(Y|X) \leq H(X)$

All follow from the "optimality" of entropy.

Theorem:

\forall pair of distributions P, Q

$$E_{x \sim P} \left[\log \frac{1}{P(x)} \right] \leq E_{x \sim P} \left[\log \frac{1}{Q(x)} \right]$$

with equality iff $P=Q$.

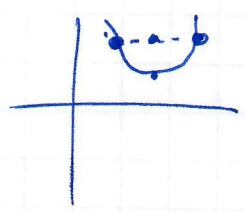
↑
 compressing according
 to "correct"
 dist.

↑
 compressing according
 to "wrong"
 distribution

Key ingredient in proof: JENSEN'S INEQUALITY

if f is convex then

$$E_{\mathbb{Z}} \left[f(\mathbb{Z}) \right] \geq f \left(E_{\mathbb{Z}} [\mathbb{Z}] \right)$$



if f is concave then

$$E_{\mathbb{Z}} [f(\mathbb{Z})] \leq f \left(E_{\mathbb{Z}} [\mathbb{Z}] \right)$$

• will apply to $f(z) = \log z$ (concave)

$$\& z = \frac{Q(x)}{P(x)} \quad x \sim P$$

$$\begin{aligned}
\text{Get: } \mathbb{E}_{x \sim P} \left[\log \frac{Q(x)}{P(x)} \right] &\leq \log \left(\mathbb{E}_{x \sim P} \left[\frac{Q(x)}{P(x)} \right] \right) \\
&= \log \left(\sum_x \frac{P(x) \cdot Q(x)}{P(x)} \right) \\
&= \log \sum_x Q(x) \\
&= \log 1 \\
&= 0
\end{aligned}$$

Conclude:

$$\mathbb{E}_{x \sim P} \left[\log \frac{1}{P(x)} \right] \leq \mathbb{E}_{x \sim P} \left[\log \frac{1}{Q(x)} \right]$$

↑
Compressing according to right dist.
↑
Compressing according to wrong dist.

Equality iff Equality in Jensen iff $Z = \text{constant}$ (for strictly concave f)
 iff $P(x) = Q(x) \forall x$.



KL Divergence

$$D(P \parallel Q) \triangleq E_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right]$$

Theorem just proved ^{says} ~~is~~ $D(P \parallel Q) \geq 0$ & $D(P \parallel Q) > 0$ iff $P \neq Q$.

Divergence \rightarrow nice : $D(P^n \parallel Q^n) = n \cdot D(P \parallel Q)$.

\rightarrow not so nice : Not a metric

$$D(P \parallel Q) \neq D(Q \parallel P)$$

$D(P \parallel Q)$ not necessarily finite

Operational meaning : loss per copy when compressing $X \sim P$ using mechanism for Q .

Using Divergence Theorem :

$$\begin{aligned} \textcircled{1} \quad H(X) &\leq \log |\Sigma| = E_{x \sim P} \left[\log |\Sigma| \right] \\ &= E_{x \sim P} \left[\log \frac{1}{Q(x)} \right] \quad Q(x) = \text{Unif}(\Sigma) \end{aligned}$$

$$\textcircled{2} \quad H(Y|X) \leq H(Y) \iff D(P_{XY} \parallel P_X \times P_Y) \geq 0$$

$$H(X, Y) \leq H(X) + H(Y)$$

\uparrow $\textcircled{3}$

$$\begin{aligned} \textcircled{2} \iff & E_{(X, Y) \sim P_{XY}} \left[\log \frac{1}{P_{XY}(x, y)} \right] \leq E_{(X, Y)} \left[\log \frac{1}{P_X \cdot P_Y} \right] \\ &= E \left[\log \frac{1}{P_X} \right] + E \left[\log \frac{1}{P_Y} \right] \end{aligned}$$

Mutual Information

$I(Y; X)$ = information in X about Y

$$\triangleq H(Y) - H(Y|X)$$

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

$$= I(X; Y) \quad \text{[symmetric]}$$

Chain Rule for information:

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

$$\begin{aligned} [I(X; Y|Z)] &\triangleq \mathbb{E}_Z [I(X|_{Z=z}; Y|_{Z=z})] \\ &\triangleq [H(X|Z) - H(X|Y, Z)] \end{aligned}$$

Data Processing Inequality

$$X - Y - Z \quad [(X \perp\!\!\!\perp Z) | Y]$$

$$\Rightarrow I(X; Z) \leq I(Y; Z).$$

Fano (In Prob 1):

$$H(X|Y) \leq h\left(\frac{e}{|S_X|}\right) + e \cdot \log |S_X|$$

$$e = \Pr[X \neq g(x)].$$