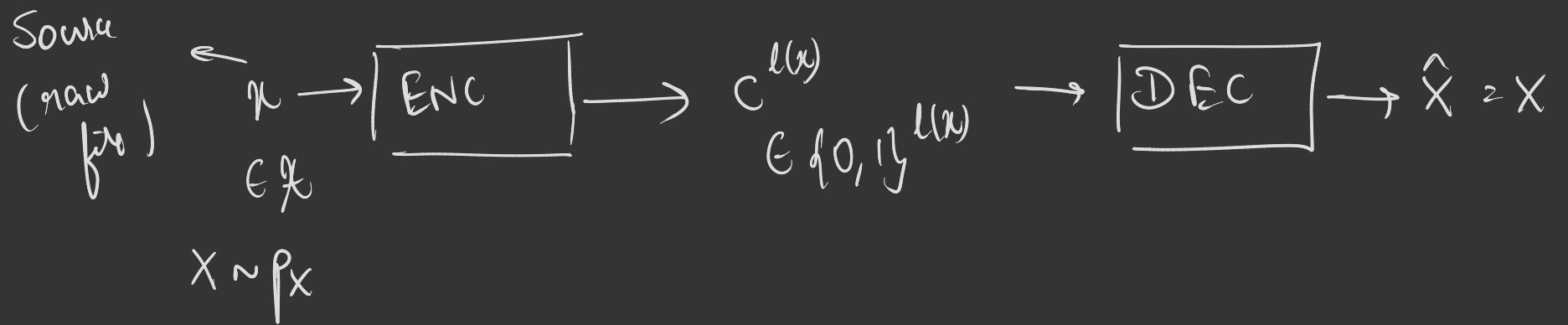


Variable Length Compression

x^n	$p(x^n)$	\mathbb{C}	$ENC(x^n)$	$L(x^n)$
000	$\frac{1}{2}$	0	<u>0</u>	1
001	$\frac{1}{4}$	10	<u>10</u>	2
010	$\frac{1}{8}$	110	110	3
011	$\frac{1}{16}$	1110	1110	4
100	$\frac{1}{32}$	11110	11110	5
101	$\frac{1}{64}$	111110	111110	6
110	$\frac{1}{128}$	1111110	1111110	7
111	$\frac{1}{128}$	1111111	1111111	7

$$L(\mathbb{C}) = \mathbb{E} L(X^n) \approx 1.984 \text{ bits}$$



Nonsingular Codes

Compression algorithm (code) is nonsingular if

$$x \neq x' \Rightarrow \text{ENC}(x) \neq \text{ENC}(x')$$

$$\forall x, x'$$

Distinct symbols are assigned distinct codewords

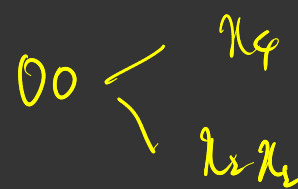
Optimal nonsingular codes

decreasing
order
of
probs

		code words		
x_1	$1/2$	1	0	—
x_2	$1/3$	01	1	0
x_3	$1/6$	001	00	—
x_4			01	00
			1	01
			1	1
			1	1

$$p(x_i) \geq p(x_{i+1})$$

x_{10}



potential
code words

- { 0, 1, 00, 10, 11, 01, 000, 001, ... }

$$\textcircled{1} \quad l(x_i) = \lfloor \log_2 i \rfloor \rightarrow \text{Claim (Check!)}$$

\downarrow
 length of
 codeword
 for symbol
 with i th largest
 prob.

$$\textcircled{2} \quad p(x_i) \leq \frac{1}{i} \rightarrow \text{Claim} \quad \underline{p(x_1) \leq 1}$$

\downarrow
 If not

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

$$p(x_2) \leq \frac{1}{2}$$

$$p(x_3) \leq \frac{1}{3}$$

$$= p(x_1) + p(x_2) + \dots + p(x_i)$$

$$> \frac{1}{1} + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{i} > 1$$

$$\Rightarrow i \leq \frac{1}{p(x_i)}$$

$$E l(x) = \sum_{i=1}^m p(x_i) l(x_i)$$

$$= \sum_{i=1}^m p(x_i) \lfloor \log_2 i \rfloor$$

$$= \sum_{i=1}^m p(x_i) \lfloor \log_2 \frac{1}{p(x_i)} \rfloor$$

$$= \sum_{i=1}^m p(x_i) \log_2 \frac{1}{p(x_i)} = H(x)$$

Entropy

Extension of a code

$$\{x_1, x_2, \dots, x_m\} = \mathcal{X}$$

$$\{a, b, c, \dots, z\} = \mathcal{Y}$$

apple

Single letter \rightarrow Multi-letter code

$$\text{ENC}(a) = 0, \quad \text{ENC}(b) = 10 \dots$$

$$\text{ENC}(ada) = \text{ENC}(a) \text{ENC}(d) \text{ENC}(a) \rightarrow \text{Extension}$$

\approx
0 111 0

$$a \text{ --- } 0$$

$$b \text{ --- } 10$$

$$c \text{ --- } 110$$

$$d \text{ --- } 111$$

Extension :

$$\text{ENC}(x(1) \ x(2) \ \dots \ x(n)) \quad \geq \quad \text{ENC}(x(1)) \ \text{ENC}(x(2)) \ \dots \ \text{ENC}(x(n))$$

$$\text{ENC}(a \ d \ b \ a \ c \ c \ a)$$

"

$$\text{ENC}(a) \ \text{ENC}(d) \ \text{ENC}(b)$$

Optimal nonsingular code

a	0
b	1
c	00
d	<u>01</u>

Extension

aa	→	00
bb	→	11
ab	→	<u>01</u>
		⋮

Extension

$$ENC: X \rightarrow \{0,1\}^*$$

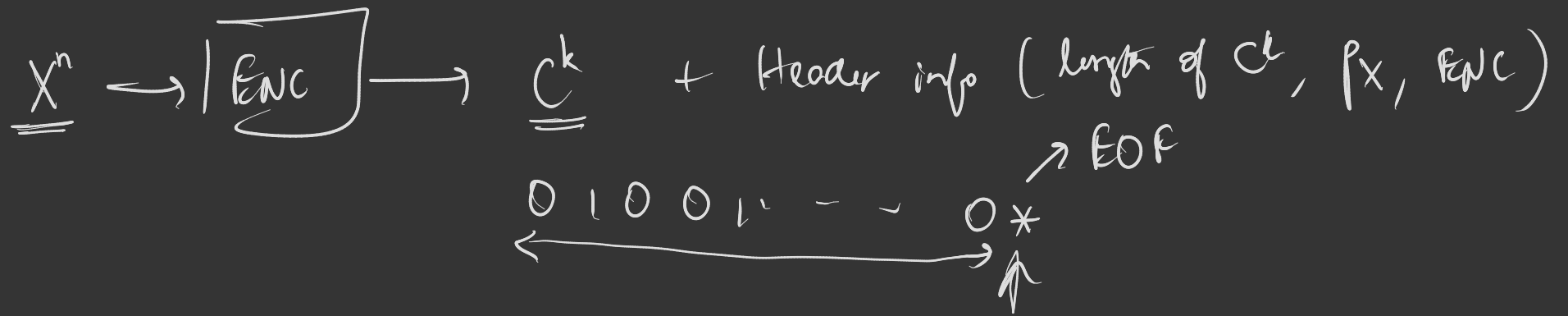
(for any set A , A^* denotes all possible
finite-length sequences from A)

$$\{0,1\}^* = \{ \epsilon, 0, 1, 00, 01, 10, 11, 000, 001, \dots \}$$

Extension of ENC:

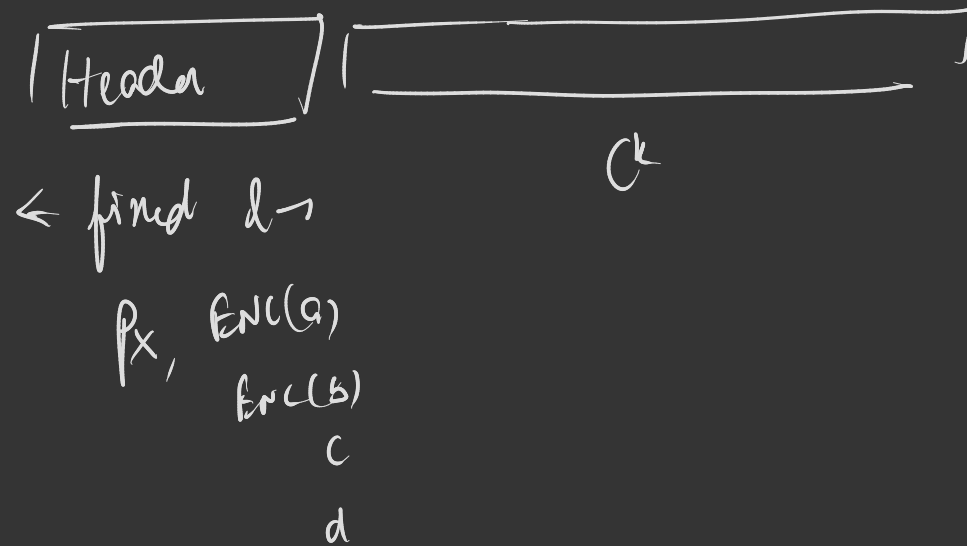
$$\{ ENC(\underline{x}) : \underline{x} \in X^* \}$$

$$\{ ENC(a), ENC(b) \dots ENC(a), ENC(aa), ENC(ab) \dots ENC(da), \\ ENC(aac) \dots \}$$



We have a mechanism for identifying end of compressed file.

abcd



Uniquely decodable code

VDC is a code for which the extension is nonsingular

$$\Rightarrow \forall x^n, \tilde{x}^m, \quad \text{ENC}(x^n) \neq \text{ENC}(\tilde{x}^m)$$
$$\forall n, \forall m, \forall x^n, \forall \tilde{x}^m$$
$$x^n \neq \tilde{x}^m$$

a	0
b	1
c	00
d	01

$$\text{ENC}(d) = \text{ENC}(ab)$$

\therefore Not uniquely dec.

a	0
b	10
c	100
d	101

$$\text{ENC}(c) = \text{ENC}(ba)$$

a 0

b 1 0

c 1 1 0

d 1 1 1



Uniquely decodable

Prefix-free or instantaneous codes (Prefix code)

A code is prefix-free if no codeword is a prefix of another

①

a	0	101
b	101	0
c	110	111
d	111	110

Prefix free

②

a	0
b	<u>10</u>
c	<u>10</u> 0
d	<u>10</u> 1

NOT P.F

③

a	<u>0</u>
b	1
c	<u>0</u> 0
d	<u>0</u> 1

NOT P.F

→
10 | 10

Is every prefix code uniquely decodable?

YES!

$x(1) x(2) \dots x(n)$

$\tilde{x}(1) \tilde{x}(2) \dots \tilde{x}(m)$

$ENC(x(1) x(2))$

$ENC(\tilde{x}(1) \tilde{x}(2) \tilde{x}(3))$

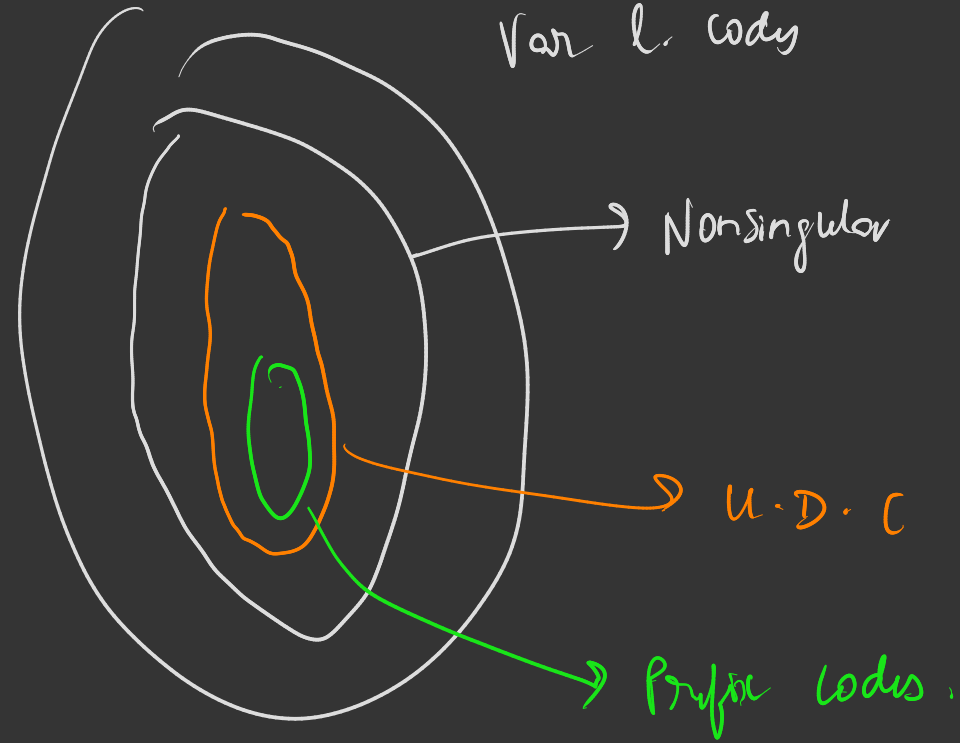
\Rightarrow Either $ENC(x(1)) = ENC(\tilde{x}(1))$

ONLY if $x(1) = \tilde{x}(1)$

OR $ENC(x(1))$ is a prefix of $ENC(\tilde{x}(1))$ X

$ENC(\tilde{x}(1))$ is a prefix of $ENC(x(1))$ X

Var L code



Summary

1. Nonsingular codes — No two symbols have the same codeword

2. Uniquely decodable codes — Every extension is nonsingular

3. Prefix codes — No codeword is a prefix of another

Example

	Singular C_1	Nonsingular Not U-D Not P-R C_2	Nonsingular Not P-R C_3	Non-singular U-D Prufin prn C_4
a	0	<u>0</u>	00	0
b	<u>00</u>	<u>010</u>	10	10
c	10	<u>01</u>	11	110
d	<u>00</u>	10	110	111

$ENC(ca)$

$= 010$

$ENC(b) = 010$

$ENC(b) = 0010$

$ENC(ac) = 0011$

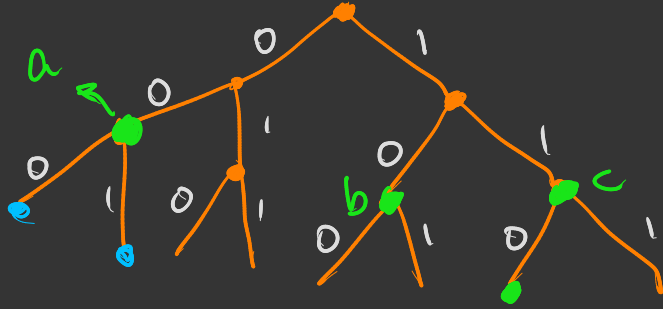
110 1110 00 110 --



a	00
b	10
c	11
d	110

① In this eg, the delay in decoding the first sym depends on # 0's (arbitrarily large)

② Decoding may be complicated.

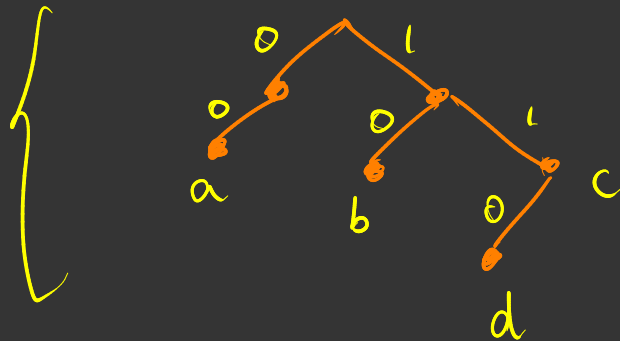


↑
 l_{max}
 ↓

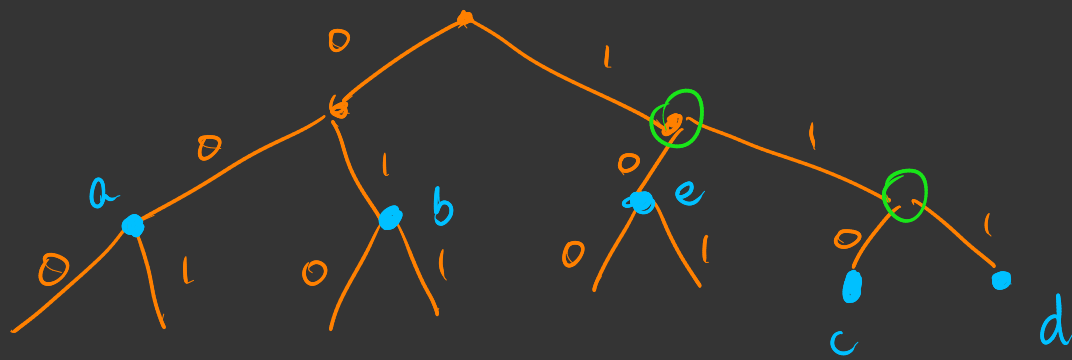
a	00
b	10
c	11
d	110
	<u> </u>
	l_{max}

Tree enumerates all binary strings of length l_{max} .

Code
 Tree



The codebook is a set of vertices on a binary tree of depth l_{max} .

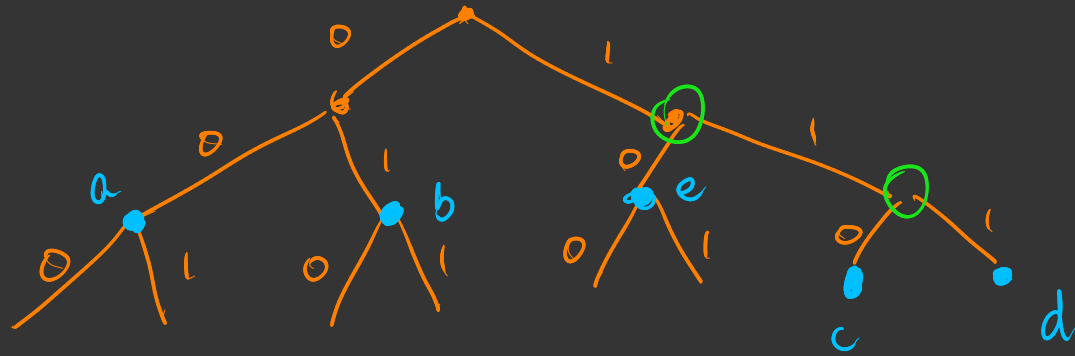


a	0	0	
b	0	1	
c	1	1	0
d	1	1	1
e	1	0	

Prefix-free

Observation: Prefix \Leftrightarrow ancestor on code tree / complete binary tree of depth l_{max}

Prefix-free code \Leftrightarrow No codeword is an ancestor of any other.



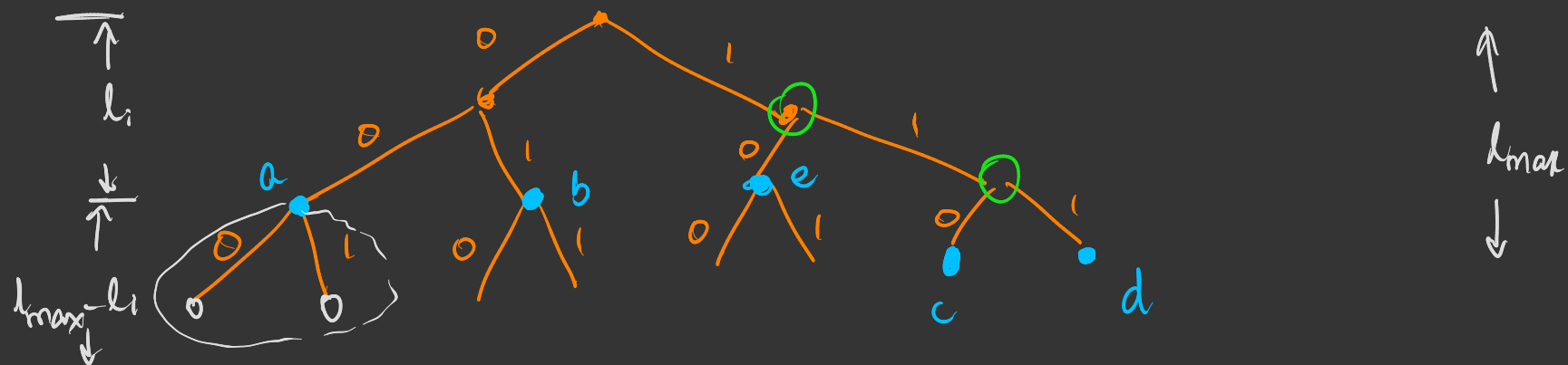
$00 \rightarrow$
 $\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow$
 a d e c

a d e c

Delay = 0

Instantaneous

Necessary and sufficient conditions on the codeword lengths of prefix codes



codeword lengths = $\{l_1, l_2, \dots, l_m\}$

$|X| = m$

$\{2, 2, 2, 3, 3\}$

Take any codeword. Count the # of leaves that are descendants.

\downarrow
 l_i

$2^{l_{\max} - l_i}$

$$C = \{0, 00, 10\}$$

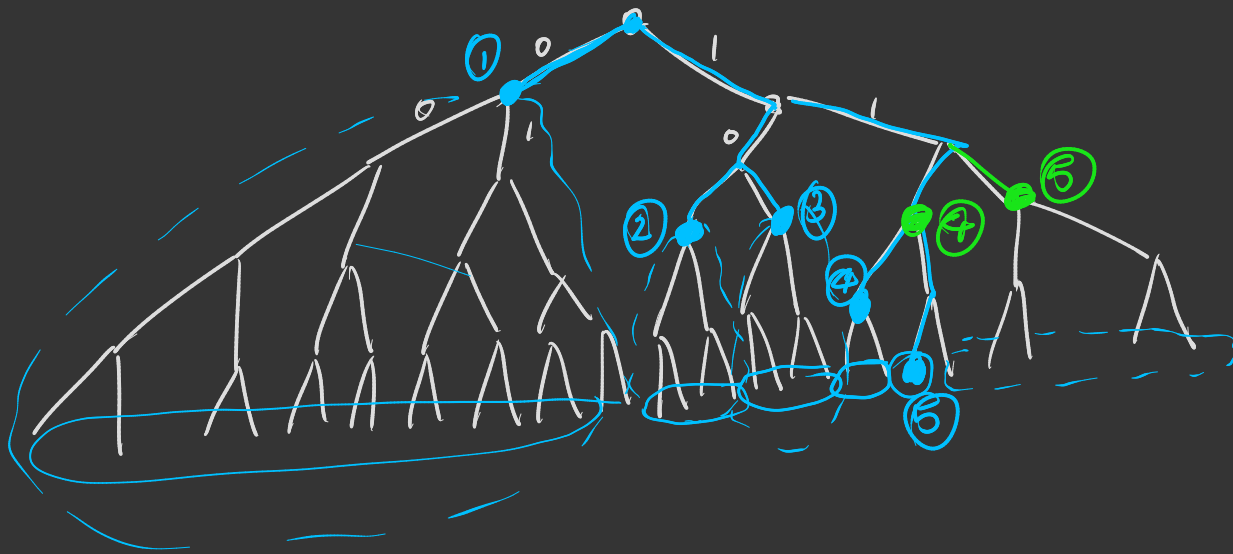
Knapsack's indep

$$\sum_{i=1}^m 2^{-l_i} = 2^{-1} + 2^{-2} + 2^{-2}$$

≈ 1

$$L = \{1, 2, 2\}$$

$$C = \{0, 10, 11\}$$



①	0
②	1 0 0
③	1 0 1
④	1 1 0 0
⑤	1 1 0 1 0

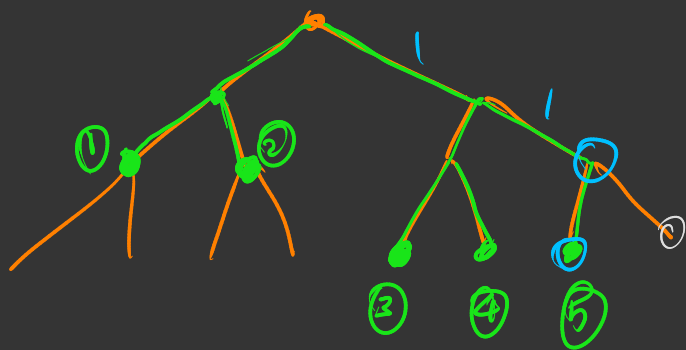
$$\{l_1, l_2, l_3, \dots, l_m\}$$

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

$$\{1, 3, 3, 4, 5\}$$

$$\sum_{i=1}^5 2^{-l_i} = 0.84 \dots < 1$$

$$\Rightarrow \sum_{i=1}^m 2^{l_{\max} - l_i} \leq 2^{l_{\max}}$$



$$\{2, 2, 3, 3, 3\}$$

$$l_{\max} = 3$$

①	00
②	01
③	100
④	101
⑤	110

$$\sum_{i=1}^m 2^{-l_i} = 0.875 < 1$$

~~~~~

Kraft sum

If Kraft sum  $< 1$ , then some of the leaves have no word ancestors.

## Kraft's inequality

For any prefix code, the codeword lengths  $\{l_1, \dots, l_m\}$  must satisfy:

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

If  $\{l_1, \dots, l_m\}$  is any set of integers satisfying

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

we can construct a prefix code with these lengths.

## Extended Kraft's inequality

Same result holds even for countably infinite  $m$ .

(A set is countable if there is a 1-1 map with  $\mathbb{Z}$ )

## Shannon Code

$\mathcal{X} \rightarrow$  alphabet  $|\mathcal{X}| = m$

Consider any pmf  $(p_1, p_2, \dots, p_m)$   $(0.2, 0.1, 0.4, 0.3)$

Goal: Construct a prefix-free code.

with small expected length

$$\sum_{i=1}^m p_i l_i$$

Know that  $\sum_{i=1}^m p_i = 1$

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

$l_1, l_2, \dots, l_m$



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

$$\left( \log_2 \frac{1}{p_i} \right) + 1 \geq l_i \geq \log_2 \frac{1}{p_i} \quad \Rightarrow \quad \begin{aligned} -l_i &\leq \log_2 p_i \\ 2^{-l_i} &\leq p_i \end{aligned}$$

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

$$\sum_{i=1}^m p_i l_i = \sum_{i=1}^m p_i \lceil \log_2 \frac{1}{p_i} \rceil \leq \sum_{i=1}^m p_i \left( \log_2 \frac{1}{p_i} + 1 \right)$$

$$\leq H(x) + 1 \text{ bit}$$