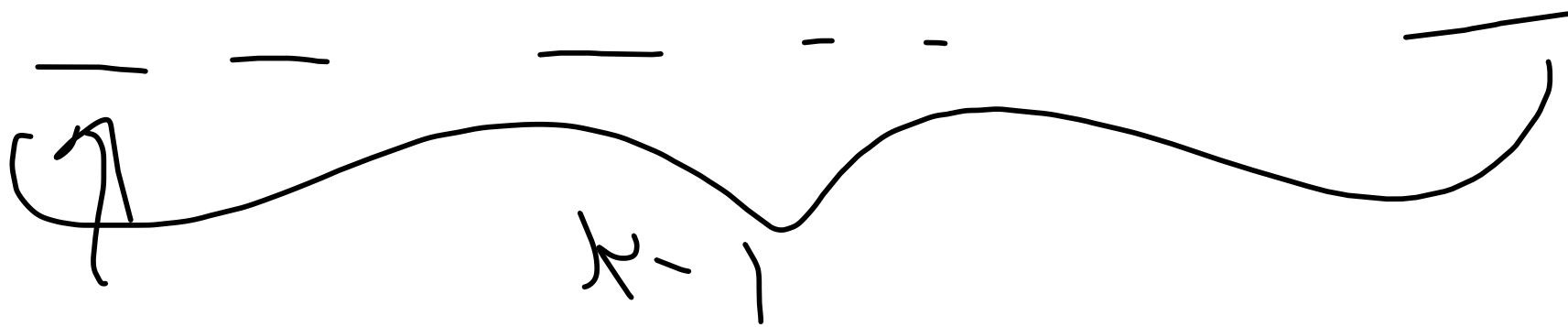


Clustering

Run Krouhkal for $n-k$ steps

Then The k -clustering obtained
by running Krouhkal's alg. for
 $n-k$ steps has maximum spacing

pf Let \mathcal{C} be the k -clusterings
consisting of C_1, C_2, \dots, C_k .
Running Kruskal for $n-k$ steps is
equivalent to running the full Kruskal
alg. & deleting the $k-1$ most expensive
edges



In the next step, Kruskal's alg.
will have picked up the k -1st
most expensive edge in the MST

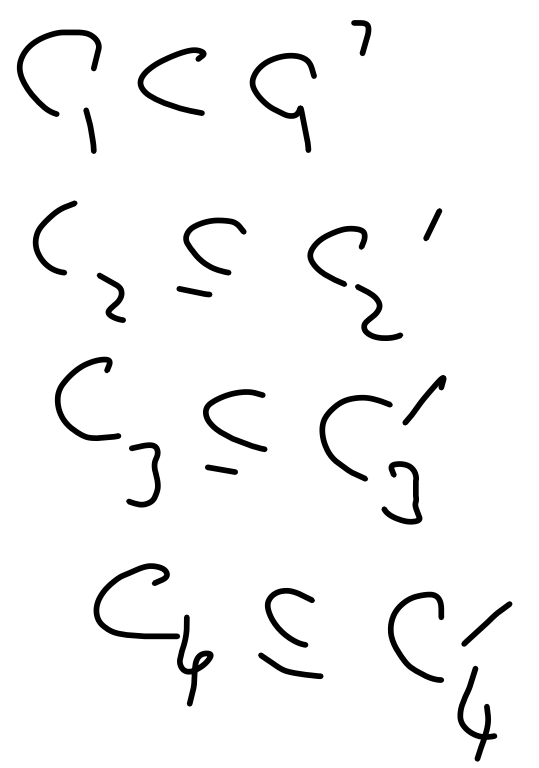
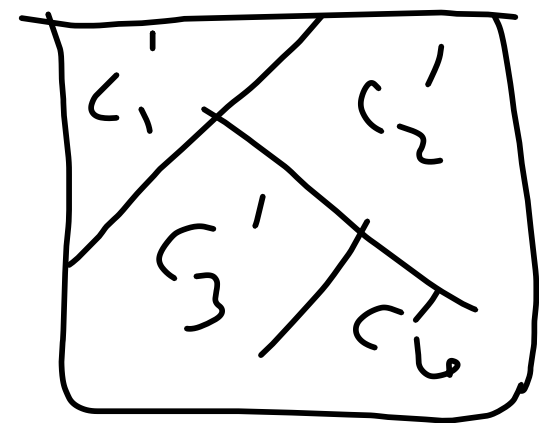
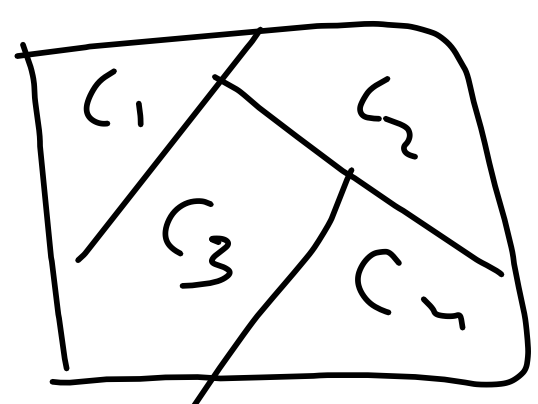
Clearly, the length of this edge
is the spacing of C . Denote it

by δ^* .

Let a be a k -clustering C consisting
of C_1, C_2, \dots, C_k

We shall show that the spacing of \mathcal{C} is at most δ^* .

If each C_i is contained in some C_j' , then $\mathcal{C} = \mathcal{C}'$.



If $C_1 \neq C_1'$, then $\exists x \in C_1 - C_1'$
 w.l.s. $\exists x \in C_1 - C_1'$

Since $\mathcal{C} \neq \mathcal{C}'$, $\exists C_\alpha$

s.t. C_α is not a subset of

any C_j' . This means \exists a pair

p_i, p_j in C_α

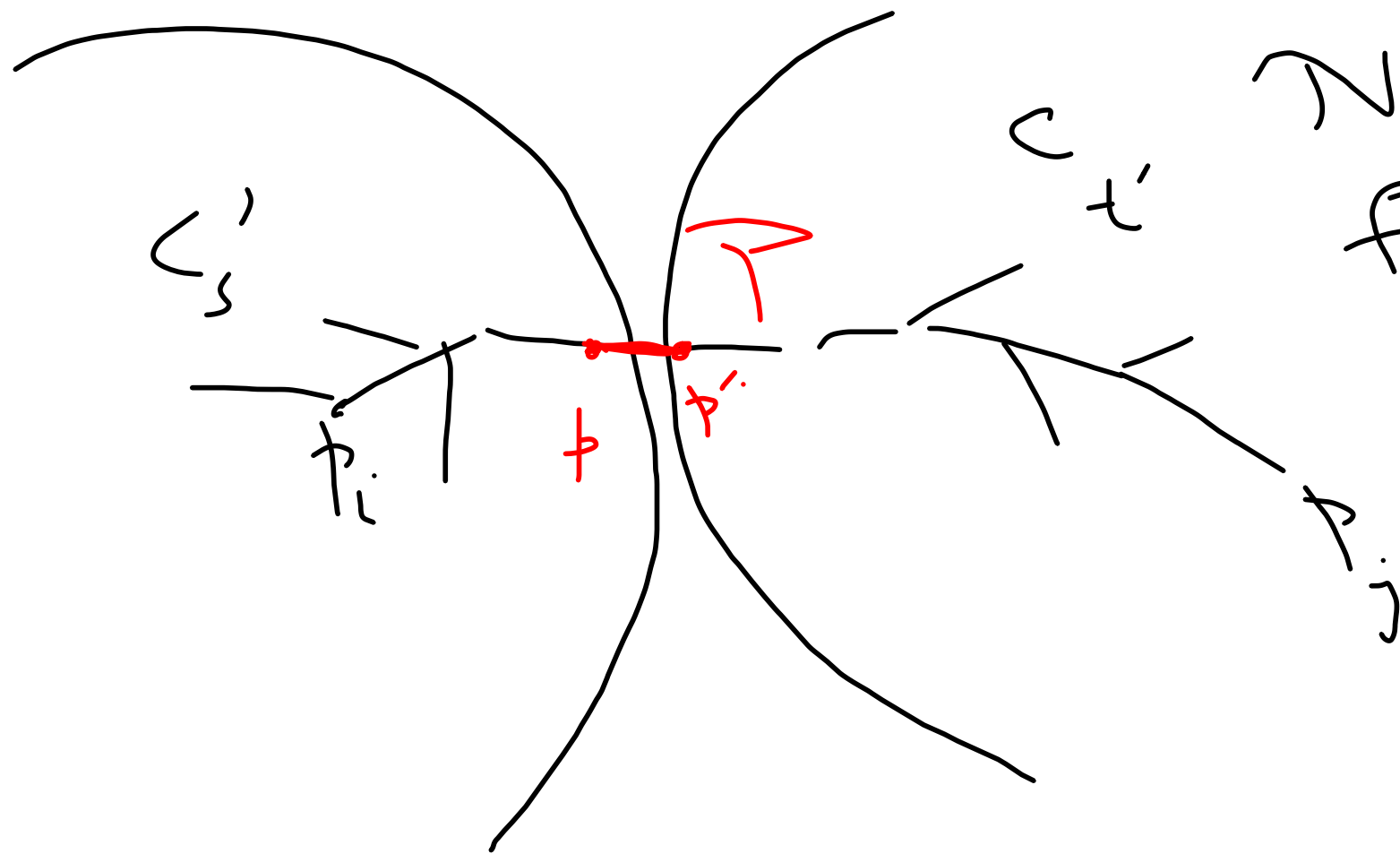
different

s.t. p_i, p_j belong to

clusters of \mathcal{C}'

Assume $p_i \in C_s'$ & $p_j \in C_t' \neq C_s'$

Consider a path P from p_i to p_j in C



Note that for each edge e in P

$$d(e) \leq \delta^*$$

Let p' be the first node on \bar{P}
that does not belong to C'_S

& let p be the previous node on \bar{P}

Clearly, $d(p, p') \leq \delta^*$.

Clearly, the spacing of $C'_S \leq d(p, p') \leq \delta^*$.

Huffman Code & Data Compression

Since computers work on bit strings
any text in a richer ~~language~~ alphabet, say
the English alphabet, requires an
encoding scheme into bit strings.
For a fixed length encoding, each English
letter, requires a 5-bit encoding.

Suppose we try to encode in
such a way that most frequently
occurring letters use shorter encodings

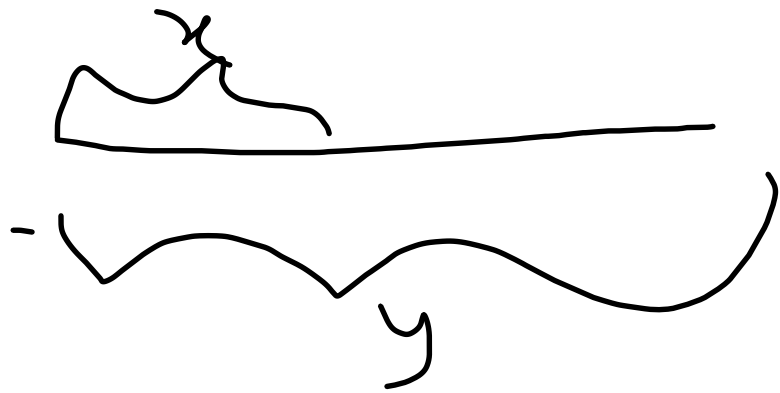
0 ← e → 25
1 ← a → 20

$$45 + 5 \times 55 \\ = \underline{320}$$

But using variable length encoding
leads to some problem while decoding

Suppose $\exists x, y$ s.t. encoding of

x is a prefix of the encoding of y .



Defⁿ A prefix code for an alphabet S is a function $\delta: S \rightarrow \{0, 1\}^*$ s.t.

for any pair $x, y \in S$, neither $\delta(x)$ nor $\delta(y)$ is a prefix of the other

Suppose δ is a prefix code of an alphabet S .
For each $x \in S$, let f_x be the fraction of the letters occurring in a text.

So if n be the no. of letters in the text,
Then $n f_x$ is no. of times the letters are equal to x

Hence the total length of the encoding

is $\sum_{x \in S} n f_x |\delta(x)|$. Hence the average

no. of bits required per letter is

$$ABL(\delta) = \sum_{x \in S} f_x |\delta(x)|$$

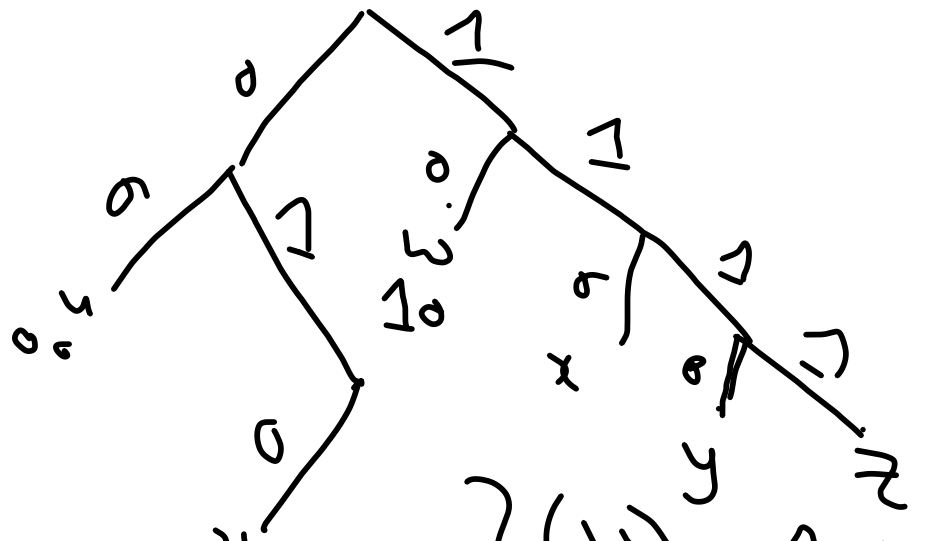
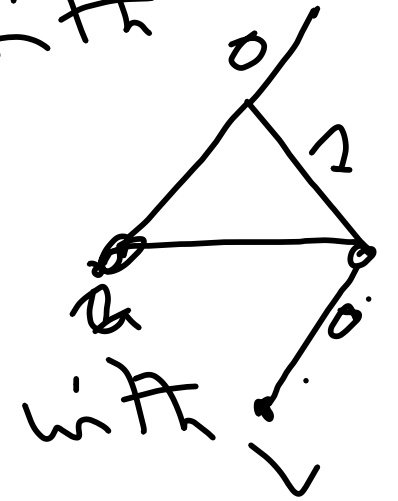
Our aim is to obtain a prefix code
of s.t. $ABL(\gamma)$ is as small as possible.

Such a code is called an Optimal
Prefix Code.

Prefix Code from binary trees.

Let T be a binary tree with
 no. of leaves equal to $|S|$

Label each leaf with
 a letter of S .



For each leaf labeled by
 x , consider the path from
 the root to the leaf

$$z(u) = 00$$

$$z(v) = 010$$

$$z(w) = 10$$

$$z(x) = 110$$

$$z(y) = 1110$$

$$z(r) = 1111$$

Each time you traverse left, write
down a 0 & each time you traverse
right, write down a 1.

The resulting string is an encoding of x .

Clearly, this is a prefix code

For each prefix code one can
construct a binary tree as
follows.

All encoding starting with 0/1 are
The encodings of the letters labelling
The leaves of the left/right subtree.
Recursively construct the left &

The right subtree.

Clearly, the length of the encoding of x
is the length of the path from the
root to the node labeled by x .

This is called the depth of v
 d is denoted by $\text{depth}_T(v)$

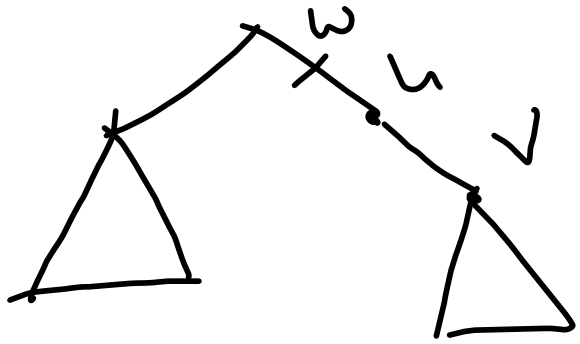
$$\text{ABL}(T) = \sum_{x \in S} f_x \text{depth}_T(x)$$

Defⁿ Call a binary tree full
is each internal node has
exactly two ~~the~~ children.

Lemma 1. The binary tree representing an optimal prefix code is full.

pf.

Suppose there is a node u with only one child v .



If u is a leaf node, remove u & make v the leaf node.

C. w. u has a parent w

Remove u & make v the child of w .

By this we obtain a tree T'
 s.t. $ABL(T') < ABL(T)$ Contradiction.

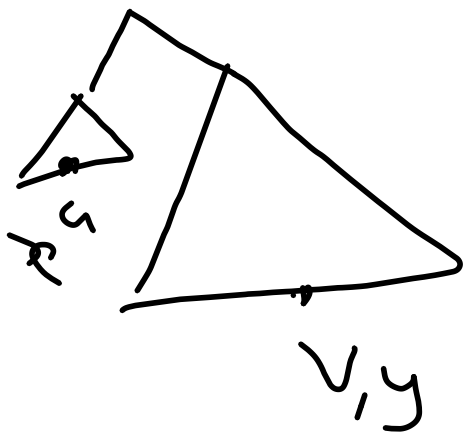
Lemma 2 For each $x \in S$, let f_x denote the fraction that x occurs. Let T be a binary tree representing an optimal prefix code.

Let u, v be two leaves of T

s.t. $depth_T(u) < depth_T(v)$

Let the labels of u, v be x, y resp.

Then $f_x \geq f_y$.



If Let T' be the tree obtained from T by interchanging the labels of x and y

$\& \checkmark$

$$ABL(T') - ABL(T)$$

$$T \implies f_x \text{depth}_T(x) + f_y \text{depth}_T(y)$$

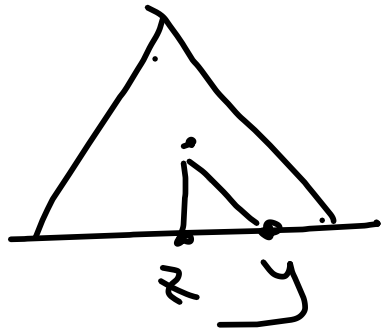
$$T' \implies f_x \text{depth}_{T'}(y) + f_y \text{depth}_{T'}(x)$$

$$= f_x \text{depth}_T(y) - f_x \text{depth}_T(x) + f_y \text{depth}_T(x) - f_y \text{depth}_T(y)$$

$$= (f_x - f_y) (\text{depth}_T(x) - \text{depth}_T(y)) \geq 0$$

Lemma 3 T^* be a binary tree dep^3
 an optimal prefix code. Let y be
 the label of a leaf with min. frequency.
 Then y has a sibling z in T^* .

pf



Since the tree is full, y has a parent w with another child z .
 If z has a child, then its depth would be greater than that of y .
 Contradiction.