# 1) Intro to information theory

## 1.1 Random Vars

Jensens inequality:
For $\mathcal{F}$ convex, ie $\mathcal{F}(\alpha x + (1-\alpha)y) \leq \alpha \mathcal{F}(x) + (1-\alpha)\mathcal{F}(y)$

we have $\mathbb{E}\,\mathcal{F}(X) \geq \mathcal{F}(\mathbb{E}X)$

e.g. $\langle x^2 \rangle \geq \langle x \rangle^2$

$\mathbb{E}\mathcal{F}$
$\mathcal{F}\mathbb{E}$

## 1.2 Entropy

Entropy $= \underset{x \sim p}{\mathbb{E}} \log \frac{1}{p(x)}$

KL Div, $D_{KL}(q\|p) = \sum_x q(x) \log \frac{q(x)}{p(x)} = \underset{q}{\mathbb{E}} \log\left(\frac{1}{p}\right)$

$D_{KL}$ is not symmetric

### From Rezende's notes:

We generally take $q_\theta$ the density of the generative model and $p$ to be the "true" density

generally want to modify $\theta$ so that $q_\theta \to p$

$KL(q\|p)$ is # of bits to communicate $q$ given that the receiver knows $p$

KL is unique divergence satisfying

i) locality:

*possible to add dependence on $\nabla p$, $\nabla q$ etc*

$$D(q\|p) = \int dx \; \mathcal{S}(q, p, x)$$

ii) invariance:
under $x \to x' = \varphi(x)$ we have $D$ invariant

$$\Rightarrow \int dx \, \mathcal{F}(q, p, x) = \int dx' \, \mathcal{F}\left(\frac{q \circ \varphi^{-1}(x')}{|\det \frac{\partial x'}{\partial x}|}, \frac{p \circ \varphi^{-1}(x')}{|\det \frac{\partial x'}{\partial x}|}, \varphi^{-1}(x')\right)$$

$\Rightarrow \mathcal{F}$ must take the form $\mathcal{F}\left(\frac{q}{p}\right)\mu(x) \Rightarrow \mathcal{F}\left(\frac{q}{p}\right)p$ or $\mathcal{F}\left(\frac{q}{p}\right)q$

<span style="color:pink">transforming as a measure</span>

iii) Subsystem independence (ie additivity of indep sub-domains)

$$\Rightarrow \quad \mathcal{F}\left(\frac{q}{p}\right) = \log\frac{q}{p}$$

## Back to Montanari

Entropy $H$ satisfies

1) $H_x \geq 0$     <span style="color:pink">proof: $E - \log p = -\log E p \geq 0$</span>

2) $H_x = 0$ only for $p(x) = \delta_x$

3) Among all distributions $p(x)$   $H$ is maximized for $p = \frac{1}{M}$

     <span style="color:pink">proof: $D_{KL}(p \mid \bar{p}) = \log_2 M - H(p) \geq 0$</span>

     <span style="color:pink">↑ uniform</span>

     <span style="color:lightblue">Q: can I get stronger bounds from other $\bar{p}$?</span>

4) For $X, Y$ indep   $H_{X,Y} = H_X + H_Y$

5) For $X, Y$ generic   $H_{X,Y} \leq H_X + H_Y$

6) For $X_1, X_2$ disjoint take   $q_{1,2} = $ Prob $x \in X_{1,2}$   resp.

     then $H_x = H(q) + \tilde{H}(q, r)$

     <span style="color:pink">$-q_1 \log q_1 - q_2 \log q_2 \quad -q_1 \sum_{x \in X_1} r^{(x)} \log r_1^{(x)} - q_2 \sum_{x \in X_2} r_2^{(x)} \log r_2^{(x)}$</span>

## <span style="color:lightblue">1.3 Sequences of random variables</span>

Def entropy rate   $h_x = \lim\limits_{N \to \infty} H[X_1 \cdots X_N]/N$

e.g. 1: $X_t$ indep $\Rightarrow P_N(X_1, \cdots, X_N) = \prod\limits_{t=1}^{N} p(x_i) \Rightarrow h_x = H(p)$

e.g. 2: <span style="color:pink">Markov Chain</span>

$\{p_t(x), x \in X\}$ on initial state
$\{w(x \to y)\}_{x,y \in X}$ are transition probabilities,   $\sum\limits_Y w(x \to y) = 1$

$\Rightarrow P_N(X_1, \cdots, X_N) = p_1(x_1) \prod\limits_{t=1}^{N-1} w(X_t \to X_{t+1})$    $\lim\limits_{t \to \infty} p_t(x) = p^*(x)$

then $h_x = -\sum\limits_x p^*(x) \sum\limits_Y w(x \to y) \log w(x \to y) = H_{Y|X}$ ← ie sum over all
     <span style="color:lightblue">$p(y|x)$ ↗</span>      letters weighted by $p(x)$
     and use entropy $H(x_{t+1|t})$

## 1.4 Correlated vars & Mutual info

**Conditional entropy**

$$H_{Y|X} := -\sum_x p(x) \sum_y p(y|x) \log p(y|x) \qquad \text{no log on } p(x)!$$

N.B. 
$$H_{X,Y} = -\sum_{x,y} p(x,y) \log(p(x,y)) = -\sum_{x,y} p(x) p(y|x) \log p(x) p(y|x)$$

$$= H_{Y|X} + \sum_x p(x) \log p(x) \sum_y p(y|x) = H_{Y|X} + H_X$$

**Mutual info:**

$$I_{X,Y} = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

$$\overset{\text{"}}{D_{KL}(p(x,y)||p(x)p(y))}$$

"reduction in uncertainty" of $X$ | knowledge of $Y$

alt: $H_{X|Y} = H_X - I_{X,Y}$

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

"the decrease in $Y$'s entropy from conditioning on $X$"

$$= H_X - H_{X|Y}$$

$$I_{X,Y} = \mathbb{E}_{x,y}\left[-\log \frac{p(x)p(y)}{p(x,y)}\right] \geq -\log \mathbb{E}_{x,y}\left[\frac{p(x)p(y)}{p(x,y)}\right] \geq 0$$

$$\uparrow \int p(x) \int p(y)$$

**Data processing inequality:** For Markov chain $X \to Y \to Z$

$$\Rightarrow p(x,y,z) = p_1(x) w_2(x \to y) v_3(y \to z)$$

Lemma: $I_{X,(YZ)} = I_{X,Z} + I_{X,Y|Z}$

here $I_{X,Y|Z} = -\mathbb{E}_{X,Y,Z} \log \frac{p(x|z)p(y|z)}{p(x,y|z)}$

$\uparrow p(x,y,z)$

$$-\mathbb{E}_{x,y,z} \log \frac{p(x) p(y,z)}{p(x,y,z)} = -\mathbb{E}_{x,z} \log \frac{p(x)p(z)}{p(x,z)} - \mathbb{E}_{x,y,z} \log \frac{p(x,z)p(y,z)}{p(x,y,z)p(z)}$$

$$= I_{X,Z} + I_{X,Y|Z} \checkmark$$

$$\Rightarrow I_{X,(YZ)} = I_{X,Z} + \overset{\geq 0}{I_{X,Y|Z}} \qquad \Rightarrow I_{X,Z} \leq I_{X,Y}$$

$$= I_{X,Y} + \overset{\geq 0}{I_{X,Z|Y}} \text{ By Markov}$$

Take $Z = f(Y) \Rightarrow I_{X,Y} \geq I_{X,f(Y)}$

Fano's inequality: Relates the info loss in a noisy channel to the probability of mischaracterization error

take $X \to Y \to \hat{X}$  w/ $\hat{X} = g(Y)$ an estimate of $X$

Let $E = \mathbb{1}_{X \neq \hat{X}}$ ,  $P_e = Pr(X \neq \hat{X}) = \mathbb{E}(E)$

$H_{X,E|Y} = H_{X|Y} + H_{E|X,Y}$  i) $H_{E|X,Y} = 0$  E is deterministic function of $X,Y$

$= H_{E|Y} + H_{X|E,Y}$  ii) $H_{E|Y} \leq H_E = \mathcal{H}(P_e)$

$\downarrow$
$\leq H_E$  iii) $H_{X|E,Y} = (1-P_e) H_{X|E=0,Y} + P_e H_{X|E=1,Y}$  $\hat{X}$ is $g(Y)$

$= P_e H_{X|E=1,Y} \leq P_e \log(|X|-1)$

$H_{X|Y} = H_{E|Y} + H_{X|E,Y} \leq H_E + P_e H_{X|E=1,Y} \leq \mathcal{H}(P_e) + P_e \log(|X|-1)$

$\uparrow$
bound on uncertainty of $X|Y$

$\leq$ Uncertainty of $X \neq \hat{X}$  ie $P_e$
$+$
$P_{error} \cdot \mathcal{H}(\text{uniform} -1)$

Exercise 1.6  $p(1) = 1-p$   For $k$ values
$p(x) = \dfrac{p}{k-1}$

take $Y$ indep of $X$ $\Rightarrow$ $H(X|Y) = H(X)$

$\Rightarrow \mathcal{H}(P_e) + P_e \log(k-1) \geq H(X)$

if $p$ small so $1-p > \dfrac{p}{k-1}$ guess 1 always

$\Rightarrow P_{error} = p$ $\Rightarrow -p \log p - (1-p) \log(1-p) + p \log(k-1) \leq H(X)$

$H(X) = -(1-p) \log(1-p) - (k-1) \cdot \dfrac{p}{k-1} \log \dfrac{p}{k-1}$

$\Rightarrow$ Equality

## 1.5 Data Compression

Sequence $\underline{X} = \{X_1 \cdots X_N\}$ for $X_i \in \mathcal{X}$ finite alphabet

assume $X_i$ are random

store a given realization $\underline{x} = \{x_1 \cdots x_N\}$ as compactly as possible

$$w: \quad \mathcal{X}^N \to \{0,1\}^*$$
$$\underline{x} \to w(\underline{x})$$

Often we take a longer stream $\to$ blocks $\underline{x}^1 \cdots \underline{x}^r$

encode each block $w(\underline{x}^1) \cdots w(\underline{x}^r)$

need concatenation of blocks to be uniquely decodable

safe if $\forall x, x'$ $w(x)$ is not prefix of $w(x')$

"instantaneous codes"

$$L(w) = \mathbb{E}_{\underline{x} \in \mathcal{X}^N} \ell_w(\underline{x}) \quad \leftarrow \text{ length of } w(\underline{x})$$

take $N=1$

$\mathcal{X} = \{1, \cdots, 8\}$

$p(i) = 2^{-i}$   $i = 1 \cdots 7$
$p(8) = 2^{-7}$   $i = 8$

| $x$ | $p(x)$ | $w_1(x)$ | $w_2(x)$ |
|---|---|---|---|
| 1 | $\frac{1}{2}$ | 000 | 0 |
| 2 | $\frac{1}{4}$ | 001 | 10 |
| 3 | $\frac{1}{8}$ | 010 | 110 |
| 4 | $\frac{1}{16}$ | 011 | 1110 |
| 5 | $\frac{1}{32}$ | 100 | 11110 |
| 6 | $\frac{1}{64}$ | 101 | 111110 |
| 7 | $\frac{1}{128}$ | 110 | 1111110 |
| 8 | $\frac{1}{128}$ | 111 | 11111110 |

$L(w_1) = 3$

both instantaneous

$$L(w_2) = \sum_{i=1}^{7} 2^{-i} i + 8 \cdot 2^{-7} \approx 2$$



← binary tree where no codeword node has ancestor codewords

What is best $w$ for a given source?

let $L_N^*$ be optimal achievable instantaneous code length. Then,

1. $\quad H_{\underline{x}} \preceq L_N^* \preceq H_{\underline{x}} + 1$

2. If the source has finite entropy rate $h = \lim_{N \to \infty} \frac{1}{N} H_{\underline{x}}$

$$\lim_{N \to \infty} \frac{1}{N} L_N^* = h$$

Lemma "Kraft's inequality"

$$\sum_{\underline{x} \in \mathcal{X}^N} 2^{-l_w(\underline{x})} \leq 1$$

Follows from "set of all leaves of binary tree sum to one"

Conversely any set of lengths $\{l_w(\underline{x})\}_{\underline{x} \in \mathcal{X}^N}$ satisfying Kraft have a code

→ start from smallest $l_w(\underline{x})$ and take first binary seq. of that length.

Goal: Find codewords $l_w^*(\underline{x})$ that minimize $L$
subject to Kraft

First, if $l$ could be real-valued

$$\min_{\ell, \, \alpha \geq 0} \sum_x p(x) \ell(x) + \alpha \left( \sum_x 2^{-\ell(x)} - 1 \right)$$

→ $p(x) - \alpha \, 2^{-\ell} \log_2 = 0$

$\ell = -\log_2 p(x) - c \qquad c = 0$ from Kraft

⟹ $\ell = \lceil -\log_2 p(x) \rceil \qquad$ also then works

$$H_x \preceq L \preceq H_x + 1 \qquad \checkmark$$

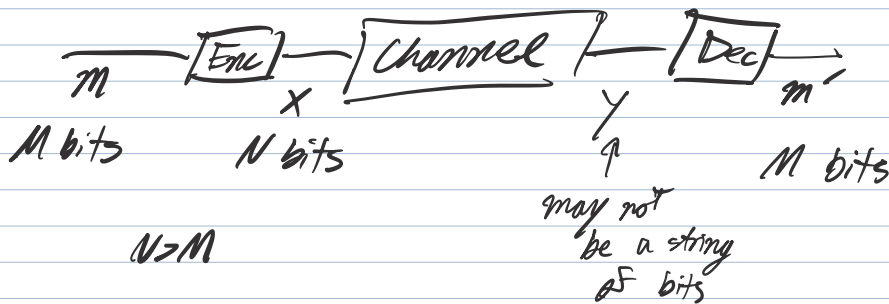"Shannon code", close to optimal for long strings

Not ideal for shorter sequences
↑ there Huffman coding is optimal
may assign super long
codeword when shorter ones
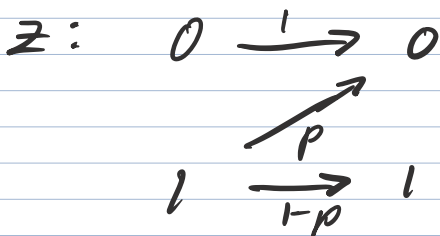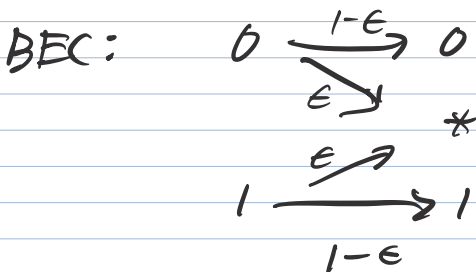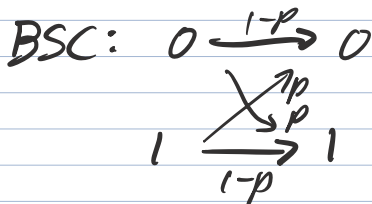are available

requires $\Theta(|X|^M)$ memory
to enumerate all $\ell(\underline{x})$

## 1.6 Data Transmission



$$M \longrightarrow \boxed{Enc} \longrightarrow \boxed{Channel} \longrightarrow \boxed{Dec} \longrightarrow m'$$

$m$  
$X$  
$Y$  
$m'$

M bits    N bits    M bits

N>M

may not
be a string
of bits

Can have a channel with insertions

Consider a memoryless channel (noise acts indep on each bit)

$$Q(\underline{y}|\underline{x}) = \prod_{i=1}^{N} Q(y_i|x_i)$$

BSC:



$$0 \xrightarrow{1-p} 0$$
$$\nearrow p \searrow p$$
$$1 \xrightarrow{1-p} 1$$

BEC:



$$0 \xrightarrow{1-\epsilon} 0$$
$$\searrow \epsilon$$
$$*$$
$$\nearrow \epsilon$$
$$1 \xrightarrow{1-\epsilon} 1$$

Z:



$$0 \xrightarrow{1} 0$$
$$\nearrow p$$
$$1 \xrightarrow{1-p} 1$$

Channel capacity C:

$$\boxed{\max_{p(x)} \; I_{X,Y}}$$

← reduction in uncertainty of $Y$ | knowledge of $X$, vice versa

We will see $C$ characterizes amount of info that can be transmitted faithfully through the channel

E.g. BSC, send a bit drawn from Bern$(q)$

$$\max_{q} \; I_{X,Y} = \sum_{x=\{0,1\}} p(x) \sum_{y=\{0,1\}} p(y|x) \log \frac{p(y|x)}{p(y)}$$

$$p(y=1) = p(y=1|x=1)p(x=1) + p(y=1|x=0)p(x=0)$$
$$= (1-p)(1-q) + p\,q$$

$$p(y=0) = (1-p)q + p(1-q)$$

$$\Rightarrow I_{X,Y} = q \cdot \left[ (1-p) \log \frac{1-p}{(1-p)q+p(1-q)} + p \log \frac{p}{(1-p)(1-q)+pq} \right]$$
$$+ (1-q) \cdot \left[ p \log \frac{p}{(1-p)q + p(1-q)} + (1-p) \log \frac{1-p}{(1-p)(1-q)+pq} \right]$$

We see $D_q I_{X,Y} = 0$ when $q = \frac{1}{2}$

Faster way
$$H(Y) - H(Y|X) = H((1-p)(1-\alpha)+p\alpha) - H(p)$$

$$\partial_\alpha = 0 \Rightarrow (2p-1)\log\frac{1-p}{p} \Rightarrow p = 0$$

$$\Rightarrow (2p-1)\alpha = p-1 \Rightarrow \alpha = \frac{1}{2}$$

$$\Rightarrow \boxed{C = H(\tfrac{1}{2}) - H(p) = 1 - H(p)}$$

N.B.
$$\partial_p H(p) = \log_2 \frac{1-p}{p}$$

**E.g.  BEC**

$$I_{XY} = q\left[(1-\epsilon)\log\frac{1-\epsilon}{q(1-\epsilon)} + \epsilon \log\frac{\epsilon}{\epsilon}\right]$$

$p(Y=0) = q(1-\epsilon)$

$p(Y=1) = (1-q)(1-\epsilon)$

$$+ (1-q)\left[(1-\epsilon)\log\frac{1-\epsilon}{(1-q)(1-\epsilon)} + \epsilon \log\frac{\epsilon}{\epsilon}\right]$$

$p(Y=*) = \epsilon$
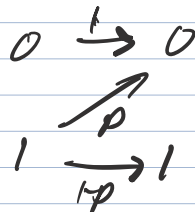
$D_q\, I_{XY} = 0$  when  $q = \frac{1}{2}$

Faster way: $H(Y) - H(Y|X) = \mathcal{H}(\epsilon) + (1-\epsilon)\mathcal{H}(\alpha) - \mathcal{H}(\epsilon)$

$H(Y) = H(Y \overset{?}{\in} *) + \sum p(*?)\mathcal{H}(Y|\overset{?}{*})$     $\partial_\alpha = 0 \Rightarrow (1-\alpha)\log_2\frac{1-\alpha}{\alpha} \Rightarrow \alpha = \frac{1}{2}$

$\quad = \mathcal{H}(\epsilon) + (1-\epsilon)\mathcal{H}(q)$

$$\Rightarrow \boxed{C = 1-\epsilon}$$

**Eg  Z-channel**

$\alpha = P(0)$   $P(Y=0) = p + \alpha(1-p)$   $P(Y=1) = (1-\alpha)(1-p)$

$\overset{\alpha + p(1-\alpha)}{}$

$0 \overset{1}{\longrightarrow} 0$   $\max_\alpha \{H(Y) - H(Y|X)\} = \max_\alpha H(Y) - \sum_x H(Y|X=x)P(x)$

$\quad\overset{p}{\nearrow}$     $= \max_\alpha \mathcal{H}((1-\alpha)(1-p)) - \alpha\cdot\mathcal{H}(Y|X=0) - (1-\alpha)\mathcal{H}(p)$

$1 \overset{}{\underset{1-p}{\longrightarrow}} 1$   $\partial_\alpha = 0 \Rightarrow -(1-p)\log\frac{1-(1-\alpha)(1-p)}{(1-\alpha)(1-p)} + H(p) = 0 \Rightarrow \frac{1}{p}-1 = 2^{H(p)/1-p}$

$$\Rightarrow \alpha = 1 - \frac{1}{(1-p)(1+2^{H(p)/1-p})}$$



$$C = \mathcal{H}\left(\frac{1}{1+2^{s(p)}}\right) - \frac{s(p)}{1+2^{s(p)}} = \log(1+2^{-s(p)}), \quad s(p) = \frac{H(p)}{1-p}$$

$$= \log(1+(1-p)p^{p/1-p})$$

Assume each bit is random —surprisingly, shannon's theorem
shows that there is no loss
in generality

$\{0,1\}^m \ni m \rightarrow \underline{x}(m) \in \{0,1\}^N$

$2^m$ codewords in $\mathbb{F}_2^N$

$Q(\underline{y}|\underline{x}) = \prod_i Q(y_i|x_i)$

$R = \frac{m}{N}$  is the rate

$$P_B(m) = \sum_{\not y} Q(\not y \mid \underline{x}(m)) \; \mathbb{I}(d(\not y) \neq m)$$

$$P_B^{max} = \max_m P_B(m) \qquad \text{``worst case''}$$

$$P_B^{av} = \frac{1}{2^m} \sum_{m \in \{0,1\}^m} P_B(m) \quad \leftarrow \text{more common}$$

Eg. 1 Repetition    k (odd) times
                         + majority

$$R = \frac{1}{k}$$

Exercise:

$$P_B^{av} = \sum_{r = \lceil \frac{k}{2} \rceil}^{k} \binom{k}{r} p^r (1-p)^{k-r}$$

Shannon, 1948

    For every rate $R < C$, there is a
sequence of codes $C_N$ of length $N$
s.t. :    $R_N \to R$    $P_B^{avg} \to 0$    as $N \to \infty$
conversely, any such sequence has $R < C$

Intuition for the role of capacity

$$H_{\underline{y} \mid \underline{x}} = N H_{y \mid x} \implies 2^{N H_{y \mid x}} \text{ outputs}$$

    need $d(y)$ to map all of them to $m$

    # possible outputs is $N H_y$

$\implies$ can distinguish $2^{N H_y} / 2^{N H_{y \mid x}}$ codewords

$$= 2^{N(H_y - H_{y \mid x})} = 2^{N I_{x,y}}$$

one needs to be able to send all $2^M$ codewords

$$\implies 2^M = 2^{NR} < 2^{N I_{x,y}}$$

$$\implies R < I_{x,y} \leq C$$

This also gives another interp of $I_{x,y}$

can distinguish $2^{N I_{x,y}}$ codewords

## Facts about channel coding:

For $p_1, p_2$ indep channels

$$(p_1 \times p_2)((y_1, y_2) | (x_1, x_2)) = p_1(y_1 | x_1) p_2(y_2 | x_2)$$

$$C(p_1 \times p_2) = \sup_{P_{x_1, x_2}} I(X_1, X_2 ; Y_1, Y_2) = \sup_{P_{x_1, x_2}} I(X_1, Y_1) + I(X_2, Y_2)$$

$$\geq C(p_1) + C(p_2)$$

Also

$$\sup \quad I(X_1, X_2 ; Y_1, Y_2) = H(Y_1, Y_2) - H(Y_1, Y_2 | X_1, X_2)$$

$$= H(Y_1, Y_2) - H(Y_1 | X_1) - H(Y_2 | X_2)$$

$$\leq H(Y_1) + H(Y_2) - \quad ''$$

$$\sup = I(X_1 ; Y_1) + I(X_2 ; Y_2)$$

$$\Rightarrow C(p_1 \times p_2) \leq C(p_1) + C(p_2)$$

$$\Rightarrow C(p_1 \times p_2) = C(p_1) \times C(p_2)$$