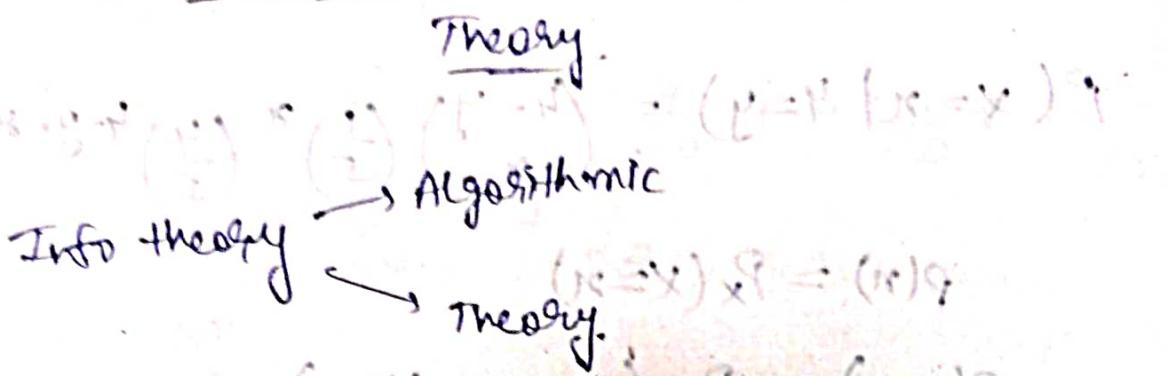


Introduction to Information Theory



→ What is information?

Information → A numerical measure of the uncertainty of an experimental outcome.

The more the random, more the information.

→ How to quantify information?

→ How can we measure the amount of info?

→ ensure correctness of info?



Information Theory

→ Information theory was invented by Shannon in 1948.

→ One of the basic postulates of information theory is that information can be treated like a measurable physical quantity, such as density or mass.

→ The information content of a message consists simply of the numbers 0's and 1's it takes to transmit it. Information can thus be measured in bits.

Basic concepts in info theory.

- Entropy
- Joint Entropy
- Conditional Entropy

Shannon Info content

→ we study info content of those events to which we can associate probability of occurrence of these events.

$$\text{Eg. } P(A) = 0.15$$

$$P(B) = 0.125$$

→ $P(B)$ provides more information

Shannon information content $h(A)$ corresponding to event A is defined as.

$$\text{Information } h(A) = \log_2 \left(\frac{1}{P(A)} \right) \text{ bits.}$$

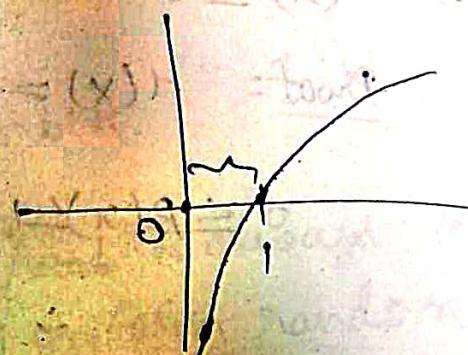
info associated with $h(A)$

$$h(A) = \log_2 \left(\frac{1}{0.15} \right)$$

$$= \log_2 2 = 1$$

$$h(B) = \log_2 \left(\frac{1}{0.125} \right)$$

$$= \log_2 8 = 3$$



we extend this definition to probability distribution function.

$$\text{Eg. } X = \{1, 2, 3\}$$

$X \rightarrow \text{discrete r.v.}$

$p(1) \swarrow p(2) \searrow p(3)$

Entropy is like the average information.

$$(q) H: \quad H(X) = \sum_{n \in X} P(n) \log \left(\frac{1}{P(n)} \right)$$

Entropy.

→ If base of logarithm is 2, then entropy is expressed in bits → default

- If base of log is e, then entropy is expressed in nats.
- If base of log is b, then we denote entropy by $H_b(x)$.

Interpretation of $H(x)$ in terms of expectation:

$$H(x) = \mathbb{E}_p \log\left(\frac{1}{p(x)}\right) \text{ where } x \sim p(x)$$

Eg. Properties of $H(x)$

$$1.) H(x) \geq 0$$

$$\text{Proof: } H(x) = - \sum_{n \in X} p(n) \log(p(n))$$

$$0 \leq p(n) \leq 1 \text{ so } \log(p(n)) \leq 0$$

$$\therefore H(x) \geq 0$$

$$2.) H_b(x) = (\log_b a) H_a(x)$$

Eg. Entropy of a Bernoulli r.v.

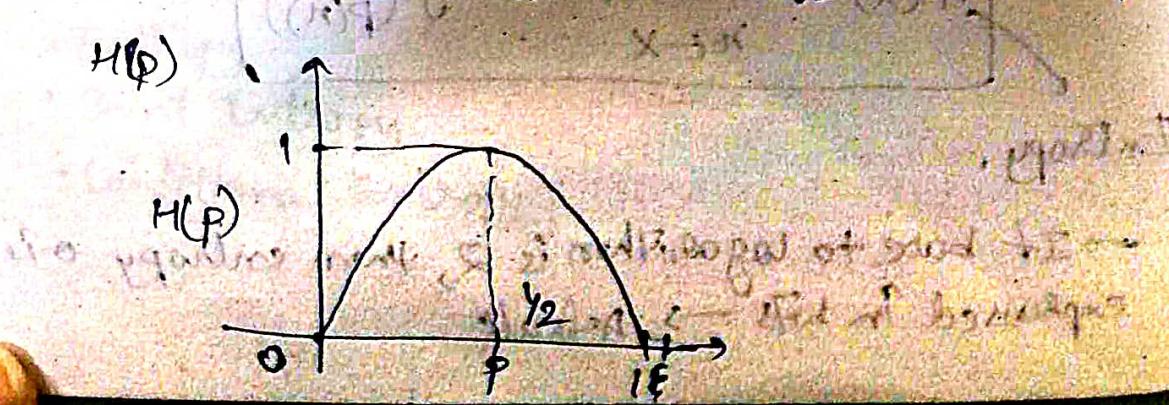
Bernoulli r.v. with parameter p :

Entropy of Bernoulli r.v.

$$= -p \log_2 p - (1-p) \log_2 (1-p) =: H(p)$$

$$H(p)$$

$$H(p)$$



$H(p)$ is always a concave function.

Entropy of a fair coin = 1 bit

Entropy of a discrete uniform r.v. $H(X) = \dots$

Eg. consider a uniform r.v. with support set

$$\{1, 2, \dots, n\}$$

$$H(X) = -\frac{1}{n} \times \log\left(\frac{1}{n}\right) - \frac{1}{n} \times \log\left(\frac{1}{n}\right) - \dots \text{ n times}$$

$$= -\frac{1}{n} \times \log\left(\frac{1}{n}\right) \times n = -\log\left(\frac{1}{n}\right) = \log n$$

$$\therefore H(X) = \log n$$

Eg. consider a fair dice.

$$H(X) = \log 6 = 2.585 \text{ bits.}$$

Entropy of any die or biased r.v. is always less than or equal to entropy of uniform random variable.

$$\begin{array}{ccccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} & 0 & 0 & 0 \end{array}$$

$$H(X) = \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{4} \times 2$$

$$= \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{3}{2} = 1.5$$

$$\therefore H(X) \leq \log(|X|)$$

Entropy of uniform r.v.

Where $|X|$ denotes the number of elements in the range of X , with equality if & only if X has a uniform distribution over X .

Joint entropy & conditional Entropy

→ Joint entropy $H(X, Y)$ of a pair of discrete s.r.s (X, Y) is defined as.

$$H(X, Y) = - \sum_{n \in X} \sum_{y \in Y} p(n, y) \log(p(n, y))$$

$$H(X, Y) = - \mathbb{E} \log p(X, Y)$$

→ If $(X, Y) \sim p(n, y)$, the conditional entropy $H(Y|X)$ is defined as

$$H(Y|X) = \sum_{n \in X} p(n) H(Y|X=n)$$

$$= \sum_{n \in X} p(n) \sum_{y \in Y} p(y|n) \log(p(y|n))$$

$$= - \sum_{n \in X} \sum_{y \in Y} p(n, y) \log p(y|n)$$

$$= - \mathbb{E} \log p(Y|X)$$

→ expectation with respect to joint s.r. and not conditional s.r. $H(Y|X) = - \mathbb{E} \log p(Y|X)$

$$\mathbb{E}(f(x)) = \sum_n f(n) \cdot p(n)$$

$$\mathbb{E}(f(x, y)) = \sum_n \sum_y f(n, y) p(n, y)$$

Entropy = uncertainty

what is $H(X, Y)$ if $X \neq Y$ are independent?

$$X \quad Y$$

$$H(X) \quad H(Y)$$

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

$H(X|Y)$ is not necessarily equal to $H(Y|X)$.

Relation between $H(X, Y)$ and $H(X|Y)$

$$\begin{aligned} H(X, Y) &= H(X) + H(Y|X) \\ &= H(Y) + H(X|Y) \end{aligned}$$

called the Chain Rule

$$\begin{aligned} H(X, Y) &= H(X) + H(Y|X) \\ &= H(Y) + H(X|Y) \end{aligned}$$

$$\text{Proof: } \text{LHS} = H(X, Y)$$

$$= -\sum_n \sum_y p(x, y) \log(p(x, y))$$

$$p(x, y) = p(x) p(y|x)$$

$$= -\sum_n \sum_y p(x, y) \log(p(x) \cdot p(y|x))$$

$$= -\sum_n \underbrace{\sum_y p(x, y) \log p(x)}$$

$$-\sum_n \sum_y p(x, y) \log(p(y|x))$$

$$-\sum_n \left(\sum_y p(x, y) \log p(x) \right)$$

$$H(Y|X)$$

$$P(x)$$

$$-\sum_n p(x) \log p(x) + H(Y|X) = \underline{H(X) + H(Y|X)}$$

Ex. Let (X, Y) have the following joint distribution.

$x \backslash y$	0	1	$H(X Y), H(Y X); H(X, Y)$
0	y_3	y_8	$p(x=0, y=0) = p(x=0) \cdot p(y=0)$
1	y_0	y_3	$p(x=1, y=1) = p(x=1) \cdot p(y=1)$

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

$H(X|Y=0) =$

$$p(x=0|y=0) = \frac{p(x=0, y=0)}{p(y=0)}$$

$x \backslash y$	0	1	$H(X Y)$
$p(x y=0)$	$\frac{1}{3}$	$\frac{2}{3}$	$\frac{1}{3} + \frac{2}{3} \cdot 1 = \frac{7}{3}$
$p(x y=1)$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2} + \frac{1}{2} \cdot 1 = \frac{3}{2}$

$$H(X|Y=0) = 0$$

$$H(X|Y=1) = 1$$

marginal

distribution

$$\therefore H(X|Y) = \frac{1}{3} \times 0 + \frac{2}{3} \times 1$$

2. $H(Y|X)$ \Rightarrow

$$H(Y|X) = \sum_n p(x) H(Y|x=n)$$

$x \backslash y$	0	1	$H(Y X)$
$p(y x=0)$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2} + \frac{1}{2} \cdot 1 = \frac{3}{2}$
$p(y x=1)$	$\frac{1}{3}$	$\frac{2}{3}$	$\frac{1}{3} + \frac{2}{3} \cdot 1 = \frac{7}{3}$

$$H(Y/X=0) = 1$$

$$H(Y/X=1) = 0$$

$$H(Y/X) = \frac{2}{3} \times 1 + \frac{1}{3} \times 0 = \textcircled{2/3}$$

$$H(X, Y) = H(X) + H(Y/X)$$

$$= \left(\frac{2}{3} \times 0 + \frac{1}{3} \times 1 \right) + \frac{2}{3}$$

$$= \frac{1}{3} + \frac{2}{3} = \textcircled{1} \text{ bits}$$

$\begin{matrix} X \\ \backslash \\ Y \end{math>$	1	2	3	4		1	2	3	4	
1	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{32}$		$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	$\rightarrow X$
2	$\frac{1}{16}$	$\frac{1}{8}$	$\frac{1}{32}$	$\frac{1}{32}$		$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\rightarrow Y$
3	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$						
4	$\frac{1}{4}$	0	0	0						

marginal distribution.

0 1 2 3 4 $\rightarrow X$

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

number of outcomes where $X=1$
total number of outcomes

$$\frac{1}{4} = P(X=1|Y=1)$$

0 1 2 3 4 $\rightarrow X$

(0, 1, 2, 3, 4) $\rightarrow Y$

$(-\infty, n] \rightarrow$ half intervals in $B(\mathbb{R})$

$A_i = (-\infty, n - \frac{1}{2}]$ for $i=1, 2, \dots, \infty$

$\bigcup_{i=1}^{\infty} A_i = (-\infty, n)$ (n is not included)

When $x \rightarrow g(x)$

where $g(x)$ is a one-one function,
then the entropy of x and $g(x)$ is same

Eg. 1) $x, y = g(x) = 2^x + x + 1$

$\therefore H(x) = H(y)$

$$\begin{matrix} p_1 & \rightarrow \\ p_2 & \rightarrow \\ 1 & \end{matrix}$$

2.) $x, y = x^2$

$$-p_1 \log p_1 - p_2 \log p_2 \geq -(p_1 + p_2) \log(p_1 + p_2)$$

Proof: $p_1 + p_2 \geq p_1 \Rightarrow p_1 + p_2 \geq p_2$

$$\Rightarrow \log(p_1 + p_2) \geq \log p_1 \quad \log(p_1 + p_2) \geq \log p_2$$

$$\Rightarrow p_1 \log(p_1 + p_2) \geq p_1 \log p_1, \quad p_2 \log(p_1 + p_2) \geq p_2 \log p_2$$

$$\Rightarrow (p_1 + p_2) \log(p_1 + p_2) \geq p_1 \log p_1 + p_2 \log p_2$$

$$\Rightarrow -(p_1 + p_2) \log(p_1 + p_2) \leq -p_1 \log p_1 - p_2 \log p_2$$

Entropy

$$H(X) = E[g(x)]$$

R.V under consideration
are discrete R.V.

$$g(x) = \log_2 \frac{1}{P(x)}$$

$$x \quad 0 \quad 1 \quad 2 \quad 3 \quad 4$$

$$g(0) \quad P(1) \quad P(2) \quad P(3) \quad P(4)$$

$P_X(x)$ itself is a R.V.

$P_X(x)$ is a function of x & hence a R.V.

$$g(x) = \log_2 \frac{1}{p_x(x)}$$

avg length = $\frac{2 \text{ bits}}{\text{symbol}}$

Eg. 1.) Source →	1	2	3	4	→	1 → 00
	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$		2 → 01
	1	2	3	3.		3 → 10
						4 → 11

2) Prefix free condition \leftarrow source

1 → 0	1	2	3	4	see. Prefix free condition
2 → 10	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	- one codeword
3 → 110	1	2	3	3.	the codeword for one source
4 → 111					(1, 2, 3, 4) do not coincide with any other.

1 2 3 4 & 1 }
 0 1 0 1 1 0 1 1 1 1 0 0. → above source
 can be obtained

from this code,

when the digits for

each source (1, 2, 3, 4),

are unique

Eg. 0 → 1

01 → X

⇒ after 0 is not

possible.

10 → 2, 101(X), 011(X), ..

each one is unique, so we are able to identify/define source from code.

$$\begin{aligned} \text{Average length} &= 1 \times \frac{1}{2} + 2 \times \frac{1}{4} + 3 \times \frac{1}{8} + 4 \times \frac{1}{8} \\ &= 2.25 \end{aligned}$$

~~Entropy~~ Average length = 2 bits/symbol

for uniform distribution.

Entropy = measure of uncertainty in the source
 or amount of info in the source.

Joint Entropy (for joint pmf)

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

$$H(X|Y) = \sum_{x \in X} \sum_{y \in Y} p_{x,y}(x, y) \log \frac{1}{p_{x,y}(x|y)}$$

Conditional entropy

$p_{x,y}$

	X	1	2	3	4
Y					
1		$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{32}$
2		$\frac{1}{16}$	$\frac{1}{8}$	$\frac{1}{32}$	$\frac{1}{32}$
3		$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$
4		$\frac{1}{4}$	0	0	0
		$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$

$$P(Y=1) = \frac{1}{4}$$

$$P(Y=2) = \frac{1}{4}$$

$$P(Y=3) = \frac{1}{4}$$

$$P(Y=4) = \frac{1}{4}$$

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

$$P_{X|Y=1}(x|y) = \sum_n p_{x,y}(x, y) = 1$$

$$P_{X|Y=1} = \frac{p_{x,y}(x, 1)}{P_Y(1)} = \frac{\frac{1}{8}}{\frac{1}{4}} = \frac{\frac{1}{16}}{\frac{1}{4}} = \frac{\frac{1}{32}}{\frac{1}{4}} = \frac{\frac{1}{32}}{\frac{1}{4}}$$

$$\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}$$

$H(X)$ = amount of uncertainty in x or amount of information in x .

$H(X|Y)$ = amount of uncertainty in x even after

observing or amount of info in X even after observing Y .

$H(X/Y) \leq H(X) \rightarrow$ amount of uncertainty of info in X

amount of uncertainty of info in X given Y

suppose if $X = g(Y)$

$H(X/Y) = H(X)$
(when X, Y are independent)

$H(g(Y)/Y) = 0$ as if Y is known, then $g(Y)$ is also known

Mutual Information between two l.v.s X & Y

$P_{X,Y}$ (Joint distribution)

~~$I(X)$~~ :

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

amount of information Y is giving about X .

Chain Rule for Entropy

$$H(X,Y) = H(X) + H(Y/X)$$

$$\boxed{I(X;Y) = I(Y;X)}$$

$$= H(Y) + H(X/Y)$$

$$\Rightarrow H(X) - H(X/Y) = H(Y) - H(Y/X)$$

Joint pmf = product of marginal pmf and conditional pmf.

$$P_{X,Y}(x,y) = P_X(x) \cdot P_{Y/X}(y/x)$$

Products in probability get converted to sums

In entropy as we are taking log.

In communication,

X (transmitter), Y (receiver).

Relative Entropy

between 2 pmfs.

is # proxy for distance between two pmfs.

$P(n), q(n), n \in X$.

$$D(p||q) = \sum_{n \in X} p(n) \log \left(\frac{p(n)}{q(n)} \right)$$

contribution $p(n)=0, q(n)\neq 0$

to sum $\rightarrow = 0$

$p(n) \neq 0, q(n)=0 \rightarrow +\infty$

$p(n) \neq 0, q(n) \neq 0 \rightarrow 0$

$$\boxed{D(p||q) \neq D(q||p)}$$

Relative entropy between joint pmfs of $X \& Y$

f product of marginals of $X \& Y$

$$I(X;Y) = D(P_{XY} || P_X P_Y)$$

relative entropy gives how far are the r.v.s from being independent

$$\text{Proof: } I(X;Y) = H(X) - H(X|Y)$$

$$D(P_{XY} || P_X P_Y) = \sum_{x,y} P_{XY}(x,y) \log \left(\frac{P_{XY}(x,y)}{P_X(x) P_Y(y)} \right)$$

$$= \sum_{x,y} P_{XY}(x,y) \log \left(\frac{P_{XY}(x,y)}{P_X(x) P_Y(y)} \right)$$

$$= \sum_{x,y} p_{x,y}(x,y) \log \frac{p_{x,y}(x,y)}{p_x(x)}$$

$$= \sum_{x,y} p_{x,y}(x,y) \log \frac{1}{p_x(x)}$$

$$= \sum_x \left(\sum_y p_{x,y}(x,y) \log \frac{1}{p_x(x)} \right)$$

$$= \sum_n p_x(n) \log \frac{1}{p_x(n)}$$

Prove that: $D(p||q) \geq 0$.

HInt: use the fact that $\forall z \in \mathbb{R}, \ln z \leq z - 1$

$$I(x; y) = H(x) - H(x|y)$$

$$I(x; y) = D(p||q)$$

$$D(p||q) \geq 0 \Rightarrow [H(x) \geq H(x|y)]$$

$$D(p||q) = \sum_n p(n) \log \frac{p(n)}{q(n)} = -\sum_n p(n) \log \frac{q(n)}{p(n)}$$

$$D(p||q) \geq \sum_n p(n) \left(\frac{q(p(n))}{p(n)} - 1 \right)$$

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

Quiz 1

$$1. (i) n(t) \xrightarrow{\text{F}} X(j\omega) * X(j\omega)$$

$$n_1(t) \cdot n_2(t) \xrightarrow{\text{F}} 2\pi X_1(j\omega) * X_2(j\omega)$$

$$\int_{-\infty}^{\infty} n(t) y(t) dt \quad \text{a hint: } \begin{cases} n(t) \\ y(t) \end{cases}$$

$$n(t) * y(t) \xrightarrow{\exists} X(j\omega) * Y(j\omega) \quad n(t) * y(t) \xrightarrow{\exists} c \cdot X(j\omega) * Y(j\omega)$$

Inverse transform of $2\pi x_1(jw) * x_2(jw)$

$$= \frac{1}{2\pi} \int_{-\infty}^{\infty} 2\pi x_1(jw) * x_2(jw) e^{jwt} dw$$

$$= \int_{w=-\infty}^{\infty} \left(\int_{w=\infty}^{\infty} x_1(jw) x_2(jw - t) dw \right) dt$$

=

$$x_1(t) * x_2(t)$$

1. $x(t) = \sum_{n=-\infty}^{\infty} x(n) \delta(t-n)$

$$(x(t))H = (x(t))\delta = (x(t))$$

$$(x(t))\delta = (x(t))$$

$$\text{(ii). } \underbrace{x(+1) + x(+4)}_{\text{Nyquist rate } = w_0} \xrightarrow{\mathcal{F}} x(jw) + x(jw) e^{-jw \cdot 1} \\ = x(jw)(1 + e^{-jw})$$

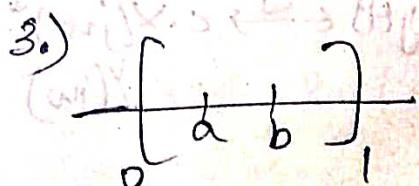
2) \mathcal{F}_1 and $\mathcal{F}_2 \rightarrow$ event space
Is $\mathcal{F}_1 \cup \mathcal{F}_2$ an event space? \rightarrow No.

$\mathcal{F}_1 = \{\emptyset, \Omega, A, A^c\}$. with $A, B \rightarrow$ disjoint

$\mathcal{F}_2 = \{\emptyset, \Omega, B, B^c\}$. with $A, B \neq \emptyset, \Omega$

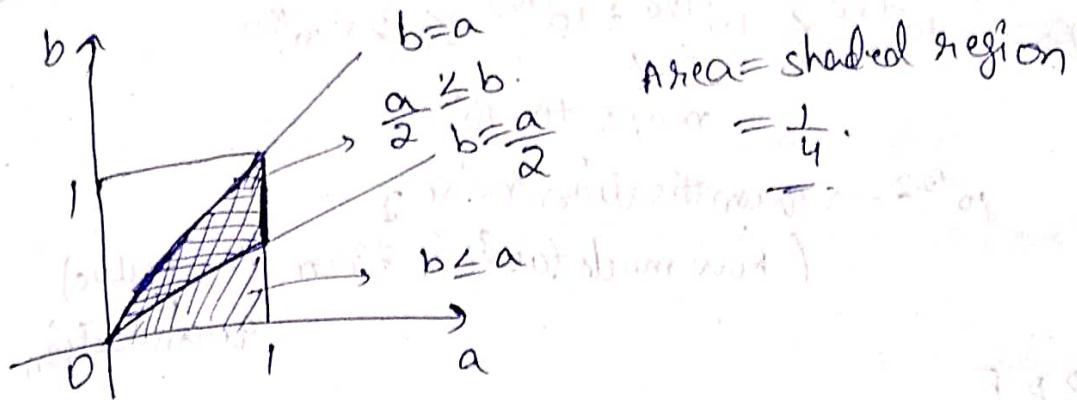
$\mathcal{F}_1 \cup \mathcal{F}_2 = \{\emptyset, \Omega, A, A^c, B, B^c\}$.

$\mathcal{F}_1 \cap \mathcal{F}_2 \rightarrow$ Yes, this is missing $A \cup B$.
an event space



what is $P(1 \leq \frac{a}{b} \leq 2)$

$$= P(b \leq a \leq 2b)$$



when there are 2 random variables independent of each other, then their joint density will be the product of the two densities.

$$F_a(a) = a, \quad F_b(b) = b. \quad \hookrightarrow \text{CDF}$$

$$f_{ab}(a, b) = f_a(a) \times f_b(b) \rightarrow \text{PDF}$$

$$= 1.1 \begin{cases} 1 & , 0 \leq a, b \leq 1 \\ 0 & , \text{otherwise} \end{cases}$$

midterm.

$$Q_1) \quad y(t) = x_1(t), \quad y_2(t)$$

$$y(t) \leftarrow \underbrace{2\pi x_1(j\omega)}_{\cdot} * \underbrace{x_2(j\omega)}_{\cdot}$$

Bandwidth = sum of bandwidths
of both x_1 and x_2

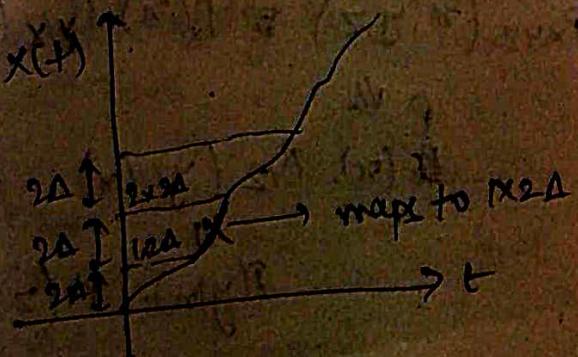
$$= \text{Bandwidth } x_1(j\omega) + \text{BW of } x_2(j\omega)$$

$$= 1000\pi + 2000\pi$$

$= 3000\pi$

$$\text{So Nyquist rate} = 2 \times 3000\pi = \underline{6000\pi}$$

Q3.) Step size = $2A$ $x(t)$



$$\text{Eg. } 10^{-120} \leq 10^{-120} + 10^{-122} \leq 2 \times 10^{-120}$$

\rightarrow maps to 10^{-120}

$10^{-122} \rightarrow$ quantisation noise

(how much far it is from the value)

quantisation.

P.D.F

$$f_e(e) = \frac{1}{2\Delta} \int_{-\Delta}^{\Delta} f_e(e) de = 1$$

$$= \int_{-\Delta}^{\Delta} f_e(e) de = \int_{-\Delta}^{\Delta} \frac{1}{2\Delta} de = \underline{\underline{1}}$$

$X \rightarrow$ uniform random variable ($[a, b]$)

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b, \\ 0, & \text{otherwise.} \end{cases}$$

Given: $e \sim$ uniform random variable ($[-A, A]$):

$$f_e(e) = \begin{cases} \frac{1}{2A}, & -A \leq e \leq A \\ 0, & \text{otherwise} \end{cases}$$

$$E(e) = \int_{-\infty}^{\infty} e \cdot f_e(e) de = \int_{-A}^{A} e \cdot \frac{1}{2A} de$$

$$E(e^2) = \frac{1}{2A} \left(\frac{e^2}{2} \right) \Big|_{-A}^{A} = \frac{1}{2A} \times \frac{1}{2} (A)^2 - (-A)^2 = 0$$

$$\text{variance} = E(X^2) - (E(X))^2$$

$$4.) P_{XYZ}(x, y, z) \equiv P_x(x) \cdot P_{YZ|X}(y/x) \cdot P_{Z|Y}(z/y)$$

\Downarrow

$$P_x(x) \cdot P_{YZ|X}(y/x)$$

$$P(y/x) \cdot P(z/y)$$

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

(chain rule)

$$= P_X(x) \cdot P_{Y|X}(y|x) \cdot P_{Z|Y}(z|y) \rightarrow \text{given } 3 \text{ in } \Omega$$

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

$$\Rightarrow \boxed{P(z|y) = P(z|y, x)} \quad \begin{cases} z \perp\!\!\!\perp X \\ (\text{conditioned on } y) \end{cases}$$

z is independent of X
conditioned on y .

$$\nabla P(yz) = P(y) \cdot P(z|y)$$

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

$$\Rightarrow \text{usually } H(X, Y, Z) = H(X) + H(Y, Z|X)$$

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

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

$$\boxed{P(z|y) = P(z|y, x)} \quad \forall z, y, x$$

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

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

Markov condition.

$$\cancel{H(Z|Y, X)} \quad \text{if } \cancel{Z \perp\!\!\!\perp X \text{ conditioned on } Y}$$

$$H(Z|Y)$$

$$= H(X) + H(Y|X) + H(Z|Y)$$

- Q.) Tossing X coins, such that exactly r_1 heads have appeared until the x th toss.

$$(a). P(X=x) = \binom{n-1}{r_1-1} (p)^{r_1} (1-p)^{n-r_1}$$

\leftarrow $= p \cdot \binom{n-1}{r_1-1} \cdot (1-p)^{n-r_1} \cdot p^{r_1-1}$

\leftarrow $= p^{r_1} \cdot \binom{n-1}{r_1-1} \cdot (1-p)^{n-r_1}$

Suppose $r=1 \rightarrow P(X=n) = p, (1-p)^{n-1}$

(this shows X is a geometric variable of p)

$$\text{Therefore } E[X] = \frac{1}{p} (1-p) + 2(1-p)^2 + \dots$$

In general for $r=r$. ($r>1$)

$$\Rightarrow E[X] = \frac{r}{p}$$

X is independent of ω

$\forall \omega \in \Omega$

$$E[X] = \sum_{\omega} p(\omega) E[X| \omega] = (\sum_{\omega} p(\omega) X(\omega))$$

$$(X(\omega_1) + X(\omega_2) + \dots + X(\omega_n)) = (X(\omega_1) + \dots + X(\omega_n))$$

$$(X(\omega_1) + \dots + X(\omega_n)) = (\sum_{\omega} p(\omega) X(\omega))$$

$$(\sum_{\omega} p(\omega) X(\omega)) + (\sum_{\omega} p(\omega) Y(\omega)) = (X + Y)$$

Source Coding (compressing source to its entropy).

operational Significance - practically what does the quantity indicate.

Let us have an information source $\rightarrow X$

X takes values from \mathcal{X} .

PMF of X is denoted by P_X .

Source code is a mapping from X to D^* .

D^* - set of all binary strings.

$$= \{ 0, 1, 00, 10, 01, 11, 000, 100, \dots \}$$

all binary strings

(countably infinite terms).

$$C: X \rightarrow D^*$$

Source code 1

$$\begin{array}{l} 1 \rightarrow 00 \\ 2 \rightarrow 01 \\ 3 \rightarrow 10 \\ 4 \rightarrow 11 \end{array} \quad \left\{ \begin{array}{l} C(1) = 00 \\ C(2) = 01 \\ C(3) = 10 \\ C(4) = 11 \end{array} \right.$$

Source code 2

$$\begin{array}{l} 1 \rightarrow 0 \\ 2 \rightarrow 10 \\ 3 \rightarrow 110 \\ 4 \rightarrow 111 \end{array} \quad \left\{ \begin{array}{l} C(1) = 0 \\ C(2) = 10 \\ C(3) = 110 \\ C(4) = 111 \end{array} \right.$$

Examples of source code

1.) ASCII - character code for English text

7-bit code.

33 control characters \Rightarrow 128 characters.

95 printing characters.

(fixed length code)
equal no. of bits for all characters.

variable length code → the no. of bits are not same
for all elements in source code

Morse code

characters more frequently used should be given less number of bits,
less frequently used characters should be assigned more no. of bits.

→ Morse code using Dots & Dashes

A. - → more frequently

Z - - - - . . → least frequently

I .. → most frequently

C - - . → less frequently

C: X → D *

symbols of the source

codewords

they are different.

Prefix-free source code

A code C is said to be prefix-free if no codeword is the prefix of any other codeword.

$$\begin{aligned} C(1) &= 1 \\ C(2) &= 10 \\ C(3) &= 01 \\ C(4) &= 110 \end{aligned}$$

$$C(2) = \overline{10}, C(4) = \underline{110}$$

C(1)

so the codeword of C(1) is a prefix of codeword of C(2), C(4)

not prefix-free code

why is prefix free property important?

Important for instant decodability

as you are passing the coded sequence, as soon as the codeword ends, you can decode

0101101011100100 1 2 3 2 2 1
↓ ↓ ↓ ↓
1 2 3 4 5 6 7 8 9 10

(Self functionality)

$X \xrightarrow{D^*}$

$x \rightarrow c(x) \quad l(n) = \text{length of codeword assigned to symbol } n.$

Average length of a source code c is denoted by $L(c)$

$$L(c) = \sum_{x \in X} l(x) \cdot P_X(x)$$

Eg-1) $X. \quad p(1) = p(2) = p(3) = p(4) = \frac{1}{4}$

$C_1 \quad 1 \rightarrow 00, 2 \rightarrow 01, 3 \rightarrow 10, 4 \rightarrow 11$

$C_2 \quad 1 \rightarrow 0, 2 \rightarrow 10, 3 \rightarrow 110, 4 \rightarrow 111$

$$\boxed{H(X) = 2 \text{ bits/symbol}}$$

$$L(C_1) = 2 \text{ bits/symbol}$$

$$H(X) = 2 \text{ bits/symbol}$$

$$L(C_2) = \frac{1+2+3+3}{4} = 2.25 \text{ bits/symbol}$$

this code is better for this source.

Eg-2) $X. \quad p(1) = \frac{1}{2}, p(2) = \frac{1}{4}, p(3) = \frac{1}{8}, p(4) = \frac{1}{8}$

$C_1 \rightarrow 1 \rightarrow 00, 2 \rightarrow 01, 3 \rightarrow 10, 4 \rightarrow 11$

$C_2 \rightarrow 1 \rightarrow 0, 2 \rightarrow 10, 3 \rightarrow 110, 4 \rightarrow 111$

$$L(C_1) = 2 \text{ bits/symbol}; L(C_2) =$$

$$L(c_2) = 1 \times \frac{1}{2} + 2 \times \frac{1}{4} + 3 \times \frac{1}{8} + 3 \times \frac{1}{8}$$

$$= 1.75 \text{ bits/symbol}$$

$$H(X) = 1.75 \text{ bits/symbol}$$

$$L(c_2) = H(X)$$

Theorem: If a code c is prefix-free for a source X , then

$$\boxed{L(c) \geq H(X)}$$

Kraft's Inequality: If a code is prefix-free, then $\{l(n)\}$ satisfy the following condition.

$$\sum_{x \in X} 2^{-l(x)} \leq 1$$

Prove the theorem using the following 2 results

1.) Kraft's Inequality

2.) $D(p||q) \geq 0$

Proof: $D(p||q) = \sum_{n \in X} p(n) \log \left(\frac{p(n)}{q(n)} \right)$

Taking $q(n) = 2^{-l(n)}$

$$\sum_{n \in X} 2^{-l(n)}$$

$$D \left\{ p(n) || \frac{2^{-l(n)}}{\sum_{n \in X} 2^{-l(n)}} \right\} \geq 0$$

$$\Rightarrow \sum_{n \in X} p(n) \log \left(\frac{p(n) \cdot \sum_{n \in X} 2^{-l(n)}}{2^{-l(n)}} \right) \geq 0.$$

$$\begin{aligned}
 & \sum p(n) \log p(n) + \sum p(n) \log (\leq 2^{-l(n)}) \geq \\
 & \quad \text{---} \quad \text{---} \quad \text{---} \quad \text{---} \\
 & \Rightarrow -H(X) + (-L(C)) \leq \sum p(n) \log (2^{-l(n)}) \\
 & \Rightarrow H(X) + K \geq L(C) \\
 & \Rightarrow \boxed{H(X) \leq L(C)}
 \end{aligned}$$

Assuming,

$$L(C) - H(X) \geq 0$$

$$\sum l(n) p(n) + \sum p(n) \log(p(n)) \geq 0$$

$$\Rightarrow + \sum p(n) \log \left(\frac{1}{2^{-l(n)}} \right) + \sum p(n) \log p(n) \geq 0.$$

\Rightarrow

$$\Rightarrow \sum p(n) \log \left(\frac{p(n)}{2^{-l(n)}} \right) \geq 0$$

\Rightarrow From $D(P||Q) \geq 0$, we know that

$$\sum p(n) \log \left(\frac{p(n)}{\frac{2^{-l(n)}}{\sum 2^{-l(n)}}} \right) \geq 0$$

$$\Rightarrow \sum_n p(n) \log \frac{p(n)}{2^{-l(n)}} + \sum_n p(n) \log \left(\sum 2^{-l(n)} \right) \geq 0.$$

$$L(C) \geq H(X)$$

$$x \quad 1 \quad 2 \quad 3 \quad 4$$

$$P(1) = \frac{1}{3}, P(2) = \frac{1}{6}, P(3) = \frac{1}{12}, P(4) = \frac{1}{12}$$

There are several sources where $L(C) \neq H(X)$
 or ~~very~~ $L(C)$ is very close to $H(X)$
 not

Only for distributions with powers of 2, i.e.

$$L(C) = H(X)$$

Designing an optimal prefix-free source code for a general source is equivalent to solving this optimization problem:

$$\min \sum_{x \in X} l(x) p_x(x)$$

such that

$$\sum_{x \in X} 2^{-l(x)} \leq 1$$

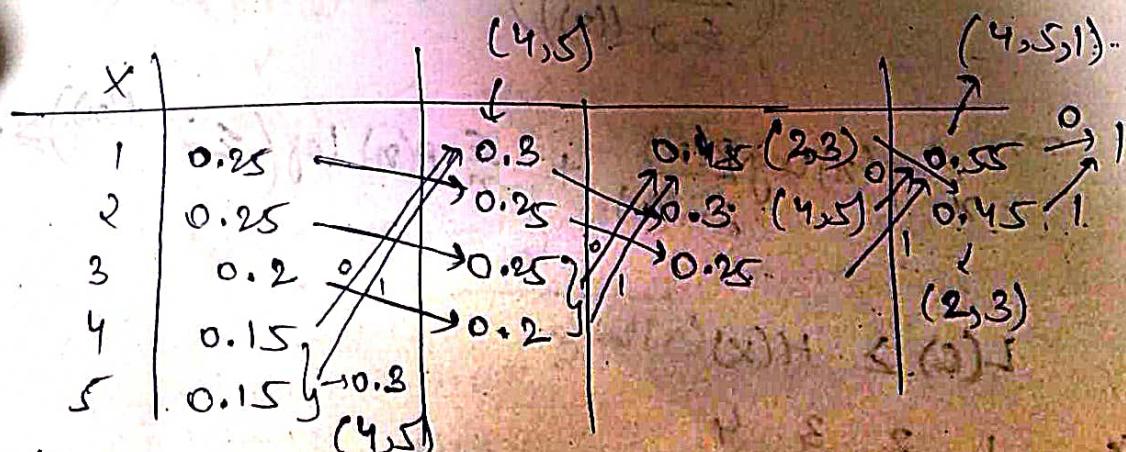
To solve the optimization problem, a technique known as Huffman Coding is used.

(JPEG).

Example of Huffman Coding — algorithm can be used for any source.

$$P(1) = 0.25, P(2) = 0.25, P(3) = 0.2$$

$$P(4) = 0.15, P(5) = 0.15$$



Step 1:

combine the least 2 probabilities into a super symbol

Step 2: Keep repeating the same process until you get

only a single super symbol

(Binary code, so we are combining 2 symbols)

- $\Rightarrow 1 \rightarrow 01$
- $2 \rightarrow 10$
- $3 \rightarrow 11$
- $4 \rightarrow 000$
- $5 \rightarrow 001$

(the 0's and 1's are taken from behind f.e. f.f starting from 1 to 0.55, 0.45 to ...)
in reverse direction

If a code is prefixfree, it satisfies Kraft's inequality.

$$\sum_{x \in X} 2^{-l(x)} \leq 1$$

If a code is prefixfree, then $L(C) \geq H(X)$

In general, we should solve the following optimization problem.

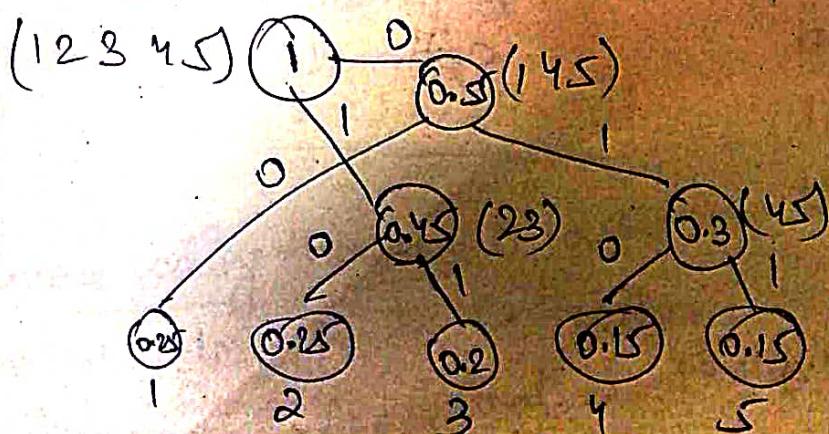
$$\min_{\{l(x)\}} \sum_{x \in X} l(x) p_x(x)$$

$$\sum_{x \in X} 2^{-l(x)} \leq 1$$

Huffman Algorithm / Huffman Code

Given $X = 1, 2, 3, 4, 5$

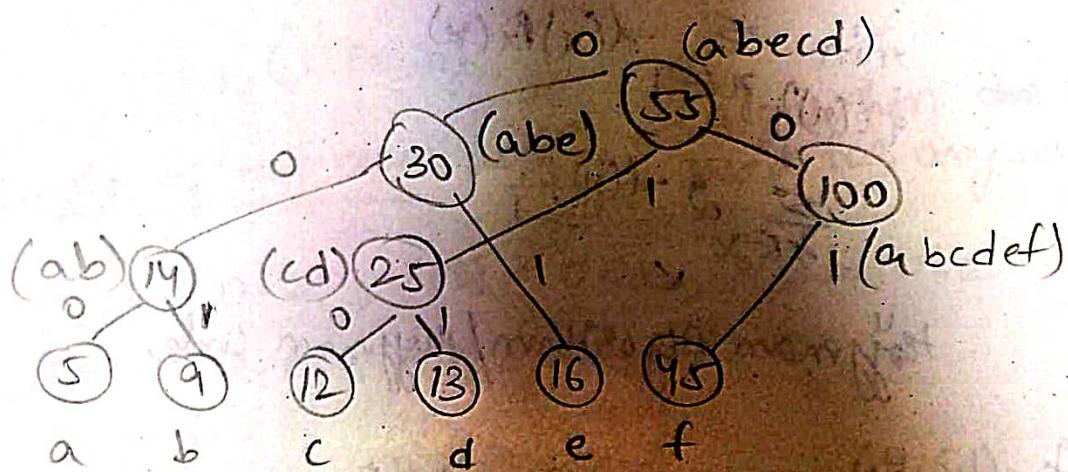
$p_x(x)$ 0.25 0.25 0.2 0.15 0.15



- 1 → 00
 2 → 10
 3 → 11
 4 → 010
 5 → 011

Eg. 2. characters frequency

a	5
b	9
c	12
d	13
e	16
f	45



- a → 0000
 b → 0001
 c → 010
 d → 011
 e → 001
 f → 1

for any source,
Huffman code is
optimal

$$L(C) \geq H(X)$$

$$L(C) \leq H(X) + 1$$

Any source code which follows Kraft's inequality, can be written in the form of prefixfree code.

$$l(n) = \left\lceil \log_2 \frac{1}{P(x)} \right\rceil$$

$$\sum_{x \in X} 2^{-l(n)} \leq \sum_{x \in X} 2^{-\log_2 \left(\frac{1}{P(x)} \right)}$$

$$\sum_{x \in X} 2^{-l(n)} \leq 1$$

$\Rightarrow \sum_{x \in X} P(x) \leq 1 \Rightarrow$ sum of all probabilities = 1

$$\boxed{\sum_{x \in X} 2^{-l(n)} \leq 1}$$

For the choice of lengths $l(n) = \left\lceil \log_2 \frac{1}{P(x)} \right\rceil$

$$L(C) = \sum_{x \in X} P(x) \left\lceil \log_2 \frac{1}{P(x)} \right\rceil$$

$$< \sum_{x \in X} P(x) \left[\log_2 \frac{1}{P(x)} + 1 \right]$$

$$\leq H(X) + 1$$

$$\boxed{L(C) \leq H(X) + 1}$$

Source coding theorem: A source can be compressed arbitrarily close to its Entropy.

IID \rightarrow independent identically distributed $H(X) + \epsilon$

IID s.v.s $x_1, x_2, \dots, x_n \sim p_x(x)$

$$p_{x_1, x_2, \dots, x_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p_{x_i}(x_i)$$

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

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

$$= n H(x)$$

distribution is similar
to $p_x(x)$ as the s.v.s are
identically distributed.

A source code can be compressed arbitrarily close
to the entropy. Designing a source code for some n .

$$n H(x) \leq L(c) \leq n H(x) + 1$$

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

by varying n , the average length can be made
as close to $H(x)$

$\frac{L(c)}{n}$ → average length per symbol / average
length per instance per symbol

If $n=2$ (from previous ex.)

$$(1,1) \rightarrow (0.25)^2 \rightarrow 0000 \quad (\text{appending both source codes})$$

$$(1,2) \rightarrow (0.25)^2 \rightarrow 0010$$

$$(1,5) \rightarrow (0.25)(0.15) \rightarrow 00011$$

$$(5,5) \rightarrow (0.15)^2$$

$$1 \rightarrow 00$$

$$2 \rightarrow 10$$

$$3 \rightarrow 11$$

Show the following inequality.

$$\frac{L_2}{2} \leq 4$$

$4 = \text{optimal length when } n=1$

$L_2 = \text{optimal length when } n=2$

$$4 = \sum_{x \in X} p(x) \cdot \log_2 f(x)$$

$$\frac{L_2}{2} \leq H(X) + \frac{1}{2}$$

$$4 \leq H(X) + 1$$

$$H(X) + \frac{1}{2} \leq H(X) + 1$$

$$\Rightarrow \boxed{\frac{L_2}{2} \leq 4}$$

$$L_2 = \sum_{x \in X} P_2(x) l_2(x)$$

\downarrow
 $2l_1(x)$

$$\frac{L_2}{2} = \sum_{x \in X} P_2(x) l_1(x)$$

$$\Rightarrow \sum_{x \in X} P_2(x) l_1(x) \leq \sum_{x \in X} P_1(x) l_1(x)$$

$$\Rightarrow \boxed{\frac{L_2}{2} \leq 4}$$

$$\frac{L_3}{3} \leq \frac{L_2}{2} \rightarrow \text{prove.}$$

$$\frac{L_3}{3} \leq 4 \quad \text{and} \quad L_2 \leq 4 + L_2 \Rightarrow \frac{L_3}{3} \leq \frac{4}{3} + \frac{L_2}{3}$$
$$\Rightarrow \frac{L_3}{3} \leq \frac{L_2}{3} \leq \frac{L_2}{2}$$

$$\Rightarrow \frac{3L_2}{2} \leq 4 + L_2 \Rightarrow L_3 \leq \frac{3L_2}{2} \Rightarrow \boxed{\frac{L_3}{3} \leq \frac{L_2}{2}}$$

Lempel-Ziv Algorithm (LZ-77) (LZ-78)

(core of zip, gzip) → makes good length = 1
 uses Dictionary based methods → no lengths = 0
 ↓
 this method does not even assume source probability.

LZ Algorithm

Ex. string → CDDCDCDDDCDDCCDCDD

1.) Parse the string from left to right

① → \subseteq ⇒ if nothing is there in the dictionary (null)
 then you add it to the dictionary.

(0, C) $\xrightarrow{②}$ D \Rightarrow (0, D)

③ $\xrightarrow{\text{D}} (1, 1, 1)$

length of the match
 (length of the no of times repeated)

which position
 from the current position is the string repeated

repetition of C so it's 1

Window

length = 4

i.e. position - curr
 ≤ 4

position \leftarrow curr
 from curr repeated position

position curr → diff = 3

④ $\xrightarrow{\text{C}}$

(1, 3, 6)

string is repeated twice
 so size of the repeated string becomes 6.

CDDCDCDD → over

position

now

CDDCDCDDDC

Tcurr

$\rightarrow (\text{DC})$

$\Rightarrow (1, 4, 2)$

CCDCDD $(1, 4, 2)$ $(1, 3, 6)$
 ↑ curr
 position. curr.
 $(1, 1, 1)$
 DCDCDD DC is repeating
 ↑ curr
 position
 $(1, 3, 2)$ (we are going back by 4 steps) C repeating
 DCCDCDD window = 4
 ↑ curr
 $(1, 2, 2)$ can move back max 4, i.e. position can be max 4.

To get back,

CDDCDCDDCDDCDD

$(0, C)$ so C un/repeated

$(0, D)$ so D un/repeated

$(1, 1, 1)$ so repeated, 1 step back, one letter to D

repeating which is D, so CDD

$(1, 3, 6)$ so repeated, 3 steps back, 6

repeating, so CDDCDD...

also length of string

Tut

2.) $(P(n) : n \in X)$

find $(l(n) : n \in X)$ such that $\sum_{n \in X} 2^{-l(n)} \leq 1$

we know that $\sum_{x \in X} P(n) = 1$

Come up with $(l(n) : n \in X)$ so that \Rightarrow a prefix code with these lengths.

$$l(n) \triangleq \left\lceil \log_2 \frac{1}{P(n)} \right\rceil \geq \log_2 \frac{1}{P(n)}$$

$\log_2 \frac{1}{P(n)} + 1$

$$\begin{aligned} \text{satisfied. } \sum_{n \in X} 2^{-l(n)} &= \sum_{n \in X} 2^{-\lceil \log_2 \frac{1}{P(n)} \rceil} \\ &\leq \sum_{n \in X} 2^{-\log_2 \frac{1}{P(n)}} = \sum_{n \in X} P(n) = 1 \end{aligned}$$

3.) $E[L] = \sum_{l(n)} P(l(n)) l(n)$

$$= \sum_{n \in X} P(n) l(n)$$

$$= \sum_{n \in X} P(n) \left[\log_2 \frac{1}{P(n)} \right]$$

$$Y = g(x)$$

$$E[Y] = \sum_y P(y) \cdot y$$

$$= \sum_x g(x) \cdot P(x)$$

as $P(n) = P(y)$

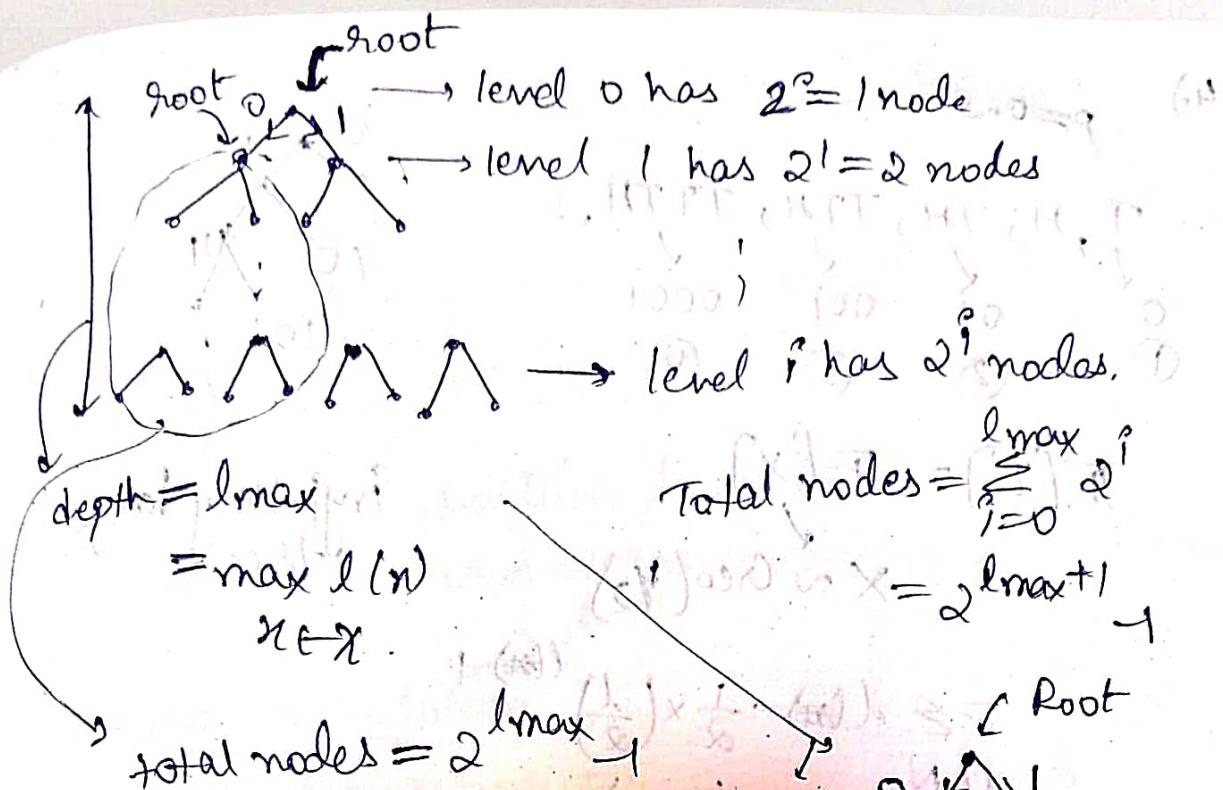
(because if we
know n , then y is also known)

$$\leq \sum_{n \in X} P(n) \left(\log_2 \frac{1}{P(n)} + 1 \right)$$

$$= H(X) + \sum_{n \in X} P(n) = H(X) + 1$$

$$H(X) \leq E[L] \leq H(X) + 1$$

corresponds to prefix code



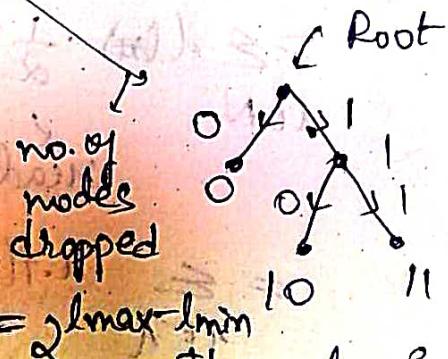
$$\text{total nodes} = 2^{l_{\max}} - 1$$

$l_{\min} = \min_{n \in X} l(n) = 1$

$$\sum_{n \in X} 2^{-l(n)} \leq 1$$

$$2^{l_{\max}+1} \leq 2^{-l(n)} \leq 2^{l_{\max}+1}$$

$$\sum_{n \in X} 2^{l_{\max}-l(n)+1} \leq 2^{l_{\max}+1}$$



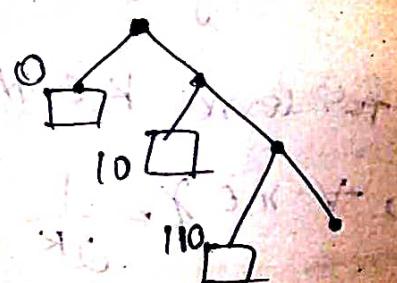
(we are dropping all the descendants)

(total no. of nodes dropped in our procedure).

we've shown

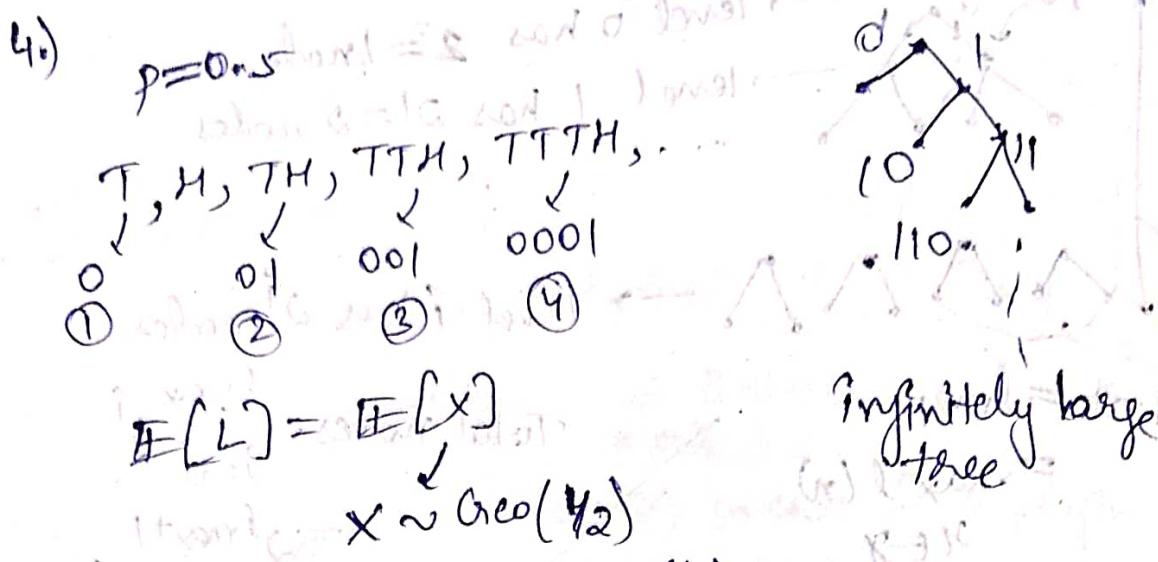
Kraft's Inequality

\Rightarrow Prefix-free code



no assigned code is assigned.

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



$$\text{with } \sum n \cdot p(n) \leq \sum n \cdot \frac{1}{2} \times \left(\frac{1}{2}\right)^n \quad (\text{why})$$

head tail.

$$\Rightarrow E[X] = \sum_{n \in \mathbb{N}} n \cdot \frac{1}{2} \times \left(\frac{1}{2}\right)^n \text{ where } p(n) = \left(\frac{1}{2}\right)^n \left(\frac{1}{2}\right)$$

$$E[L] = 2.$$

$$H(X) \leq E[L] \Rightarrow H(X) = E[L]$$

So the expected scheme
we have is optimum
coding scheme

5.) $X = \{1, 2, 3, \dots, 2^k\}$ for some $k \in \mathbb{N}$.

$$p(n) = \frac{1}{|X|}, \forall n \in X = \frac{1}{2^k}$$

Goal: Find $n \in X$, using the least no. of questions.

Perform binary search! \Rightarrow exactly k questions.

$$H(X) = \log \left(\frac{1}{p(n)} \right) \geq \log |X| = k$$

Since $p(n) = \frac{1}{2^k}$; $X \sim \text{uniformly distributed}$