

LECTURE 3

Note Title

2/1/2016

TODAY

- Basic Quantities in Information Theory
 - Entropy, Conditional Entropy (already known)
 - Information, KL Divergence (new)
- Basic Relationships
 - Positivity of Information
 - Upper bound on Entropy
 - Convexity of Entropy
 - Data Processing Inequalities
 - Positivity of Divergence
 - Fano's Inequality ...

ERRATA: Proof of Kraft's Inequality (\Leftarrow)

l_1, \dots, l_n pattern of prefix-free C

$$\Leftarrow \sum_{i=1}^n 2^{-l_i} \leq 1$$

———— ✗ ————

Proof Claimed: Can just pick C greedily...

But clearly wrong

Example: $l_1 = l_2 = n$; $l_3 = 1$

$C(1) = 000000$; $C(2) = 111111$

$C(3) = ?$ Can't start with 0, Can't start with 1.

———— ✗ ————

Fix: Assume $l_1 \leq l_2 \leq l_3 \dots \leq l_n$

Now greedy works.

After assigning $C(i)$, claim # ^{alive} nodes at level l_i or deeper $= 1 - \sum_{j=1}^i 2^{-l_j} > 0$



Review

• $X \sim P_x = (p_1, \dots, p_k)$ $\Omega = \{1, \dots, k\}$

• $H(X) \triangleq \sum_{i=1}^k p_i \log \frac{1}{p_i} \triangleq \mathbb{E}_{x \sim P_x} \left[\log \frac{1}{P_x} \right]$

• $H(X|Y) \triangleq \mathbb{E}_{y \sim P_y} \left[H(X|Y=y) \right]$

• Axioms

(?) • $H(X) \leq \log k$ (not yet proved!)

• $H(U_\Omega) = \log k$ ($U_\Omega \triangleq$ uniform dist. on Ω)

(?) • $H(X|Y) \leq H(X)$ (not yet proved!)

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

New Concepts

(Mutual) Information:

Suppose X, Y jointly distributed

$I(X; Y) =$ Amount of information that Y contains about X

Intuitive Examples:

$$X \perp Y : I(X, Y) = 0$$

$$X = Y : I(X, Y) = H(X)$$

$$X = g(Y) : I(X, Y) = H(X)$$

$$Y = f(X) : I(X, Y) = H(Y)$$

$$\left. \begin{array}{l} X = (A, B) \\ Y = (A, C) \\ A \perp B \perp C \end{array} \right\} : I(X, Y) = H(A)$$

Formal Definition $I(X; Y) \triangleq H(X) - H(X|Y)$

$$I(X; Y|Z) \triangleq \mathbb{E}_{z \sim P_z} [I(X; Y|Z=z)]$$

Questions: • Is Information (as defined)
non-negative?

(follows from unproven axiom (3))

- Is there some monotonicity to Information wrt conditioning?

Example: $X \perp Y \in U\{0,1\}$
 $Z = X \oplus Y$

$$I(X; Y) = 0 \quad ; \quad I(X; Y | Z) = 1$$

$$I(X; Y | Z) > I(X; Y)$$

Example: $X = Y = Z \in U\{0,1\}$

$$I(X; Y) = 1$$

$$I(X; Y | Z) = 0$$

$$I(X; Y | Z) < I(X; Y)$$

Conclusion: No monotonicity.

Chain Rule For Information

$$\begin{aligned} I(x_1 \dots x_n; Y) \\ &= I(x_1; Y) + I(x_2; Y | x_1) \\ &\quad + \dots + I(x_n; Y | x_1 \dots x_{n-1}) \end{aligned}$$

Now: for some proofs.

Mother of all inequalities

• Let P, Q be distributions of Ω

then
$$\mathbb{E}_{x \sim P} \left[-\log \frac{Q(x)}{P(x)} \right] \geq 0$$

What? Why? How?

↑
Divergence Inequality

Intuition

$$\mathbb{E}_{x \sim P} \left[-\log \frac{1}{P(x)} \right] \approx \begin{array}{l} \text{(optimal)} \\ \wedge \text{ length of compressing} \\ X \sim P \end{array}$$

$$\mathbb{E}_{x \sim P} \left[-\log \frac{1}{Q(x)} \right] \approx \begin{array}{l} \text{length of some} \\ \text{compression scheme} \\ \text{for } X \sim P, \\ \text{one that pretends } X \sim Q \end{array}$$

Inequality asserts

"Some" Compression \geq "Minimum" Compression

Food for thought

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

Proof of Divergence Inequality

Key Ingredient: Jensen's Inequality

if $f: \mathbb{R} \rightarrow \mathbb{R}$ is nice & convex

$$E_x[f(x)] \geq f(E_x[x])$$

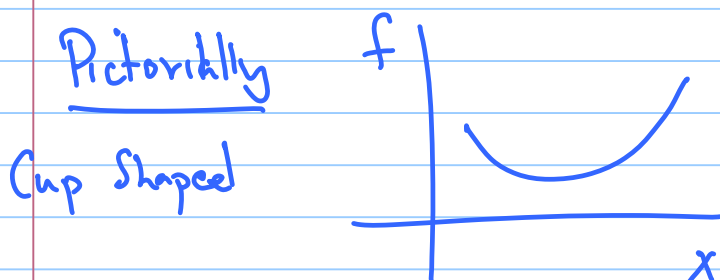
(measure on x s.t. $E_x[x]$ exists;
 $E_x[f(x)]$ exists etc.)

for us all true: x will have
finite support
 $f(x)$ bounded



Convex: $\forall x, y, \lambda \in [0, 1]$

$$f(\lambda x + (1-\lambda)y) \geq \lambda f(x) + (1-\lambda)f(y)$$



For us:

$-\log x$ is convex.

Divergence Inequality.

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

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

$$= -\log \left[\sum_x Q(x) \right]$$

$$= -\log 1 = 0$$



Using Divergence Inequality

$$\textcircled{1} \quad H(X) \leq \log \Omega$$

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

$$\text{Let } U(x) = \frac{1}{\Omega} \quad \forall x \in \Omega$$

$$\Leftrightarrow \mathbb{E}_{x \sim P} \left[-\log \frac{U(x)}{P(x)} \right] \geq 0 \quad \boxtimes$$

(Divergence between P & Uniform)

$$\textcircled{2} \quad H(X|Y) \leq H(X)$$

$$\Leftrightarrow H(X, Y) \leq H(X) + H(Y) \quad [\text{Chain Rule}]$$

$$\Leftrightarrow \mathbb{E}_{(x,y) \sim P_{xy}} \left[-\log \frac{1}{P_{xy}(x,y)} \right] \leq \mathbb{E}_{(x,y)} \left[-\log \frac{1}{P_x(x)} \right] + \mathbb{E}_{(x,y)} \left[-\log \frac{1}{P_y(y)} \right]$$

$$\Leftrightarrow \mathbb{E}_{(x,y) \sim P_{xy}} \left[-\log \frac{P_x(x) \cdot P_y(y)}{P_{xy}(x,y)} \right]$$

Apply Divergence Inequality to

$$P = P_{x,y} \quad \& \quad Q = P_x \times P_y$$

$$\Leftrightarrow \mathbb{E}_{x,y \sim P} \left[-\log \frac{Q(x,y)}{P(x,y)} \right] \geq 0 \quad \boxtimes$$

Kullback-Leibler Divergence aka Relative Entropy

$$D(P \parallel Q) \triangleq \mathbb{E}_{x \sim P} \left[-\log \frac{Q(x)}{P(x)} \right]$$

Key Measure of "distance between distributions"

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

$$D(P \parallel R) \leq D(P \parallel Q) + D(Q \parallel R) ?$$

(I think not ... exercise)

But key, key quantity! Why?

① Divergence Inequality: Miracle?

② Chain Rule:

$$D(P_{xy} \parallel Q_{xy}) = D(P_x \parallel Q_x) + D(P_{y|x} \parallel Q_{y|x})$$

$$\stackrel{\text{IID}}{\mathbb{E}}_{x \sim X} \left[D(P_{y|x} \parallel Q_{y|x}) \right]$$

in general cumbersome.

But for i.i.d variable

X_1, \dots, X_n

$X_i \sim P$?

or $X_i \sim Q$?

$$D(P^n \parallel Q^n) = n \cdot D(P \parallel Q).$$

$$P_x [x \sim P \text{ typical for } Q] \leq \exp(-D(P \parallel Q))$$

[Will see later]

Markov Chains & Data Processing Inequality:

Notation: $X \rightarrow Y \rightarrow Z$

implies $(X \perp Z) | Y$

conditioned on Y , X & Z independent

Data Processing Inequality

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

Fano's Inequality (Entropy \Rightarrow Unpredictability)

let $X \rightarrow Y \rightarrow \hat{X}$ & $P_e = \Pr[X \neq \hat{X}]$

then $H(P_e) + P_e \log |\Omega| \geq H(X|Y)$

$$\Rightarrow P_e \geq \frac{H(X|Y) - 1}{\log |\Omega|}$$

Proof: $E = 1$ if $\hat{X} \neq X$
 $= 0$ o.w.

$$\bullet H(E, X | \hat{X}) = H(X | \hat{X}) + \underbrace{H(E | X, \hat{X})}_0$$

$$\geq H(X | Y)$$

$$\bullet H(E, X | \hat{X}) = \underbrace{H(E | \hat{X})}_{\leq H(P_e)} + \underbrace{H(X | E, \hat{X})}_{P_e \cdot \log \Omega}$$

□