## Lecture 3 (Jan 22)

Exercise. $H(X|Y) = 0 \iff X = g(Y)$ for some func. $g(\cdot)$

## Chain rule for MI.

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

Proof:

$$\text{LHS} = H(X_1, \dots \to X_n) - H(X_1, \dots \to X_n | Y)$$

$$= \sum_{i=1}^{n} H(X_i | X_1, \to X_{i-1})$$

$$- \sum_{i=1}^{n} H(X_i | X_1, \to X_{i-1}, Y) = \checkmark.$$

M.I. under conditioning.

$$I(X \, ; \, Y) \quad vs. \quad I(X \, ; \, Y | Z) \quad ?$$

① if $Z = Y \Rightarrow I(X; Y|Z) = 0 \leq I(X; Y)$.

② $Z = X \oplus Y$ & $X \perp Y$ $\Rightarrow I(X; Y) = 0 < I(X; Y|Z) = 1$.

$$H(X_1, \dots, X_n) \le H(X_1) + \dots + H(X_n)$$
$$\& = \text{ if } X_i \text{ are independent.}$$

__Lemma:__ If $X_1, \dots, X_n$ are independent,

$$I(X_1, \dots, X_n ; Y) \ge \sum_{i=1}^{n} I(X_i ; Y)$$

LHS: $= H(X_1, \dots X_n) - H(X_1, \dots X_n | Y)$

$\overset{indep}{=} \sum_{i=1}^{n} H(X_i) \cancel{+ H(Y) -}$

$\cancel{(H(Y))} + H(X_1 | Y) + H(X_2 | X_1, Y) + \dots$
$$\qquad\qquad + H(X_n | X_1, \dots, X_{n-1}, Y))$$

$\overset{regroup}{=} \sum_{j=1}^{n} H(X_i) - \underbrace{H(X_i | X_1, \dots X_{n-1}, Y)}_{\le H(X_i | Y)}$

$\ge \sum_{i=1}^{n} I(X_i ; Y)$.

## Def.

### Relative Entropy / Information Divergence/

#### Kullback-Liebler Divergence:
#### (KL)‾

Let $p, q$ be distributions on $\mathcal{U}$, Then,

$$D(p \| q) \overset{def}{:=} \sum_{x \in \mathcal{U}} p(x) \log \frac{p(x)}{q(x)}.$$

⚠ Not symmetric!

\* Clearly, for $p = q \Rightarrow D(p \| q) = 0$.

\* $D(p \| q)$ ~~finite~~ $\Rightarrow Supp(p) \subseteq Supp(q)$.

\* <u>Gibbs inequality</u>: $D(p \| q) \geq 0$

$$\& = 0 \text{ iff } p = q.$$

<u>Proof.</u> (Concavity) ~~…~~

$$\left[ Z = \frac{q(x)}{p(x)} \text{ w.p. } p(x) \right]$$

$$-D(p \| q) = \sum_x p(x) \log \frac{q(x)}{p(x)} = \underset{x}{\mathbb{E}} \log Z$$

Jensen
$$\leq \log \mathbb{E}(Z) = \log 1 = 0.$$

$(\text{equality iff } Z = \mathbb{E}(Z) \text{ everywhere.})$

**Theorem.** $I(X;Y) = \mathcal{D}\left( p(x,y) \| p(x)p(y) \right)$

$\begin{cases} \text{in particular, if } X \perp Y \Rightarrow \mathcal{D}(\cdot \| \cdot) = 0. \checkmark \\ \\ \text{on the other hand: if } X = Y \Rightarrow \mathcal{D} \text{ is maximal.} \end{cases}$

**Proof.**

$$\underline{RHS.} = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

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

$$= \underbrace{\sum_{x,y} p(x,y) \log \frac{1}{p(x)}}_{H(X)} - \underbrace{\sum_{x,y} p(x,y) \log \frac{1}{p(x/y)}}_{H(X/Y)}$$

$$= I(X;Y). \qquad \square$$

Some viewpoints of on KL − divergence.

① Source Coding interpretation:

" Increase in expected encoding length
due to incorrect knowledge of the
distribution. "

② Rejection Sampling

$p, q \leftarrow$ Universe $U$.

<u>Setting</u>: Sequence $X_1, X_2, \ldots$   i.i.d. $\leftarrow q$.

<u>Goal</u>: Output $i^*$ s.t. $X_{i^*}$ is

distributed according to $p$.
(with no modification)

\* Say $p$ is uniform on $\{1, \to 50\}$
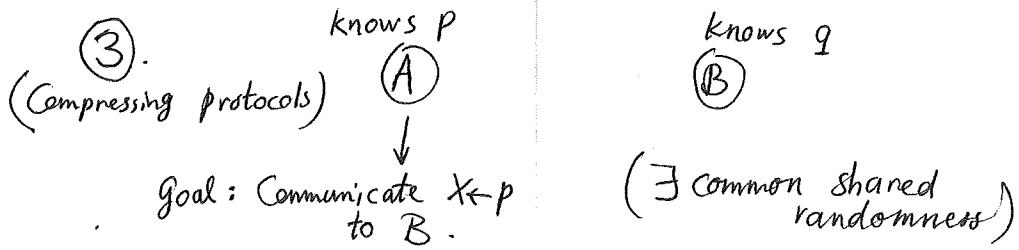
& $q$   "   "   " $\{1, \to 100\}$.

Then, Sample from $q$ until you hit $\{1, \to 50\}$.
but the inverse is not possible.

$\downarrow$

<u>Solution</u>: Output the first $i^*$ s.t. $X_i \leq 50$

$\Rightarrow E(\cancel{\text{max}} \, i^*) = 2, \quad E(\text{length } i) = 1.$

<u>Thm.</u> $\exists$ strategy that achieves

$E(\text{length } i^*) \leq D(p \| q).$

$\downarrow$ over $p$ & internal randomness of strategy.
(No proof now!)

③ (Compressing protocols)

knows p
Ⓐ

knows q
Ⓑ

↓

Goal: Communicate $x \leftarrow p$ to B.

($\exists$ common shared randomness)

**Thm.** $\exists$ interactive protocol between

A & B with Expected "$D(p \| q)$ + small"

bits of communication st. in the end,

① A outputs ~~answer~~ $a \leftarrow p$

② B outputs $b$ st.

$$\forall x, \quad Pr(b=x \mid a=x) \geq 1-\varepsilon.$$

(i.e, $D(\|)$ is the complexity of "agreement").

---

② & ③ are useful for Compression of protocols.

## Data processing inequality:

Suppose:
$$X, \quad Y, \quad g(Y) \text{ deterministic func.}$$

Then:
$$I(X; Y) \geq I(X; g(Y)).$$

Proof next time!