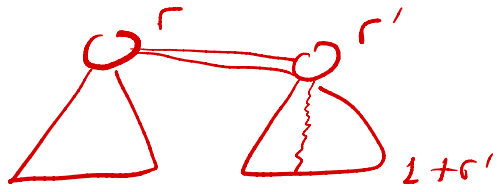
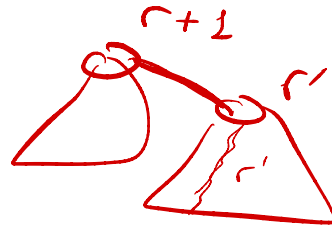


Data Compression



$$r > r'$$

$$r \geq 1+r'$$



$$r = r'$$

if - root of a tree has rank r
 \rightarrow number of vertices in tree $\geq 2^r$

Minimum Connected Subgraph

Given $G = (V, E)$

undirected

connected

$w(u, v)$ for each $(u, v) \in E$

Want to find $S \subseteq E$

such that

- graph (V, S) is connected
- $\sum_{(a, b) \in S} w(a, b)$ minimized

Given

sequence of elements from a set V

$$V = \{a, b, \dots, z, \sqcup\}$$

given "we hold these truths to be self."

use $\lceil \log_2 |V| \rceil$ bits for each symbol from V

file of n symbols $n \cdot \lceil \log_2 |V| \rceil$

want to use k bits for each symbol of V

2^k k -bit strings

$$2^k \geq |V|$$

$$k \geq \lceil \log_2 |V| \rceil$$

$$V = \{a, b, c\}$$

$k=1$

0

1

Variable length encodings

a 0
b 000
c 010
d 011
e 1
P
:

c 010
aea 010

a → 01
b → 10
c → 001
d → 11
e → 000

10011000001

prefix-free

An injective mapping

$$V \rightarrow \{0, 1\}^k$$

is always prefix free

┌ two different binary
└ strings of same length

impossible one is prefix of
other

$\begin{bmatrix} 0110 \\ 011011 \end{bmatrix}$ not prefix-free

V fixed-length encoding
 $\lceil \log_2 V \rceil$ per symbol

V sequence of length n
for each v
 $f(v) = \#$ times that v
appears in the sequence

suppose I find a variable-length
encoding of V where each $v \in V$
is mapped to a bit string of
length $l(v)$

Then sequence has encoding

that uses $\sum_v f(v) \cdot l(v)$ bits

Given V , frequencies $f(v)$ for each $v \in V$
find a prefix-free encoding of V
such that

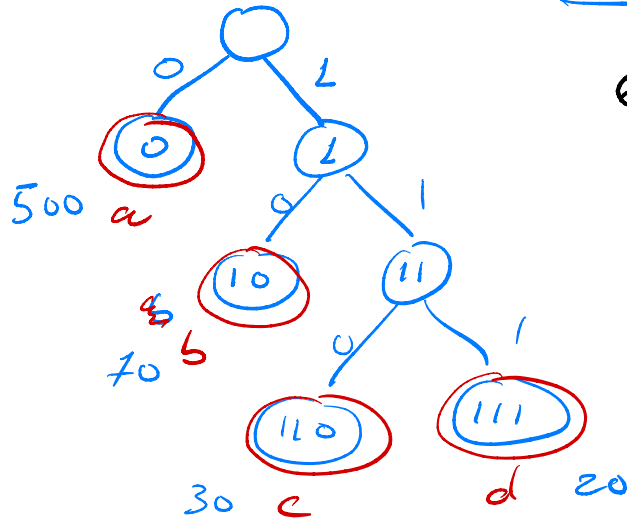
$$\sum_v f(v) l(v)$$

is minimized, where $l(v)$ is
the number of bits in the encoding
of v

V	f(V)	encoding	
a	500	0	00
b	70	10	01
c	30	110	10
d	20	111	11

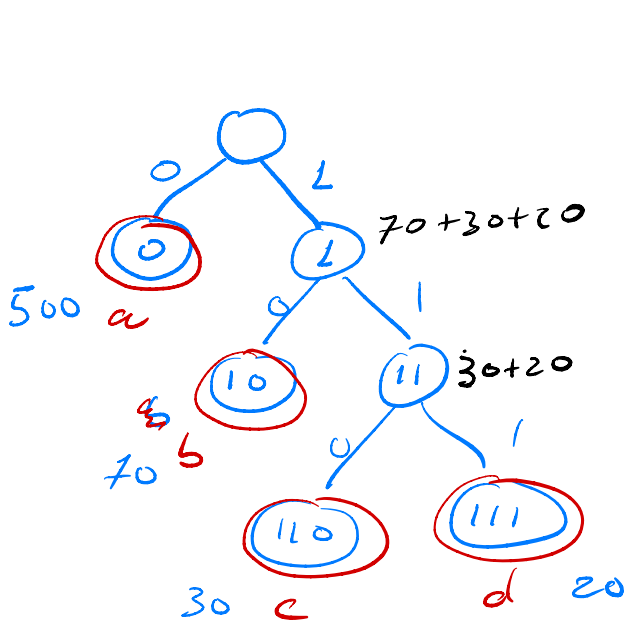
sits in encoding

$$500 \cdot 1 + 70 \cdot 2 + 30 \cdot 3 + 20 \cdot 3 = 790$$



$$620 \cdot 2 = 1240$$

1001101110101100
 b a c d a b c d



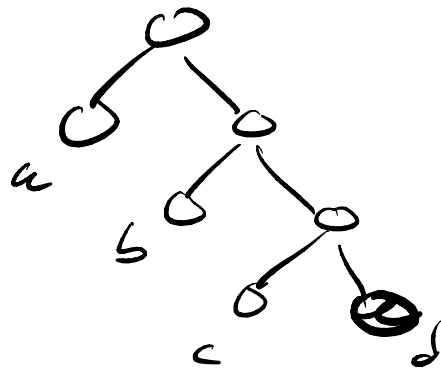
$$\begin{aligned}
 &500 \\
 &+ 70 + 70 \\
 &+ 30 + 30 + 30 \\
 &+ 20 + 20 + 20
 \end{aligned}$$

$$\begin{aligned}
 &500 + 120 + 70 + 50 + 30 + 20 \\
 &= 790
 \end{aligned}$$

given $V = \{a, b, c, d\}$

f $f(a) = 500$ $f(b) = \dots$

want to construct a binary tree and associate elements of V to the leaves

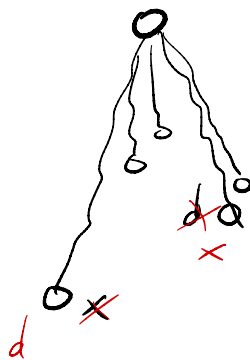


associate to each node sum of $f(v)$ for all descendants of node

minimize sum of above over all non-root nodes

a	500
b	70
c	30
d	20

An optimal tree exists in which
 d is a deepest leaf (has an encoding
 \geq all other encodings)



$$f(x) \geq f(d)$$

$$l(x) > l(d)$$

switch x, d

before $(\sum_{v \neq x, d} f(v) l(v)) + f(x)l(x) + f(d)l(d)$

after $(\sum_{v \neq x, d} f(v) l(v)) + f(x)l(d) + f(d)l(x)$

before - after

$$= f(x)l(x) + f(d)l(d)$$

$$- f(x)l(d) - f(d)l(x)$$

$$= f(x)(l(x) - l(d))$$

$$- f(d)(l(x) - l(d))$$

$$= \underbrace{(f(x) - f(d))}_{>0} \underbrace{(l(x) - l(d))}_{>0}$$

≈ 0

before \approx after

V set of symbols

f frequencies

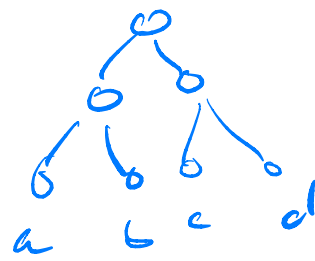
z is an element of V with
smallest $f(\cdot)$

minimizes $\sum_{v \in V} f(v)l(v)$

Then there is a minimal prefix-free
encoding of V such that z has a
longest encoding, that is

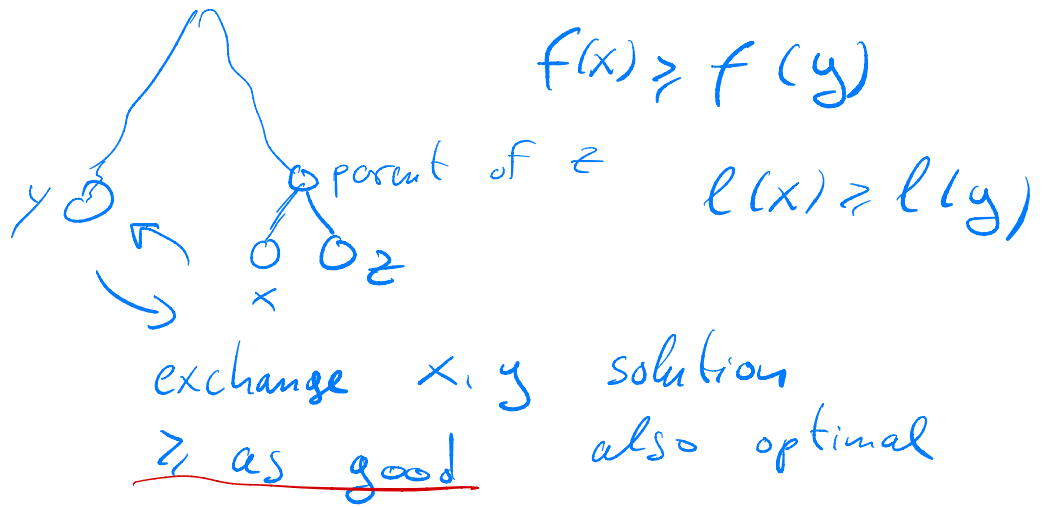
$$\forall v \in V \quad l(z) \geq l(v)$$

a 100
b 100
c 100
d 100

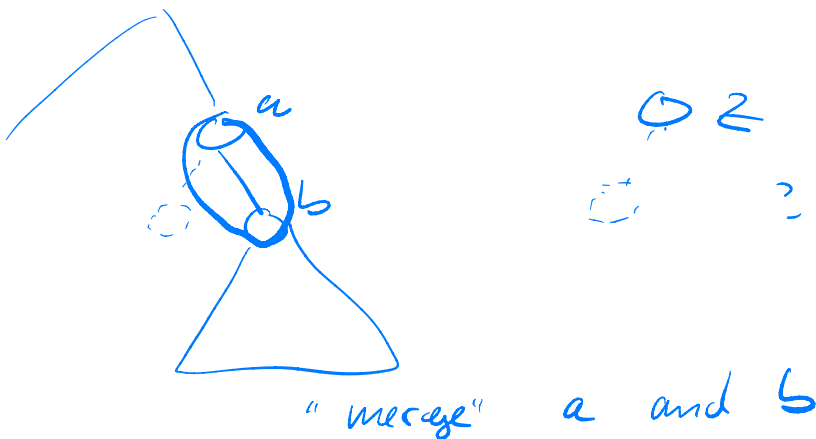


Suppose y, z two least frequent elements of V [$\forall x \in V - \{y, z\} f(x) \geq f(y)$ and $f(x) \geq f(z)$]
 There is an optimal solution in which y, z are siblings and longest

From before, there is optimal solution in which z deepest



In optimal solution every non-leaf node has two children



a 500
 b 70
 c 30
 d 20

a 500
 b 70
 x 50

a 500
 y 120

