Today

- Properties of Information & Entropy.

- Tool: Convexity & Jensen's Inequality;
  Informational Divergence

- Results: - Positivity of Mutual Information

   - Data Processing theorem

   - Fano's Theorem (Prob. Error high
      if Conditional Entropy high)

Review of last time

## Entropy

$$X \in \{1 \ldots N\}$$

$$Pr[X=i] = P_i$$

$$H(X) = H(P_1 \ldots P_N) = \sum_{i=1}^{N} -P_i \log P_i$$

## Conditional Entropy

$$H(X|Y) = \sum_y P_Y(y) \sum_x P_r(X=x|Y=y) \cdot \log(\quad)$$

$$= \sum_y \sum_x P(x,y) \log \frac{P_Y(y)}{P(x,y)}$$

## Mutual Information

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

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

## Chain Rule of Entropy

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

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

Today: Some more complex properties

Two objectives:

- $I(x;y) \geq 0$ ?

- Conditional Entropy $\Rightarrow$ Unpredictability ?

$$\xrightarrow{\hspace{3cm}} \times \xrightarrow{\hspace{3cm}}$$

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

Question: $I(x;y) \geq 0$ ?

Why is $\sum_{x,y} P(x,y) \log \frac{P(x,y)}{P_x(x) \cdot P_y(y)} \geq 0$

$$\Updownarrow$$

$$E\left[ \log \frac{P(x,y)}{P_x(x), P_y(y)} \right] \geq 0$$

Expression is of the form

$$E_n\left[f(x)\right] \geq 0 \ ?$$

"Pattern matching" points to Jensen's Inequality

**Jensen's Inequality**

for every "Convex function" $f: R^n \to R$

$$E\left[f(x)\right] \geq f\left(E[x]\right)$$

What is a convex function?

$= x^2$ good example; Also $e^x$;

$\qquad \qquad \qquad \qquad \qquad \qquad -\log x;$

$= $ Calculus defn: $f: R \to R$ is convex

if $f''(x) \geq 0$

— More generally

$f: \mathbb{R}^n \to \mathbb{R}$ is convex if

$$f(\lambda \cdot x + (1-\lambda) \cdot y) \leq \lambda f(x) + (1-\lambda) f(y)$$

———————————⊱———————————

Jensen's Inequality: Trivial for distribution on
finite random-variables from defn.

But Calculus definition $\Rightarrow$ general defn ...
takes a little work      Won't do here

————————⊱————————

Implication
———————

$$\mathbb{E}[-\log z] \geq -\log \mathbb{E}[z]$$

for any distot $p, q$

Apply to      $Z = \dfrac{q(x)}{p(x)}$      ← r.v. $X$ with
                                              distl $p$

$$E\left[\log \frac{P(n)}{q(n)}\right] = E\left[-\log Z\right] \geqslant -\log E[Z]$$

$$= -\log \sum_{n} P(n) \cdot \frac{q(x)}{P(x)}$$

$$= \log 1$$

$$= 0$$

Conclude: $\forall$ pair of distributions $P, q$

$$E\left[\log \frac{q(x)}{P(x)}\right] \geqslant 0$$

When is $E\left[\log \frac{q(x)}{P(x)}\right] = 0$ ?

iff $\forall x, \quad \frac{q(x)}{P(x)} = \quad = 1$

i.e. if $q(x) = P(x)$

Motivates "Relative Entropy"

"(Informational / KL) Divergence"

Kullback-Liebler

$$D(p \| q) = E\left[ \log \frac{q(x)}{p(x)} \right]$$

- $D(p \| q) \geq 0$

- $D(p \| q) = 0 \quad \Rightarrow \quad p = q.$

- $I(x; y) = D\left( p(x, y) \| P_x(x) \cdot P_y(y) \right)$

- $H(x) = \log |x| - D(P \| U)$

(Introduction a little premature ----)

Later we will see that the
"best" compression of $x$ should
take $\approx \left| \log \frac{1}{P_x} \right|$ bits.

So if we compress $X$ thinking it
comes from dist $q$ will use

$$- \sum p(x) \log \frac{1}{q(x)} \quad \text{bits}$$

— Should have taken

$$\sum p(x) \log \frac{1}{p(x)} \quad \text{bits}$$

$$\text{Inefficiency} = \sum p(x) \log \frac{p(x)}{q(x)} = D(p \| q).$$

More on $D(p \| q)$

$$D(p \| q) \neq D(q \| p)$$

$D(p \| q)$ could be infinite

## Example

$X = 0$    w.p. 1        $X = 0$   w.p. $\frac{1}{2}$

$\phantom{X} = 1$    w.p. 0          $= 1$   w.p. $\frac{1}{2}$

$\underbrace{\hphantom{XXXXXXXX}}_{p}$           $\underbrace{\hphantom{XXXXXX}}_{q}$

one of $D(p \| q)$ / $D(q \| p)$ is infinite

$p$

Back to consequences :

①   $H(x) \leq \log |\mathcal{L}_x|$

since   $H(x) = \log |\mathcal{L}_x| - D(p \| u)$

① $H(X|Y) \leq H(X)$ [Conditioning reduces uncertainty]

② $H(X_1 \cdots X_n) \leq \sum_{i=1}^{n} H(X_i)$

s.t. $H(X_1 \cdots X_n) = \sum H(X_i | X_1 \cdots X_{i-1})$

$\leq \sum H(X_i)$

---⟑---

## Concavity of Entropy

$H(\lambda \cdot p + (1-\lambda) \cdot q)$

$\geq \lambda \cdot H(p) + (1-\lambda) H(q)$

$X \sim p$ ; $Y \sim q$ ; $b = 0$ w.p. $\lambda$
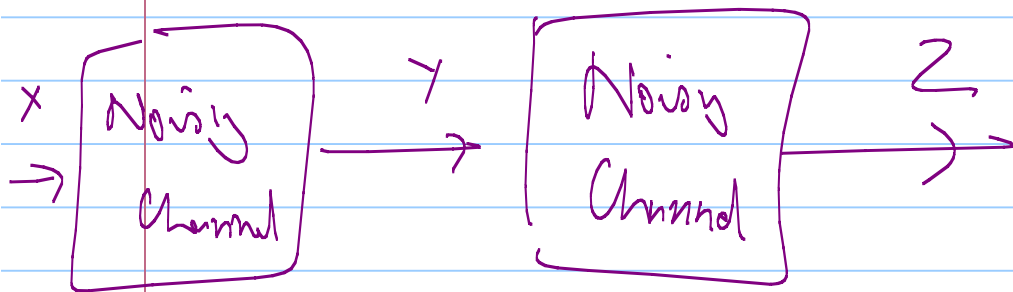$\phantom{X \sim p ; Y \sim q ; b = 0} 1$ o.w.

$Z = X$ if $b = 0$ $\&$ $Y$ if $b = 1$

$$H(Z) = H(\lambda \rho + (-\lambda) \ell) \geq H(Z|b)$$

$$= \lambda H(x) + (1-\lambda) H(\ell)$$

## Data Processing Theorem



would expect

$$I(x;z) \leq I(x;y)$$

To formulize : " Markov Chains

$X, Y, Z$     form    M.C    $(X \to Y \to Z)$

if

$$P_{Z|(X,Y)}(z \mid (x,y)) = P_{Z|Y}(z \mid y)$$

$$\underrightarrow{\qquad\qquad}_{\phantom{P}} \underline{\qquad}$$

**Equivalently**

$$P_{(X,Z)|Y}((x,z) \mid y) = P_{X|Y} \cdot P_{Z|Y}$$

$$X \to Y \to Z \qquad (\Leftrightarrow) \qquad Z \to Y \to X$$

$$\underline{\qquad\qquad}_Y \underline{\qquad}$$

$$I(X; (Y,Z)) = I(X; Z) + I(X; Y|Z)$$

$$= I(X; Y) + \underbrace{I(X; Z|Y)}_{\ddot\smile}$$

$$\Rightarrow \quad I(x; y) = I(x; z) + I(x; y|z)$$

$$\geq I(x; z)$$

Finally .... Fano's Lemma



$$P_e = P_r\left[X \neq \tilde{X}\right] \quad ; \quad I : \text{indicator of}$$
$$\text{event}$$
$$\tilde{X} \neq X$$

Fano's Lemma:

$$H(I) + P_e \cdot \log\left(|\Omega_x| - 1\right) \geq H(X|Y)$$

**Corollary**   $P_e \geq \dfrac{H(x|y) - 1}{\log |\Omega_x|}$

**Proof:**

$$H((I, X)|Y) = H(x|y) + H(I|(x,y))$$
$$= H(I|y) + H(x|(I,y))$$
$$\leq H(I) + H(x|(I,y))$$

$$H(x|(I,y))$$

$$= P_e \cdot H(x|(I=1, y))$$
$$+ (1-P_e) \cdot \underbrace{H(x|(I=0,y))}_{=0}$$

$$= P_e \cdot H(x|(I=1, y))$$

$$\boxed{H(X \mid (I=1, Y)) \leq \log(|\Omega_x| - 1)}$$

$$\leq P_e \cdot \log(|\Omega_x| - 1)$$

---

## Application :

$$X = n \quad \text{random bits} = X_1 \ldots X_n$$

$$Y_i = X_i \quad \text{w.p} \quad 1-p$$
$$= \bar{X_i} \quad \text{w.p.} \quad p$$

$$H(X \mid Y) = n \cdot H(p)$$

Fano's Lemma :     $P_e \geq \dfrac{n \cdot H(p) - 1}{n}$

$$\simeq H(p)$$