# Class no 9 :- FIXED LENGTH SOURCE CODE (output of encoder is some binary tuple of a fixed length)
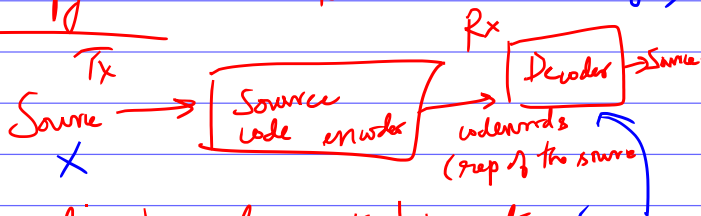
Kings Significance of entropy & other terms.

① Idea no 1:
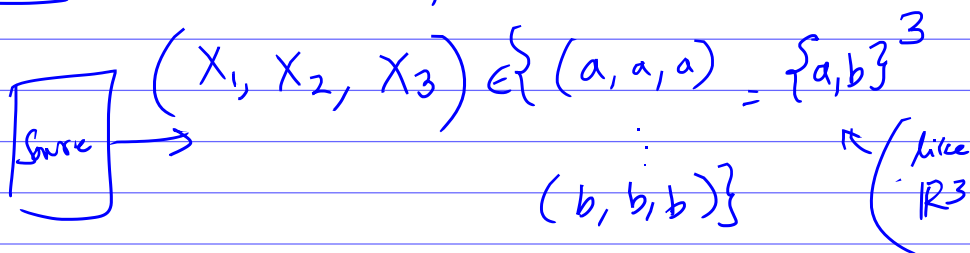


If we are willing to live with / tolerate some small probability of decoding error, we can compress the source better, => we can have smaller length for representing the source.

(decoder is ordered to know the $P_x$)

→ ignoring means, not encoding (or) mans encoding

by 'ignoring' source symbols which have very low probability of occurrence.

② Idea no 2: 'club' multiple source RV instances.



$$(X_1, X_2, X_3) \in \{ (a, a, a) = \{a, b\}^3$$
$$\vdots$$
$$(b, b, b) \} \qquad \leftarrow \left( \text{like } \mathbb{R}^3 \right)$$

( Assume that $X_1, X_2, X_3$ are all independent RVs

$$P_{X_1, X_2, X_3}(x_1, x_2, x_3) = P_{X_1}(x_1) P_{X_2}(x_2) P_{X_3}(x_3), \quad \forall \; x_1, x_2, x_3 \in \{a, b\}. )$$

We know joint distribution from the individual distributions

marginal distributions

$$\underbrace{(x_1, x_2, x_3)}_{=} \qquad (P_{X_1}(x_1), P_{X_2}(x_2), P_{X_3}(x_3))$$

$$P_{X_1, X_2, X_3}\left( \underline{x} \right) = P^{n_a(\underline{x})} (1-P)^{n_b(\underline{x})}$$

$$P_X(a) = P$$
$$P_X(b) = 1-P.$$

$n_a(\underline{x})$: no of times 'a' occurs in $\underline{x}$

$n_b(\underline{x})$ = no of times 'b' occurs in $\underline{x}$.

Suppose we have some 'compression scheme' (a mapping)

from $C_s : \{a, b\} \longrightarrow \{0, 1\}$.   $\left( C_s(x) \in \{0, 1\} \right)$

Can we use this to get a scheme for $\{a, b\}^3$ ?

Yes: $C_s' : (x_1, x_2, x_3) \longrightarrow \left( C_s(x_1), C_s(x_2), C_s(x_3) \right)$.

$: \{a, b\}^3 \longrightarrow \{0, 1\}^3$   (Fixed length code)

$C_{s'}(a, b, a) \longrightarrow (0, 1, 0)$   $\left( \begin{array}{l} \text{suppose} \\ C_s(a) = 0 \\ C_s(b) = 1 \end{array} \right)$

length of this code: 3 bits. to 3 source symbols

$\downarrow$

This code $C_s'$ is as good as the original code $C_s$.


$\longrightarrow$ In the case of encoding only one source symbol,

our possible code lengths were either 0 or 1 (only)

$\longrightarrow$ Here we have more choices 0, 1, 2, 3 length
binary strings (vectors or tuples) can be used.

$\boxed{C_{s'}} : \{a, b\}^3 \longrightarrow \{0, 1\}$   $\left[ \begin{array}{l} \text{length} = 1 \\ \text{Normalized length} = 1/3 \end{array} \right]$

This is a good code if   $(a, a, a) \Rightarrow$ very high prob

$\left\{ 7 \text{ vectors} \right\} \Rightarrow$ totally has a small prob.

In this case we can map

$C_{s'}\left( (a, a, a) \right) = 0$, & $C_{s'}(\underline{x}) = 1$ $\forall$ $\underline{x} \in \{a, b\}^3 \Big|_{(a, a, a)}$

$(\underline{x}) \longrightarrow \boxed{\text{Encoder}} \longrightarrow 0 \Rightarrow \boxed{R_x \text{ decodes}} \longrightarrow (a, a, a)$

knows the code, knows $P_x$

<u>Idea no.2</u> Suppose we are allowed to combine multiple source symbols & encode them together into some fixed length binary string, then this gives a more 'efficient' source code ( smaller (normalized length))

<u>Example:-</u>

Source $X \sim P_X$ ; $P_X(a) = p$
$P_X(b) = 1-p$

Remember : ① We are allowing for encoding long source strings & we can tolerate some small prob of error.

② We have to a <u>fixed length source code</u>

( every 'n' length source string is to be encoded into a 'l' length binary vector / string / tuple .

<u>Fixed length</u> means $l$ doesn't change with the source string. )

Assumption : 'n' - length Random Source vector is represented by

$(X_1, \ldots \ldots, X_n)$, where

$X_i$ is the RV representing $i^{th}$ output of source $\in \{a, b\}$

[ $X_i$s have same distribution ] $\Leftarrow$ $P_{X_i} = P_X$ ($P_{X_i}(a) = P_X(a)$
[ $X_i$s are independent ] $P_{X_i}(b) = P_X(b)$ )

$\longrightarrow$ In the language of Communications , $X_i : i \in 1. --, n$ are said to be 'independent and identically distributed" [ i.i.d ]

**Question:**

Suppose $n$ is 'very large', how many $a$'s and $b$'s do we expect to 'see' in the random source sequence $(X_1, \ldots, X_n)$?

No of $a$'s $\sim np$

$\ldots$ "$b$'s $\sim n(1-p)$

No of such sequences with such a distribution of $a$'s & $b$'s

Set of
Typical sequences $\sim \binom{n}{np}$ [Notice that this is only a subset of the $2^n$ sequences]

Idea of the ~~source~~ efficient code we want to use is that we will encode only these $\binom{n}{np}$ sequences with unique codewords

→ [for all other sequences we will use a single codeword]

Class no 10:

Atypical sequences := not typical $\Rightarrow$ no of $a$'s is much different from $np$.

& no of $b$'s is very deff from $n(1-p)$

Intuitively easy to say that

w·h·p (with high probability), any $n$-length sequence obtained as output of the source (running $n$ times) is going to be a __typical__ sequence.

& $P$(we will get a atypical sequence as output) is very small.

Source code: (Fixed length source code) $\rightarrow$ [ "Block $\rightarrow$ Block Source coding" (book terminology) ]

fixed length (n) - sequences

Assign to each typical sequence a unique codeword of some length. L

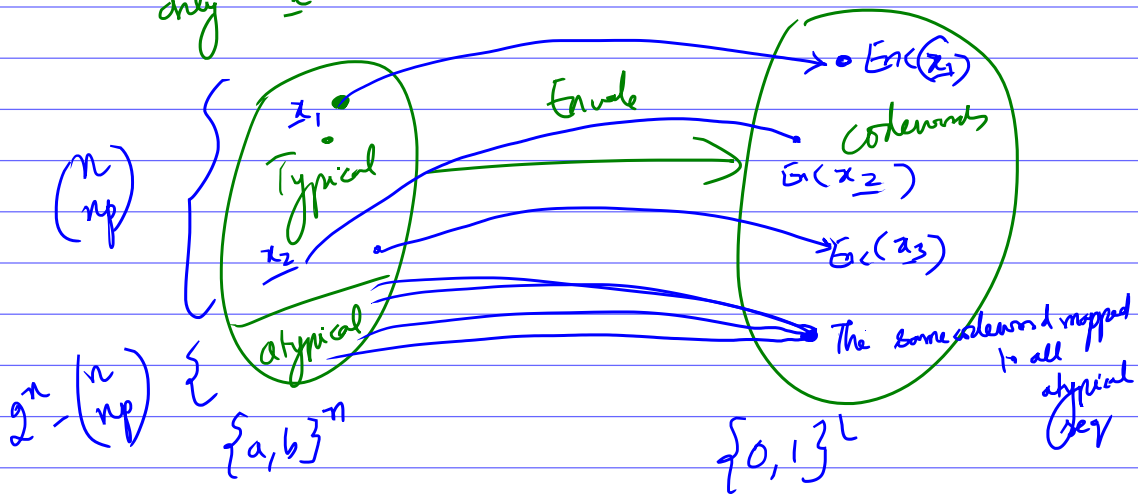[ Typical sequences $\rightarrow$ Codewords is a One-one map. ]

Assign to each atypical sequence, assign some same codeword of length L

different from those assigned to typical sequences

$\underline{x} \in$ Source output. (n length seq)

$Enc(\underline{x}) \rightarrow$ codeword associated with $\underline{x}$

If $\underline{x} \in$ Typical seq, then decoder/Rx will be able to identify $\underline{x}$ without errors as $Enc(\underline{x})$ will be uniquely associated with only $\underline{x}$



$\binom{n}{np}$ { Typical $x_1$ $x_2$

$2^n - \binom{n}{np}$ { atypical

$\{a,b\}^n$

Encode

$\rightarrow Enc(x_1)$ codewords
$Enc(x_2)$
$\rightarrow Enc(x_3)$
The same codeword mapped to all atypical seq

$\{0,1\}^L$

for defining such a map, what should be min length of L?

L is atleast $\log_2 \binom{n}{np} + 1$

$\underbrace{\phantom{\log_2 \binom{n}{np}}}_{typical}$ $\rightarrow$ for atypical.

Now,

$$\log_2\binom{n}{np} = \log_2\left(\frac{n!}{(np)!\,(n-np)!}\right)$$

$$= \log_2(n!) - \log((np)!)$$
$$\quad - \log((n(1-p))!)$$

rough
$$\approx n\log_2 n - np\log_2 np$$
$$\quad - n(1-p)\log_2(n(1-p))$$
$$\quad - o(n)$$

$\hookrightarrow$ much smaller than other terms in the summation.

$$n! = \textcircled{n}(n-1)(n-2)\ldots(n-(n-1))^+$$
$$\approx n^n - \begin{bmatrix}\text{Some poly in } n \\ \text{of degree} < n.\end{bmatrix}$$

(Better approx are there, we are doing rough analysis) $\rightarrow$ Stirlings approximation

$$= \underline{n\log n} - \underline{np\log p} - \underline{np\log n}$$
$$\quad - \underline{n(1-p)\log n} - \underline{n(1-p)\log(1-p)} - \text{small term}$$

$$= n\left[\underline{p\log_2\frac{1}{p} + (1-p)\log_2\frac{1}{1-p}} - \frac{\underline{\text{small term}}}{n}\right]$$

Remains constant as $n$ grows          $\hookrightarrow$ goes to 0 as $n\to\infty$.

$$\approx n\,H(X), \text{ where } X \text{ is the source RV.}$$

$\Rightarrow$ It is sufficient to have length of codewords $\approx \underline{n\,H(X)+1}$

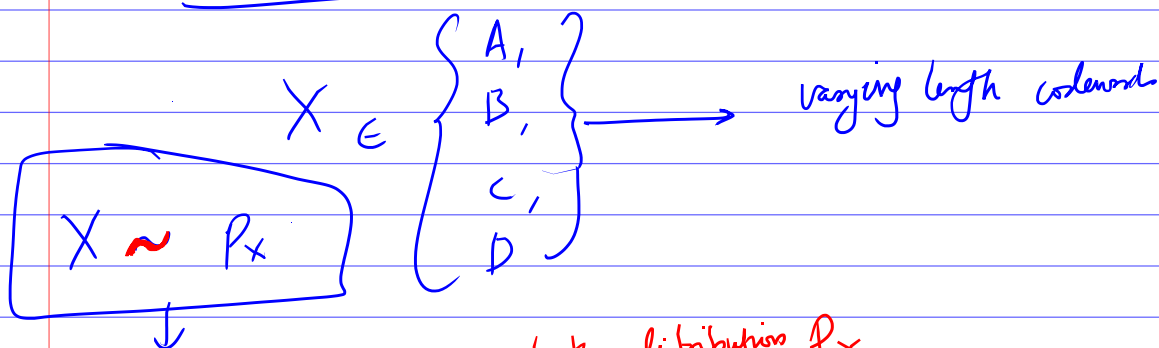$n$ length source sequences $\longrightarrow$ $n\,H(X)+1$ length source code.

Per source symbol, what is the length of the codeword?
$$= H(X) + \frac{1}{n} \quad \left(\rightarrow H(X) \text{ as } n\to\infty.\right)$$
bits

Per source symbol, we "are using" $H(X)$ bits. [sufficient]
Infact, we "need" $H(X)$ bits otherwise we will have large $P(error)$.

Fixed length source sequences to variable length codewords

$$X \in \left\{ \begin{matrix} A, \\ B, \\ C, \\ D \end{matrix} \right\} \longrightarrow \text{varying length codewords}$$

$$\boxed{X \sim P_X}$$

$\downarrow$

X RV is distributed according to the distribution $P_X$

eg
Source code 1:

A $\longrightarrow$ 0

B $\longrightarrow$ 1

C $\Longrightarrow$ 10

D $\longrightarrow$ 11

$\longrightarrow$ We expect to do 'better' than previous fixed-fixed length scenario because we have the freedom here to set varying length codewords.

$\longrightarrow$ We will therefore demand zero probability of error

$\downarrow$

In prev fixed-fixed scenario we cannot do any compression if we wanted P(error)=0.

However we have a problem :—

Suppose we use Source code1 above

Then if source generates BA. $\longrightarrow$ 10    some codeword

or    C $\longrightarrow$ 10

Confusion at decoder