

Class 12

General comment about engineering quantities [like length
compression or rate of compression]

in info theory :-

Achievability : ("a certain 'length' is achievable of comp " rate R")



There exists some scheme in which we can show that the length of compression (or) 'rate' (or) engg quantity of interest, happens to be equal to L (or) R (or) that value which is achievability

Converse:

No scheme exists which can improve upon some value.

Max

10m/s speed is 'Achievable': (Optimal) Speed \geq 10m/s.
among humans

Running rate: ① If some person who can run 10m/s.

Converse: Speed $<$ 10m/s.

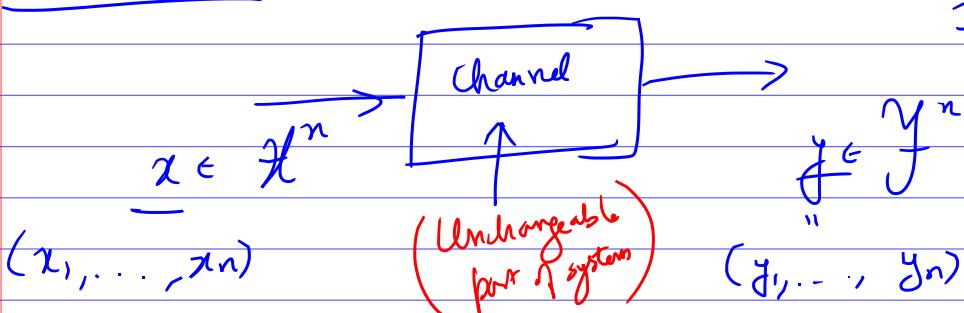
② If no person who can run at 10m/s.

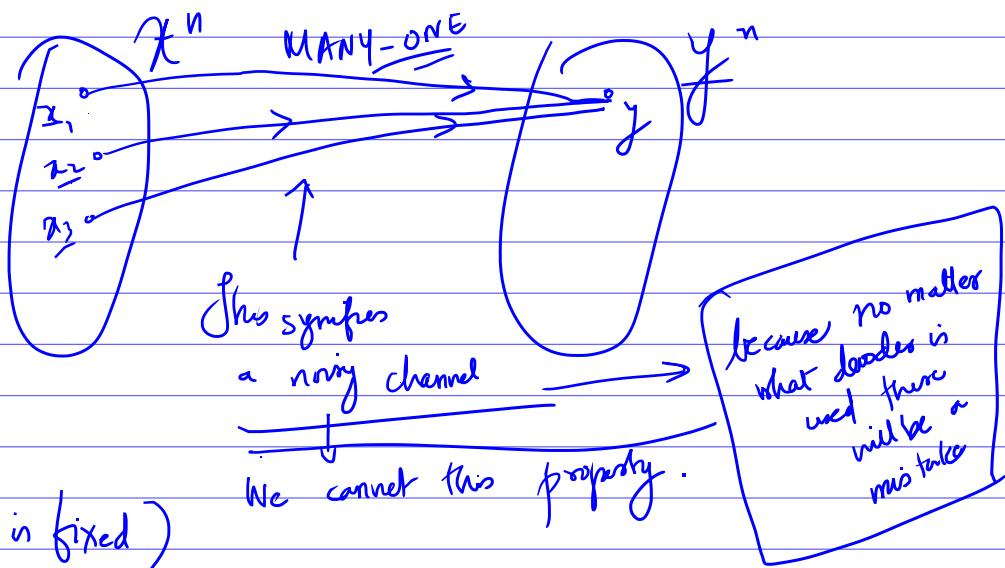
Matching converse: means "No human being can run at a speed $10 + \epsilon$, for any $\epsilon > 0$ ".

$\forall \epsilon \in (0, \infty)$

Channel Coding:

$y \in \mathcal{Y} \text{ (output)}$





(Assuming n is fixed)

To make this channel one-one (& therefore ensure correct decoding),

we emit some sequences (n -length vectors from \mathcal{X}^n)

from set of all possible transmittable sequences.

This subset of transmittable sequences is called Channel Code (or simply Code). Denoted generally by \mathcal{C}

Note that $\mathcal{C} \subseteq \mathcal{X}^n$. Each seq/vecs in \mathcal{C} is called codeword.

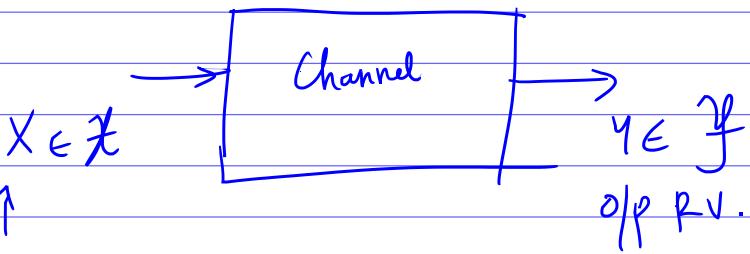
No of bits that is required to represent $|\mathcal{C}|$ codewords
 $= \log_2 |\mathcal{C}|$ bits.

Rate of the code $\mathcal{C} = \frac{\log_2 |\mathcal{C}|}{n}$ bits per channel use

bpcu
 b/cu.

Intuitively: Higher the rate, more the chance for many-one kind of system, & thus more the chance of error.

Probabilistically Noisy Channel (or) Random Channel (or) Random Noise:



$\xrightarrow{\text{If } P\text{ RV}}$ For $X = x \in \mathcal{E}$, there will be a probability distribution on the output RV. Y .

$$P_{Y|X=x} = \{ P(Y=y | X=x) : y \in \mathcal{Y} \}$$

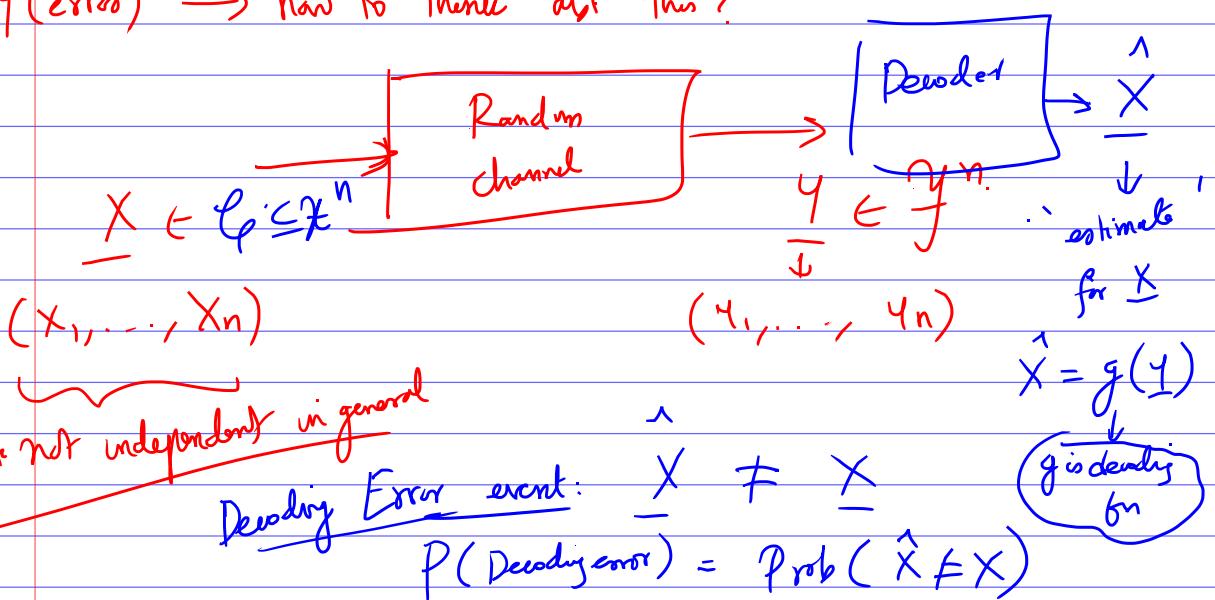
\hookrightarrow represents the cond distribution on Y gn $X=x$

For every $x \in \mathcal{E}$, we will have (kind of) distribution $P(y|x=x)$

Those distributions $P_{Y|X}(y|x)$ fix completely

characterize or describe the random channel.

$P(\text{error}) \rightarrow$ how to think abt this?



Class no 13

$$x \in \mathcal{X}^3 \rightarrow \boxed{P_{Y|x}(y|x)} \rightarrow y \in \mathcal{Y}^3.$$

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

$\mathcal{C} = \begin{Bmatrix} 000 \\ 111 \end{Bmatrix}$ instead of all eight sequences.

Intuitively $P(\text{error})$ decreases, but Rate of $C = \frac{\log_2 |\mathcal{C}|}{n}$

(Intuitively,
Prob of error can thus
also not be reduced) \leftarrow [Non zero rate code with
smaller rate is not possible for
 $n=3$] $= 1/3$.

To reduce $P(\text{error})$ further we have to increase n ,

So if we pick $n=4$, then we can get codes of rate $< \frac{1}{3}$,

& $P(\text{error})$ will be (intuitively) smaller

Intuitively, it seems like, if we want $P(\text{error}) = \epsilon$, for some extremely small ϵ , then we should expect rate to be very close to 0.

Surprisingly this is NOT the case (at least in the ideal sense)

(^{assuming} we have freedom in choosing n) :

For ANY $\epsilon > 0$, there exists a code C with $P(\text{error}) \leq \epsilon$, & rate of the code $R(\epsilon) = \max_{P_X} I(X; Y) - f(\epsilon)$ \rightarrow small for ϵ .

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

(b) Distribution of P_X

Note:

X is not a 'natural source' upon which we have no control, but it is the output of some encoder which encodes the 'real source'.

$$\begin{aligned} P_Y \text{ also, but this is dependent on } (a) P_{Y|X=x} \\ \hookrightarrow P_Y(y) &= \sum_{x \in X} P(x,y) \\ &= \sum_{x \in X} p(x) P_{Y|X}(y|x) \end{aligned}$$

So $p_X(x)$ is generally assumed to be 'controllable' in the mathematical framework of Info. Theory.

The quantity $\max_{P_X} I(X;Y)$ is called the "Channel Capacity". Generally denoted by C .

Channel Coding Theorem: Achievability: Given in prev page.

Converse: No matter what we do, we cannot get a code with rate $> C$ (channel capacity) & expect small probability of error.

Note: To make the rate very close to C , we have to have very high value of the code length n .

Eg: Binary Symmetric Channel: (Bit-flip channel)

Binary Input /Binary output : BSC(p)

$$X = \{0, 1\} = Y.$$

$$P_{Y|X}(y|x=0) = \begin{cases} p & , \text{ if } y=1 \\ 1-p & , \text{ if } y=0 \end{cases}$$

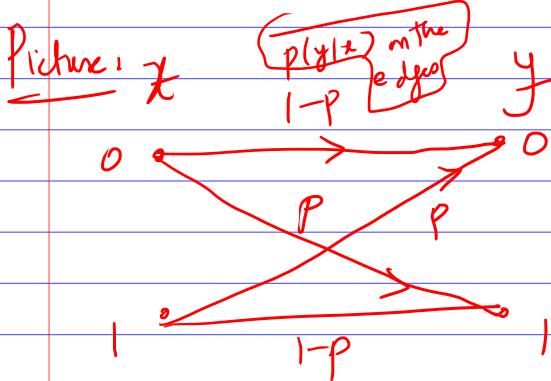
$$P_{Y|X=1}(y|x=1) = \begin{cases} 1-p & , \text{ if } y=1 \\ p & , \text{ if } y=0 \end{cases}$$

$$\begin{array}{c} X \\ \downarrow 0 \\ \left[\begin{array}{c|c} 1-p & p \\ p & 1-p \end{array} \right] \\ \downarrow 1 \end{array} \quad \left. \begin{array}{l} \xrightarrow{\text{Symmetric matrix}} \\ \text{This is a 'symmetric' channel.} \end{array} \right\}$$

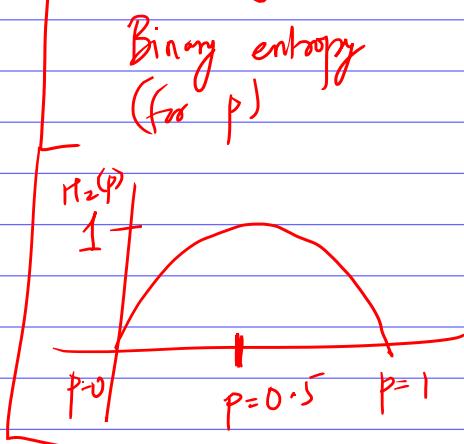
$P_{Y|X}$ matrix

Note that this means $H(Y|X=0) = H(Y|X=1)$

$$\Rightarrow H(Y|X) = H_2(p) = H_2(p) \quad \left\{ \begin{array}{l} H_2(p) \triangleq \\ -p \log_2 p - (1-p) \log_2 (1-p) \end{array} \right.$$



We want to calculate $\max_{P_X} I(X; Y)$



$$I(X; Y) = H(Y) - H(Y|X) \quad \left[\begin{array}{l} \max_{P_X} \\ = H(Y) - H(X|Y) \end{array} \right]$$

Exercise: computation

$$H(Y) = -\sum_{y \in Y} p(y) \log_2 p(y)$$

$$\max_{P_X} \left(H(Y) - H(Y|X) \right) = \left(\max_{P_X} H(Y) \right) - H_2(p)$$

$\hookrightarrow H_2(p)$ which is not affected by choice of P_X as p is

only dependent on the channel,
which cannot be changed.

[it is given statement that
channel transmission prob = p]

$$\max_{P_X} H(Y) \leq 1 \quad (\text{this is already known to us})$$

Question : Can it ever attain 1 for any P_X ?

Answer : Yes, $H(Y) = 1$ for P_X being the uniform distribution

Channel Capacity of BSC(p)

$$C_{BSC(p)} = 1 - H_2(p).$$

\hookrightarrow Exercise : (Please try !)

Show that
If P_X is uniform
then P_Y is also
uniformly distributed

(Class no 14 :-)

We will show a converse for the BSC(p) capacity.
i.e.,

★ we will show, no matter what we do,

we cannot get a rate to be $> C_{BSC(p)} = 1 - H_2(p)$

& simultaneously have small $P(\text{error})$

Let \mathcal{C} be the code which has rate R , & P_{err} very small.

We want to show that $R < C$.

$$R = \frac{\log_2 |\mathcal{C}|}{n} \Rightarrow |\mathcal{C}| = 2^{nR}$$

(n is very large : assumption)

Suppose that $\underline{c} \in \mathcal{C}$ is transmitted,

then

$$\underline{c} = (c_1, \dots, c_n) \rightarrow \text{BSC}(p) \rightarrow (y_1, \dots, y_n)$$

→ There are $\approx np$ positions in \underline{c} which are flipped

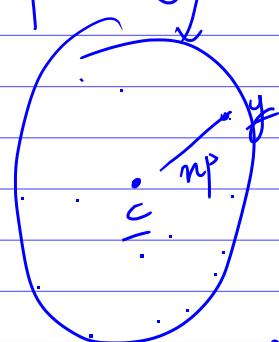
(approx)
to get \underline{y} . [channel is independently acting on
each input bit].

→ but decoder doesn't know which positions are flipped.

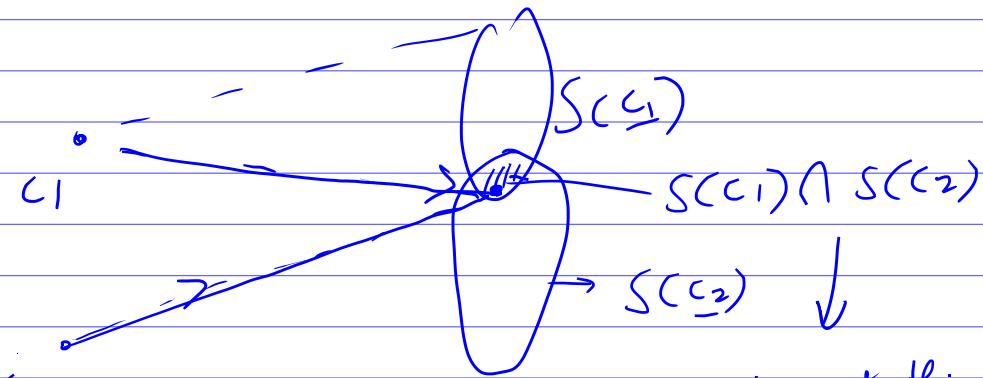
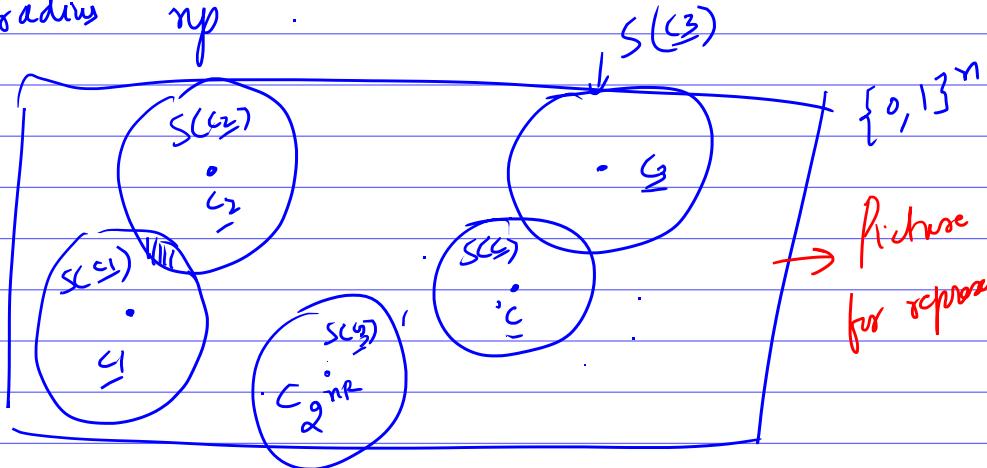
→ We can expect any seq in the set $S(\underline{c})$
as the opp response, with high probability.

$$S(\underline{c}) = \{ \underline{y} \in \{0, 1\}^n : d_H(\underline{c}, \underline{y}) = np \}$$

↓
Hamming distance bw \underline{y} & \underline{c} \triangleq no of positions
in \underline{c} which we
have to flip to get
 \underline{y} .



Now around every codeword we draw this 'Hamming ball' of radius r_{hp} .



Otherwise, if intersection is non-empty, then we want this intersection to be \sim empty.

there will be a many-one mapping from $C \rightarrow \mathbb{R}^n$, which will mean decoding error ($\text{err}(P)$) will be high.

Because the code C has small P-err, this means that the balls around the codewords in C are non-interacting.

$$\Rightarrow |C| \leq \frac{2^n}{\text{[No of vectors in any ball]}}$$

$$\hookrightarrow |S(c)| = \text{some}$$

But $|S(c)| = |\text{No of vectors in my ball}|$ for any $c \in C$.

$$= \binom{n}{r_{hp}}$$

$$|G| \leq \frac{2^n}{\binom{n}{np}}$$

Recall from
Source coding

$$\log_2 |G| \leq n - \log_2 \left(\frac{n}{np} \right) \approx n - nH_2(p)$$

$$\Rightarrow R = \frac{\log_2 |G|}{n} \leq 1 - H_2(p) \rightarrow \text{Converse is over.}$$

$$R \leq C_{BSC}(p)$$

Achievability to be done later.

Claude Shannon : "A mathematical theory of communication"
Bell's System Journal, 1948.

For a class of channels called Discrete Memoryless Channels (without feedback) Theory of error-correcting codes =

Shannon's achievability result for Channel Coding was not constructive [not identifying a specific code which works]

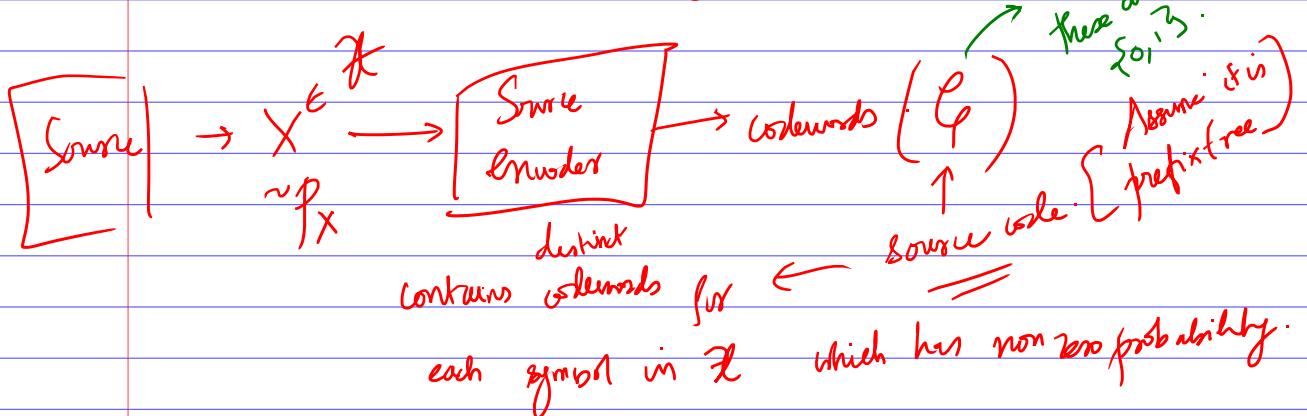
it was 'existence' proof ['We show that such a code exists but we don't know what it is precisely']

It took about 40-50 years to come up with candidate constructions which have rate close to capacity, & very small prob of error, which could actually be implemented in the real world

- ① Info theory = Used to model the comm system
 & find out limits of communication
 (without actual constructions)
- ② Coding theory = Theory ^{Synthesis & Analysis} & Practice of codes which are constructible/ explicit constructions.
- ③ Communication theory :— We have limits & rules from the above ① & ② →
- CT enables us to actually realize the comm system in practice.

Class notes :-

Go back to Source Coding:



$c(x)$ be the codeword assigned to $x \in X$.

$l(x)$ be the length of the codeword assigned to x .

(Fixed - Variable length source coding) $\left| \begin{array}{l} Q = \{ c(x) : x \in X \} \\ \uparrow \text{source code} \end{array} \right.$

length
length of some seq (here = 1)

Average length $L_Q = \sum_{x \in \Sigma} p(x)l(x)$ prefix-free

Goal: is to design a code \mathcal{L} which has minimum L_Q .

Optimal length

$$L^* = \min_{\mathcal{L}} L_Q$$

Lemma:

$$L^* \geq H(X) \cdot \text{for } \mathcal{L}^*$$

\Rightarrow Any prefix free code \mathcal{L} has average length at least $H(X)$

Proof:

To prove this we need the following claim

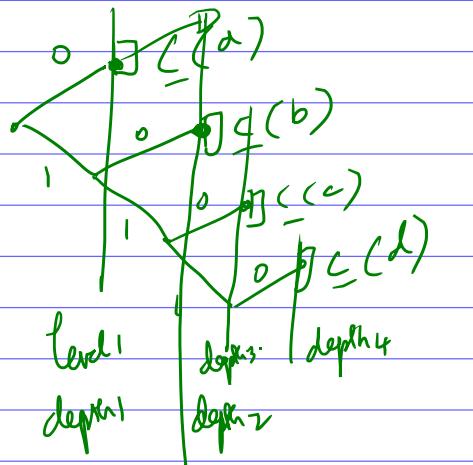
Claim: Let \mathcal{L} be any prefix-free code. Then

[Kraft's inequality] $\sum_{x \in \Sigma} 2^{-l(x)} \leq 1.$

Proof:

We know that any prefix-free code can be represented using a binary tree which has the 'leaves' as the codewords.

Σ	$C(x)$
a	0
b	1 0
c	1 1 0
d	1 1 1 0
e	1 1 1 1



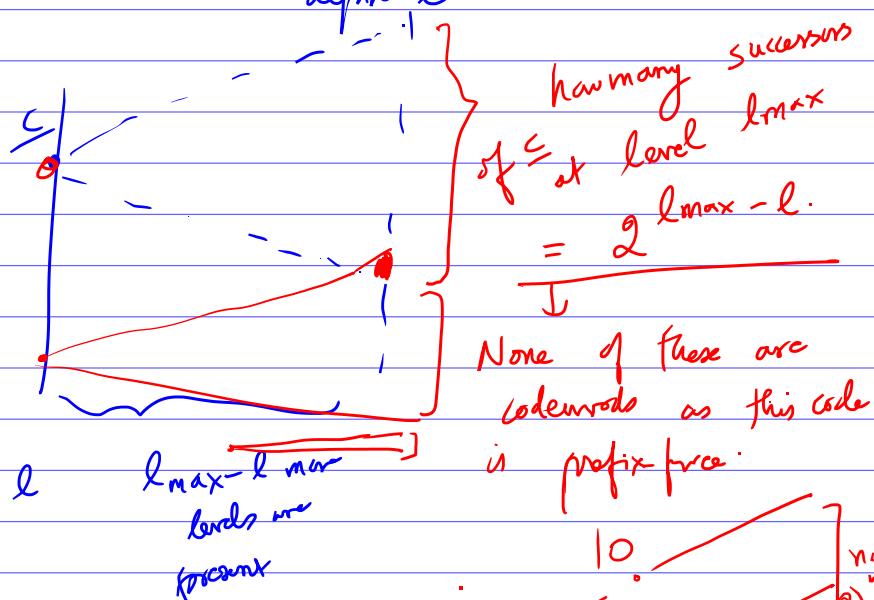
For the binary tree corresponding to the code,

What is the depth of the tree? $\rightarrow = \text{Length of the largest codeword in the code}$

$$(l_{\max} = \max_{x \in X} l(x))$$

Suppose there is a codeword \underline{c} with length l . ($l \leq l_{\max}$)

represented by a node in the tree at depth l .



$$\sum_{x \in X} 2^{l_{\max}-l(x)} \leq 2^{l_{\max}} \quad ?(A)$$

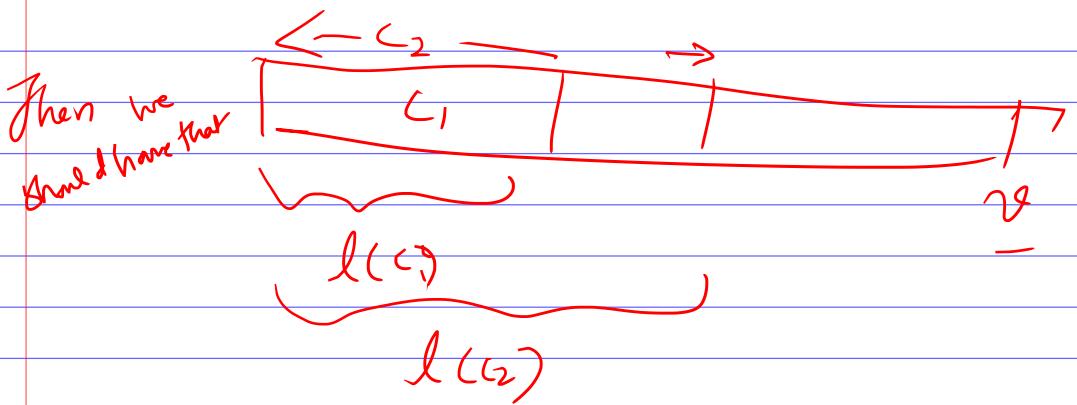
\rightarrow This is true if distinct codewords $\underline{c}_1, \underline{c}_2$ don't have common successors at l_{\max} level.

\rightarrow Any successor of \underline{c}_1 at l_{\max} level has

\underline{c}_1 as a prefix.

\rightarrow \underline{c}_1 by for \underline{c}_2 as well.

So suppose $l(\underline{c}_1) \leq l(\underline{c}_2)$ there can be a common successor \underline{c} at l_{\max} level



$\Rightarrow c_1$ should be a prefix of c_2 , which is not true as ℓ is a prefix free code.

\Rightarrow No pair of the codewords in ℓ have any common suffixes.

\Rightarrow (A) is true

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

\rightarrow This proves the claim (Kraft inequality for p-f code).

Now back to lemma prob. For any p-f code ℓ we must have $L_\ell - H(X) \geq 0$.

$$\sum_{x \in \mathcal{X}} p(x) l(x) - \sum_{x \in \mathcal{X}} p(x) \log \frac{1}{p(x)} \geq 0$$

Recall $D_C(p || q) = \sum_{x \in \text{Supp}(p_x)} p(x) \log \frac{p(x)}{q(x)}$.

Then this is a valid prob dist

Let $q_{|X}(x) \triangleq$

$$\left(\frac{2^{-l(x)}}{\sum_{x' \in \mathcal{X}} 2^{-l(x')}} \right) \geq 0$$

$\left(\sum_{x \in \mathcal{X}} q_{|X}(x) = 1 \right)$
Note that this satisfies

$$\begin{aligned}
 \text{Now } D(p_x \parallel q_x) &= \sum_{x \in \text{supp}(p_x)} p(x) \log \frac{p(x)}{\underset{x \in X}{\sum} 2^{-l(x)}} \\
 &= -\sum_{x \in \text{supp}(p_x)} p(x) \log \frac{1}{p(x)} + \sum_{x \in \text{supp}(p_x)} p(x) \log \left(\frac{1}{\underset{x \in X}{\sum} 2^{-l(x)}} \right) \\
 &= -H(X) + \sum_{x \in \text{supp}(p_x)} p(x) \log \underset{x \in \text{supp}(p_x)}{2^{-l(x)}} \\
 &\quad + \sum_{x \in \text{supp}(p_x)} p(x) \log \left(\sum_{x \in \text{supp}} 2^{-l(x)} \right) \xrightarrow{\text{Kraft.}} \leq 1 \text{ by} \\
 &\leq -H(X) + L_q
 \end{aligned}$$

$$0 \leq D(p_x \parallel q_x) \leq -H(X) + L_q$$

known result

$$\Rightarrow L_q - H(X) \geq 0$$

$$\Rightarrow L_q \geq H(X) //$$

Equality happens for $p_x = q_x$ & this term being 0.

Class 1b

Kraft's inequality: For a given prefix-free code with length function $l(\cdot)$,

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

✓

Lemma: Now suppose that we have a RV $X \in \{x_1, \dots, x_k\}$ & ^k positive integers $\sum_{i=1}^k 2^{-l_i} \leq 1$.

Then there exists a p-f code for X with codeword lengths l_1, \dots, l_k .

Proof:

We show that we can construct a binary tree

with nodes at depths (levels) l_1, \dots, l_k

such that none of those nodes are successors of each other

(\Rightarrow they are the leaves of some binary tree
 \Rightarrow they represent a valid p-f code).

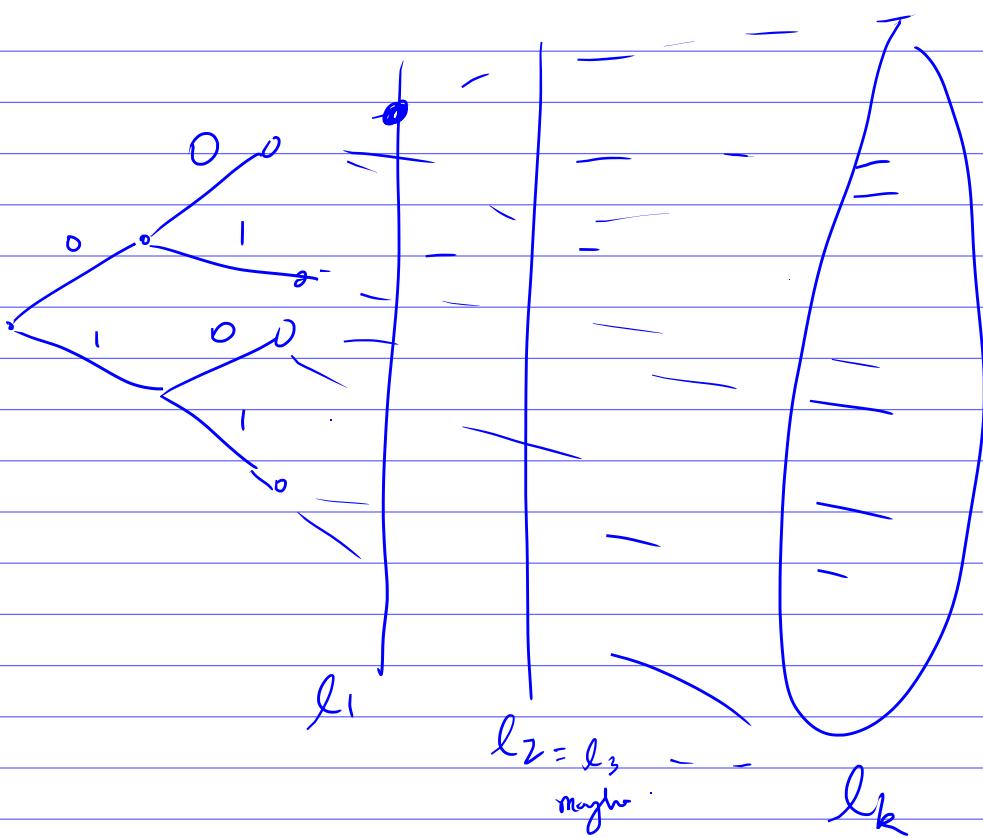
Assume that "Without loss of generality"
 $l_1 \leq l_2 \leq \dots \leq l_k$.

For any $i \leq k$, $i-1$

$$\sum_{j=1}^{i-1} 2^{-l_j} < 1 \rightarrow \textcircled{A}$$

Observation
(by given statement)

Imagine that we take the full binary tree upto level l_k .

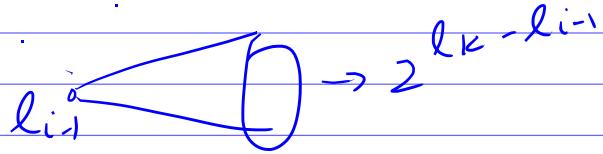
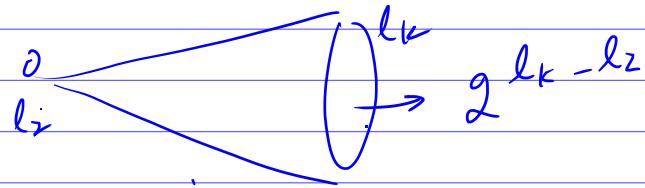
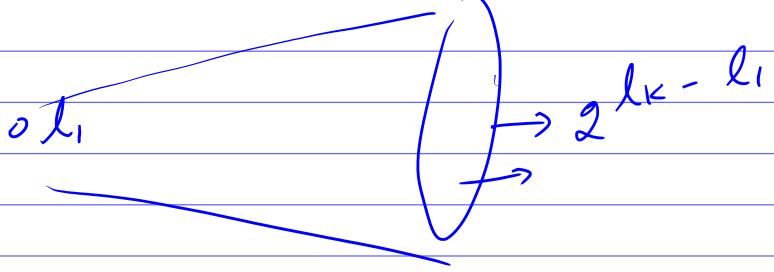


At each step i of the algorithm, we 'intend' to pick one available (undeleted) node from the above tree at level l_i & delete all its successors from the tree. Repeat this process for $i=1, \dots, k$. We will then have a 'p-f tree'.

We have to show that at each step $i=1, \dots, k$, there is at least one node left undeleted at level (depth) l_i . For this we will use Observation (A).

Base: Clearly at step 1, there is a node at level l_1 , $(l_1 > 1)$.

Induction: For any $i \leq k$, after $(i-1)$ steps, assume that we have picked nodes at levels l_1, \dots, l_{i-1} & appropriately deleted some. We want to show there is node at level i .



Total nodes at level l_k which are not in the tree after $(i-1)$ th step = $\sum_{j=1}^{i-1} 2^{l_k - l_j}$

$$\text{No. of Nodes remaining at level } l_k = 2^{l_k} \left(1 - \sum_{j=1}^{i-1} 2^{-l_j} \right) > 1$$

\Rightarrow No. of nodes remaining in tree > 1 after $(i-1)$ th step $\stackrel{\text{at } l_k}{\longleftarrow} < 1$ by obs(A)

\Rightarrow At least one survivor node should be present at level l_i also. So we can pick a node for the i th step from level l_i also. //

This completes the proof.

Remark: We construct the tree from smallest length to largest length backwards. Contrast this with

Optimal Source code construction that will follow later (Huffman Codes).

Now, suppose that the source RV $X \sim P_X$. ($\{x\} = K$)
 +ve
 [Supp(P_X) = \mathcal{X}]

We want to obtain a collection of integers l_1, \dots, l_K

such that $\sum_{i=1}^K 2^{-l_i} \leq 1$. Then we know how to

obtain a code for X .

We want small avg length
 $\sum_{i=1}^K p_i l_i$

Then choose small l_i for larger p_i .

$$\sum_{i=1}^K 2^{-l_i} \leq 1$$

Suppose all codewords are of some length

$$K 2^{-l} \leq 1$$

$$2^l \geq K$$

$$l \geq \log_2 K$$

\rightarrow Pick $l = \lceil \log_2 K \rceil$ will work ("ceil")

\rightarrow But there is no guarantee that this code is "good" i.e., it may not have small avg length L

Idea: (behind Shannon-Fano code as in McEliece book)

For $i=1, \dots, K$

$$\text{Fix } l_i = \lceil \log_2 \frac{1}{p_i} \rceil \rightarrow ①$$

where p_i is the prob of X taking the i^{th} value in \mathcal{X} .

Clearly $l_i > 1$.

Checking K -inequality:

$$\begin{aligned} \sum_{i=1}^K 2^{-l_i} &= \sum_{i=1}^K 2^{-\lceil \log_2 \frac{1}{p_i} \rceil} \leq \sum_{i=1}^K 2^{-\log_2 \frac{1}{p_i}} \\ &= \sum_{i=1}^K p_i = 1 \end{aligned}$$

\Rightarrow The lengths given by ① satisfy K -inequality.

\Rightarrow We can use the tree-fanning algorithms to get a p-f code.

This ~~Pf~~ code obtained is called as The Shannon - Fano code for X .

Q: What is the L value for Shannon Fano Code?

$$L_{\text{Shannon-Fano}} = \sum_{i=1}^k p_i \lceil \log_2 \frac{1}{p_i} \rceil \\ < \sum_{i=1}^k p_i \left(\log_2 \frac{1}{p_i} + 1 \right)$$

$$L_{\text{SF}} < H(X) + 1$$

$\overbrace{\quad}^{\text{great!}}$

→ But SF code is not always an optimal length p-f code
 L_{SF} is not ^{always} the smallest length among any p-f code for X .