

Info Lect. 1

Measures of Information

ex:

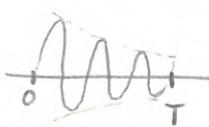


chirp \equiv modulated signal

$$s(t) = a(t) \cdot \cos(2\pi f_0 t)$$

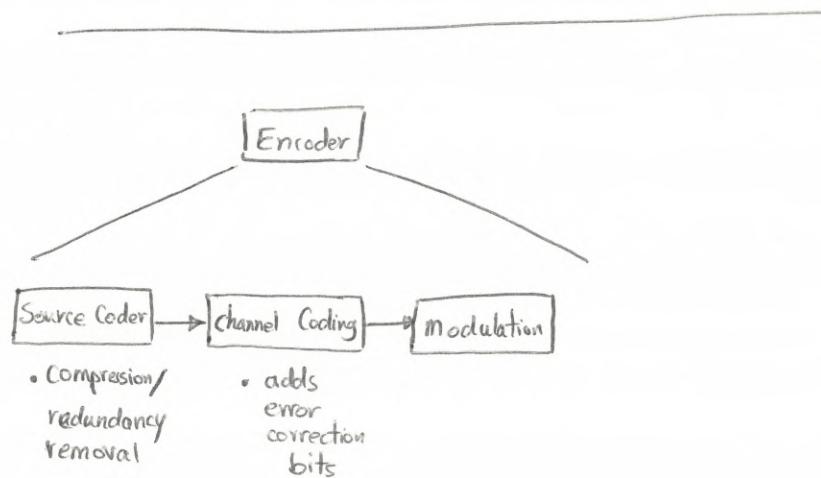
Received signal

$$r(t) = C \{ a(t-t_0) \cos(2\pi(f_0 + f_D)(t-t_0)) \} + w(t)$$



receiver
noise

⊕ radio
interferences



* Joint source-channel coding is an OPEN problem.

3 INFORMATION Measures

1. Entropy limit of lossless compression

2. Mutual Information distortions in the channel

3. Relative Entropy cost of using wrong pmf in encoding, related to Hypothesis Testing

Entropy:

x is a discrete R.V.

\mathcal{X} : all possible outcomes, finite number of elements

$p(x)$

$-\log_2 p(x)$:
+ uncertainty of observing x
+ Description Length of outcome x (in bits)

$$H(X) = \sum_{x \in \mathcal{X}} p(x) \log \frac{1}{p(x)}$$

ex. $X \sim \text{Bernoulli}(p)$

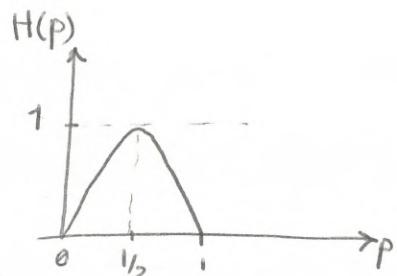
$$\begin{cases} 1, & p \\ 0, & 1-p \end{cases} \Rightarrow H(X) = \underbrace{H(p)}_{\substack{\text{Binary} \\ \text{Entropy}}} = p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p}$$

Properties:

1. $H(x) \geq 0$
2. If $X \sim \text{Bernoulli}(p) \Rightarrow 0 \leq H(p) \leq 1$

3. $H_b(x) = \sum p(x) \log_b p(x)$
 $= (\log b) H_a(x)$

$0 \leq H(x) \leq \log_2 n$ number of outcomes



4. Entropy is a functional. (Function of Function: $H(p(x))$)

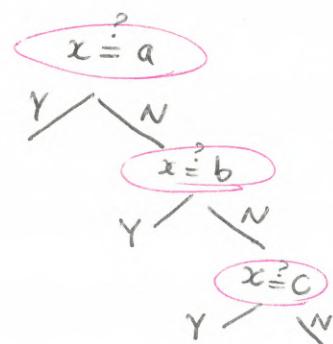
5. $H(X)$ is concave in $p(x)$.

6. n letter alphabet (power of 2)

Ex.

a	b	c	d
$1/2$	$1/4$	$1/8$	$1/8$

$$\begin{aligned} H(X) &= 1 \times \frac{1}{2} + 2 \times \frac{1}{4} \\ &\quad + 3 \times \frac{1}{8} + 3 \times \frac{1}{8} \\ &= \frac{7}{4} \text{ [bits]} \end{aligned}$$



3 questions
needed!

Joint Entropy of $(X, Y) \sim P(x, y)$

- $-\log P(x, y)$

- $H(X, Y) = \sum_{x, y} P(x, y) \log \frac{1}{P(x, y)}$

$$H(X|Y=y) = - \sum_x p(x|y) \log p(x|y)$$

$$H(X|Y) = \sum_y H(X|Y=y) p(y)$$

$$= - \sum_y \sum_x p(x,y) \log \frac{p(x|y)}{\frac{p(x,y)}{p(y)}}$$

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

$$\Rightarrow \boxed{H(X,Y) = H(Y) + H(X|Y)}$$

OR

$$\boxed{H(X,Y) = H(X) + H(Y|X)}$$

2. Relative Entropy

$$X \quad p_X(x)$$

$$q_X(x)$$

$$D(p||q) \triangleq \sum_x p(x) \log \frac{p(x)}{q(x)}$$

"distance" between p and q

- not symmetric $D(p||q) \neq D(q||p)$

- not satisfy triangular inequality $D(p||q) \not\leq D(p||r) + D(r||q)$

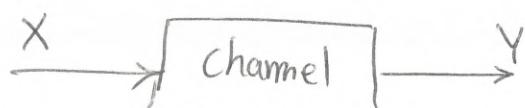
3. Mutual Information (Description of Channel)

X, Y are 2 R.V.'s

$$p(x, y)$$

$$I(X;Y) \triangleq \sum_x \sum_y p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

- How Much Information Y has about X or X has about Y



- $I(X;Y) = D(p(x,y)||p(x)p(y))$

$$\bullet \quad I(X;Y) = \sum_y \sum_x p(x,y) \log \frac{p(y|x)p(x)}{\cancel{p(x)p(y)}}$$

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

$$\bullet \quad D > 0 \Rightarrow I(X;Y) > 0 \Rightarrow H(Y) \geq H(Y|X)$$

Info. Let 3

Entropy: limit of lossless compression

Review:

X : discrete R.V

X : alphabet

$$H(X) = - \sum_{x \in X} p(x) \cdot \log_2 p(x) = E_{x \sim P(X)} \left[\log_2 \frac{1}{p(x)} \right]$$

- $H(X) \geq 0$
- concave in P

Relative Entropy

X

$p(x)$

$q(x)$



symmetry



△ ineq.

$$D(p \parallel q) = \sum_x p(x) \log_2 \frac{p(x)}{q(x)}$$



symmetry

$$D(p \parallel q) + D(q \parallel p) \equiv \text{Kullback - Liebler Divergence}$$



△ ineq.

symmetrical (✓)

$$\cdot D(p \parallel q) \geq 0$$

- concave in p and also in q

Mutual Information: description of channel / noise and distortion

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

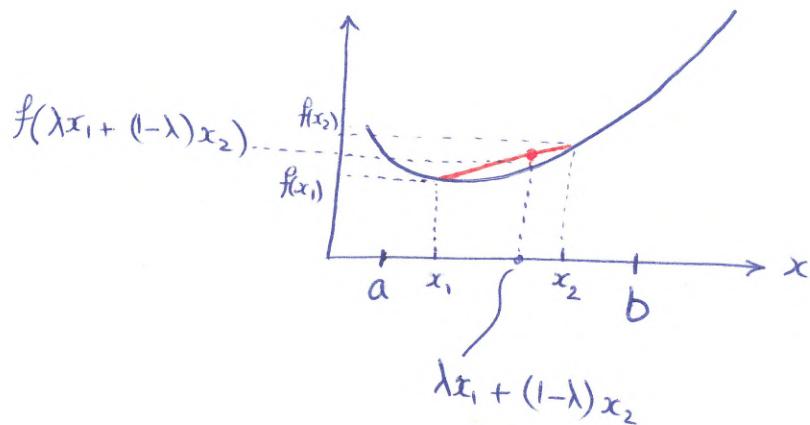
$$I(X;Y) = D(p(x,y) \parallel p(x)p(y))$$

- $I(X; Y) \geq 0$
- $I(X; Y) = H(X) - H(X|Y) \geq 0$
 ↳ conditioning reduces entropy

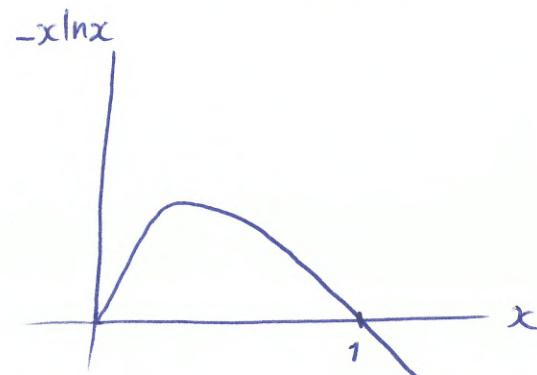
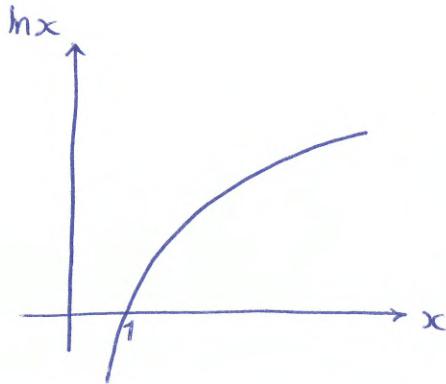
Convexity

Defn. $f(x)$ is convex \vee on (a, b) if for $x_1, x_2 \in (a, b)$ and for $\lambda \in (0, 1)$

$$f(x_1 \lambda + x_2 (1-\lambda)) \leq \lambda f(x_1) + (1-\lambda) f(x_2)$$



ex.



Thm. $f(x)$ is twice continuous and differentiable,

then if $\frac{d^2 f(x)}{dx^2} \geq 0$, then $f(x)$ is convex \checkmark

Proof.

Taylor Expansion

$$f(x) = f(x_0) + f'(x_0) \cdot (x - x_0)$$

$$+ \frac{1}{2} \frac{d^2 f(x^*)}{dx^2} \cdot (x - x_0)^2$$

$$\Rightarrow f(x) \geq f(x_0) + f'(x_0) \cdot (x - x_0)$$

$$x_1 \rightarrow x \quad \cancel{x\lambda} \quad f(x_1) \geq f(x_0) + f'(x_0) (x_1 - x_0)$$

$$x_2 \rightarrow x \quad \cancel{x(1-\lambda)} \quad f(x_2) \geq f(x_0) + f'(x_0) (x_2 - x_0) \quad \lambda \in (0, 1)$$

$$\lambda f(x_1) + (1-\lambda) f(x_2) \geq \underline{\lambda f(x_0)} + \underline{(1-\lambda) f(x_0)}$$

$$+ \cancel{\lambda f'(x_0) (x_1 - x_0)} + \cancel{(1-\lambda) f'(x_0) (x_2 - x_0)}$$

replace x_0 with $\lambda x_1 + (1-\lambda)x_2$:

$$\lambda f(x_1) + (1-\lambda) f(x_2) \geq f(x_0) + \lambda f'(x_0) x_1 + f'(x_0) x_2 - f'(x_0) x_0 - \lambda f'(x_0) x_2$$

$$\geq f(x_0) + f'(x_0) (\cancel{\lambda x_1 + x_2} - \cancel{\lambda x_1 - x_2 + \lambda x_2} - \cancel{\lambda x_2})$$

$$\geq f(x_0)$$

$$\cancel{f(x_0)} <$$

$$\boxed{f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda) f(x_2)} \quad \checkmark$$

(P 2.2, C.T'06)

| X
H(X)

(a) $Y = 2^X$

(b) $Y = \cos(X)$

?
H(Y) $\begin{cases} = \\ \leqslant \\ \geqslant \end{cases}$ H(X)

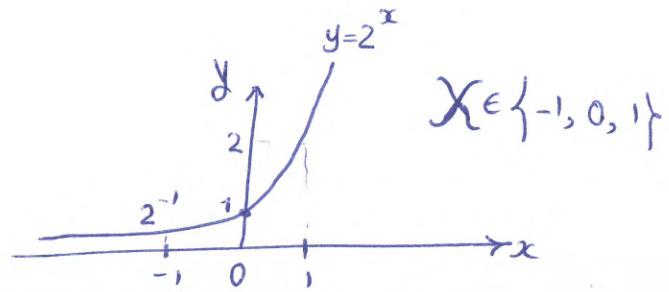
(a) $y = 2^x$ is monotonically increasing \uparrow

one-to-one mapping

$P(Y = \frac{1}{2}) = P(X = -1)$

$P(Y = 1) = P(X = 0)$

$P(Y = 2) = P(X = 1)$



$\mathcal{Y} \in \{\frac{1}{2}, 1, 2\}$

$\Rightarrow H(Y) = H(X)$ (Y : monotonic of X)

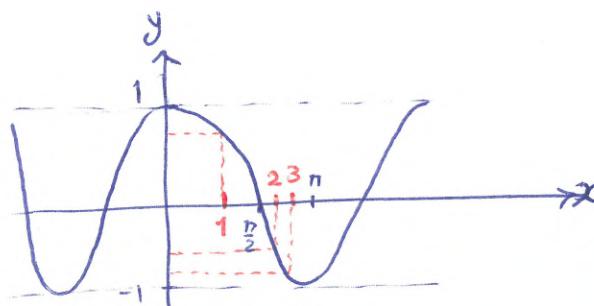
(b)

$X \in \{1, 2, 3\} \xrightarrow{\text{monotonic}} H(X) = H(Y)$

$X \in \{0, \cancel{\pi/2}, \pi, 2\pi\}$

M.E.
 $P(Y = 1) = P(X = 0) + P(X = 2\pi)$

$P(Y = -1) = P(X = \pi)$



\Rightarrow so in general
for many-to-one $\therefore H(Y) \leq H(X)$

[P, 2.30, CT '06]

$$\begin{array}{l} X \\ \times \\ H(X) \\ E[X] = \sum_{n=0}^{+\infty} np(n) = A \text{ (known)} \end{array}$$

$$\max H(X)$$

$$P(x): \frac{E[X]=A}{1}$$

$$\frac{\sum_{n=0}^{+\infty} p(n) = 1}{2}$$

Lagrange Multipliers (λ, μ)

$$L(\lambda, \mu) = \sum_{n=0}^{+\infty} p(x) \ln p(x) + \lambda \left(\sum_{n=0}^{+\infty} np(n) - A \right) + \mu \left(\sum_{n=0}^{+\infty} p(n) - 1 \right)$$

$$\frac{\partial L(\lambda, \mu)}{\partial p(n)} = 0$$

$$-\log p(n) - \frac{p(n)}{p(n)} + \lambda n - \mu = 0$$

$$p(n) = e^{\mu - 1 + \lambda n}$$

HW 1 \rightarrow P 7.4 : starting point

$$\ln x \leq x - 1$$

(came from Taylor series)

Info. Lect. 4

P. 2.30

$X \sim \text{discrete R.V.}$

integer-valued X

\Rightarrow maximize $H(X)$

$$P(n), E(X) = A$$

$$\sum_{n=0}^{+\infty} P(n) = 1$$

$$L(p, \lambda, \mu) = \underbrace{E_p \left[\log \frac{1}{P(X)} \right]}_{-\sum_{n=0}^{+\infty} p(n) \log p(n)} + \lambda \left(\sum_{n=0}^{+\infty} n p(n) - A \right)$$

$$+ \mu \left(\sum_{n=0}^{+\infty} p(n) - 1 \right)$$

$$\frac{\partial L(p, \lambda, \mu)}{\partial p(n)} = -\ln p(n) - 1 + \lambda n + \mu \Big|_{p^*(n)} = 0$$

$$p(n) = e^{\lambda n + (\mu - 1)} = e^{\lambda n} \cdot e^{(\mu - 1)}$$

to find μ :

$$\sum_{n=0}^{+\infty} e^{\lambda n} \cdot e^{\mu - 1} = 1 \Rightarrow e^{\mu - 1} = 1 - e^\lambda$$

assume $e^\lambda < 1 \Rightarrow \lambda < 0$

To find λ :

$$P(n) = e^{\lambda n} (1 - e^\lambda)$$

$$\sum_{n=0}^{+\infty} n e^{\lambda n} (1 - e^\lambda) = A$$

$$\sum_{n=0}^{+\infty} n e^{\lambda n} = \frac{d}{d\lambda} \left(\sum_{n=0}^{+\infty} e^{\lambda n} \right) = \frac{d}{d\lambda} \cancel{\left(\frac{1}{1-e^\lambda} \right)} = \frac{e^\lambda}{(1-e^\lambda)^2}$$

$$\Rightarrow \frac{e^\lambda}{(1-e^\lambda)^2} \cancel{(1-e^\lambda)} = A$$

$$e^\lambda = \frac{A}{1+A}$$

$$\Rightarrow P(n) = \left(\frac{A}{1+A} \right)^n \left(\frac{1}{A+1} \right)$$

$$P(n) = \left(\frac{A}{A+1} \right)^n \frac{1}{A+1}$$

Jensen's Inequality

- f is convex over (a, b)
- X is a discrete R.V. with $\frac{\text{alphabet}}{\text{support}} \mathcal{X}$ and outcomes in (a, b)

$$E[f(X)] \geq f(E(X))$$

- Concave $f \Rightarrow$ invert inequality

ex.1

$$D(p \parallel q) \geq 0$$

proof:

$$\begin{aligned} \sum_x p(x) \log \frac{p(x)}{q(x)} &= E_{p(x)} \left[\log \frac{p(x)}{q(x)} \right] \\ &= -E_{p(x)} \left[\log \frac{q(x)}{p(x)} \right] \end{aligned}$$

• \log is concave.

$$E_{p(x)} \left[\log \frac{q(x)}{p(x)} \right] \leq \log \left[E_{p(x)} \left(\frac{q(x)}{p(x)} \right) \right]$$

$$\Rightarrow -E_{P(x)} \left[\log \frac{q(x)}{p(x)} \right] \geq -\log \underbrace{\left(E_{P(x)} \left[\frac{q(x)}{p(x)} \right] \right)}_{1} \underbrace{\left[\frac{1}{\sum_{x \in \mathcal{X}} p(x)} \right]}_{0} \quad \blacksquare$$

ex. 2

$$H(X) \leq \log |\mathcal{X}|$$

$|\mathcal{X}|$: Cardinality = size of alphabet \mathcal{X}

$$H(X) = E_{P(x)} \left[\log \frac{1}{P(x)} \right]$$

$$\leq \log \underbrace{\left(E_{P(x)} \left[\frac{1}{P(x)} \right] \right)}_{1} = \log |\mathcal{X}|$$

equality
holds

$$\text{when } p(x) = \frac{1}{|\mathcal{X}|}$$

(uniform distribution)
(like white noise)

$$\sum_{x \in \mathcal{X}} P(x) \underbrace{\frac{1}{P(x)}}_1 = |\mathcal{X}|$$

Jensen's
Inequality Proof:

$$|X|=2$$

○ $X = \begin{cases} x_1 & \text{with prob. } P(1) \\ x_2 & \text{with prob. } P(2) \end{cases}, \quad P(1) + P(2) = 1$

$$f(x_1) \underbrace{\cdot P(1)}_P + f(x_2) \underbrace{\cdot P(2)}_{1-P} \geq f(p x_1 + (1-p)x_2)$$

Convexity

○ By Induction

$|X|=k$ and assume that the inequality is proved
for alphabet of size $k-1$.

$$\mathbb{E}_{p(x)}[f(X)] = \sum_{n=1}^k f(x_n) \cdot p(n)$$

a PMF over $k-1$
size alphabet

$$= f(x_k) P(k) + \sum_{i=1}^{k-1} f(x_i) \cdot \underbrace{\frac{P(i)}{\sum_{j=1}^{k-1} P(j)}}_{\text{a PMF over } k-1 \text{ size alphabet}} \times \sum_{j=1}^{k-1} P(j)$$

our proved for $k-1$
assumed

$$\geq f(x_k) P(k) + \left(\sum_{j=1}^{k-1} P(j) \right) f\left(\sum_{i=1}^{k-1} x_i \cdot \frac{P(i)}{\sum_{j=1}^{k-1} P(j)} \right)$$

Sum to ONE!

Convexity

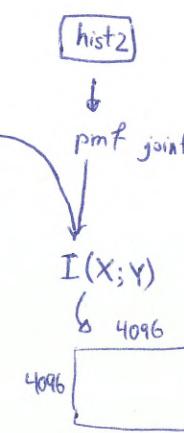
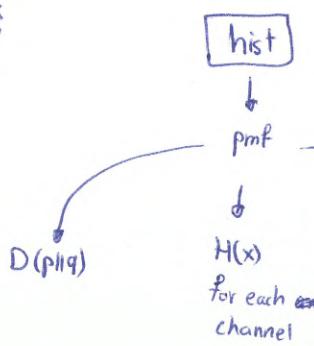
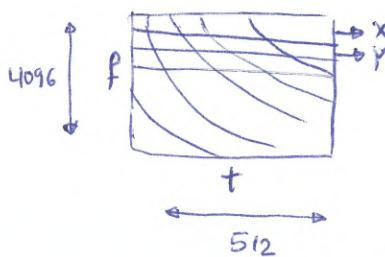
$$\geq f(x_k) P(k) + \cancel{\sum_{j=1}^{k-1} P(j)} \sum_{i=1}^{k-1} x_i \frac{P(i)}{\cancel{\sum_{j=1}^{k-1} P(j)}}$$

$$f(E(X))$$



project assignment:

4096×512



- fit Gaussians to each channel?

[info. lect.]

Project due: 2 weeks from now.

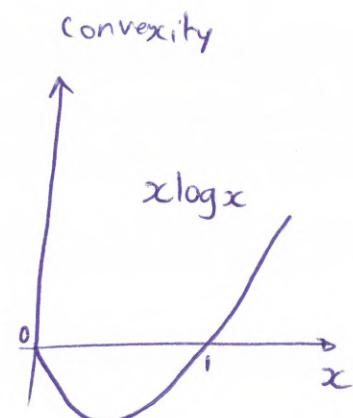
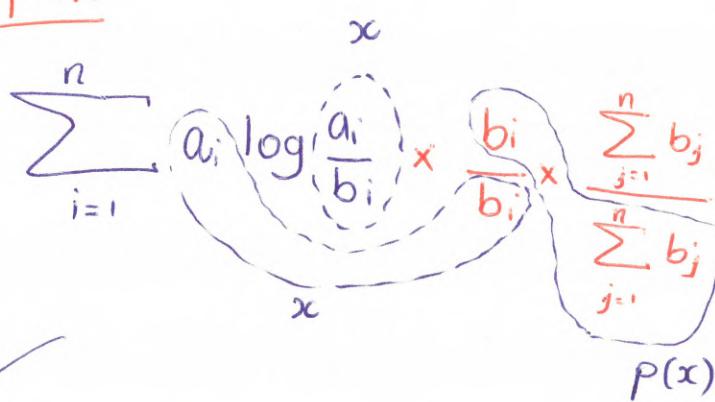
Log-Sum Inequality

$$\underline{a} = (a_1, \dots, a_n) \quad a_i > 0$$

$$\underline{b} = (b_1, \dots, b_n) \quad b_i > 0$$

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \sum_{i=1}^n a_i \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i}$$

proof:



$$\begin{aligned} &\geq \left(\sum_{j=1}^n b_j \right) \cdot \sum_{i=1}^n \frac{a_i}{b_i} \cdot \frac{b_i}{\sum_{j=1}^n b_j} \log \sum_{i=1}^n \frac{a_i}{b_i} \cdot \frac{b_i}{\sum_{j=1}^n b_j} \\ &= \left(\sum_{i=1}^n a_i \right) \log \frac{\sum a_i}{\sum b_j} \end{aligned}$$

prove it myself again!

9.1.22-1

Remarks:

1) $H_p(x)$ is concave in p .

$$2) I(X; Y) = \sum \sum p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

$$= \sum \sum q(y|x) p(x) \log \frac{q(y|x)p(x)}{p(x)p(y)} \rightarrow \sum_x p(x)q(y|x)$$

$\Rightarrow M.I.$ is concave in $p(x)$ and convex in $q(y|x)$.

3) $D(p||q)$ is convex in both p and q .

$$H(\lambda p(x) + (1-\lambda)q(x)) = H_{\lambda p(x) + (1-\lambda)q(x)}(x)$$

$$= - \sum_x \left[\underbrace{\lambda p(x)}_{a_1} + \underbrace{(1-\lambda)q(x)}_{a_2} \right] \log \frac{\lambda p(x) + (1-\lambda)q(x)}{\underbrace{\lambda + (1-\lambda)}_{b_1 + b_2}} \underbrace{\log}_{\text{Convex}}$$

$$\leq a_1 \log \frac{a_1}{b_1} + a_2 \log \frac{a_2}{b_2}$$

Just forget
about the sum
for a second:

Considering
minus

$$\geq -\lambda p(x) \log \frac{\lambda p(x)}{\lambda} - (1-\lambda)q(x) \log \frac{(1-\lambda)q(x)}{1-\lambda}$$

$$\rightarrow = \lambda H_p(x) + (1-\lambda) H_q(x)$$

brought
back the sum:

$$\Rightarrow H_{\lambda p + (1-\lambda)q}(x) \geq \lambda H_p(x) + (1-\lambda) H_q(x)$$

Data Processing Inequality

$$X \rightarrow Y \rightarrow Z$$

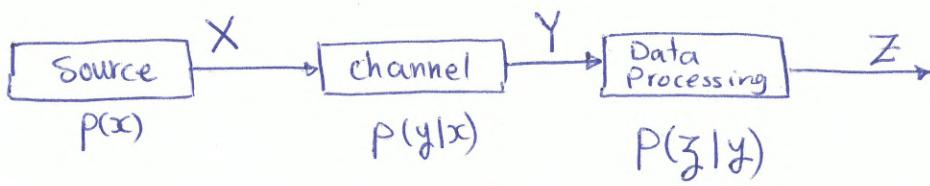
markov chain (no feedback)

$$P(x, y, z) P(x|y) = \frac{P(x, y, z)}{P(y)} = \frac{P(z|x, y)}{P(x, y)} \cdot \frac{P(x|y) P(y)}{P(z)}$$

$$P(x, z|y) = P(z|y) \cdot P(x|y)$$

Markov

* Given present (y), past (x) and future (z) are independent.



- $I(X; Y) \geq I(X; Z)$ \rightarrow equality: when data processing is perfect.

proof

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

$\underbrace{I(X; Z | Y)}_{\geq 0} + \underbrace{I(X; Y)}_{\geq 0}$

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

$\underbrace{I(X; Y | Z)}_{\geq 0} + \underbrace{I(X; Z)}_{\geq 0}$

$$I(X; Z | Y) = \sum_{x,y,z} p(x,y,z) \log \frac{\overbrace{p(x,y,z)}^{\text{markov } p(x|y) \cdot p(y|z)}}{\overbrace{p(x,y,z)}^{p(x|y) \cdot p(y|z)}} = 0$$

$$\Rightarrow I(X; Y) \geq I(X; Z)$$

$$\cancel{H(X) - H(X|Y)} \geq \cancel{H(X) - H(X|Z)}$$

$$\Rightarrow \boxed{H(X|Z) \geq H(X|Y)}$$

info. lect.

Telescoping Properties

$$1. \quad H(x_1, \dots, x_n) = - \sum_{\underline{x} \in X^n} p(\underline{x}) \log p(\underline{x})$$

notation
 $\underline{x}^n = \underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = (x_1, \dots, x_n)$

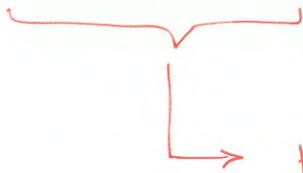
$$\begin{aligned} p(\underline{x}) &= p(x_1, \dots, x_n) = p(x_n | x_{n-1}, \dots, x_1) \underbrace{p(x_{n-1}, \dots, x_1)}_{p(x_{n-1} | x_{n-2}, \dots, x_1) p(x_{n-2}, \dots, x_1)} \\ &= \left[\prod_{k=2}^n p(x_k | x_{k-1}, \dots, x_1) \right] p(x_1) \end{aligned}$$

$$\Rightarrow H(x_1, \dots, x_n) = - \sum_{\underline{x}} p(\underline{x}) \left\{ \sum_{k=2}^n \log p(x_k | x_{k-1}, \dots, x_1) + \log p(x_1) \right\}$$

$$= - \sum_{k=2}^n \left\{ \sum_{\underline{x}} p(\underline{x}) \log p(x_k | x_{k-1}, \dots, x_1) + \sum_{\underline{x}} p(\underline{x}) \log p(x_1) \right\}$$

$H(x_1, \dots, x_n) = \sum_{k=2}^n H(x_k | x_{k-1}, \dots, x_1) + H(x_1)$

$$2. \ I(\overbrace{X, Y}^{\text{joint}}; Z) = I(X; Z|Y) + I(Y; Z)$$



$$H(X, Y) - H(X, Y|Z) = [H(X|Y) + H(Y) - H(X|Y, Z)] - [H(Y|Z) - I(X; Z|Y) - I(Y; Z)]$$

- AEP: fancy title for usage of law of large numbers.
- chapter 3 of Cover & Thomas

X_1, \dots, X_n iid

$$\left\{ \begin{array}{l} S_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{a sequence of sample means} \\ V_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - S_n)^2 \quad \text{a sequence of sample variances} \end{array} \right.$$

X_i : random variable, S_n : R.V., V_n : R.V.

- Do they converge?
- In what sense?
- What is the limit?

Different kinds of convergence:

- ✓ 1. Convergence in probability (WLLN)
- ✓ 2. Convergence in distribution (CLThm)
- ✓ 3. Mean-square convergence (continuity and 2nd order statistics of R.P.'s)
- ✗ 4. Probability 1 convergence (almost surely)

$$\begin{aligned} & \xrightarrow{\text{P}} \lim_{i \rightarrow +\infty} X_i = a \\ & \quad \xrightarrow{\text{number}} \end{aligned}$$

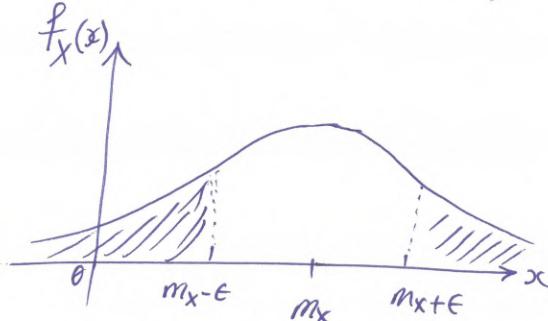
related to sandwich theorem

we are going to use WLLN, so 2 results need to be recalled: Chebyshov and WLLN.

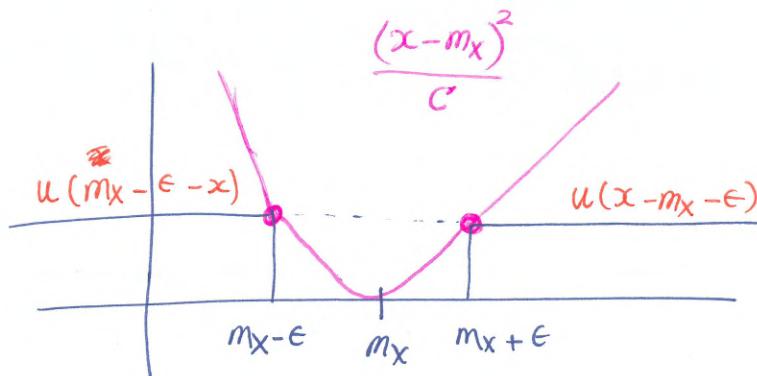
Chebyshov Inequality

X is a R.V. with mean m_X and variance σ_X^2

$$P(|X - m_X| > \epsilon) \leq \frac{\sigma_X^2}{\epsilon^2}$$



$$P(|X - m_x| > \epsilon) = \int_{m_x - \epsilon}^{m_x + \epsilon} f_X(x) dx + \int_{m_x + \epsilon}^{+\infty} f_X(x) dx$$



$$= \int_{-\infty}^{+\infty} \left\{ u(m_x - \epsilon - x) + u(-m_x - \epsilon + x) \right\} x f_X(x) dx$$

2 points chosen

to tighten the bound

$$\frac{(m_x + \epsilon - m_x)^2}{C} = 1 \Rightarrow C = \epsilon^2$$

tightest bound

$$\Rightarrow \leq \int_{-\infty}^{+\infty} \frac{(x - m_x)^2}{\epsilon^2} f_X(x) dx$$



X_1, \dots, X_n, \dots

iid

$m_x \rightarrow \sigma_x^2$: mean and variance of each X_i

RV.

$$S_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$n = 1, 2, 3, \dots$

$$E[S_n] = \frac{1}{n} \sum E[X_i] = m_x$$

$$\text{Var}[S_n] = \frac{1}{n} \sigma_x^2$$

(think how to prove)

$$\text{Var}[S_n] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right]$$

$$= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sigma_x^2$$

$$= \frac{n}{n^2} \sigma_x^2 = \frac{1}{n} \sigma_x^2$$

$$P(|S_n - E[S_n]| > \epsilon) \leq \frac{\text{Var}[S_n]}{\epsilon^2} = \frac{\sigma_x^2}{n \epsilon^2}$$

as

$$n \rightarrow +\infty$$

$$P(|S_n - E[S_n]| > \epsilon) \rightarrow 0$$

WLLN



WLLN (iid data)

$$S_n \xrightarrow{P} \underbrace{E[S_n]}_{m_X} \quad \text{as } n \rightarrow +\infty$$

Sample entropy

$$\frac{1}{n} \log \frac{1}{P(x_1, \dots, x_n)}$$

$$\frac{1}{n} \log \frac{1}{P(x_1, x_2, \dots, x_n)} \underset{iid}{=} \frac{1}{n} \log \frac{1}{\cancel{P(x_1) \cdots P(x_n)}} \\ = -\frac{1}{n} \sum_{i=1}^n \log P(x_i) \triangleq S_n^e$$

$$S_n^e \xrightarrow{P} ? \quad H(X)$$

as $n \rightarrow +\infty$

AEP

$$\text{Let } Y_i = -\log \underbrace{P(x_i)}_{iid} \Rightarrow Y_i: iid$$

$$S_n = \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow{P} E[Y_i]$$

$$E_{P_{Y_i}}[Y_i] = E_{P_{X_i}}[-\log P(x_i)] = -\sum_x p(x) \log p(x) = H(X)$$

info. lect.

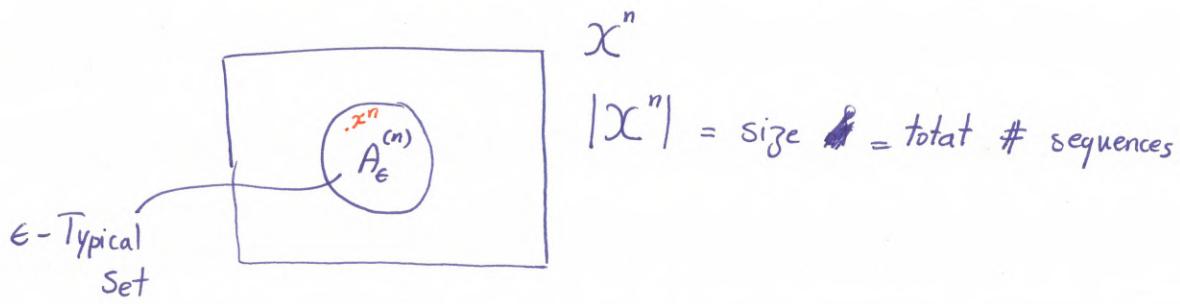
$\mathcal{X}: X_1, \dots, X_n \quad \text{iid RV}$

AEP

$$\frac{1}{n} \log \frac{1}{P(X^n)} \xrightarrow[n \rightarrow +\infty]{P} H(X)$$

$$P(\left| \frac{1}{n} \log \frac{1}{P(X^n)} - H(X) \right| > \epsilon) \rightarrow 0$$

as $n \rightarrow +\infty$



$$\left. \begin{array}{l} \text{for } x^n \in A_\epsilon^{(n)} \\ |A_\epsilon^{(n)}| \cong 2^{nH(X)} \\ P(x^n) \cong 2^{-nH(X)} \\ P(A_\epsilon^{(n)}) \approx 1 \end{array} \right\}$$

* typical set sequences have their empirical entropy close to $H(X)$.

* $A_\epsilon^{(n)}$ is a small set of high probability with uniformly distributed sequences in it.

Def:

$A_\epsilon^{(n)}$ is typical set

$$A_\epsilon^{(n)} = \left\{ x^n : \left| \frac{1}{n} \log \frac{1}{P(x^n)} - H(X) \right| < \epsilon \right\}$$

ex. (ternary alphabet)

$$\mathcal{X} = \{1, 2, 3\}$$

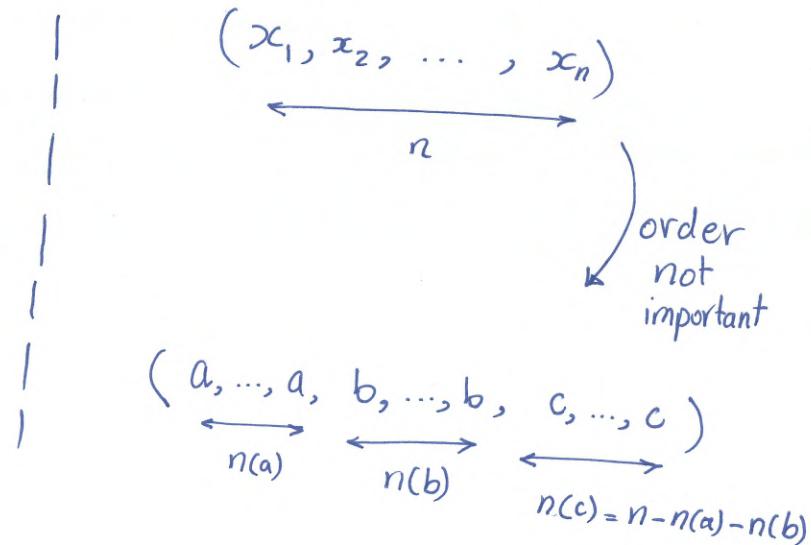
~~$\{x_1, x_2, x_3\}$~~

~~$\{a, b, c\}$~~

$p(a)$

$p(b)$

$p(c) = 1 - p(a) - p(b)$



$$P(x^n) = P(a)^{n(a)} \cdot P(b)^{n(b)} \cdot P(c)^{n(c)}$$

\underline{x}

$n \cdot \frac{n(a)}{n} \rightarrow f_n(a)$: frequency of observing a

~~WLLN~~ **WLLN**

$$\frac{n(a)}{n} \xrightarrow[\text{as } n \rightarrow +\infty]{P} P(a)$$

$$\underset{\text{as } n \rightarrow +\infty}{\approx} \frac{n P(a)}{P(a)} = \frac{n P(b)}{P(b)} = \frac{n P(c)}{P(c)}$$

$$= 2^{\log(P(a)^{np(a)})} \cdot 2^{\log(\cdot)} \cdot 2^{\log(\cdot)}$$

$$= 2^{np(a) \log P(a)} \cdot 2^{np(b) \log P(b)} \cdot 2^{np(c) \log P(c)}$$

$$\Rightarrow P(x^n) \approx 2^{-nH(X)}$$

Corollaries

1 If $x^n \in A_\epsilon^{(n)}$:

$$\left| \frac{1}{n} \log \frac{1}{P(x^n)} - H(X) \right| < \epsilon$$

$$H(X) - \epsilon < -\frac{1}{n} \log P(x^n) < H(X) + \epsilon$$

$$-nH(X) - n\epsilon < \log P(x^n) < -nH(X) + n\epsilon$$

) $\times (-n)$

$$2^{-n(H(X)+\epsilon)} < P(x^n) < 2^{-n(H(X)-\epsilon)}$$

) $\log: \text{monotonic}$

$$2/ \quad |A_{\epsilon}^{(n)}| \leq 2^{-n}(H(x) + \epsilon)$$

$$\begin{aligned} 1 &= \sum_{x^n \in \mathcal{X}^n} p(x^n) \geq \sum_{x^n \in A_{\epsilon}^{(n)}} p(x^n) \\ &\geq \sum_{x^n \in A_{\epsilon}^{(n)}} 2^{-n}(H(x) + \epsilon) \\ &= 2^{-n}(H(x) + \epsilon) |A_{\epsilon}^{(n)}| \end{aligned}$$

Cardinality of $A_{\epsilon}^{(n)}$

$$3/ \quad P(A_{\epsilon}^{(n)}) > 1 - \epsilon$$

From AEP property:

$$\underbrace{\frac{1}{n} \log \frac{1}{P(x^n)}}_{\text{empirical entropy}} \xrightarrow{P} H(x)$$

$$P\left(\left|\frac{1}{n} \log \frac{1}{P(x^n)} - H(x)\right| > \epsilon\right) \rightarrow 0$$

$\epsilon \xrightarrow{n \rightarrow +\infty} 0$

We can define $n(\epsilon)$ and find the n that $n > n(\epsilon)$.

4/

$$1 - \epsilon \leq P(A_\epsilon^{(n)}) = \sum_{x^n \in A_\epsilon^{(n)}} p(x^n) \leq \sum_{x^n \in A_\epsilon^{(n)}} 2^{-n} (H(x) - \epsilon) = 2^{-n} (H(x) - \epsilon) |A_\epsilon^{(n)}|$$

~~$$\Rightarrow (1 - \epsilon) 2^{-n} (H(x) - \epsilon) \leq |A_\epsilon^{(n)}|$$~~

$$|A_\epsilon^{(n)}| \approx 2^{-n H(X)}$$

Worst case
largest typical set $H(X) = \log |X|$

$$\Rightarrow |A_\epsilon^{(n)}|_{\max} = 2^{n \log |X|}$$

$B_{\delta}^{(n)}$ sets of smallest size and largest possible probability
 n can vary

$$P(B_{\delta}^{(n)}) > 1 - \delta$$

$$0 < \delta < 0.5$$

ex.

$$X \sim \text{Bernoulli}(0.8)$$

↳ $P(X=1)$

$$\delta = 0.4$$

$$n=1$$

$$B_{0.4}^{(1)} = \{1\}$$

$$P(B_{0.4}^{(1)}) = 0.8 \quad \text{greater than } \overbrace{1-\delta}^{0.6} \Rightarrow \textcircled{O}$$

$$n=2$$

$$B_{0.4}^{(2)} = \{11\}$$

$$P(B_{0.4}^{(2)}) = 0.64 > 0.6 \quad \textcircled{O}$$

$n=3$

$$B_{0.4}^{(3)} = \{ 111, 110 \}$$

$$P(X=1) = 0.512 < 0.6$$

!

We need to add more to the typical set.

$$P(110) = 0.128$$

Now $0.512 + 0.128 > 0.6$ (1)

satisfied

$n=4$

$$B_{0.4}^{(4)} = \{ 1111, 1110, 1101 \}$$

$$P(1111) = (0.64)^2 = 0.4096$$

$$P(1110) = 0.512 \times 0.2 = 0.1024$$

$$H(X) = -0.8 \log 0.8 - 0.2 \log 0.2 = 0.7219 \left[\frac{\text{bits}}{\text{Symbol}} \right]$$

Defn:

$$A_{\epsilon}^{(n)} = \left\{ x^n : \left| \frac{1}{n} \log \frac{1}{P(x^n)} - H(X) \right| < \epsilon \right\}$$

$$H(X) - \epsilon < \frac{1}{n} \log \frac{1}{P(x^n)} < H(X) + \epsilon$$

(wavy line)

$n=1$

$$\epsilon = 0.4 \quad \text{pick } \{1\}$$

$$0.7219 - 0.4 < \frac{1}{1} \log \frac{1}{0.8} < 0.7219 + 0.4$$

$$0.3219 \cancel{<} 0.2231 < 1.1219$$

$n=2$

$$\cancel{\cancel{<}} \quad \underbrace{\frac{1}{2} \log \frac{1}{(0.8)^2}}_{\log \frac{1}{0.8} = 0.2231} <$$

Sequences of all 1 will not make a typical set.

$n=2$

$\{10\}$ ↘

$$0.3219 \quad \leftarrow \frac{1}{2} \log \underbrace{\frac{1}{0.16}}_{1.3219} \quad \cancel{1.1219}$$

$n=3$

111 is not typical.

let's try 110:

$$\frac{1}{3} \log \underbrace{\frac{1}{0.64 \times 0.2}}_{0.9886} = \hat{H}(X)$$

$$A_{0.4}^{(3)} = \{110, 101, 011\}$$

$n=4$

$$A_{0.4}^{(4)} = \{1110, 1101, 1011, 0111\}$$

$$\hat{H}(X) = 0.6616 \quad \textcircled{d}$$

also tried: 1100 → $\hat{H}(X) \approx 1.3219 \quad \times$

info. lect.

Constrained Optimization

- P(error)
 - MAP
 - $\max P_{\text{Detection}}$ given $P_{\text{False Alarm}} \leq d$
 - ML estimate
 - MAP estimate, MMSE
-

$$\text{minimize } f(\underline{x}) \quad \underline{x} \in \mathbb{R}^n$$

assume m equality constraints

$$h_1(\underline{x}) = 0, \dots, h_m(\underline{x}) = 0$$

and p inequality constraints

$$g_1(\underline{x}) \leq 0, \dots, g_p(\underline{x}) \leq 0$$

$$\min f(\underline{x})$$

$$\underline{x} \in \mathbb{R}^n : h_1(\underline{x}) = 0, \dots, h_m(\underline{x}) = 0$$

$$g_1(\underline{x}) \leq 0, \dots, g_p(\underline{x}) \leq 0$$

ex.

$$\min f(\underline{x}_1, \underline{x}_2) \quad \overbrace{3x_1^2 + 2x_1x_2 + 3x_2^2 - 20x_1 + 4x_2}$$

$$(\underline{x}_1, \underline{x}_2) : \quad x_1 \geq 0 \\ x_2 \geq 0$$

$$\left. \begin{array}{l} \frac{df}{dx_1} = 6x_1 + 2x_2 - 20 = 0 \\ \frac{df}{dx_2} = 2x_1 + 6x_2 + 4 = 0 \end{array} \right\} \rightarrow \left\{ \begin{array}{l} x_2 = -2 \\ x_1 = 4 \end{array} \right.$$

$$(x_1, x_2) = (4, -2)$$

$$f(4, -2) = -44 \quad (\text{Global minimum})$$

→ violating the constraints!

If both constraints are ON, then set $x_2 = 0$

Cause this is the first point satisfying the constraint on x_2

$$\Rightarrow \left\{ \begin{array}{l} x_1 = \frac{10}{3} \\ \text{OR} \\ \cancel{x_1 = -2} \end{array} \right.$$

When set to equality,
the constraint is active.

Constrained Solution

$$(x_1^*, x_2^*) = \left(\frac{10}{3}, 0\right)$$

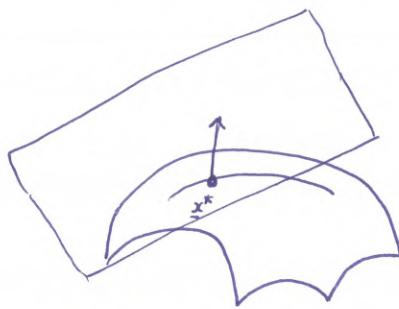
$$f\left(\frac{10}{3}, 0\right) = -33.33$$

$$f(-2, 0) = 46$$

Geometric Point of View

○ ∀i: $\nabla f(\underline{x}^*) = \lambda_i \nabla h_i(\underline{x}^*)$

$$\Rightarrow \nabla f(\underline{x}^*) = \sum_{i=1}^m \lambda_i \nabla h_i(\underline{x}^*)$$



Lagrangian: $L(\underline{x}, \lambda_1, \dots, \lambda_m) = f(\underline{x}) + \sum_{i=1}^m \lambda_i h_i(\underline{x})$

f(x) λ_i h_i(x)
Lagrange multipliers

ex. 1



$$x + 2y = R$$

$$\text{area} = xy = f(x, y)$$

$$\max f$$

$$(x, y) : x + 2y = R$$

$$L(x, y, \lambda) = xy + \lambda(x + 2y - R)$$

$$\begin{array}{l} y + \lambda = 0 \\ x + 2\lambda = 0 \end{array} \quad \left. \begin{array}{l} y = -\lambda \\ x = -2\lambda \end{array} \right\} \Rightarrow -2\lambda + (-2\lambda) = R$$

$$S^* = 1250$$

$$\begin{array}{l} R = 100 \\ y = 25 \\ x = 50 \end{array}$$

$$\left. \begin{array}{l} \\ \\ \end{array} \right\} \quad \leftarrow$$

$$\begin{array}{l} y = \frac{R}{4} \\ x = \frac{R}{2} \end{array} \quad \left. \begin{array}{l} \lambda = -\frac{R}{4} \\ \end{array} \right\} \quad \leftarrow$$

ex. $f(\underline{x}) = x_1^2 + x_1x_2 + x_2x_3$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

maximize such that $x_1 + x_2 = 4$

$$L(\underline{x}, \lambda) = x_1^2 + x_1x_2 + x_2x_3 + \lambda(x_1 + x_2 - 4)$$

$$\frac{\partial L}{\partial x_1} = 0 \rightarrow 2x_1 + x_2 + \lambda = 0 \quad x_1 = -\frac{\lambda}{2}$$

$$\frac{\partial L}{\partial x_2} = 0 \rightarrow x_1 + x_3 + \lambda = 0 \quad x_3 = -\frac{\lambda}{2}$$

$$\frac{\partial L}{\partial x_3} = 0 \rightarrow x_2 = 0$$

$$\frac{\partial L}{\partial \lambda} = 0 \rightarrow x_1 + x_2 = 4 \quad x_1 = 4$$

$$x_3 = 4$$

Constrained Solution:

$$(x_1^*, x_2^*, x_3^*) = (4, 0, 4)$$

ex.3 (maximum entropy)

x_1, \dots, x_n outcomes

$p(x_1), \dots, p(x_n)$ prob. assignments

$\max H(X)$ over choices for $p(x_i)$ such that $\sum_{i=1}^n p(x_i) = 1$

$$L(p(x_1), p(x_2), \dots, p(x_n)) = -\sum_{i=1}^n p(x_i) \ln p(x_i) + \lambda \left(\sum_{i=1}^n p(x_i) - 1 \right)$$

$\underbrace{\text{n equations}}_{\text{eq.}} \left\{ \begin{array}{l} \frac{\partial L}{\partial p(x_i)} = -\ln p(x_i) - 1 + \lambda = 0 \\ i=1, \dots, n \end{array} \right.$

$$\frac{\partial L}{\partial \lambda} \stackrel{0}{\Rightarrow} \sum_{i=1}^n p(x_i) = 1$$

$$\underbrace{p(x_i)}_{=} = e^{-1+\cancel{\lambda}}$$

$$\sum_{i=1}^n e^{-\lambda+1} = 1 \Rightarrow$$

$$e^{-\lambda+1} = \frac{1}{n}$$

UNIFORM Distribution

ex.4 $X = \{1, 2, \dots, n\}$

$$p(1), p(2), \dots, p(n)$$

$$\max H(X) \text{ over } p. \quad \sum_{i=1}^n p(i) = 1, \quad \sum_{i=1}^n i p(i) = m$$

$$L(p_1, \dots, p_n, \lambda, \mu) = -\sum_{i=1}^n p(i) \ln p(i) + \lambda \left(\sum_{i=1}^n p(i) - 1 \right) + \mu \left(\sum_{i=1}^n i p(i) - m \right)$$

$$-\ln p(i) - 1 + \lambda + \mu \cdot i = 0$$

$$p(i) = e^{-1+\lambda} \cdot e^{\mu \cdot i}$$

$$\sum_{i=1}^n p(i) = 1$$

$$\sum_{i=1}^n i p(i) = m$$

$$e^{-1+\lambda} \sum_{i=1}^n e^{\mu i} = 1 \rightarrow e^{-1+\lambda} = \frac{1}{\sum_{i=1}^n e^{\mu i}} q^i$$

$$\Rightarrow \sum q^i - 1 = \frac{1 - q^{i+1}}{1 - q} - 1$$

$$\sum_{j=1}^n j \frac{e^{\mu_j}}{\sum_{i=1}^n e^{\mu_i}} = m$$

$$\frac{d}{d\mu_j} \left(\underbrace{\sum_{j=1}^n e^{\mu_j}} \right) = \sum_{j=1}^n j e^{\mu_j}$$

$$\frac{e^\mu (1 - e^{\mu n})}{1 - e^\mu}$$

Inequality Constraints

ex1

$$\min f(x, y)$$

$$f(x, y) = 3x^2 + 2xy + 3y^2 - 20x + 4y$$

$$(x, y) : x \geq 0, y \geq 0$$

1. check if any local min or max

$$\frac{\partial f}{\partial x} = 6x + 2y - 20 = 0 \quad (1)$$

$$-3x \left(\frac{\partial f}{\partial y} = 2x + 6y + 4 = 0 \right) \quad (2)$$

$$-16y - 32 = 0$$

$$\boxed{y = -2} \Rightarrow \boxed{x = 4}$$

* NOT A SOLUTION!

$$f(4, -2) = 48 - 16 + 12 - 80 - 8 = 0$$

2. set $x = 0, y = -\frac{2}{3}$ (from (2))

* NOT A SOLUTION!

3. set $y = 0, x = \frac{10}{3}$ (from (1))

* POTENTIAL SOLUTION! $\Rightarrow f(\frac{10}{3}, 0) = -\frac{100}{3}$

4. $x = 0, y = 0$

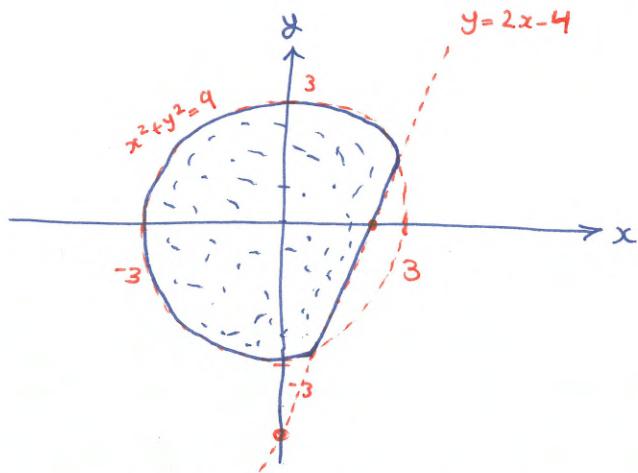
$$f(0, 0) = 0$$

solution $\begin{cases} (x, y) = (\frac{10}{3}, 0) \\ f(x, y) = -\frac{100}{3} \end{cases}$

ex.2 $\min f(x, y)$

$$(x, y) : \begin{cases} x^2 + y^2 \leq 9 \\ 2x - y \leq 4 \end{cases}$$

$$f(x, y) = 3x^2 + 4y^2 + 6xy - 6x - 8y$$



continuous f and closed region \Rightarrow Must have local min/max.

$$\underline{1.} \quad \frac{\partial f}{\partial x} = 6x + 6y - 6 = 0 \quad (3)$$

$$\underline{-1} \times \left(\frac{\partial f}{\partial y} = 8y + 6x - 8 = 0 \right) \quad (4)$$

$$-2y + 2 = 0$$

$$\boxed{y = 1} \quad \Rightarrow \quad \boxed{x = 0}$$

* Possible Solution $(0, 1)$: $\underline{f(0, 1) = -4}$

2. Both constraints are active.

involve $(3, 4)$

$$6x + 6y - 6 = \lambda(2x) + \mu(2)$$

$$8y + 6x - 8 = \lambda(2y) + \mu(-1)$$

$$\nabla f(x, y) = \lambda \nabla g_1 + \mu \nabla g_2$$

$$\text{and} \quad \begin{aligned} x^2 + y^2 &= 9 \\ 2x - y &= 4 \end{aligned}$$

4 equations
4 unknowns

$$\Rightarrow (x_1 = 2.6770, \quad y_1 = 1.3541)$$

both ON the
boundary!

$$(x_2 = 0.5230, \quad y_2 = -2.9541)$$

$$f(x_1, y_1) = 23.6881$$

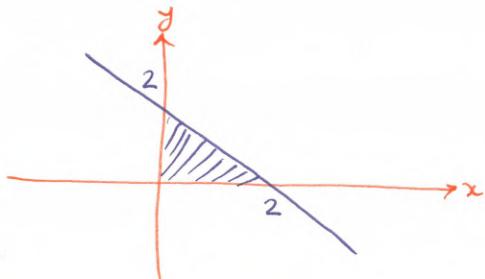
$$f(x_2, y_2) = 46.9522$$

\Rightarrow Final solution: $(0, 1)$

ex. 3 $\max f(x, y)$

$$\begin{cases} x + y \leq 2 \\ x \geq 0, \quad y \geq 0 \end{cases}$$

$$f(x, y) = \ln\left(1 + \frac{x}{5}\right) + \ln\left(1 + \frac{y}{2}\right)$$



1. $\frac{\partial f}{\partial x} = \frac{1}{5+x}$

$$\frac{\partial f}{\partial y} = \frac{1}{2+y}$$

$\} \rightarrow$ no solution!

options:

$x+y=2$, then check $x \geq 0, y \geq 0$.

$x=0$ or $y=0$ or both

2. set $x+y=2$ (use Lagrange multipliers with equality constraints)

$$\frac{1}{5+x} = \lambda \rightarrow x = \frac{1-5\lambda}{\lambda}$$
$$\frac{1}{2+y} = \lambda \quad = \underbrace{\frac{1}{\lambda}}_{\mu} - 5$$

↓

$$y = \underbrace{\frac{1}{\lambda}}_{\mu} - 2$$

$$\Rightarrow \mu - 5 + \mu - 2 = 2 \Rightarrow \boxed{\mu = \frac{9}{2}}$$

$$\Rightarrow x = -\frac{1}{2}, y = \frac{5}{2} \quad \text{NOT A SOLUTION!}$$

3.

set $x=0$

$f(x,y)$ increasing due to $y \}$ $\Rightarrow y=2$ (max of $f(0,y)$)

$$f(0,2) = \ln(2)$$

4.

set $y=0$ $\xrightarrow[\text{reasoning}]{\text{the same}} f(2,0) = \ln(\frac{7}{5})$

5.

$$f(0,0) = 0$$

Final Solution
 $\Rightarrow (x,y) = (0,2)$

if $x + y \leq 5$

$\Rightarrow \underline{2.} \quad f(1, 4) = \ln\left(\frac{6}{5}\right) + \ln(3) = 1.28$

$\underline{3.} \quad f(0, 5) = 1.2528$

$\underline{4.} \quad f(5, 0) = 0.6931$

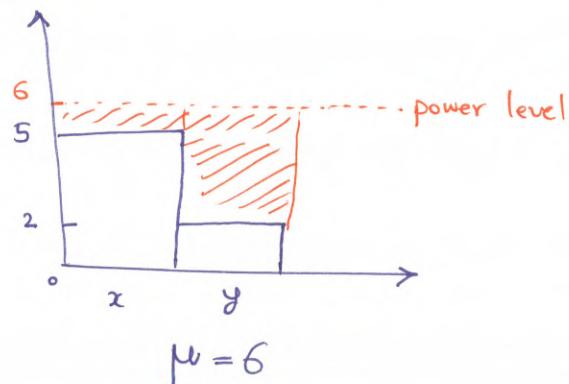
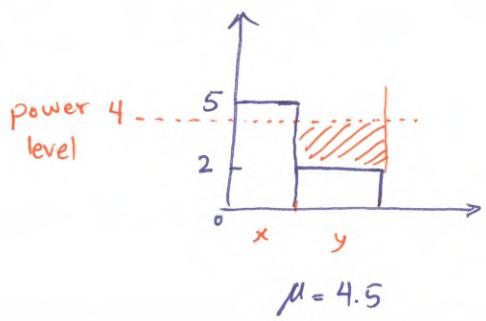
\Rightarrow Final
solution $(x, y) = (1, 4)$

$$\max f(x,y) = \ln\left(1 + \frac{x}{5}\right) + \ln\left(1 + \frac{y}{2}\right)$$

$$(x,y) : x \geq 0, y \geq 0$$

$$x+y \leq 5$$

$$x+y \leq 2$$



Kuhn - Tucker Conditions

$$\min f(\underline{x})$$

$$\underline{x} \in \mathbb{R}^n$$

$$\underline{x} : h_1(\underline{x}) = 0$$

$$g_1(\underline{x}) \leq 0$$

$$\vdots$$

$$h_m(\underline{x}) = 0$$

$$g_p(\underline{x}) \leq 0$$

Def'n

\underline{x}^* is regular if $\nabla_{\underline{x}} h_i(\underline{x})$ and $\nabla_{\underline{x}} g_j(\underline{x})$ for all i, j are linearly independent.

Theorem: Let \underline{x}^* is a local minimum.

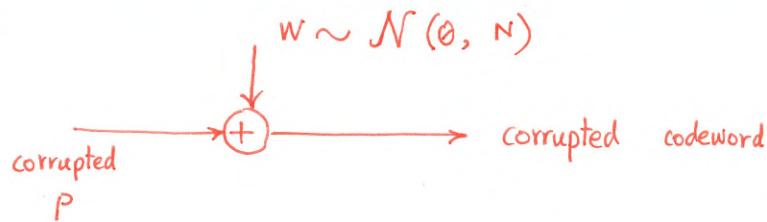
and it is a regular point.

then there exist $\underline{\lambda}$ and $\underline{\mu}$ such that:

$$\nabla_{\underline{x}} f(\underline{x}^*) + \nabla_{\underline{x}} h^T(\underline{x}^*) \cdot \underline{\lambda} + \nabla_{\underline{x}} g^T(\underline{x}^*) \underline{\mu} = \underline{0} \quad (1)$$

and $g^T(\underline{x}^*) \cdot \underline{\mu} = 0 \quad \text{for } \underline{\mu} \leq \underline{0}$ (2)

ex.



capacity $\leftarrow C = \frac{1}{2} \ln \left(1 + \frac{P}{N} \right)$

n channels

$$\sum_{i=1}^n p_i = P$$

$$p_i \geq 0$$

$$\max_{\underline{p}} \sum_i \frac{1}{2} \ln \left(1 + \frac{p_i}{N_i} \right)$$

$$P : \sum p_i = P$$

$$p_i \geq 0 \quad i=1, \dots, n$$

$$\nabla_{\underline{P}} \left\{ \sum_i \frac{1}{2} \ln \left(1 + \frac{P_i}{N_i} \right) + \lambda \nabla_{\underline{P}} \left\{ \sum_i P_i - p \right\} \right\}$$

$$+ \nabla_{\underline{P}} \left\{ \sum_i P_i \cdot \mu_i \right\} = 0$$

$$\sum_i \mu_i P_i = 0$$

1) $P_i > 0$ ($\mu_i = 0$)

$$\frac{1}{2} \underbrace{\frac{1}{N_i + P_i}} + \lambda = 0$$

$$P_i = \frac{-\frac{1}{2\lambda}}{N_i} \quad \text{kappa} \rightsquigarrow (\kappa)$$

2) $P_i = 0$ ($\mu_i \leq 0$)

$$\underbrace{\frac{1}{2} \frac{1}{N_i + P_i} + \lambda}_{-\frac{1}{2} \frac{1}{N_i + P_i} - \lambda} + \mu_i = 0$$

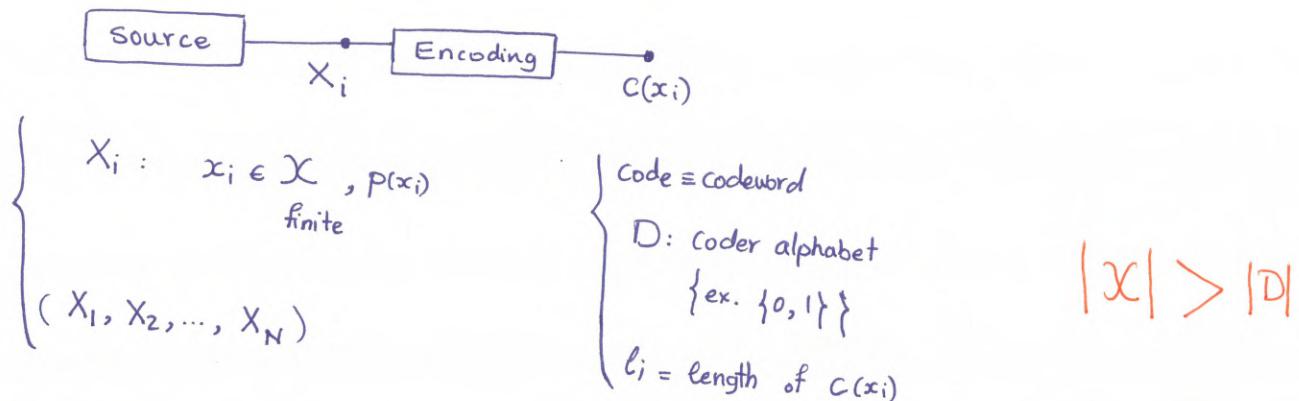
$$-\frac{1}{2} \frac{1}{N_i + P_i} - \lambda \leq 0$$

$$P_i > \frac{1}{-2\lambda} - N_i \rightarrow P_i > x^{\kappa} - N_i$$

(||) assumption

Data Compression (Compaction)

- remove redundancies



- memoryless source : X_i are independent

ex. MPEG : encodes difference between two frames

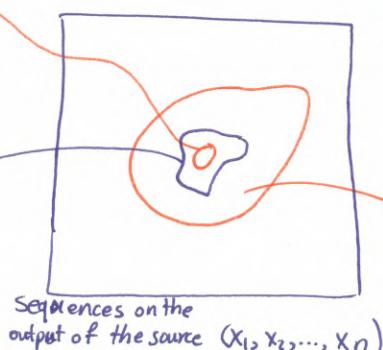
↗ a codeword that is a prefix for another codeword

* Prefix Codes

(Instantaneous Codes)

Uniquely Decodable

Codes
have non-singular
extension



Non singular codes : assign individual code $c(x_i)$ to every x_i (ONE-to-ONE)

ex. $\mathcal{X} = \{1, 2, 3\}_4$

Extension of Code.

x	$C(x)$
1	0
2	010
3	01
4	10

Code for sequence $(1, 2, 3, 4)$.

decoding:
0 010 01 10
1 1 4 3 4
1 2 3 4

NO ONE-to-ONE Decoding!

$\xrightarrow{C(1). C(2) C(3) C(4)}$

concatenated symbols 1, 2, 3, 4
encoded

ex. Morse Code:

$$|X| = 28$$

$$|D| = 4$$

\downarrow
 $\{\cdot, -, \text{letter space}, \text{word space}\}$

$\rightarrow "e"$ is assigned $\cdot \cdot$.

$\rightarrow "q"$ is assigned $\cdot \cdot - -$.

Kraft Inequality

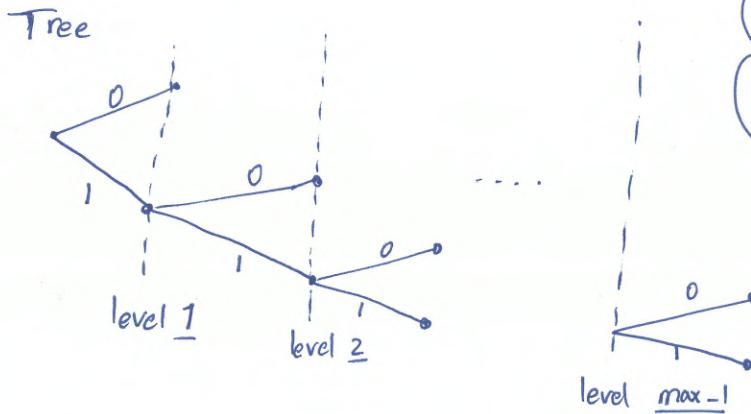
- ① A prefix code with alphabet D and code words length's $\ell_1, \dots, \ell_{\max}$, then ~~the~~ $\{\ell_i\}$ satisfy:

$$\sum_{i=1}^{\max} |D|^{-\ell_i} \leq 1$$

|| (choose $|D|=2$)

- ② If codeword lengths $\ell_1, \dots, \ell_{\max}$ satisfy the inequality, then a prefix code can be constructed with codewords of lengths $\ell_1, \dots, \ell_{\max}$.

Proof: ① choose $D=2$



$2^{\ell_{\max}} = \text{Max number of Codewords}$

not necessarily prefix

codewords removed (at each level)

$$2^{\ell_{\max} - \ell_1} + 2^{\ell_{\max} - \ell_2} + \dots + 2^{\ell_{\max} - \ell_k} + \dots$$

$$= \sum_{k=1}^{\ell_{\max}} 2^{\ell_{\max} - \ell_k} \leq 2^{\ell_{\max}}$$

$$\Rightarrow \sum_{k=1}^{\max} 2^{-\ell_k} \leq 1$$

$$\mathcal{L} = \sum_i l_i p_i$$

\downarrow
 $p(x_i)$

Average codeword length (description length)

$$\begin{array}{ll} \text{minimize}_{l_i} & \sum_{i=1}^n l_i p_i \\ \text{subject to} & \sum_i 2^{-l_i} \leq 1 \end{array} \quad (\text{constrained optimization})$$

$$\nabla_{l_i} \left\{ \sum_{i=1}^n l_i p_i \right\} = \lambda \nabla_{l_i} \left\{ \sum_{i=1}^n 2^{-l_i} - 1 \right\}$$

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

\downarrow
 $e^{-l_i \ln 2}$

$$P_i = \lambda (-\ln 2) \underbrace{e^{-\ell_i \ln 2}}_{2^{-\ell_i}}$$

$$\Rightarrow 2^{-\ell_i} = -\frac{P_i}{\lambda \ln 2}, \quad \frac{1}{-\lambda \ln 2} \sum_{i=1}^n P_i = 1$$

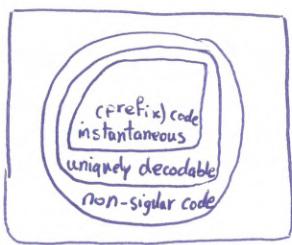
$$\Rightarrow \lambda = -\frac{1}{\ln 2}$$

$$, \quad P_i = 2^{-\ell_i}$$

$$\ell_i = -\log_2 P_i$$

Model Selection Based on MDL (Minimum Description Length)

info. lect. |



$|X|^n$: # of all possible sequences of length n consisting of alphabet X

$$(x_1, \dots, x_n) \rightarrow C(x_1, \dots, x_n)$$

\doteq block code \doteq

Kraft Inequality

- necessary and sufficient conditions to assign prefix codes

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

$$\min E[\ell(X)]$$

$$\ell_i: \sum_{i=1}^n 2^{-\ell_i} = 1 \quad \Rightarrow \quad \ell_i = -\log_2 p_i$$

expected value \Rightarrow

$$\sum_{i=1}^n p_i \cdot \ell_i = -\sum p_i \log p_i = H(X)$$

Limit of Compression

If p_i have base 2: $\{\frac{1}{2}, \frac{1}{4}, \dots\} \rightarrow \ell_i$: integer-valued

In general $\ell_i = \lceil -\log p_i \rceil$

Yes, it's possible!

Is it still prefix code? check Kraft Inequality!

$$\sum_{i=1}^n 2^{-\ell_i} = \sum_{i=1}^n 2^{-\lceil -\log p_i \rceil} \leq \sum_{i=1}^n 2^{-\log \frac{1}{p_i}} = \sum_{i=1}^n p_i = 1$$

- Prefix codes have $L = E[\ell(X)] \geq H(X)$.

proof.

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

$$\log \sum_{i=1}^n 2^{-\ell_i} \times \frac{p_i}{p_i} \leq 0$$

$$\log (E_{p_i, y_i}) \leq 0$$

\log is concave

$$E_{p_i}(\underbrace{\log y_i}_{\text{concave}}) \leq \log (E_{p_i}(y_i)) \leq 0$$

$$-\ell_i - \log p_i$$

$$H(X) - E[\ell(X)]$$

$$\Rightarrow E[\ell(X)] \geq H(X)$$

$$P_i \times \left(-\log p_i \leq \lceil -\log p_i \rceil \leq -\log p_i + 1 \right)$$

$$\underbrace{-\sum p_i \log p_i}_{H(X)} \leq L \leq H(X) + 1$$

for symbol-by-symbol coding.

block coding

$$l(x_1, \dots, x_n) = \lceil -\log p(x_1, \dots, x_n) \rceil$$

$$H(x^n) \leq \sum_{x^n \in \mathcal{X}^n} p(x^n) \lceil -\log p(x_1, \dots, x_n) \rceil \leq H(x^n) + 1$$

$L(x^n)$

$$\frac{1}{n} H(x^n) \leq \frac{1}{n} L(x^n) \leq \frac{1}{n} H(x^n) + \frac{1}{n}$$

by increasing n , we can tighten
the sandwich bound as far as
we need.

Huffman Coding

$$X \quad \mathcal{X} = \{x_1, \dots, x_n\}$$

$$P_1 \geq P_2 \geq \dots \geq P_n$$

- combine two outcomes with smallest probabilities

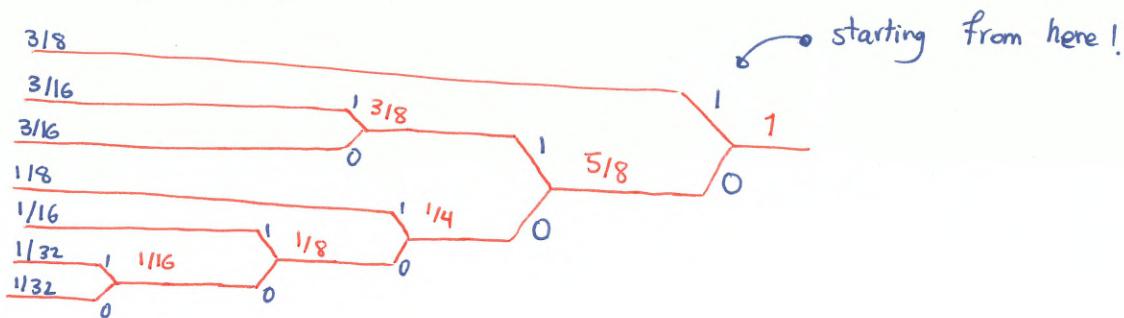
$$x_{n-1}^* \quad P_{n-1} + P_n = P(x_{n-1}^*)$$

- this will lead to a new C^* with $L^* \leq L$

ex.1

$\mathcal{X} \in \{A, \dots, G\}$

1	A	3/8
011	B	3/16
010	C	3/16
001	D	1/8
0001	E	1/16
00001	F	1/32
00000	G	1/32



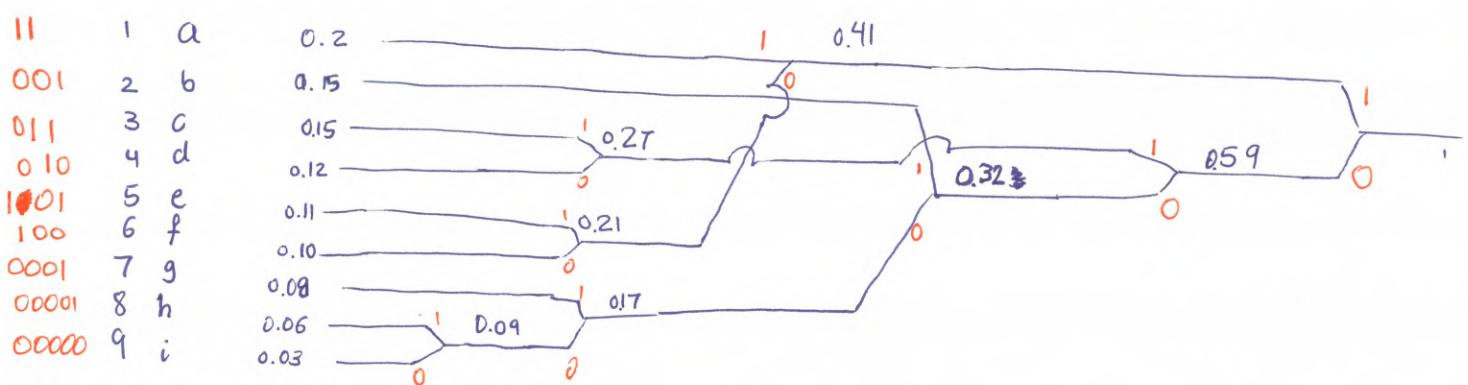
$$H(x) = 2.28 \left[\frac{\text{bits}}{\text{symbol}} \right]$$

$$\begin{aligned} \text{average codeword length} &= L(x) = 1 \cdot \frac{3}{8} + 3 \left(\frac{3}{16} + \frac{3}{16} + \frac{1}{8} \right) + 4 \cdot \frac{1}{16} + 5 \left(\frac{1}{32} + \frac{1}{32} \right) = \\ &= \frac{3}{8} + \frac{12}{8} + \frac{2}{8} + \frac{5}{16} = \frac{39}{16} = 2.4375 \left[\frac{\text{bits}}{\text{symbol}} \right] \end{aligned}$$

Info. lect.]

- if distribution changes, you have to run Huffman from scratch!
but block codes/stream codes like arithmetic coding can deal with this issue easier.

ex. 2 Huffman code



$$L = E[\ell(X)] = 2(0.2) + 3(0.15) + 4(0.12) + 5(0.08) + 5(0.09) = 3.06$$

[bit/symbol]

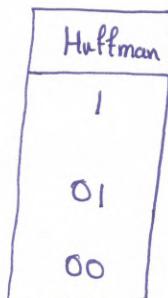
$$H(X) = 3.02 \text{ bits/symbol}$$

ex. 3 (block code)

A $\frac{3}{4}$

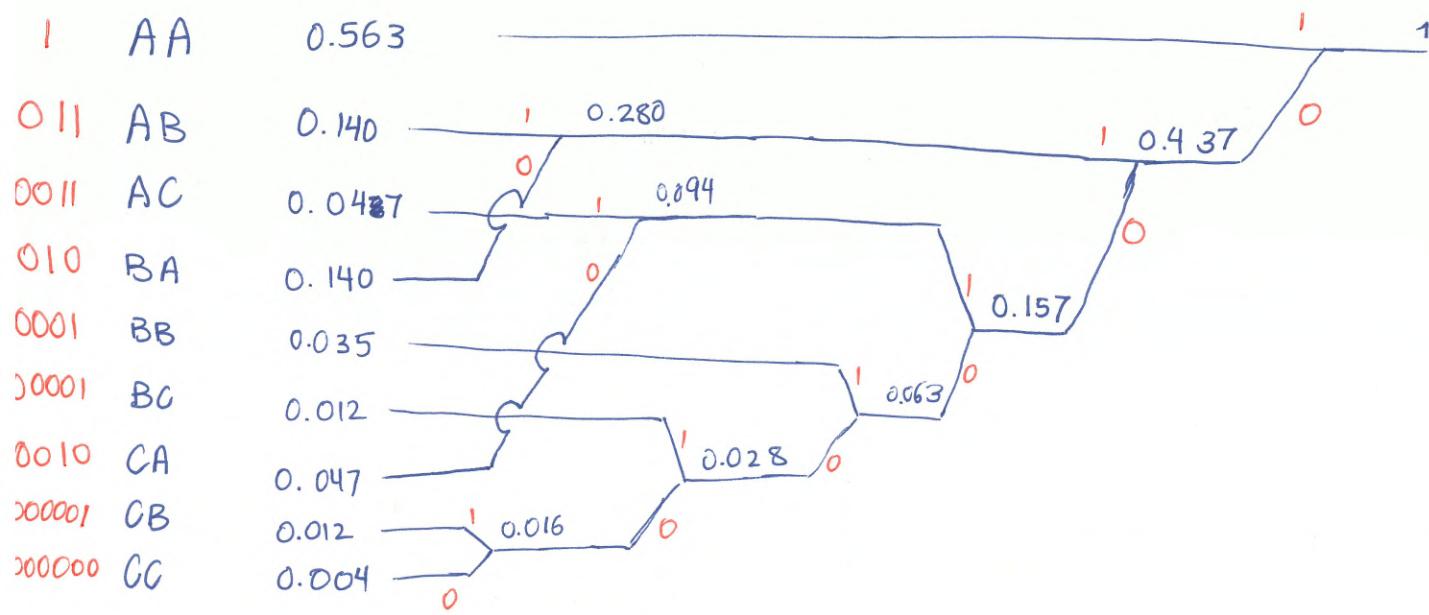
B $\frac{3}{16}$

C $\frac{1}{16}$



$$L(X) = E[\ell(X)] = 1.25, \quad H(X) = 1.012$$

what if we plan to code blocks of length two?



$$\frac{L}{2} = \frac{E[\ell(x)]}{2} = 1.037 \text{ [bits]}_{\text{Symbol}}$$

↳ closer to the Entropy!

$$\text{length of alphabet } x_i : \ell(x_i) = \lceil -\log p_i \rceil$$

$$E_{\sim q_i} \text{ or } \sum q_i \Rightarrow -\log p_i \leq \ell(x_i) \leq -\log p_i + 1$$

q_i : true pmf
 p_i : model pmf

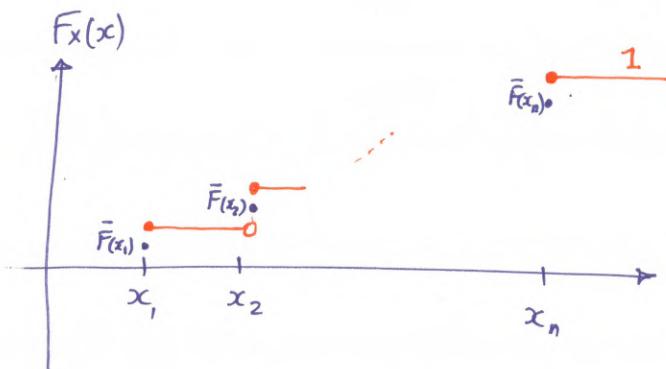
$$D(q||p) + H(x) \leq L \leq D(q||p) + H(x) + 1$$

Shannon Codes

$$\ell(x_i) = \lceil -\log p_i \rceil + 1$$

$$H(X) < L < H(X) + 2$$

$$x_1, \dots, x_n \\ p_1, \dots, p_n$$

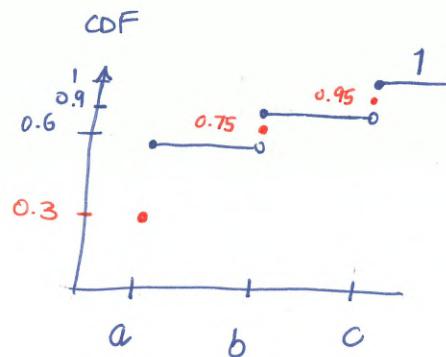


$$\bar{F}(x_i) = F(x_{i-1}) + \frac{p_i}{2}$$

in general: real-valued.

$$\begin{aligned} & \lfloor \bar{F}(x_i) \rfloor_{\ell(x_i)} && \sum_{k=1}^{\ell(x_i)} b_k (\frac{1}{2})^k && b_k \in \{0, 1\} \\ & \bar{F}(x_i) - \lfloor \bar{F}(x_i) \rfloor_{\ell(x_i)} && 0.b_1 b_2 b_3 \dots & = \sum_{k=1}^{+\infty} b_k (\frac{1}{2})^k \\ & \sum_{k=\ell(x_i)+1}^{+\infty} b_k \frac{1}{2^k} &\leq & \sum_{m=0}^{+\infty} \frac{1}{2^{m+1-\ell(x_i)}} &= \frac{1}{2^{\ell(x_i)}} = \frac{1}{2^{\lceil -\log p_i \rceil + 1}} \\ & m = k - \ell(x_i) - 1 && \leq \frac{1}{2^{-\log p_i}} \cdot \frac{1}{2} = \frac{p_i}{2} \end{aligned}$$

1	a	0.6
01	b	0.3
00	c	0.1



$$H(X) = 1.30 \quad \text{bit/sym}$$

$$\underset{\text{Huffman}}{E(\ell(x))} = 1.40 \quad \text{bit/sym}$$

$$\bar{F}(a) = 0.3 = 0. \frac{1}{2} + 1 \cdot \frac{1}{2^2} + 0 \cdot \frac{1}{2^4} + \dots$$

$$\bar{F}(b) = 0.75 = 1 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2^2} \dots$$

$$\bar{F}(c) = 0.95 = 1 \frac{1}{2} + 1 \frac{1}{2^2} + 1 \frac{1}{2^3} + \dots$$

$$\ell(a) = \left\lceil \underbrace{-\log 0.6}_{0.7370} \right\rceil + 1 = \cancel{1.7370} \quad 2$$

$$\ell(b) = \left\lceil \underbrace{-\log 0.3}_{1.7370} \right\rceil + 1 = 3$$

$$\ell(c) = \left\lceil \underbrace{-\log 0.1}_{3.3219} \right\rceil + 1 = 5$$

Shannon codes:	a	01
	b	110
	c	11110

$$\mathcal{L} = \cancel{2.6} 2(0.6) + 3(0.3) + 5(0.1) = \underline{\underline{2.6}} \quad \text{bit/sym} \quad < H(X) + 2$$

Arithmetic Coding

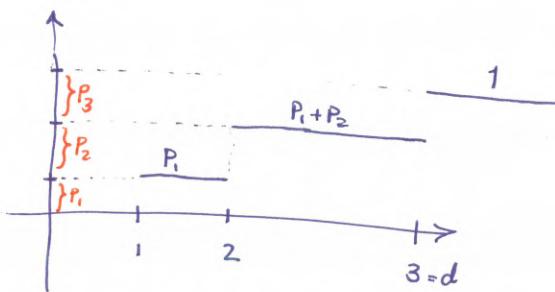
- recursive
- inspired by Shannon Coding

X

$$\mathcal{X} = \{a_1, a_2, \dots, a_d\}$$

$$P_1 \quad P_2 \quad \dots \quad P_d$$

$$d=3:$$



to represent 1 any value between $[0, P_1]$ can be selected.

to represent 2 " " " $[P_1, P_1 + P_2]$ " " "

to represent 3 " " " $[1 - P_3, 1]$ " " "

$$1 \rightarrow 0.1$$

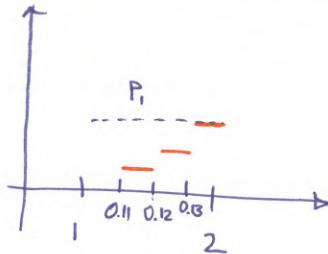
$$2 \rightarrow 0.2$$

$$3 \rightarrow 0.3$$

pair coding:

$$(x_1, x_2) \rightsquigarrow 3^2 = 9 \text{ possible combinations.}$$

0.11	$P_1 \cdot P_1$
0.12	$P_1 \cdot P_2$
0.13	$P_1 \cdot P_3$
0.21	$P_2 \cdot P_1$
0.22	$P_2 \cdot P_2$
0.23	$P_2 \cdot P_3$
0.31	$P_3 \cdot P_1$
0.32	$P_3 \cdot P_2$
0.33	$P_3 \cdot P_3$



to respresent 0.12 any value between $[P_1 P_1, P_1 P_1 + P_1 P_2]$ can be chosen.

$$F_{x^n}(x^n) = P(X_1 < x_1)$$

$$+ P(X_1 = x_1, X_2 < x_2)$$

X^n : sequence of length n .

$$+ \dots + P(X_1 = x_1, X_2 = x_2, \dots, X_n < x_n)$$

1. n -length string

$$(x_1, x_2, \dots, x_n)$$

$$\ell(x_n) = \lceil -\log p(x^n) \rceil + 1$$

2. Given an interval to encode

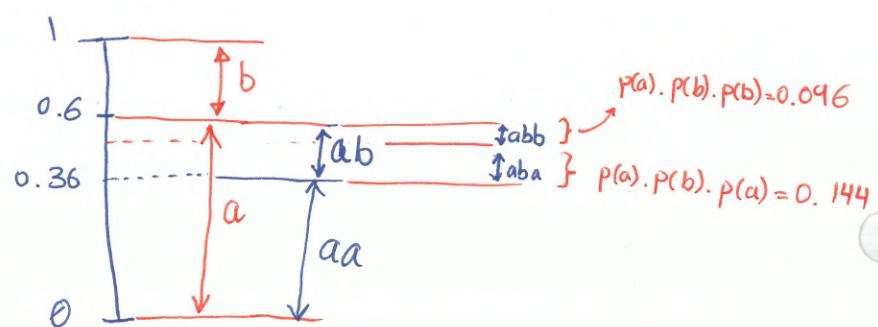
assign midpoint to x^n , encode the midpoint.

ex.

$$a < b \quad X = \begin{cases} a, & P(a) = 0.6 \\ b, & P(b) = 0.4 \end{cases}$$

$$n=3$$

$$(a, b, b)$$



$$\ell(a, b, b) = \lceil -\log 0.096 \rceil + 1 = 5 \text{ bits}$$

midpoint of $[0.504, 0.6] = 0.552$

$$0.552 = 1 \frac{1}{2} + 0 \frac{1}{4} + 0 \frac{1}{8} + 0 \frac{1}{16} + 1 \frac{1}{32}$$

length is 5 bits from previous page.

Final code of abo : 10001

Decoding: $(10001)_2 = (0.53125)_{10}$

question 1: Is 0.53125 in $[0, 0.6]$ or $[0.6, 1]$?
 $\Rightarrow \underline{a}$

question 2: Is 0.53125 in $[0, 0.36)$ or $[0.36, 0.6]$?

$\Rightarrow \underline{ab}$

question 3: Is 0.53125 in $[0.36, 0.504]$ or $[0.504, 0.6]$

$\Rightarrow \underline{abb}$

$$F_{X^3}(a, b, b) = \underbrace{P(X_1 < a)}_0 + \underbrace{P(X_1 = a, X_2 < b)}_{P(X_1=a) \cdot P(X_2 < b)} + P(X_1 = a, X_2 = b, X_3 < b)$$

memoryless source

$$+ P(X_1 = a, X_2 = b, X_3 < b)$$

$$0.6 \times 0.4 \times 0.6 = 0.504$$

* Mackay Laplace / Dirichle model
to have adaptive conditional probability model.

arithmetic:

$$L \leq H(X^n) + 2$$

as $n \rightarrow \infty$ arithmetic buys us a lot in comparison to Huffman.

~~by having n. of limits arithmetic better.~~

Project.

}

2. texts

2 p and q

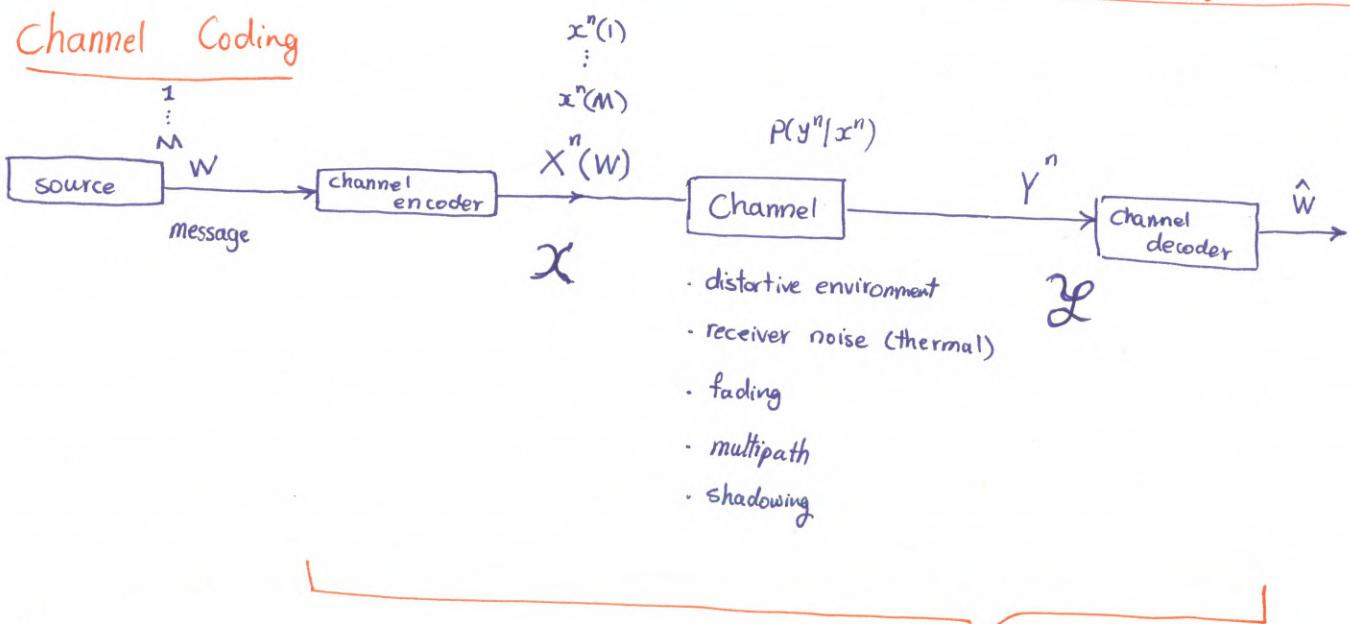
chapters 5 and 6.

Info. lect.

project { 2 texts: Alice in Wonderland + sonnet 8 of Shakespeare
 Huffman + Arithmetic + Lempel Ziv (universal coding)

area of research: Joint Source/Channel Coding.

10, 11, 22 - 1



the whole block is ideally an error-free channel!

* if channel is memoryless: $P(y^n|x^n) = \prod_{i=1}^n P(y_i|x_i)$

ex, $n=1000$

$$P(\text{flip}) = P(Y \neq X) = 0.1$$

intuition: about 900 bits will be transmitted error-free.

capacity: $1 - H_b(0.1) = * 0.531$ bits/channel use

measure of interest: Mutual Information

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

Concave in $P_X(x)$

* we can maximize $I(X; Y)$.

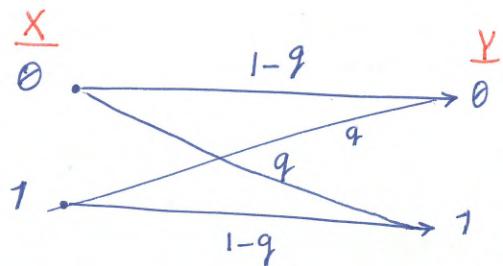
10, 13, 22 - 1

2 important channels : {

- Discrete channel
- Gaussian channel

Discrete channels

1. Binary Symmetric Channel



$$\begin{aligned} P_X(0) & \quad P_Y(y) = \sum_{x_i \in \mathcal{X}} P(y|x_i) p(x_i) \\ P_X(1) & \end{aligned}$$

$$\begin{bmatrix} P_Y(0) \\ P_Y(1) \end{bmatrix} = \begin{bmatrix} P_{Y|X}(0|0) & P_{Y|X}(0|1) \\ P_{Y|X}(1|0) & P_{Y|X}(1|1) \end{bmatrix} \begin{bmatrix} P_X(0) \\ P_X(1) \end{bmatrix}$$

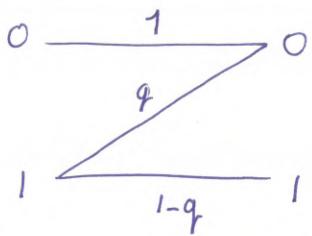
Transition
Matrix

$$\begin{bmatrix} 1-q & q \\ q & 1-q \end{bmatrix}$$

Transition Matrix

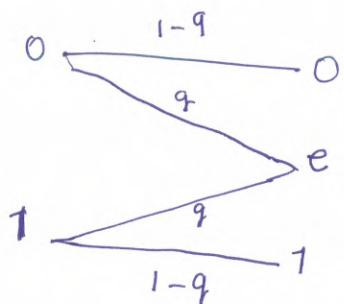
of Binary Symmetric
Channel

2. \mathbb{Z} - channel



$$\begin{bmatrix} 1 & 0 \\ q & 1-q \end{bmatrix}$$

3. Erasure Channel

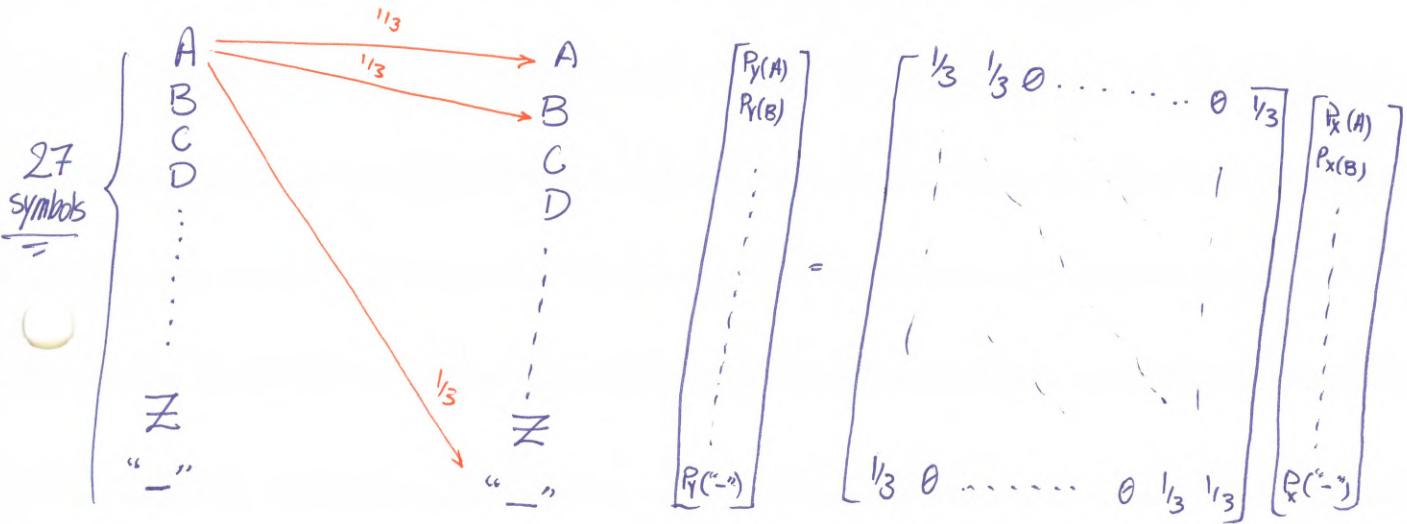


$$\begin{bmatrix} P_Y(0) \\ P_Y(e) \\ P_Y(1) \end{bmatrix} = \begin{bmatrix} 1-q & 0 & 0 \\ q & q & q \\ 0 & 1-q & 1-q \end{bmatrix} \begin{bmatrix} P_X(0) \\ P_X(1) \end{bmatrix}$$

3×2

4. Type Writer Channel

A- \mathbb{Z} but any other alphabet is valid.



Information in the Channel

Def'n: $C = \max_{P_X(x)} I(X; Y)$

ex.1 binary symmetric

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

$$H(Y|X) = H(Y|X=0) P(X=0) + H(Y|X=1) P_X(1)$$

$$= \underbrace{\left\{ -\log(1-q) \cdot (1-q) - q \log q \right\}}_{H_b(q)} P_X(X=0) + H_b(q) P_X(X=1)$$

$$= H_b(q)$$

$$I(X; Y) \leq 1 - H_b(q)$$

equality achieved by $P_Y(0) = P_Y(1) = \frac{1}{2}$

ex.2 (\mathbb{Z}_c -channel)

○ $I(X;Y) = H(Y) - H(Y|X)$

$$\begin{aligned} H(Y|X) &= H(Y|X=0) \underbrace{P(X=0)}_P + H(Y|X=1) \underbrace{P(X=1)}_{1-P} \\ &= (- (1-q) \log(1-q) - q \log q) P \\ &= H_b(q) \cdot P \end{aligned}$$

$P_Y(0) = 1 - P + qP$

$P_Y(1) = (1-q)P$

$$H(Y) = -(1-P+qP) \log_e(1-P+qP) - (1-q)P \log_e(1-q)P$$

$$\begin{aligned} \frac{\partial I}{\partial P} &= (1-q) \cancel{\log(1-P+qP)} + (1-q) \\ &\quad - (1-q) \log(1-q)P + q - H_b(q) \end{aligned}$$

$$1-P+qP = P(1-q) \left(e^{(H_b(q)-1)/(1-q)} \right)$$

$$\Rightarrow P^* = \frac{1}{(1-q) \left(1 + e^{\frac{H_b(q)-1}{1-q}} \right)}$$

$$\text{for } q = \frac{1}{2} \Rightarrow C(q = \frac{1}{2}) = \log 5 - 2 \text{ bits/channel use}$$

look at Cover & Thomas!

Absent on 10, 18, 22.

info. lect.

Ex. from Mackay

$$X \sim \text{iid Bernoulli}(0.1)$$

B.S.C. $q = 0.2$

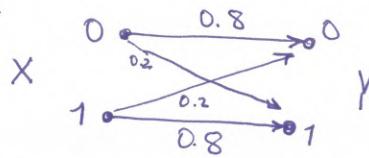
$n = 100$

$X^{100}:$

$\underbrace{11\dots11}_{10} \underbrace{0\dots0}_{T_2} \underbrace{0\dots0}_{18}$

\downarrow

$\underbrace{11\dots11}_{8} \underbrace{00}_{2} \underbrace{0\dots0}_{T_2} \underbrace{1\dots1}_{18}$



$$P(Y=0) = 0.1 \times 0.2 + 0.9 \times 0.8 = 0.74, P(Y=1) = 0.26$$

70 ones out of 100 symbols: ϵ -typical!
(empirical dist. in agreement with true dist.)

ϵ -typical! (empirical dist. in agree. with true dist.)



Jointly Typical Sequence! (chapters 9, 10 of Mackay!)

Channel Coding Theorem

Thm. All rates below capacity, $R < C$ are achievable.

If there is a sequence of $\frac{(M, n)}{2^{nR}}$ codes with $\lambda_{\max}^{(n)} \rightarrow 0$, as $n \rightarrow \infty$,

then the rate $R < C$.

Proof.

Random Coding

1. Pick $p(x)$

$2^{nR} = M$
↓
 M messages/labels

Codeword \xrightarrow{C}

$$C = \begin{bmatrix} x_1(1) & \dots & x_n(1) \\ \vdots & \ddots & \vdots \\ x_1(2^{nR}) & \dots & x_n(2^{nR}) \end{bmatrix}$$

iid sequences

2. Share C with decoder.

3. Pick at random $W=w$

$$P(W=w) = \frac{1}{2^{nR}} = \frac{1}{M}$$

transmit $X^n(w)$.

4. At the decoder, Y^n

$$P(Y^n = y^n | X^n(w)) = \prod_{i=1}^n P(Y_i = y_i | X(w))$$

5. Decode by joint typicality

$$g(Y^n | y^n) = \hat{\omega}$$

$(X^n(\hat{\omega}), Y^n) \in A_\epsilon^{(n)}$ this sequences are jointly Typical.

Analysis

Generate many codes, iid

$$\text{Error} = \{ \hat{W} \neq W \}$$

$$P(\text{Error}) = \sum_{\substack{C \\ \text{all codes}}} P(C) P_e^{(n)}(C) \xrightarrow{\text{ }} \frac{1}{2^{nR}} \sum_{w=1}^{2^{nR}} \lambda_w(C)$$

$$= \frac{1}{2^{nR}} \sum_{i=1}^{2^{nR}} \sum_c P(c) \lambda_w(c)$$

due to symmetry

$$= \sum_c P(c) \lambda_1(c)$$

$$= P(\text{Error} | W=1)$$

Event

$$E_i = \left\{ (X^n(i), Y^n) \in A_\epsilon^{(n)} \right\}$$

$$\begin{aligned} \{\text{Error } | W=1\} &= E_1^c \cup E_2 \cup E_3 \cup \dots \cup E_{(2^{nR})=M} \\ &\quad \xrightarrow{\text{complement}} \\ P(\text{Error} | W=1) &= P(E_1^c) \xrightarrow{\text{Union Bound}} \\ &\leq P(E_1^c) + \sum_{i=2}^{M=2^{nR}} P(E_i) \\ &\quad \xrightarrow{\text{very small}} \\ &\leq 2^{-n(I(X;Y)-3\epsilon)} \end{aligned}$$

$$P(\text{Error} | W=1) \leq \epsilon + 2^{nR-n(I(X;Y)-3\epsilon)}$$

$$\text{If } R < I(X;Y) - 3\epsilon \text{ then } 2^{-n(I(X;Y)-3\epsilon)} \leq \epsilon$$

$$\Rightarrow P(\text{Error} | W=1) \leq 2\epsilon$$

$$\lambda_{\max}^{(n)} \leq 4\epsilon$$

Pick $p(x) = p^*(x)$

$$R < C - 3\epsilon$$

$$\frac{M}{2} = 2^{nR-1}$$

$$R - \frac{1}{n} < C - 3\epsilon$$

Lemma 1 For the n^{th} extension (using n times) of channel

$$I(X^n; Y^n) \leq nC$$

(for any $p(x^n)$)

Lemma 2
(Fano's Inequality)

$$H(X^n | Y^n) \leq 1 + P_e^{(n)} \cdot nR$$

$$nR = H(w) = H(w|Y^n) + I(w; Y^n)$$

$$I(w; Y^n) = H(w) + H(w|Y^n)$$

data processing
inequality

$$\leq \underbrace{H(W|X^n)}_{H(X^n(w)|Y^n)} + I(X^n(w); Y^n)$$

Fano's
Inequ.

$$\leq 1 + P_e^{(n)} \cdot nR + nC$$

$\Rightarrow R \leq \frac{1}{n} + P_e^{(n)} \cdot R + C$

 \downarrow
 $\hookrightarrow 0$
 as $n \rightarrow \infty$

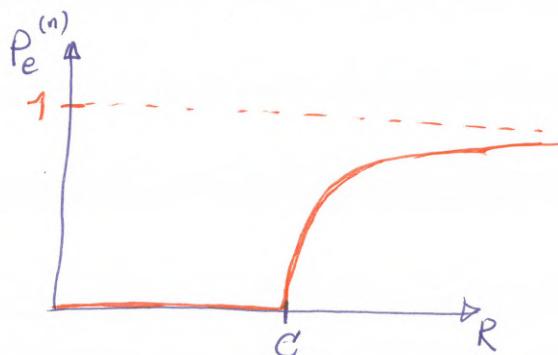
$\Rightarrow R \leq \frac{1}{n} + C$

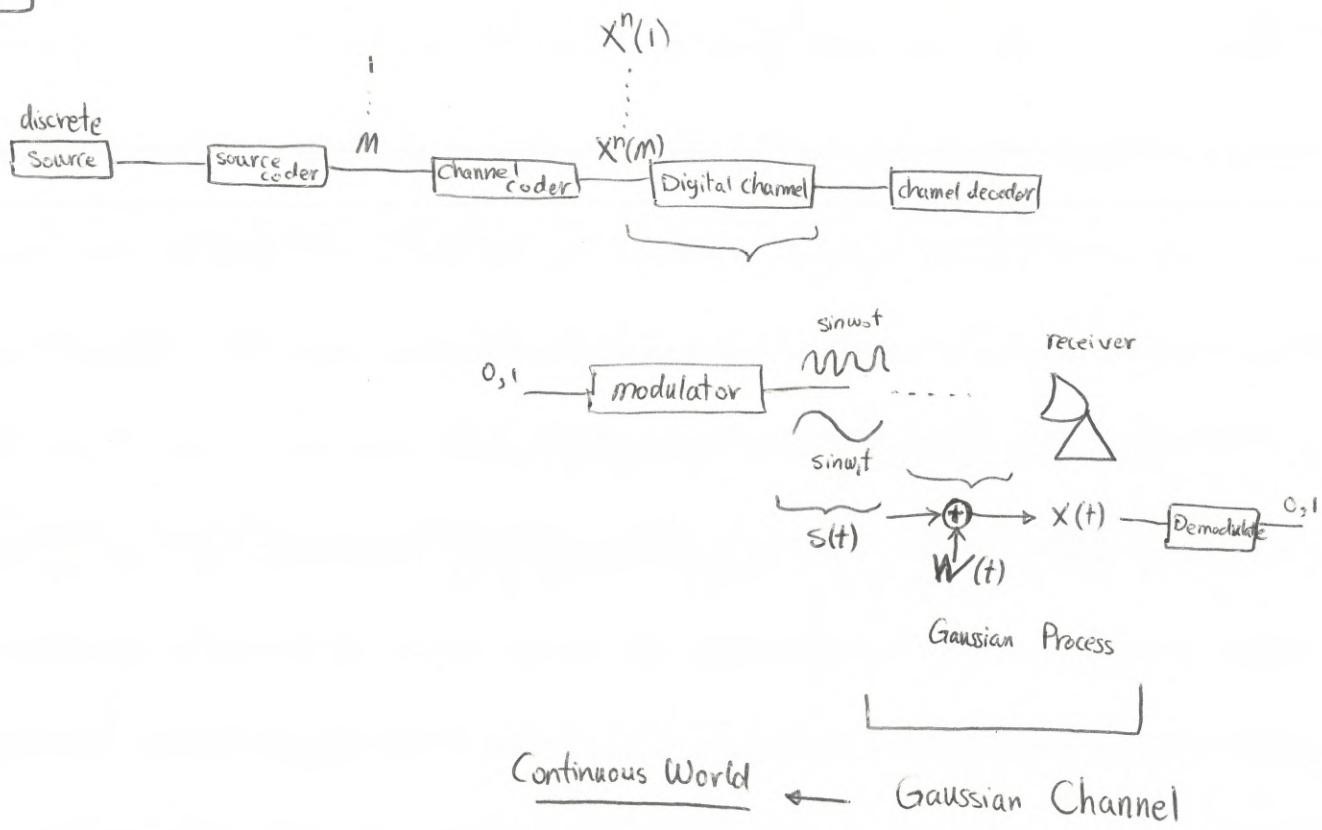
Another Way to Write

$$P_e^{(n)} \geq 1 - \frac{1}{nR} - \frac{C}{R}$$

if $R > C$ then $P_e^{(n)} \rightarrow 1$.

look at Mackay
for more illustrations
+ intuitions!





Differential Entropy

X is continuous R.V. with pdf $f_X(x)$

$$f_X(x) = \frac{d}{dx} F_X(x)$$

\hookrightarrow cdf

Defn

$$h(X) = E \left[\log_2 \frac{1}{f_X(x)} \right] = \int_{\substack{x: f_X(x) > 0}} f_X(x) \log_2 \frac{1}{f_X(x)} dx$$

ex.1

$$X \sim \text{unif}(\theta, a)$$

$$h(x) = + \int_{\theta}^a \log_2 a \times \frac{1}{a} dx = \log a$$

- [If $a > 1$, then $h(x) > 0$]
- [If $0 < a < 1$, then $h(x) < 0$]

\in typical sets:

$$A_{\epsilon}^{(n)} \xrightarrow[\text{size of set}]{} 2^n H(x)$$

continuous set $n=1$: $A_{\epsilon}^{(1)} \xrightarrow[\text{size}]{} 2^{h(x)}$

Volume of set! (not size)
(always positive)

ex.2

$$X \sim N(\mu, \sigma^2)$$

$$h(x) = -E[\ln f_X(x)] = E\left[-\frac{1}{2\sigma^2}(x-\mu)^2 - \frac{1}{2}\ln(2\pi\sigma^2)\right]$$

$$= \frac{1}{2} \cancel{\ln e} + \frac{1}{2} \ln(2\pi\sigma^2) = \frac{1}{2} \ln(2\pi e \sigma^2) \quad \text{nats/symbol}$$

AEP

$X_1, \dots, X_n \sim \text{iid from } f_X(x)$,

then $-\frac{1}{n} \log f_{X^n}(x^n) \xrightarrow[\substack{\text{converge} \\ \text{in} \\ \text{probability}}]{P} h(x)$

(proof by using WLLN)

Typical Set

$$A_\epsilon^{(n)} = \left\{ x^n : \left| -\frac{1}{n} \log f_{X^n}(x^n) - h(x) \right| < \epsilon \right\}$$

$$1) P(A_\epsilon^{(n)}) \rightarrow 1 \quad \text{as } n \rightarrow \infty$$

$$2) \text{vol}(A_\epsilon^{(n)}) \leq 2^{n(h(x) + \epsilon)}$$

$$3) \text{vol}(A_\epsilon^{(n)}) \geq (1 - \epsilon) 2^{n(h(x) - \epsilon)}$$

set

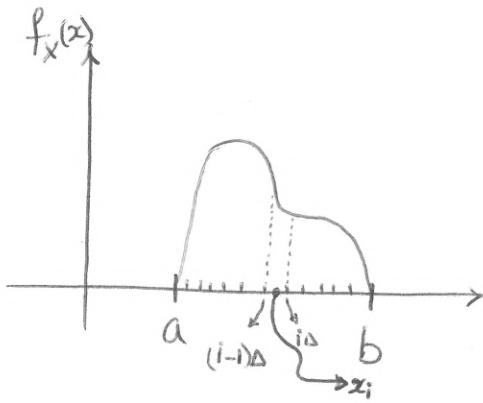
★ $\text{vol}(A_\epsilon^{(n)}) \underset{\text{OR}}{\approx} 2^{nh(x)}$

(≈)

Volume of a cube $\rightsquigarrow \ell^n \approx 2^{nh(x)}$

$$2^{n \log e} \approx 2^{nh(x)}$$

$$\ell = 2^{h(x)} \stackrel{\text{uniform } X}{=} 2^{\log a} = a$$



n intervals of length Δ

mean value theorem

$$\int_{(i-1)\Delta}^{i\Delta} f_X(x) dx = f(x_i) \cdot \Delta$$

$$X^\Delta = \left\{ x_i \quad \text{with} \quad p(x_i) = f(x_i) \cdot \Delta \right.$$

$$H(X^\Delta) = - \sum_i f_X(x_i) \Delta \times \underbrace{\log(f_X(x_i) \Delta)}_{\log f_X(x_i) + \log \Delta}$$

If $\Delta \rightarrow 0$
 $n \rightarrow +\infty$

$$\Rightarrow - \sum_i f_X(x_i) \overbrace{\log f_X(x_i)}^dx \xrightarrow{\text{red}} h(x)$$

$$-\log \Delta \sum f_X(x_i) \Delta \rightarrow -\log \Delta$$

as $n \rightarrow +\infty$
 $\Delta \rightarrow 0$

$$H(X^\Delta) \cancel{+} + \log \underbrace{\Delta}_{2^n} \xrightarrow{\text{number of divisions}} h(x)$$

ex. $X \sim \text{unif}(0, 1/2)$

$$h(x) = \log a = -1$$

division step : $\Delta = 2^{-4}$

$$H(X^\Delta) = -1 + 4 = 3$$



* as we decrease Δ , number of bits/symbol goes to ∞ !!!

Defn Relative Entropy $D(f||g) = \int f_X(x) \log \frac{f_X(x)}{g_X(x)} dx$

Jensen's Inequality $D(f||g) \geq 0$

$$I(X;Y) = \iint f_{XY}(x,y) \log \frac{f_{XY}(x,y)}{f_X(x)f_Y(y)} dx dy \Rightarrow I(X;Y) \geq 0$$

$$I(X;Y) = h(Y) - h(Y|X) = h(Y) + h(X) - h(X,Y)$$

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

lect.
info.

ex.

$$\underline{x} = \underline{X}^n \sim N(\underline{\mu}, K_{n \times n})$$

$$h(\underline{x}) = E[-\log_e f_{\underline{x}}(\underline{x}^n)]$$

$$f_{\underline{x}^n}(\underline{x}^n) = \frac{1}{\sqrt{(2\pi)^n \det(K)}} e^{-\frac{1}{2} (\underline{x}^n - \underline{\mu})^T K^{-1} (\underline{x}^n - \underline{\mu})}$$

$$= -E\left[-\frac{1}{2} (\underline{x}^n - \underline{\mu})^T K^{-1} (\underline{x}^n - \underline{\mu})\right] - \frac{E}{2} \left[-\ln(2\pi)^n \det(K)\right]$$

$$= \frac{1}{2} \underbrace{\text{Tr}(K K^{-1})}_{\parallel} + \frac{1}{2} \underbrace{\ln(2\pi)^n \det(K)}_{\text{I}}$$

$$\underbrace{\frac{n}{2} = \frac{1}{2} \ln e^n}_{\text{II}}$$

$$= \frac{1}{2} \ln \left((2\pi e)^n \det(K) \right)$$

Proposition

$$E[\underline{x}^n] = \emptyset$$

$$E[\underline{x}^n (\underline{x}^n)^T] = K$$

show: $h(\underline{x}^n) \leq \frac{1}{2} \ln ((2\pi e)^n \det(K))$

if $\underline{x}^n \sim N(\emptyset, K)$

Proof: $D(f_{x^n} \| g_{x^n}^{\text{pdf of jointly Gaussian}})$

$$= \mathbb{E}_{f_{x^n}} \left[\ln \frac{f_{x^n}}{g_{x^n}} \right] \geq 0$$

~~skipped~~

$$-h_{f_{x^n}}(x^n) + \underbrace{\mathbb{E}_{f_{x^n}} \left(\ln \frac{1}{g_{x^n}} \right)}_{\frac{1}{2} \ln [(2\pi e)^n \det(K)]} \geq 0 \quad \textcircled{D}$$

$$\mathbb{E}_{f_{x^n}} \left[\frac{1}{2} x^{n^T} K^{-1} x^n + \frac{1}{2} \ln ((2\pi)^n \det(K)) \right]$$

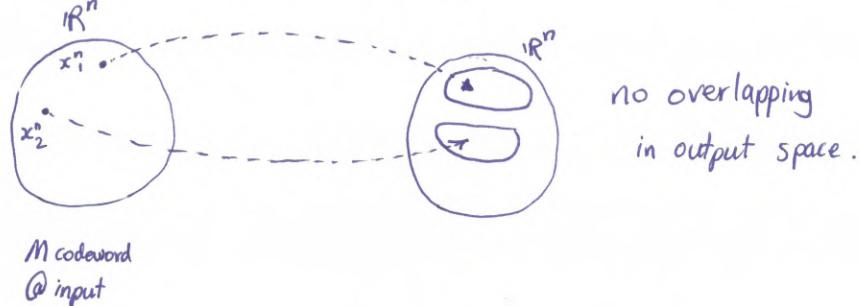
$$= \frac{1}{2} \text{Tr}(K K^{-1}) + \frac{1}{2} \ln (2\pi)^n \det(K)$$

Gaussian Channel



• power constraint: $E[X^2] \leq P$

$$Z \sim N(0, N)$$



Defn The Capacity of
 (Theoretical Capacity)

$$C = \max_{\substack{f_X(x): \\ E[X^2] \leq P}} I(X; Y) \quad (\text{assumption: Gaussian Channel})$$

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

$$\downarrow \int f_X(x) \left(\int f_{Y|X}(y|x) \log \frac{1}{f_{Y|X}(y|x)} dy \right) dx$$

for $X = x$

$$Y = x + Z \quad \xrightarrow{\quad} \quad \frac{1}{2} \ln (2\pi e N)$$

\Downarrow

$Y \sim N \Leftarrow x \text{ is fixed}$

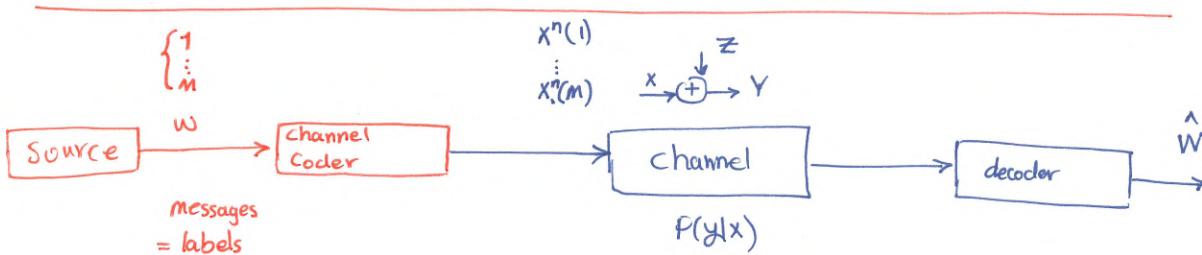
$$\Rightarrow I(X; Y) = h(Y) - \frac{1}{2} \ln (2\pi e N)$$

$$h(Y) - \frac{1}{2} \ln(2\pi e N) \leq \frac{1}{2} \ln(2\pi e (N+P)) - \frac{1}{2} \ln(2\pi e N)$$

to maximize
if $X \sim N(\theta, P)$

$$\frac{1}{2} \ln \left(1 + \frac{P}{N} \right)$$

$$\Rightarrow C = \frac{1}{2} \ln \left(1 + \frac{P}{N} \right)$$



Defn: ~~(M, n)~~ power constrained code. consist of: ~~set of messages~~

1. a set messages ~~labels~~
2. an encoder $f: w \rightarrow x^n(w)$

$$w: 1 \dots M$$

3. ensure power constraint ~~labels~~

$$\frac{1}{n} \sum_{i=1}^n x_i^2(w) \leq P$$

4. an encoder $g: Y^n \rightarrow \hat{W}$

Rate of code:

R is achievable ~~with error probability~~

$$\left\{ \begin{array}{l} R = \frac{\log M}{n} \\ M = 2^{nR} \end{array} \right.$$

if there is a sequence of (M, n) power constrained codes with $\lambda^{(n)} \rightarrow 0$ as $n \rightarrow \infty$
 max. Probability of error

operational capacity is the maximum of achievable R .

ϵ -typical set on the output of the Gaussian channel

$$A_\epsilon^{(n)} = \left\{ y^n : \left| -\frac{1}{n} \sum \log f_{Y^n}(y^n) - h(Y) \right| < \epsilon \right\}$$

$$Y_i \sim N(\theta, P+N)$$

$$\sum_{i=1}^n \ln f_{Y_i}(y_i)$$

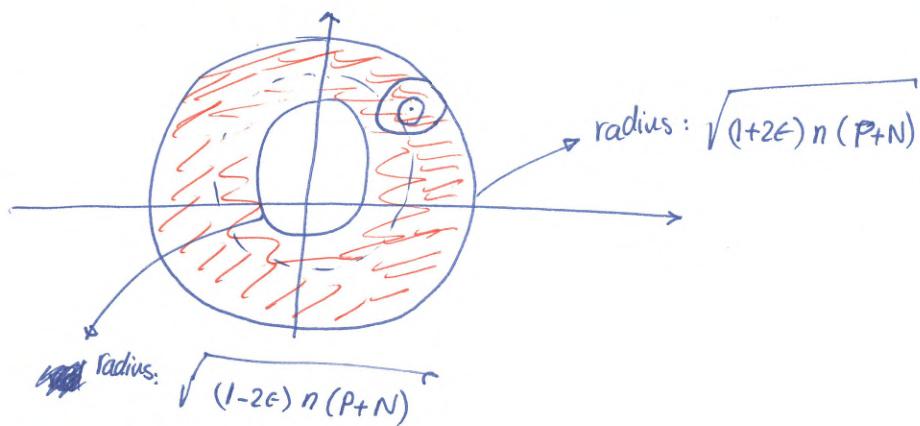
$$-\sum_{i=1}^n \left[\frac{1}{2} \frac{(y_i)^2}{(P+N)} + \frac{1}{2} \ln (2\pi(P+N)) \right]$$

$$\left| \frac{1}{2n(P+N)} \sum_{i=1}^n y_i^2 + \underbrace{\frac{1}{2} \ln(2\pi(P+N)) - \frac{1}{2} \ln(2\pi e(P+N))}_{-\frac{1}{2}} \right| \leq \epsilon$$

$$(1-2\epsilon)n(P+N) \leq \underbrace{\sum_{i=1}^n y_i^2}_{r^2} \leq (2\epsilon + 1)n(P+N)$$

Sphere in n -dim space.
with radius r^2

~~2ε n(P+N)~~



For a given $\underline{x}(w)$, the output $y_i \sim N(x_i(w), N)$

↳ much shallower!!!

info. lect.

$$z_i \sim N(0, N)$$

$$x_i \xrightarrow{+} y_i$$

$$E[x^2] \leq P$$

$$\text{entropy} \rightarrow h(Y)$$

this will enforce the partition on the output, unless, this problem would be unconstrained with no partitions.

$$h(Y) = \frac{1}{2} \ln(2\pi(N+P))$$

capacity:

$$\max_{\substack{\text{f}_X: E[x^2] \leq P}} I(X; Y) = \frac{1}{2} \ln \left(1 + \frac{P}{N} \right)$$

(M, n) - code

$$\text{with } \frac{1}{n} \sum_{i=1}^n x_i^2(\omega) \leq P$$

encoding function

$$f: \omega \xrightarrow{\text{label}} X^n(\omega)$$

• decoding function *by joint typicality*

$$\begin{array}{c} Y^n \\ \xrightarrow{g} \\ \theta(Y^n) = \hat{\omega} \end{array}$$

$$R = \frac{\log M}{n}$$

$$A^{(n)}(Y^n) = \left\{ Y^n : \left| \frac{1}{n} \ln f_{Y^n}(Y^n) - h(Y) \right| < \epsilon \right\}$$

$$Y_i \sim N(0, P+N)$$

$$-\frac{1}{n} \left(-\frac{1}{2(P+N)} \sum y_i^2 - \frac{n}{2} \ln(2\pi(P+N)) - h(Y) \right)$$

$$\frac{1}{2} \ln(2\pi e(P+N))$$

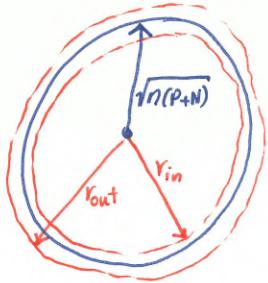
$$\Leftrightarrow \frac{1}{2n} \frac{1}{P+N} \sum_{i=1}^n y_i^2 + \frac{1}{2} \ln(2\pi(P+N)) - \frac{1}{2} \ln(2\pi e(P+N))$$

$$-\frac{1}{2} \ln e = -\frac{1}{2}$$

$$(P+N)2n\left(\frac{1}{2} - \epsilon\right) < \sum y_i^2 < (\epsilon + \frac{1}{2})2n(P+N)$$

$$\sum_1^n y_i^2 = r^2$$

Sphere of n -dimensional with radius r



$$r_{out} = \sqrt{n(1+2\epsilon)(N+P)}$$

$$r_{in} = \sqrt{n(1-2\epsilon)(N+P)}$$

Fix $X^n = x^n(\omega)$

$$Y^n = x^n(\omega) + Z^n$$

$$Y_i = x_i(\omega) + Z_i$$

$\sum_{i=1}^N \mathcal{N}(x_i(\omega), N)$

$$A_\epsilon^{(n)}(y^n \mid x^n(\omega)) = \left\{ y^n : \left| -\frac{1}{n} \log \left(f_{Y^n \mid X^n}(y^n \mid x^n(\omega)) \right) - h(Y^n \mid X^n(\omega)) \right| < \epsilon \right\}$$

empirical entropy

$$\left| \frac{1}{2n} \sum_{i=1}^N (y_i - x_i)^2 + \frac{1}{2} \ln(2\pi N) - \frac{1}{2} \ln(2\pi e N) \right| < \epsilon$$

$$(\frac{1}{2} - \epsilon)2nN < \sum (y_i - x_i)^2 < (\epsilon + \frac{1}{2})2nN$$

$r = \sqrt{nN}$

$r_{out} = \sqrt{(2\epsilon + 1)nN}$

$r_{in} = \sqrt{(1-2\epsilon)nN}$

$$\text{Volume}_{\text{Sphere}}(r^n) = C_n \cdot r^n$$

↙
just a positive number

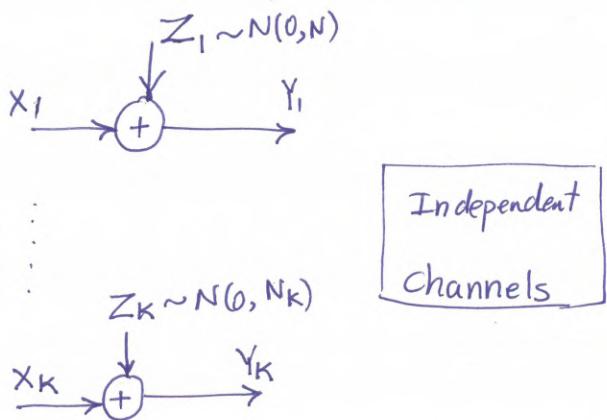
(sphere-packing)

$$\frac{C_n \cdot \{n(P+N)\}^{n/2}}{C_n \cdot \{nN\}^{n/2}} = \left(1 + \frac{P}{N}\right)^{n/2} = M = 2^{nR} = 2^{nC}$$

Capacity ↘

$$\Rightarrow C = \frac{1}{2} \ln \left(1 + \frac{P}{N}\right)$$

Parallel Gaussian Channels (sth like Eigenvalue Decomposition on Images)



$$\sum_{i=1}^n P_i \leq P$$

↓ all power to
be distributed over
the channels

$E[X_i^2]$

$$\max I(X^k, Y^K) = ?$$

$$\int$$

$$h(Y^K) - \underbrace{h(Y^K | X^K)}_{h(Z^K)}$$

$$= \underbrace{h(Y^K)}_{\sum_{i=1}^K h(z_i)} - \sum_{i=1}^K \underbrace{h(z_i)}_{\frac{1}{2} \ln(2\pi e N_i)}$$

$$\leq \sum_{i=1}^K h(y_i) \quad \hookrightarrow N(0, P_i + N_i)$$

$$= \sum_{i=1}^K \frac{1}{2} \ln(2\pi e(N_i + P_i)) - \sum_{i=1}^K \frac{1}{2} \ln(2\pi e N_i).$$

$$= \frac{1}{2} \sum_{i=1}^K \ln\left(1 + \frac{P_i}{N_i}\right)$$

we will use the maximum available power.

$$\sum_{i=1}^K P_i = P$$

$$\rightarrow \max_{P_i: \sum_{i=1}^K P_i = P, P_i > 0} \frac{1}{2} \sum \ln\left(1 + \frac{P_i}{N_i}\right) \quad \text{Lagrange Multiplier}$$

$$J(P_1, \dots, P_K, \lambda, \mu_1, \dots, \mu_K)$$

$$= \frac{1}{2} \sum_{i=1}^K \ln\left(1 + \frac{P_i}{N_i}\right) + \lambda \left(\sum_{i=1}^K P_i - P \right) + \sum_{i=1}^K \mu_i P_i$$

Kuhn Tucker
2 cases

$$(\sum \mu_i P_i = 0)$$

$P_i > 0, \mu_i = 0$

$P_i = 0, \mu_i < 0$

why?

(I) $P_i > 0, \mu_i = 0$

$$\frac{\partial J}{\partial P_i} = 0 \Rightarrow \frac{1}{2} \frac{1}{P_i + N_i} + \lambda = 0 \Rightarrow P_i + N_i = \frac{1}{-2\lambda}$$

hard to describe

$$\lambda > 0 \Rightarrow (\lambda < 0)$$

$$\Rightarrow P_i = \lambda - N_i$$

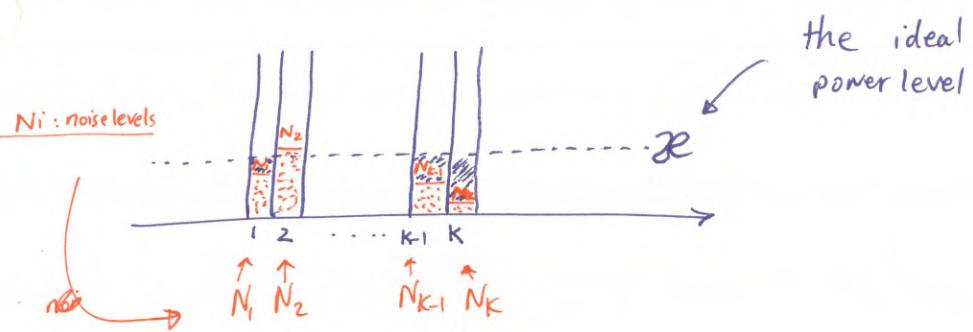
(II) $P_i = 0, \mu_i < 0$

$$\frac{\partial J}{\partial P_i} = 0 \Rightarrow \frac{1}{2} \frac{1}{P_i + N_i} + \lambda + \mu_i \Big|_{P_i=0} = 0 \Rightarrow \mu_i = -\underbrace{\frac{1}{2(P_i + N_i)} \Big|_{P_i=0}}_{< 0} - \lambda < 0$$

$$\Rightarrow \mu_i = -\frac{1}{2N_i} - \lambda < 0 \rightarrow N_i > \frac{-1}{2\lambda}$$

So if noise level is above λ we don't put any power into that channel.

$$\sum_{i=1}^K (x - N_i)^+ = P$$



Bounds on Population of Binary Codes

$$M = 2^{nR}$$

↳ # of distinct codewords

$$\rightarrow \mathcal{X} = \{0, 1\}$$

→ length n codewords.

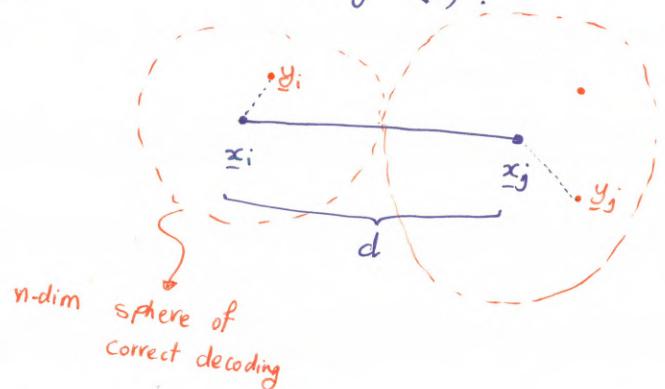
→ 2^n is total number of sequences.

→ Form a code C s.t. $d(\underline{x}_i, \underline{x}_j) \geq d$ (i)

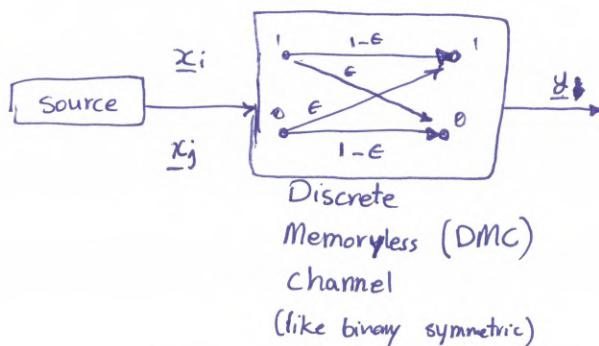
distance
between
a pair of
codewords

Goal: To Find the

max population of
a binary code C ,
with code words $\{\underline{x}_i\}$
satisfying (i).



→ channel may introduce distortion.



- choose $d(\underline{x}_i, \underline{x}_j)$ as Hamming distance.

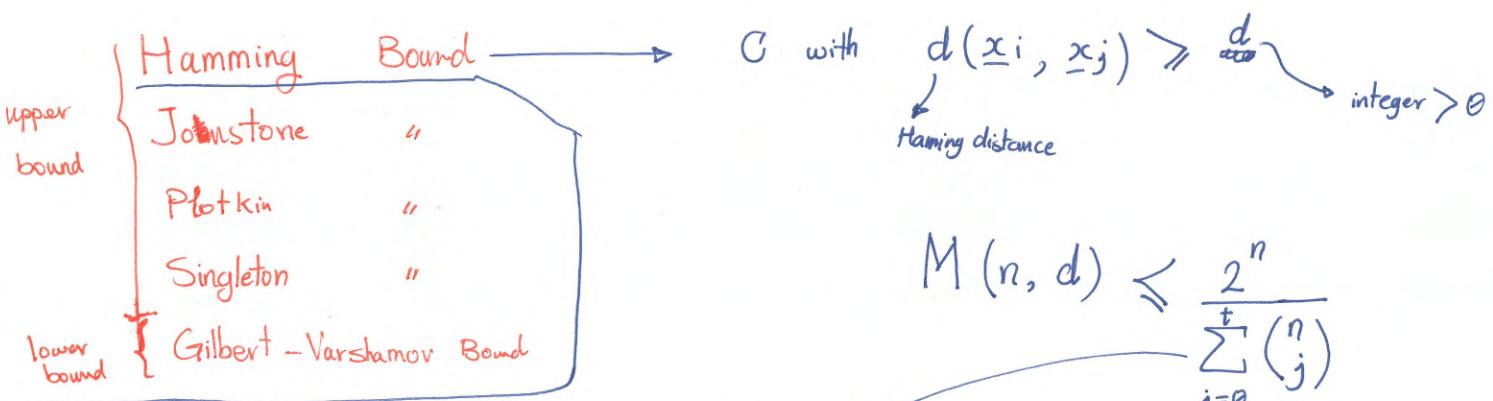
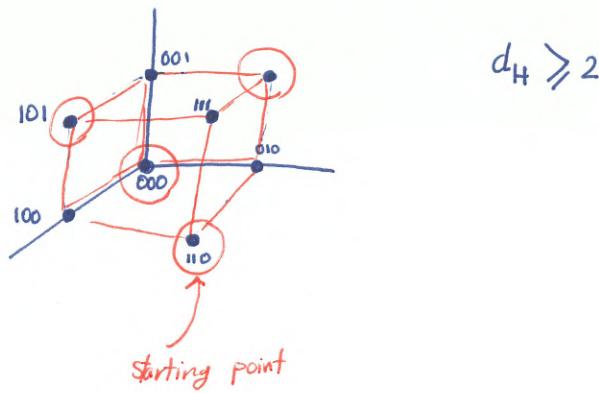
ex. $\begin{array}{r} 1100 \\ 1000 \end{array} \rightarrow d_H = 1$

modular addition

$$d_H = \sum_k x_{i,k} \oplus x_{j,k}$$

- Decode \underline{y} as \underline{x}_i if $d_H(\underline{x}_i, \underline{y}) < \frac{d}{2}$

$N=3$ → 8 possible codewords.

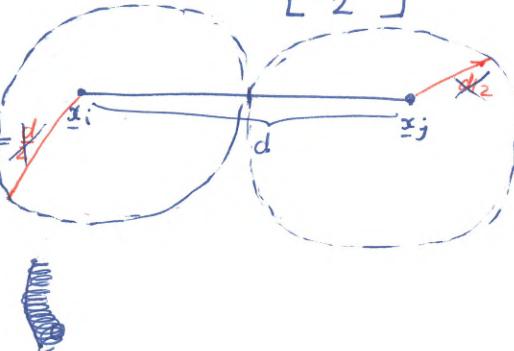


binary sequences that differ from \underline{x}_i by j bits.

$$\therefore \binom{n}{0} + \binom{n}{1} + \binom{n}{2} + \dots + \binom{n}{t}$$

$\xrightarrow{\substack{\text{1 bit difference to } \underline{x}_i \\ \underline{x}_i \text{ itself}}}$

$$\text{with } t = \left\lfloor \frac{d-1}{2} \right\rfloor$$



Sphere Packing Bound

Plotkin Bound

1) If d is even and $2d > n$, then

$$M(n, d) \leq 2 \left\lfloor \frac{d}{2d-n} \right\rfloor$$

2) If d is odd and $2d+1 > n$, then

$$M(n, d) \leq 2 \left\lfloor \frac{d+1}{2d+1-n} \right\rfloor$$

3) If d is even and $n = 2d$, then

$$M(n, d) \leq 4d$$

4) If d is odd and $n = 2d+1$, then

$$M(n, d) \leq 4d$$

proof of 1)

$$\sum_{\substack{i \neq j \\ x_i \in C}} d_H(x_i, x_j) \geq M \cdot (M-1) \cdot d \quad (*)$$

(M-comparison)

Si zeros
M-si ones}

Codewords

x_{11}		x_{1n}
x_{21}		x_{2n}
\vdots		
x_{m1}		x_{mn}

ith column

$$2 \underbrace{s_i(M-s_i)}$$

it is maximized

if $s_i = \frac{M}{2}$ (equal number of zeros and ones)

$$\sum_{\substack{i+j \\ \underline{x}_i \in C}} d_H(\underline{x}_i, \underline{x}_j) = \sum_{i=1}^n 2 s_i (M-s_i) \leq \sum_{i=1}^n \underbrace{\frac{M^2/2}{2 \frac{M}{2}(M-\frac{M}{2})}}_{(*)} \frac{n M^2}{2}$$

combining (*) and (**): ~~$M(M-1)d$~~ $\leq \frac{n M^2}{2}$

$$M(d - \frac{n}{2}) \leq d$$



~~scribble~~

Singleton Bound

C with $d(\underline{x}_i, \underline{x}_j) \geq d$

$$M(n, d) \leq 2^{n-d+1}$$

- remove $d-1$ first bits from each codeword, then the new codewords are still distinct.

- The number of new ~~codewords~~ sequences of lengths $n-d+1$ is 2^{n-d+1} thus it is ^{upper} ~~bound~~ a bound on M .

C with $d(\underline{x}_i, \underline{x}_j) \geq d$

$$M(n, d) \geq \frac{2^n}{\sum_{j=0}^{d-1} \binom{n}{j}}$$

There is at least ONE sequence \underline{z}

among 2^n sequences, ~~one~~ with property

$$d_H(\underline{x}_i, \underline{z}) \leq d-1$$

Then the entire 2^n set of sequences ~~one~~ can be covered by $\bigcup_{\underline{x}_i \in C} B(\underline{x}_i, d-1)$ ball with center \underline{x}_i with radius $(d-1)$.

Then $2^n = \left| \bigcup_{\underline{x}_i \in C} B(\underline{x}_i, d-1) \right|$

↑ cardinality (size)

$$\leq \sum_{\underline{x}_i \in C} |B(\underline{x}_i, d-1)|$$

$$\binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{d-1}$$

$$= \left(\sum_{j=0}^{d-1} \binom{n}{j} \right) M$$

$$M \geq \frac{2^n}{\sum_{j=0}^{d-1} \binom{n}{j}}$$

~~scribble~~

- thomas & cover tried to develop an RD theory which covers both source quantization and channel quantization.

ex.

$$X_i \sim \text{Bernoulli}(\frac{1}{2})$$

with chosen Hamming Distortion.

$$\mathbb{E}[d(x, \hat{x})] = D$$

$$R(D) = \min I(x; \hat{x})$$

$$P(\hat{x}|x): \quad \mathbb{E}[x \oplus \hat{x}] \leq D$$

$$\mathbb{E}[x \oplus \hat{x}] = D$$

$$I(x; \hat{x}) = \underbrace{H(x)}_{H_b(p)=1} - \underbrace{H(x|\hat{x})}_{\begin{array}{l} \text{channel is} \\ \text{symmetric} \\ \text{and it} \\ \text{can be} \\ \text{reversed.} \end{array}}$$

$$H_b(D)$$

$$\Rightarrow R(D) = 1 - H_b(D)$$

$\hat{x} \oplus x \sim \begin{cases} 1 \\ 0 \end{cases}$

$P[x \neq \hat{x}]$

$P[x = \hat{x}]$

$D = 1 \times \underbrace{P[x \neq \hat{x}]}_{P(x=0, \hat{x}=1) + P(x=1, \hat{x}=0)} + 0 \times P[x = \hat{x}]$

$\underbrace{P(x=0, \hat{x}=1)}_{D} + \underbrace{P(x=1, \hat{x}=0)}_{D}$

$\underbrace{P(\hat{x}=0|x=1)P(x=1)}_{1/2} + \underbrace{P(\hat{x}=1|x=0)P(x=0)}_{1/2}$

$1 \xrightarrow{D} 1 - D \xrightarrow{D} 1 \quad P(\hat{x}=1) = \frac{1}{2}$

$0 \xrightarrow{1-D} 0 \xrightarrow{D} 1 \quad P(\hat{x}=0) = \frac{1}{2}$

symmetry

$$X_i \sim \text{Bernoulli}(p)$$

$$H(x|\hat{x}) = H(x \oplus \hat{x}|\hat{x}) = - \sum_{\hat{x}=0}^1 p(\hat{x}) \cdot \underbrace{\sum_{x=0}^1 p(x|\hat{x}=\hat{x}) \log p(x|\hat{x}=\hat{x})}_{-\overbrace{P(x=0|\hat{x}=0) \log p(x=0|\hat{x}=0)} + \dots}$$

$$H(x \oplus \hat{x} | \hat{x}) = - \sum_{i=0}^1 p(\hat{x}_i) \underbrace{\sum_{\hat{x}}}_{\hat{x}=0} P(x \oplus \hat{x}_i | \hat{x}_i) \log P(x \oplus \hat{x}_i | \hat{x}_i)$$

- $P(\underbrace{x=0 \oplus \hat{x}=0}_{x=0} | \hat{x}=0) \cdot \log (\cdot)$
 - $P(\underbrace{x=1 \oplus \hat{x}=0}_{x=1} | \hat{x}=0) \cdot \log$

We Knew: $H(x \oplus \hat{x} | \hat{x}) \leq H(x \oplus \hat{x})$

so we have lowerbound on M.I.:

$$I(x; \hat{x}) \geq H_b(p) - H(x \oplus \hat{x})$$

as previous page.
 $H_b(D)$

ex.2 $X \sim N(0, \sigma^2)$

$$E[(x - \hat{x})^2] \leq D$$

$$I(x; \hat{x}) = h(x) - h(x | \hat{x})$$

$$= \frac{1}{2} \log(2\pi e \sigma^2) - h(x - \hat{x} | \hat{x})$$

$$\Rightarrow I(x; \hat{x}) \geq \frac{1}{2} \log \frac{\sigma^2}{D}$$

Variance is given.

$\log N(0, D)$ given variance, maximized when Gaussian

$$\leq h(x - \hat{x})$$

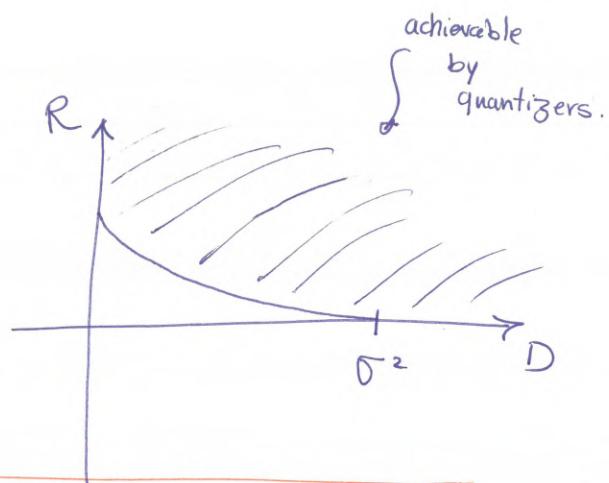
$$\leq \frac{1}{2} \ln(2\pi e D)$$

when is it achieved?

let's go for the test channel.

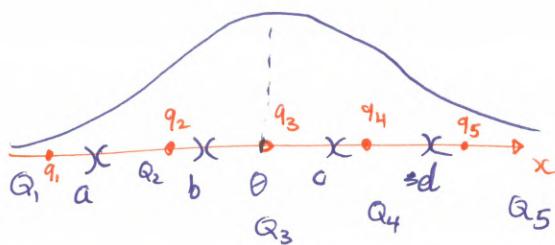
Gaussian channel makes sense as test channel.

$$\begin{array}{c} \text{W} \sim N(0, D) \\ \xrightarrow{+} X \sim N(0, \sigma^2) \\ \downarrow \\ \tilde{X} \sim N(0, \sigma^2 - D) \\ \downarrow \\ \text{variance decreased!} \end{array}$$



RV $X \sim N(0, \sigma^2)$

Loyd Quantizer



Least Squares Quantizing
(K Means)

$$\min_{Q_i, q_i} \sum_{i=1}^5 E[(X - q_i)^2]$$

q_i : quanta

Q_i : partition

Alternating minimization

1. Given Q_i , $i = 1, 2, \dots, 5$

calculate

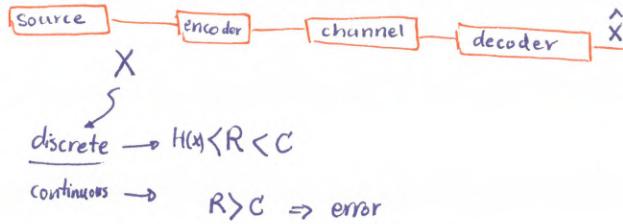
$$q_i = \frac{\int x \cdot f_X(x) dx}{\int_{x \in Q_i} f_X(x) dx}$$

conditional expectation

2. Update region

if $|x - q_i|^2 > |x - q_j|^2$

then move x to ~~Q_i~~ Q_j .



rate-distortion
 how to quantize
 how to do joint source-channel coding

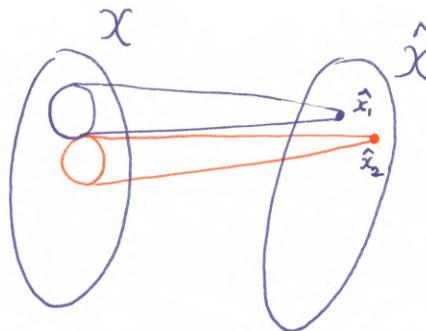
Rate-Distortion (RD)

$$\begin{cases} R(D) \\ D(R) \end{cases}$$

- Theoretical RD function

$$R(D) = \min_{P(\hat{x}|x)} I(x; \hat{x})$$

$$P(\hat{x}|x): E[d(x, \hat{x})] \leq D$$

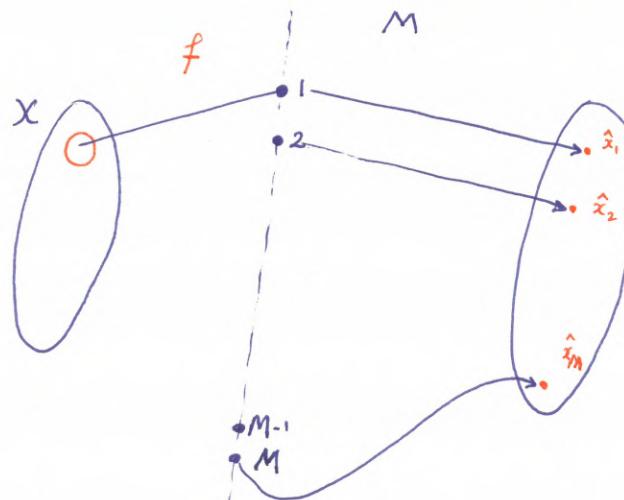


- Operational RD function

(n, M) codes

$$f: x^n \rightarrow \{1, \dots, M\}$$

$$R = \frac{\log M}{n}$$



$$\lim E[\bar{d}(x, \hat{x})] = D$$

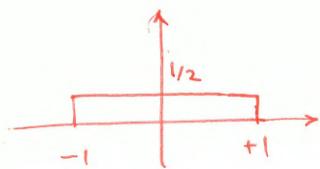
$$d \cancel{\bar{d}} = \frac{1}{n} \sum_{i=1}^n d(x_i, \hat{x}_i)$$

ex.

$$X \sim \text{unif}(-1, 1)$$

Quantize to 1 bit representation.

$$R = 1$$



encode sign!

$$M = 2^R = 2 \rightarrow \{1, 2\}$$

$$(-1, 0) \xrightarrow{\text{encode}} 1 \xrightarrow{\text{decode}} x_1$$

$$(0, 1) \xrightarrow{\text{encode}} 2 \xrightarrow{\text{decode}} x_2$$

$$\mathbb{E}[d(x, \hat{x})] = \mathbb{E}[(x - \hat{x})^2]$$

$$\min_{x_1, x_2} \mathbb{E}[(x - \hat{x})^2]$$

$$\int_{-1}^0 \frac{1}{2} (x - x_1)^2 dx + \int_0^1 \frac{1}{2} (x - x_2)^2 dx$$

$d/dx_1 ($

$$\int_{-1}^0 (x - x_1) dx = 0$$

$$\left. \frac{x^2}{2} - x_1 x \right|_{-1}^0 = 0$$

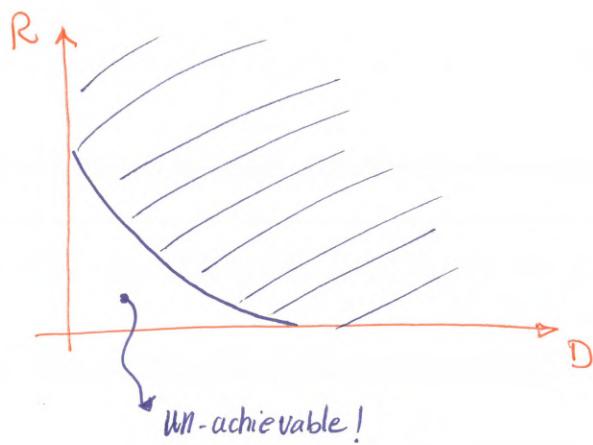
$$-\frac{1}{2} + x_1 = 0$$

$$x_1 = \frac{1}{2}$$

symmetry

$$x_2 = \frac{1}{2}$$

Reproduction
Points



ex.2 Binary Source

$$X_i \sim \text{Bernoulli}(p)$$

$$R(D) = \begin{cases} H(p) - H(D), & \text{if } D \leq \min\{p, 1-p\} \\ 0, & \text{else} \end{cases}$$

(for binary source, distortion is always Hamming distance)

$$d(x, \hat{x}) = \begin{cases} 1, & \text{if } x \neq \hat{x} \\ 0, & \text{else} \end{cases}$$

$$R(D) = \min I(x; \hat{x})$$

$$d(x, \hat{x}) = x \oplus \hat{x}$$

↓
modulo 2 addition

$$P(\hat{x}|x) : E[d(x, \hat{x})] \leq D$$

$$I(x; \hat{x}) = \underbrace{H(x)}_{H_b(p)} - \underbrace{\sum_{\hat{x}} H(x|\hat{x})}_{\downarrow} \underbrace{P(\hat{x})}_{\downarrow} \underbrace{\sum_{\hat{x}} H(x|\hat{x}) P(\hat{x})}_{\downarrow}$$

adding constant to RV. doesn't change entropy $\rightarrow H(x \oplus \hat{x} | \hat{x} = \hat{x})$
 we did the same in channel coding notes. 11, 15, 22 - 2

$$H(x \oplus \hat{x} | \hat{x}) \leq H(x \oplus \hat{x})$$

$$\Rightarrow I(x; \hat{x}) \geq H_b(p) - \underbrace{H(x \oplus \hat{x})}_{H_b(D)}$$

$$x \oplus \hat{x} = \begin{cases} 1 & \text{if } x \neq \hat{x} \\ 0 & \text{if } x = \hat{x} \end{cases}$$

$$\cancel{E[x \oplus \hat{x}] = 1 \cdot P[x \neq \hat{x}] + 0 \cdot P[x = \hat{x}]}$$

$$= D$$

we used
the upperbound
for $E[d(x, \hat{x})]$

Conditions
for the equality
to be hold

when $H(x|\hat{x}) = ?$ $H_b(D) ?$

what if $D > \min\{p, 1-p\}$

pick $\hat{x} = 0$, the $P(\hat{x} = 0) = 1$

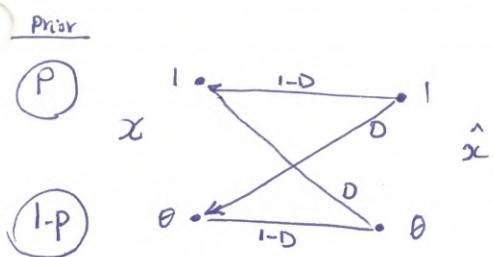
$$P[x \oplus \hat{x} | \hat{x} = 0] = P[x] = p$$

$$I(x; \hat{x}) = 0$$

If $D < p$

Design a test channel

BSC



$$P = q(1-D) + (1-q)D$$

$$P+D = q(1-2D)$$

$$q_h = \frac{P+D}{1-2D}$$

$$P(\hat{x}|x) = \frac{P(x|\hat{x})p(\hat{x})}{P(x)}$$

Introduction to Large Deviations (LD)

Consider a sequence $\{S_n\}$

$$S_n \xrightarrow{P} c, \text{ as } n \rightarrow \infty$$

Same as

$$P(|S_n - c| > \epsilon) \rightarrow 0, \text{ as } n \rightarrow \infty$$

ex.

$$\{Y_i\}$$

form sample mean.

$$S_n = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$P(|S_n - c| > \epsilon) = k(n, \epsilon, c) e^{-nI(c, \epsilon)}$$

slowly varying
compared to $e^{-nI(c, \epsilon)}$



Now we can say S_n satisfies LD principle.

2 results1. Cramer's Theorem

$X_i \sim \text{iid R.V.'s}$

$$S_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Assume $E[X_i] = m$

By the LLN

$$S_n \xrightarrow{P} m$$

5th
like
chernoff Bound!

Consider function $I(\theta) = \sup_{\theta} [\theta x - \log M(\theta)]$

(Sup is like max but applies to both open and closed sets)

$$M(\theta) = E_{f_x}[e^{\theta x}] = \int_{-\infty}^{+\infty} e^{\theta x} f_x(x) dx$$

Moment Generating Function \rightarrow It is CONVEX. $\Rightarrow \log M(\theta)$ is also CONVEX.

$I(x)$: Large Deviation Rate Function

properties of LD rate function

P.1

$I(x)$ is a convex (U) function.

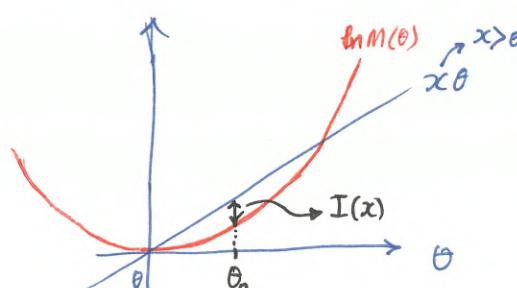
P.2

$I(x)$ has minimum at $E[X_i] = m$.

and furthermore $I(m) = 0$.

P.3

If $M(\theta) < \infty$ $I(x) = x\theta_0 - \ln M(\theta_0)$



θ_0 : where the difference is maximum.

ex.1

$$X \sim N(0, 1)$$

$$M(\theta) = E[e^{\theta X}] = e^{\frac{1}{2}\theta^2}$$

$$I(x) = \sup_{\theta} x\theta - \frac{\theta^2}{2}$$

$$\frac{\partial}{\partial \theta} I(x) = 0 \Rightarrow x = \theta_0$$

$$\Rightarrow I(x) = \frac{x^2}{2}$$

ex.2

$$X \sim \text{Bernoulli}(\frac{1}{2})$$

$$M(\theta) = E[e^{\theta X}] = e^{\theta} \frac{1}{2} + \frac{1}{2}$$

$$I(x) = \sup_{\theta} (x\theta + \log_2 - \log 1 + e^{\theta})$$

$$x - \frac{e^{\theta_0}}{1 + e^{\theta_0}} = \theta$$

$$x + e^{\theta_0}(x-1) = 0$$

$$\theta_0 = \log \frac{x}{1-x}$$

$$I(x) = x \log \frac{x}{1-x} + \log_2 - \log \frac{1}{1-x} = \overbrace{x \log x + (1-x) \log 1-x}^{\text{very like binary entropy...}} + \log_2 11, 29, 22 - 2$$

Now → Cramer's Theorem

Assume $M(\theta) < \infty$ for all θ . Then for every closed set $F \subseteq \mathbb{R}^1$

$$\overline{\lim}_{n \rightarrow +\infty} \frac{1}{n} \log P(S_n \in F) \leq -\inf_{x \in F} I(x)$$

approaching
from above

and for every open set $G \subseteq \mathbb{R}^1$

Sandwich

$$\underline{\lim}_{n \rightarrow +\infty} \frac{1}{n} \log P(S_n \in G) \geq -\inf_{x \in G} I(x)$$

approaching
from below

If F is $[a, b]$ and $G = (a, b)$, then

$$\overline{\lim} (\dots) = \underline{\lim} (\dots) = -\inf_{x \in [a, b]} I(x)$$

ex.3

$$X_i \sim \text{iid } N(0, 1)$$

$$S_n = \frac{1}{n} \sum_{i=1}^n X_i^2$$

$S_n \xrightarrow{\text{P}} 1$, (since mean is zero, S_n will be sample variance)
as $n \rightarrow +\infty$

$$P(S_n > \gamma)$$

γ should be larger than 1
($\gamma > 1$)

$$\frac{S_n - 1}{\sqrt{2/n}} \xrightarrow{\text{in dist.}} Z \sim N(0, 1)$$

from chi-squared distribution

$$\Rightarrow S_n \sim N(1, 2/n)$$

The central-limit theorem.

$$P(S_n > \gamma) \approx \int_{\gamma}^{+\infty} \frac{1}{\sqrt{2\pi^2/n}} e^{-\frac{(t-1)^2}{2^{2/n}}} dt = Q\left(\frac{\gamma-1}{\sqrt{2/n}}\right)$$

$$\approx k e^{-\frac{n}{4}(\gamma-1)^2} \quad \begin{array}{l} \text{rate function} = \frac{(\gamma-1)^2}{4} \\ \text{Constant} \end{array}$$

Now using Cramer's Thm.

$$P(S_n > \gamma) \approx k(\gamma, n, 1) \cdot e^{-n \inf_{x>\gamma} I(x)}$$

mean

$$I(x) = \sup_{\theta} (x\theta - \log M(\theta))$$

$$M(\theta) = E \left[e^{\frac{1}{n} \sum_{i=1}^n x_i^2 \cdot \theta} \right]$$

$$= \left(E \left[e^{\frac{\theta}{n} x_i^2} \right] \right)^n$$

$$\int_{-\infty}^{+\infty} e^{\frac{\theta}{n} x^2} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$$

$$\underbrace{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(1-\frac{2\theta}{n})x^2}}_{\text{make it a Gaussian.}} \times \underbrace{\sqrt{\frac{1}{1-2\theta/n}}}_{\sqrt{\frac{1}{1-2\theta/n}}}$$

$$= \left(\frac{1}{1-\frac{2\theta}{n}} \right)^{\frac{n}{2}}$$

$$I(x) = \sup_{\theta} (x\theta + \frac{n}{2} \log (1 - \frac{2\theta}{n}))$$

$$\frac{\partial}{\partial \theta} I(x) = x + \frac{n}{2} \frac{-2/n}{1 - \frac{2\theta_0}{n}} = \theta \rightarrow x = \frac{1}{1 - \frac{2\theta_0}{n}}$$

$$\rightarrow x - \frac{2\theta_0}{n} x = 1 \rightarrow \boxed{\theta_0 = \frac{x-1}{x} \frac{n}{2}}$$

$$\Rightarrow I(x) = n \left(\frac{x-1}{2} \right) + \frac{n}{2} \log \left(1 - \frac{x-1}{x} \right) = \frac{n}{2} \left(x - \log x - 1 \right) \Big|_{x=\gamma}$$

$$I(\gamma) = \frac{1}{2} (\gamma - \log \gamma - 1)$$

Compare

set $\gamma = 3$.

$$\frac{(\gamma-1)^2}{4}$$

The CLT
rate func. = 1

The LD rate

function $I(3) = 0.451$

Cramer's Theorem

$$X_1, \dots, X_n, \dots \text{ iid}$$

$$\mathbb{E}[X_i] = m_X$$

$$S_n = \sum_{i=1}^n X_i \quad \text{asymptotic approximation, } n \rightarrow \infty$$

$$P\left(\underbrace{\frac{S_n}{n} - m_X}_{\frac{S_n}{n} > \gamma} > \epsilon\right) \stackrel{\circ}{=} e^{-n \inf_{x>\gamma} I(x)}$$

$$\begin{aligned} \log & \rightarrow x \frac{1}{n} \\ & \downarrow \lim_{n \rightarrow \infty} \end{aligned}$$

$$I(x) = \sup_{\theta} \left\{ \theta x - \log M_x(\theta) \right\}$$

the asymptotic will become equality.

ex. $X_i \sim N(0, 1)$

$$Y_n = \sum_{i=1}^n X_i^2$$

$$\frac{Y_n}{n} \xrightarrow[n \rightarrow \infty]{P} 1 \quad (\text{LLN})$$

$$P\left(\frac{Y_n}{n} > \underbrace{\gamma + \epsilon}_{\gamma}\right) \doteq e^{-n I(\gamma)}$$

$$\begin{aligned} \Rightarrow M_{X^2}(\theta) &= E[e^{\theta X^2}] \\ &= \int_{-\infty}^{+\infty} e^{\theta x^2} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx \times \frac{1}{\sqrt{1-2\theta}} \\ &= \frac{1}{\sqrt{1-2\theta}} e^{-\frac{1}{2}(1-2\theta)x^2} \end{aligned}$$

$$\Rightarrow I(x) = \sup_{\theta} \left\{ \theta x + \frac{1}{2} \ln(1-2\theta) \right\}$$

$$x - \frac{1}{1-2\theta} = 0$$

Gartner and Ellis (1984)

they wanted to remove assumptions.

$\Rightarrow \{Y_n\}$ possibly dependent

$$\text{Define } \varphi_n(\theta) = \frac{1}{n} \log E[e^{\theta Y_n}]$$

Assumption 1

$$\lim_{n \rightarrow +\infty} \varphi_n(\theta) = \varphi(\theta)$$

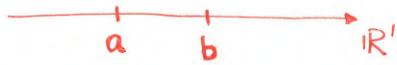
$$D_\varphi = \{ \theta : \varphi(\theta) < \infty \}$$

Assumption 2

$\varphi_n(\theta)$ is differentiable on D_φ .

$$I(x) = \sup_{\theta} \{ \theta x - \varphi(\theta) \}$$

Theorem (Gartner - Ellis)



$$\overline{\lim}_{n \rightarrow +\infty} \frac{1}{n} \log P(Y_n \in [a, b]) \leq \inf_{x \in [a, b]} I(x)$$

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \log P(Y_n \in (a, b)) \leq \inf_{x \in (a, b)} I(x)$$

ex. (AutoRegressive Process of order 1)

$$X_i = \alpha X_{i-1} + W_i \quad (\alpha < 1)$$

$$X_0 = 0$$

This R.P. is stationary
the sandwich will be equal
from both sides.

W_i are iid over $[-\frac{1}{2}, \frac{1}{2}]$ ↳ can have any kind of distribution.

$$\mathbb{E}[W_i] = 0$$

Form $Y_n = \sum_{i=1}^n X_i$

NOT IID \rightarrow No Cramer Usage!

Approximate / Analyze $P(\frac{Y_n}{n} > \gamma)$

$$X_1 = W_1$$

$$X_2 = \alpha W_1 + W_2$$

$$X_3 = \alpha^2 W_1 + \alpha W_2 + W_3$$

⋮

$$X_n = \alpha^{(n-1)} W_1 + \alpha^{(n-2)} W_2 + \dots + \alpha W_{n-1} + W_n = \sum_{i=1}^n \alpha^{n-i} W_i$$

$$Y_n = \sum_{i=1}^n X_i = (1 + \alpha + \dots + \alpha^{n-1}) W_1 + (1 + \alpha + \dots + \alpha^{n-2}) W_2 + \dots + W_n$$

$$= \frac{1 - \alpha^n}{1 - \alpha} W_1 + \frac{1 - \alpha^{n-1}}{1 - \alpha} W_2 + \dots + W_n$$

$$\varphi_n(\theta) = \frac{1}{n} \log E[e^{\theta Y_n}] \quad Y_n \text{ is a sum of independent R.V.'s.}$$

$$= \frac{1}{n} \log \prod_{i=1}^n M_{W_i} \left(\theta - \frac{(1-\alpha)^{n-i}}{1-\alpha} \right)$$

as $n \rightarrow \infty$
 $\alpha^{n-i} \rightarrow 0$

$$\Rightarrow = \frac{1}{n} \cdot n \log M_W \left(\frac{\theta}{1-\alpha} \right)$$

$$= \log M_W \left(\frac{\theta}{1-\alpha} \right)$$

Assume $W_i = \begin{cases} \frac{1}{2}, \text{ with prob. } \frac{1}{2} \\ -\frac{1}{2}, \text{ " " " } \end{cases}$



$$= \log e^{\frac{\theta}{1-\alpha} \frac{1}{2} \frac{1}{2}} + e^{-\frac{\theta}{1-\alpha} \frac{1}{2} \frac{1}{2}}$$

$$= \log \cosh \left(\frac{\theta}{1-\alpha} \right)$$

$$I(x) = \sup_{\theta} \left\{ \theta x - \log \cosh \left(\frac{\theta}{1-\alpha} \right) \right\}$$

$$x = \frac{\frac{1}{2} e^{\frac{\theta}{1-\alpha} \frac{1}{2} \left(\frac{1}{2(1-\alpha)} \right)} - \frac{1}{2} \left(\frac{1}{2(1-\alpha)} \right) e^{-\frac{1}{2} \frac{\theta}{1-\alpha}}}{\cosh \left(\frac{\theta}{1-\alpha} \right)}$$

$$= \frac{1}{2(1-\alpha)} \frac{\sinh \left(\frac{\theta}{1-\alpha} \right)}{\cosh \left(\frac{\theta}{1-\alpha} \right)}$$

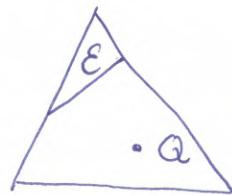
$\rightarrow \tanh \left(\frac{\theta}{1-\alpha} \right)$

Sanov's Theorem

Let X_1, \dots, X_n iid with cdf/pdf $Q(x)$

and \mathcal{E} is a set of cdfs $P(x)$

\downarrow
does not
include Q .



Form empirical cdf:

$$L_n(x) = \frac{1}{n} \sum \delta_{X_i}(x)$$

like what
the objective of
MLE is!

$$P(L_n \in \mathcal{E}) \doteq e^{-n \inf_{P \in \mathcal{E}} D(P \parallel Q)}$$

if $D_{KL} \downarrow \Rightarrow$ possibility that L_n lies in $\mathcal{E} \uparrow$

Ex. 3 (Cover, p. 365)

A fair coin toss $X_i \in \{0, 1\}$

$n=1000$

$$\text{Find } P\left(\sum_{i=1}^{1000} X_i > 700\right) = P\left(\frac{\sum X_i}{1000} > 0.7\right)$$

converges
 $\xrightarrow{} \frac{1}{2}$

$$P(X_i=1) = P$$

Set \mathcal{E} consist of all $P(x)$ which $(\mathbb{E}_{x \sim P(x)} [X_i]) > 0.7$.

$$e^{-n \inf_{P \geq 0.7} D(\varphi, P) \| (0.5, 0.5)}$$

$$\underbrace{D(0.7, 0.3) \| (0.5, 0.5)}_{0.119}$$

0.119

$$\Rightarrow e^{-1000 \times 0.119} = e^{-119} \rightsquigarrow \text{such low probability!!} \quad \text{☺}$$

ex.4 $X_i \sim \text{iid } N(0, 1)$

sample average power
isn't it energy?

$$Y_n = \frac{1}{n} \sum_{i=1}^n X_i^2 \rightarrow 1$$

Approximate $P(Y_n/n > 5)$

first define \mathcal{E} .

average power

$$\mathcal{E} = \left\{ P(X) : \widehat{E[X_i^2]} \geq 5 \right\} \quad \text{①}$$

$$\inf_{P_X(x) \in \mathcal{E}} D(P_X \| N(0, 1))$$

we have no idea about the family of P_X . \rightarrow Using Lagrange Multipliers.

$$\mathcal{L}(P_X, \lambda, \mu) = \int_{-\infty}^{+\infty} P_X(x) \log \frac{P_X(x)}{q_X(x)} dx$$

$$+ \lambda \left[\int_{-\infty}^{+\infty} x^2 P_X(x) dx - 5 \right]$$

$$+ \mu \left[\int_{-\infty}^{+\infty} P_X(x) dx - 1 \right]$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial P_X} = \log \frac{P_X(x)}{q_X(x)} + 1 + \lambda x^2 + \mu$$

project • just several bounds on an application

$$P(T(x^n) > \gamma) = e^{-n D(P^* \| Q)}$$



Hypothesis Testing

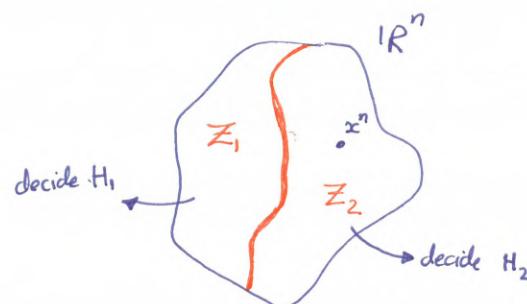
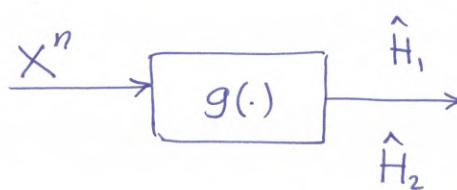
X_1, \dots, X_n iid

$$H_1: X_i \sim P_1(x)$$

$$H_2: X_i \sim P_2(x)$$

decision rule

$$g(X^n)$$



4 outcomes

decide H_1	H_1 is true	correct decision
H_2	H_2	
H_1	H_2	wrong decision
H_2	H_1	

in pattern recognition
 P_{FA} and P_{miss}
 are trade-off for each other.

$$P_{FA} = P(\text{decide } H_1 \mid H_2 \text{ is true}) = P_2(X \in Z_1)$$

↳ actually there is no signal out there, but we wrongly decide that there is!

$$P_{miss} = P(\text{decide } H_2 \mid H_1 \text{ is true}) = P_1(X \in Z_2)$$

↳ Actually there were a signal but we missed it.

Neyman - Pearson Criterion

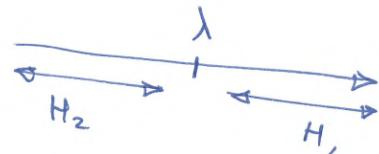
$$\min P_{\text{miss}}$$

$$(\underline{z}_1, \underline{z}_2) : P_{\text{FA}} \leq \epsilon \\ \epsilon > 0$$

$$\begin{aligned} J(\underline{z}_1, \underline{z}_2, \lambda) &= P_1(\underline{x} \in \underline{z}_2) + \lambda (P_2(\underline{x} \in \underline{z}_1) - \epsilon) \\ &= \sum_{x \in \underline{z}_2} P_1(x) + \lambda \left(\sum_{x \in \underline{z}_1} P_2(x) - \epsilon \right) \\ &\quad \underbrace{\qquad\qquad\qquad}_{1 - \sum_{x \in \underline{z}_1} P_1(x)} \end{aligned}$$

Decide H_1 if $P_1(\underline{x}) > \lambda P_2(\underline{x})$

" H_2 if $P_1(\underline{x}) < \lambda P_2(\underline{x})$



$$P_1(\underline{x}) \begin{cases} > \lambda P_2(\underline{x}) \\ < \lambda P_2(\underline{x}) \end{cases} \quad \lambda > 0$$

a R.V.

$$\Lambda \triangleq \frac{P_1(\underline{x})}{P_2(\underline{x})} \quad \begin{matrix} > \lambda \\ < \lambda \end{matrix} \rightarrow \begin{matrix} H_1 \\ H_2 \end{matrix}$$

Lagrange Multiplier
= Threshold

Likelihood Ratio Λ $\xrightarrow{\text{LRT (LR Test)}}$ use the constant $P_2(\underline{x} \in \underline{z}_1) = \epsilon$

If X_i 's are iid:

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \log \frac{P_1(x_i)}{P_2(x_i)}}_{\frac{1}{n} \mathcal{L}(x)} \gtrless \begin{cases} H_1 & \\ H_2 & \end{cases} \frac{\log \lambda}{n}$$

$$-D(P_2||P_1)$$



$$D(P_1||P_2)$$

If X_i are from H_1 , then

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \log \frac{P_1(x_i)}{P_2(x_i)}}_{\text{Converge in } P \text{ on long run}} \rightarrow D_{KL}(P_1||P_2)$$

If X_i are from H_2 , then

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \log \frac{P_1(x_i)}{P_2(x_i)}}_{\text{in } P \text{ on long run}} \rightarrow -D_{KL}(P_2||P_1)$$

Using
AEP
(LLN
in
Information
Theory)

$$P_{FA} = P\left(\frac{1}{n} \mathcal{L}(x) > \frac{1}{n} \log \lambda\right) = e^{-n?} \xrightarrow{\text{asymptotically}} \min_{P \in \mathcal{E}} D(P||P_2)$$

true probability

Using Sanov's Theorem

$$\mathcal{E} = \left\{ P(x) : E_P \left[\frac{1}{n} \mathcal{L}(x) \right] > \frac{\ln \lambda}{n} \right\}$$

$\frac{1}{n} \sum_{i=1}^n \log \frac{P_1(x_i)}{P_2(x_i)}$

$$\sum_x P(x) \log \frac{P(x)}{P_2(x)} - P(x) \log \frac{P(x)}{P_1(x)} > \frac{\log \lambda}{n}$$

↑
↑
i.i.d
n will
be committed

$$D(P||P_2) - D(P||P_1) > T$$

$$\mathcal{E} = \left\{ P(x) : D(P||P_2) - D(P||P_1) > T \right\}$$

So:

$$\min D(P||P_2)$$

$$P: D(P||P_2) - D(P||P_1) > T$$

$$\Rightarrow J(p(x), \lambda, \mu) = \sum_x p(x) \log \frac{p(x)}{P_2(x)} + \lambda \left(\sum_x p(x) \log \frac{P_1(x)}{P_2(x)} - T \right) + \mu \left(\sum_x p(x) - 1 \right)$$

$$\frac{\partial J}{\partial p(x)} = \log \frac{p(x)}{P_2(x)} + 1 + \lambda \log \frac{P_1(x)}{P_2(x)} + \mu = 0$$

$$p(x) = P_2(x) e^{-1 - \mu - \lambda \log \frac{P_1(x)}{P_2(x)}}$$

$$\sum P(x) = 1$$

$$\begin{aligned} & p \log p + q \log q \\ & \text{Key: } p \log p - p \log p + q \log q - q \log p \end{aligned}$$

Q

$$P^*(x) = \frac{P_2(x)^{1+\lambda} P_1(x)^{-\lambda}}{\sum_t P_2(t)^{1+\lambda} P_1(t)^{-\lambda}}$$

- λ is found from the constraint.

$$D(P \parallel P_2) - D(P \parallel P_1) = T$$

Hypothesis Testing Problems

$$X_i \sim \text{iid} \quad i=1 \dots n$$

$$H_1: P_1(x)$$

$$H_2: P_2(x)$$

Neyman Pearson Criterion $\rightarrow \frac{1}{n} \log \frac{P_1(\underline{x})}{P_2(\underline{x})} \stackrel{H_1}{>} \frac{1}{n} \log \frac{\lambda}{\gamma}$

$$P_{FA} = P_2(\underline{x} \in \mathcal{Z}_1) = P_2\left(\frac{1}{n} \log \frac{P_1(\underline{x})}{P_2(\underline{x})} > \frac{\log \lambda}{n}\right) = e^{-n D(P_\lambda^* \| P_2)}$$

$$P_{\text{miss}} = P_1(\underline{x} \in \mathcal{Z}_2) = P_1\left(\ell(\underline{x}) < \frac{\log \lambda}{n}\right) = e^{-n D(P_\lambda^* \| P_1)}$$

Sanov's Theorem

where $P_\lambda^*(x) = \frac{P_1^\lambda(x) P_2^{1-\lambda}(x)}{\sum_x P_1^\lambda(t) P_2^{1-\lambda}(t)}$

with $\lambda \checkmark D(P_\lambda^* \| P_1) - D(P_\lambda^* \| P_2) = \frac{\log \lambda}{n}$

$$D(P_\lambda^* \| P_1) - D(P_\lambda^* \| P_2) = \frac{\log \lambda}{n}$$

$$\text{if } \lambda = \theta \Rightarrow P_{\lambda}^*(x) = P_2(x) \Rightarrow \begin{cases} P_{FA} = 1 \\ P_{miss} = e^{-nD(P_2||P_1)} \end{cases}$$

$$\text{if } \lambda = 1 \Rightarrow P_{\lambda}^*(x) = P_1(x) \Rightarrow \begin{cases} P_{FA} = e^{-nD(P_1||P_2)} \\ P_{miss} = 1 \end{cases}$$

Chernoff - Stein Lemma

(Large Deviation extension of Neyman Pearson Lemma)

$$P_{FA} = P_2 \left(\ell(\underline{x}) > \frac{\log \lambda}{n} \right) \leq \epsilon, \quad 0 < \epsilon < \frac{1}{2}$$

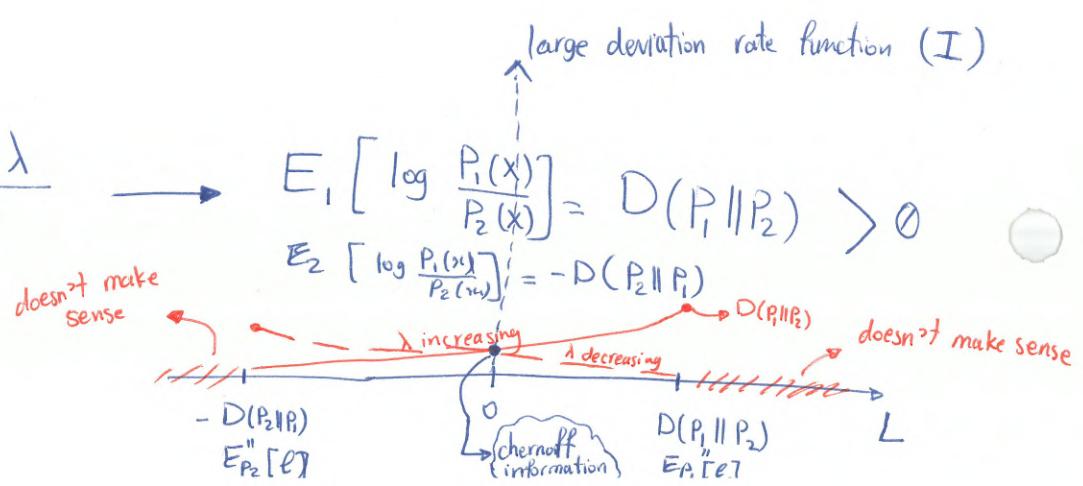
Chernoff-Stein $\rightarrow P_{miss} = P \left(\ell(\underline{x}) < \frac{\log \lambda}{n} \right) = e^{-nD(P_2||P_1)}$

Likelihood Ratio Test

$$\frac{1}{n} \log \frac{P_1(\underline{x})}{P_2(\underline{x})} \stackrel{H_1}{\geq} \frac{\log \lambda}{n} \stackrel{H_2}{<} \frac{\log \lambda}{n}$$

$$H_1: \ell \sim P_1(L)$$

$$H_2: \ell \sim P_2(L)$$



Bayesian criterion

π_1, π_2 priors

$$P(\text{error}) = P_2 \left(\ell(\underline{x}) > \frac{\log \lambda}{n} \right) \pi_2 + \pi_1 P_1 \left(\ell(\underline{x}) < \frac{\log \lambda}{n} \right)$$

$$\stackrel{*}{=} \pi_2 e^{-n D(P_\lambda^* \| P_2)} + \pi_1 e^{-n D(P_\lambda^* \| P_1)}$$

dominate? ↗ dominate? ↘

$$\stackrel{*}{=} e^{-n \min \{ D(P_\lambda^* \| P_2), D(P_\lambda^* \| P_1) \}}$$

to maximize minimum of the two : @ Chernoff Information

$$\stackrel{*}{=} e^{-n C(P_1, P_2)}$$

$$\text{where } D(P_\lambda^* \| P_i) = D(P_\lambda \| P_i)$$

Fully Bayesian Criterion

$$\underbrace{\log \frac{P_1(\underline{x})}{P_2(\underline{x})}}_{\ell(\underline{x})} \begin{cases} \geq \log \frac{\pi_2}{\pi_1} \\ \leq \log \frac{\pi_1}{\pi_2} \end{cases} \quad \left. \begin{array}{c} \rightarrow \\ \text{Chernoff Bound} \end{array} \right\} \quad P_2 \left(\ell(\underline{x}) > \underbrace{\log \frac{\pi_2}{\pi_1}}_{\text{RV.}} \right)$$

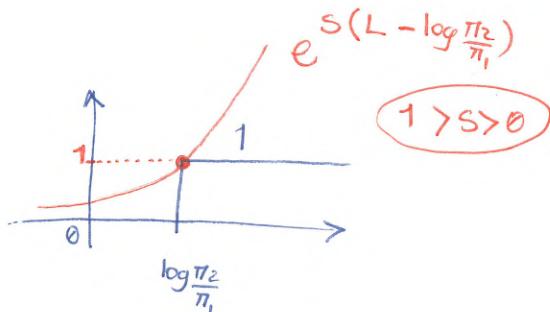
$$H_1: \ell \sim P_1(L)$$

$$H_2: \ell \sim P_2(L)$$

$$P_2(L > \log \frac{\pi_2}{\pi_1}) = \sum P_2(L)$$

$$L > \log \frac{\pi_2}{\pi_1}$$

$$= \sum_{L \in \mathbb{R}'} P_2(L) \cdot u(L - \log \frac{\pi_2}{\pi_1})$$



$$\leq \sum P_2(L) e^{s(L - \log \frac{\pi_2}{\pi_1})}$$

$$= e^{-s \log \frac{\pi_2}{\pi_1}} \underbrace{\sum e^{sL} P_2(L)}_{\text{MGF } \varphi_2(s)}$$

Laplace Transform of $P_2(\cdot)$

$$= e^{-s \log \frac{\pi_2}{\pi_1}} + \varphi_2(s)$$

$$\log \varphi_2(s)$$

$$\varphi_2(s) = E_{\pi_2}[e^{sL}]$$

$$\sum_{i=1}^n \log \frac{P_1(x_i)}{P_2(x_i)}$$

$$\sum P_2 \left(\frac{P_1}{P_2} \right)^s =$$

$$= E_{P_2(x)} \left[e^{s \log \frac{P_1(x)}{P_2(x)}} \right]$$

$$0 \leq s \leq 1$$

$$\stackrel{x_i \sim iid}{=} \left\{ E_{P_2(x)} \left[e^{s \log \frac{P_1(x)}{P_2(x)}} \right] \right\}^n = \left\{ \sum_x P_1^s(x) P_2^{1-s}(x) \right\}^n$$

$$= e^{-s \log \frac{\pi_2}{\pi_1} + \underbrace{\varphi_2(s)}_{\log \varphi_2(s)}}$$

$$= e^{-s \log \frac{\pi_2}{\pi_1} + n \log \left\{ \sum p_i^s(x) p_2^{1-s}(x) \right\}}$$

minimize wrt to $0 \leq s \leq 1$

Let's Tighten the Bound!

$$P_2 \left(\ell(X) > \log \frac{\pi_2}{\pi_1} \right) \leq e^{-n \max_{0 \leq s \leq 1} \left\{ \frac{s \log \frac{\pi_2}{\pi_1}}{n} - \log \left\{ \sum p_i^s(x) p_2^{1-s}(x) \right\} \right\}}$$

if $n \rightarrow \infty$
goes to zero!

by letting #samples to ∞ , priors effects wash away!

$$\rightarrow e^{-n \left\{ -\min_{0 \leq s \leq 1} \log \left(\sum p_i^s(x) p_2^{1-s}(x) \right) \right\}}$$

Chernoff Information

$$D(P_\lambda^* || P_1) = D(P_\lambda^* || P_2)$$

agrees with

