

## lec 9: Data Compression

Sep 20<sup>th</sup> 2017

Defn-

A source code  $\mathcal{L}$  for a random variable  $X$  is a mapping from

$$x \rightarrow \mathbb{D}^* \text{ (strings of finite length from alphabet } \mathbb{D}\text{)}$$

$$\forall x \in \mathcal{X}_0, \Rightarrow c(x) \in \mathbb{D}^*$$

$$l(x) = \text{length of } c(x) \\ (\# \text{ } \mathbb{D}\text{-symbols in } c(x))$$

Defn-

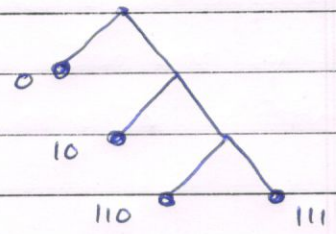
The expected length  $L(\mathcal{L})$  of a source code  $\mathcal{L}$  for a random variable  $X$  having pmf.  $p(x)$  is given by,

$$L(\mathcal{L}) = \sum_{x \in \mathcal{X}_0} p(x) l(x)$$

we assume wlog,  $\mathbb{D} = \{0, 1, \dots, D-1\}$ .

- Note that any string of encoded bits (without commas) can be uniquely decoded, if the forward mapping is 1:1. (Eg. 5.1.1 in Cover & Thomas)

The code can be visualised as a tree



Defn- A code is said to be nonsingular if every element in the range of  $c$  maps onto a different string in  $\mathbb{D}^*$ ,

ie-  $x \neq x' \Rightarrow c(x) \neq c(x')$

Defn- The extension  $c^*$  of a code  $c$  is the mapping from finite length strings of  $X_0$  to finite length strings of  $\mathbb{D}$ , defined by:

$$c(x_1 \dots x_n) = c(x_1) c(x_2) \dots c(x_n)$$

A code is said to be uniquely decodable if its extension is non-singular (however you may need to look at the entire string before decoding even the 1st source symbol).

Defn- A code is said to be a prefix code or an instantaneous code if no codeword is a prefix of any other codeword.

singular

non-singular

uniquely decodable

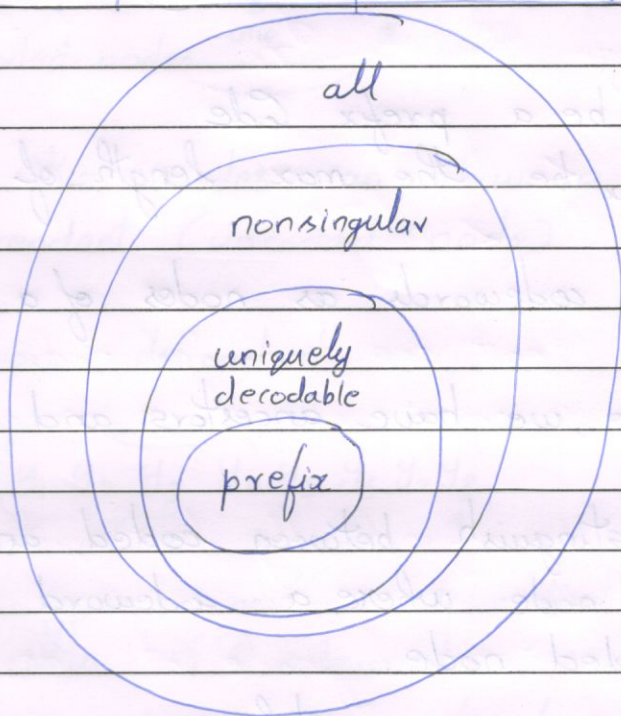
prefix



Code



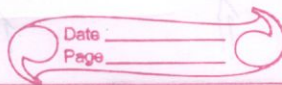
X	$\mathcal{L}_1$	$\mathcal{L}_2$	$\mathcal{L}_3$	$\mathcal{L}_4$
1	0	0	10	0
2	0	010	00	10
3	1	01	11	110
4	1	10	110	111



Nesting of different Code classes

Note -  $\mathcal{L}_3$  is uniquely decodable, but must wait till end v for decoding, correctly. of sequence

## Kraft's Inequality -



Theorem- A prefix code over an alphabet  $\mathcal{D}$  of size  $D$  and with codeword  $c(x_i)$  having length  $l_i$ ,  $1 \leq i \leq n$  exists iff,

$$\sum_{i=1}^n D^{-l_i} \leq 1 \quad \dots \quad (1)$$

Proof- Let  $\mathcal{C}$  be a prefix code  
Let  $l_{\max}$  be the max length of a codeword

Place the codewords as nodes of a  $D$ -ary tree.

In the tree, we have ancestors and descendants

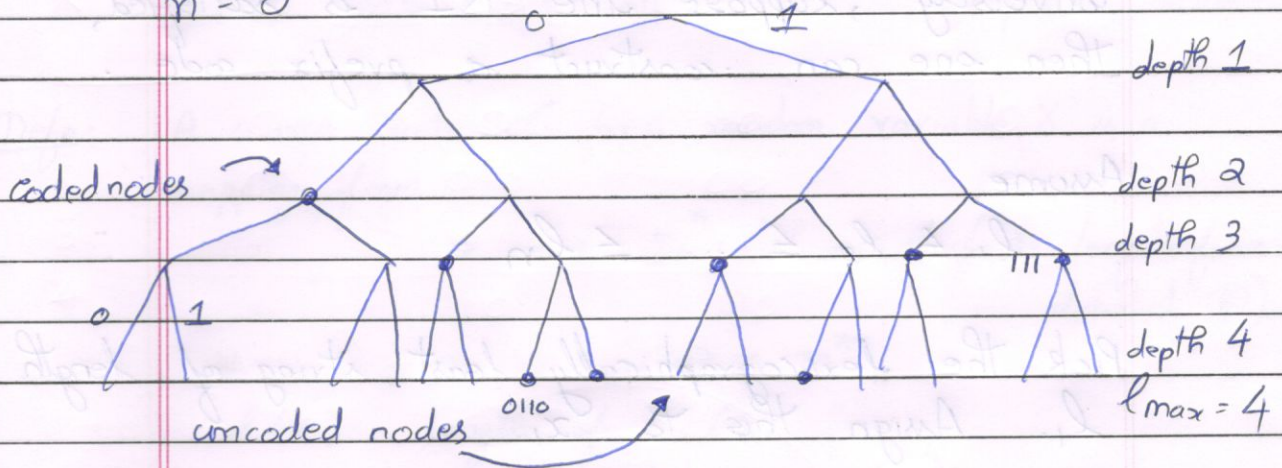
Let us distinguish between coded and uncoded nodes. A node where a codeword is located is a coded node.

The # of descendants, <sup>at  $l_{\max}$</sup>  (coded or uncoded) of a coded node where the code is of length  $l_i$  is

$$l_i = D^{l_{\max} - l_i}$$

$$D = 2$$

$$n = 8$$



Claim- No two nodes can have a common descendant (uncoded node)

Pf- Common descendant

$y_1 y_2 y_3 y_4 y_5 y_6 y_7 y_8$  ← fig A.

Say  $y_8$  is common descendant

But then  $\exists$  2 codeword  $y_1 y_2 y_3 y_4$  &  $y_1 y_2 y_3$  (say)

But  $y_1 y_2 y_3$  is prefix of  $y_1 y_2 y_3 y_4$

It can be verified that if the common descendant has two coded nodes as ancestors, then one of the coded nodes must be an ancestor of the other (see figure A)

$$\therefore \sum_{i=1}^n D^{l_{\max} - l_i} \leq D^{l_{\max}} \Rightarrow \sum_{i=1}^n D^{-l_i} \leq 1$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

Conversely, suppose the KI is satisfied, then one can construct a prefix code.

Assume,

$$l_1 \leq l_2 \leq \dots \leq l_n$$

Pick the lexicographically least string of length  $l_1$ . Assign this to  $x_1$ .

$$\begin{aligned} \therefore D^{-l_1} &< 1 \\ \therefore D^{l_{\max} - l_1} &< D^{l_{\max}} \end{aligned}$$

$\therefore \exists$  at least one uncoded node at depth  $l_{\max}$  where parent is not a coded node.

Trace the lexicographically least path to the node to depth  $l_2$ .

Assign the terminal node at depth  $l_2$  to codeword  $x_2$ . Then

$$D^{-l_1} + D^{-l_2} < 1 \quad \text{Continue in this fashion}$$

Until all codewords are determined.