# Binary Codes

$\Sigma = \{a, b, c, d\}$

def: Given an alphabet $\Sigma$, a binary code is a function $f : \Sigma \to \{0,1\}^*$

ex: Morse code, ASCII, Braille

Suppose you have a message where the letter "a" occurs 50% of the time, "b" 30%, and "c" 20%. Which is the best binary encoding of $\Sigma = \{a, b, c\}$?

A): $f(a) = 00$
$f(b) = 01$
$f(c) = 10$

B) $f(a) = 0$
$f(b) = 1$
$f(c) = 01$

C) $f(a) = 0$
$f(b) = 10$
$f(c) = 11$

**A)** $f(a) = 00$
$f(b) = 01$
$f(c) = 10$

Clear for decoding, but doesn't take advantage of different prob.

**B)** $f(a) = 0$
$f(b) = 1$
$f(c) = 01$

Ambiguous for decoding

$01 \rightarrow c$
$\rightarrow ab$

**C)** $f(a) = 0$
$f(b) = 10$
$f(c) = 11$

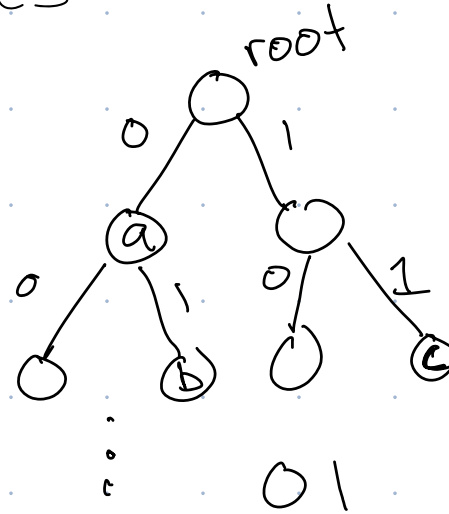No ambiguity for decoding, average bits per is small

# of bits in $f(i)$

Average letter length: $L(f) = \sum_{i \in \Sigma} |f(i)| \cdot p(i)$

ex:

$$|f(a)| \cdot p(a) + |f(b)| \cdot p(b) + |f(c)| \cdot p(c) =$$

$$2 \cdot .5 + 2 \cdot .3 + 2 \cdot .2 = 2$$

$$1 \cdot .5 + 2 \cdot .3 + 2 \cdot .2 = 1.5$$

# Binary Trees & Binary Codes

$f(a) = 0$

$f(b) = 01$

$f(c) = 11$

$\longleftrightarrow$



root

0    1

(a)

0    1    0    1

(b)    (c)

$\vdots$

01

$|f(i)| = d(i) =$ depth of node $i$

a $\rightarrow$ [○]

b $\rightarrow$ [○ 1]

Ambiguity arises if multiple letters lie on same path.

<u>def</u>: A code is "prefix free" if all letters are at leaves in corresponding binary tree

# Merge Trees to Create Prefix Free Codes

**1**
- .2 (a)
- .25 (b)
- .1 (c)
- .15 (d)
- .3 (e)

**2**
- .2 (a)
- .25 (b)
- .25 (node with children c, d)
- .3 (e)

**3**
- tree: root with children a and node(c, d)
- (b)
- (e)

**4**
- tree: root with children a and node(c, d)
- tree: node with children b, e

**5**
- tree: root with two subtrees: left node(a, node(c, d)), right node(b, e)

# Optimal Binary Encoding Problem

Input:
- $\Sigma$ (alphabet of symbols)
- $p : \Sigma \to \mathbb{R}$ (probabilities / frequencies for each symbol)

Output:

$f : \Sigma \to \{0,1\}^*$ such that

- $f$ is prefix free    ← constraint
- minimize average letter length $L(f)$    ↑ Objective function

# Huffman's Algorithm   $O(n^2)$

For each $i \in \Sigma$:   $O(n)$
- Create a tree with node labelled $i$   $O(i)$
- Give tree weight $p(i)$
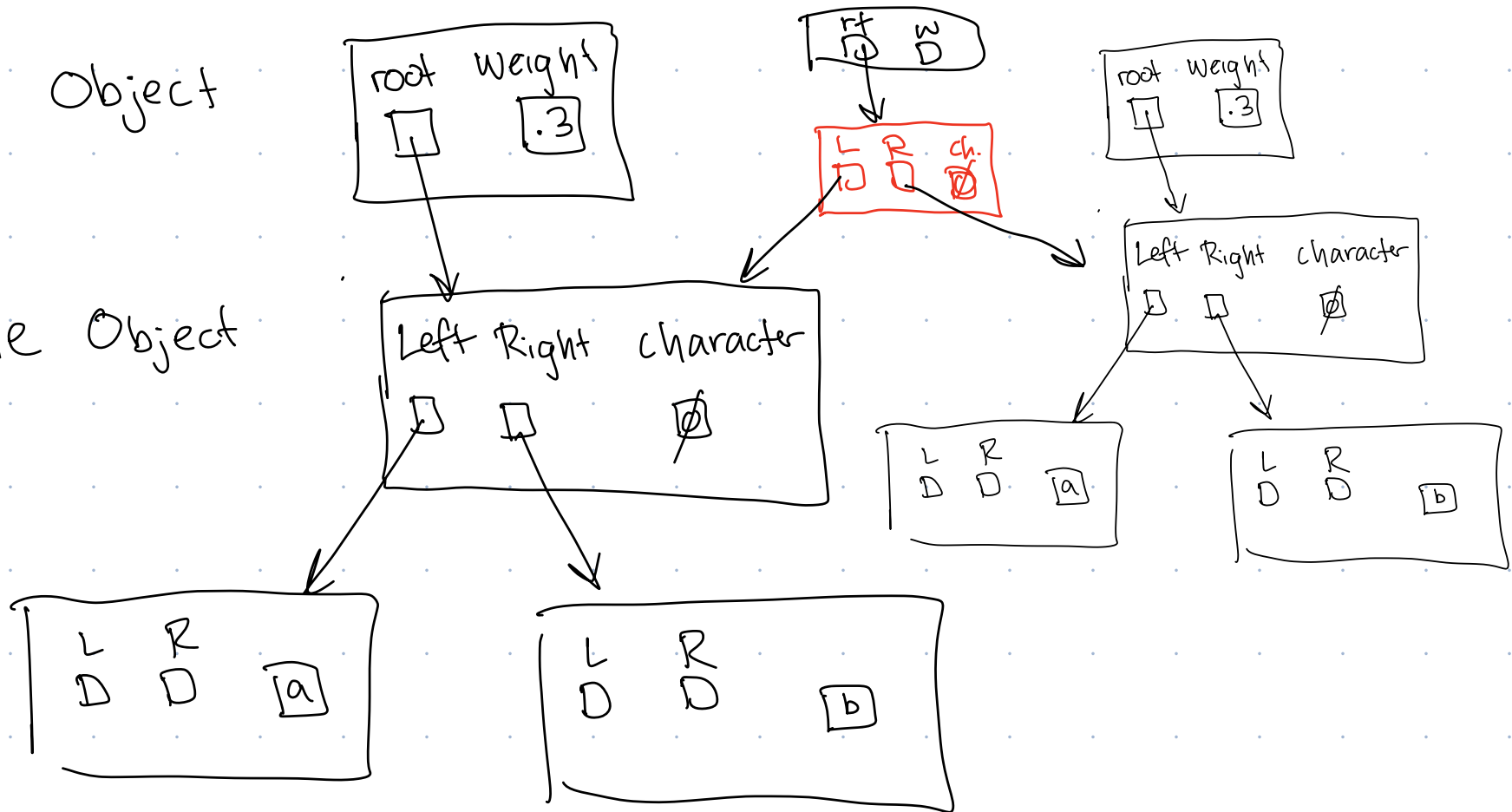
While there is more than one tree in our forest:   $O(n)$
- Merge 2 trees $O(i)$ with smallest weight   $O(n)$
- Set weight of merged tree to be sum of 2 weights   $O(1)$

---

| $i$ | $p(i)$ |
|-----|--------|
| a | .3 |
| b | .25 |
| c | .2 |
| d | .15 |
| e | .1 |

- Use Huffman's algorithm to create a binary code
- What is the average letter length of your code     2.25
- What is the runtime of Huffman's in terms of $|\Sigma| = n$? Ideas to improve?
- Why greedy?     • Midd Bucket list

# Tree Object

# Node Object

# Huffman's Algorithm

For each $i \in \Sigma$:

| • Create a tree with one node, label i $\}$ $O(n\log n)$
| • Give tree weight P(i)

While there is more than one tree: $O(n)$

$O(i) \rightarrow$ | • [Merge] [two trees with smallest weights] $\rightarrow$ 2 pops $O(\log n)$

| • Set weight of merged tree to be sum of weights

• Reinsert into heap $\leftarrow O(\log n)$

$O(n\log n)$

## Min Heap to store trees

- Initialize n items in heap $\rightarrow O(n\log n)$
- Pop min-value off heap $\rightarrow O(\log n)$
- Push new item into heap $\rightarrow O(\log n)$

# Why greedy: Order using a simple score function and take the object(s) with the best score(s)

# Huffman's Algorithm

For each $i \in \Sigma$:

- Create a tree with one node, label $i$
- Give tree weight $P(i)$

While there is more than one tree:                                    $n$

- Merge two trees with smallest weights
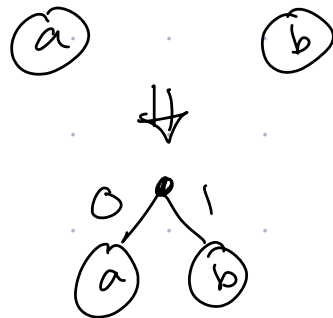- Set weight of merged tree to be sum of weights

---

Improve runtime?

**Thm**: Huffman's Algorithm produces a prefix free code that minimizes average letter length.

**Pf**: We will prove correctness by induction on $n = |\Sigma|$.

**Base case**: If $n = 2$, there are two letters, $a, b$.

Huffman

(a)    (b)

$\Downarrow$

0 / \ 1

(a)  (b)

This is optimal because there is no code with less than 1 bit per letter.

Inductive step: Assume for induction that Huffman's alg is optimal for any alphabet with $k$ characters. Consider an alphabet $\Sigma$ s.t. $|\Sigma| = k+1$.

 Let $a, b$ be letters with smallest weight in $\Sigma$.
Define $\Sigma^- = \Sigma - \{a, b\} \cup \{a/b\}$. Note $|\Sigma^-| = k$
Set $P(a/b) = P(a) + P(b)$              (a/b) is a single letter

ex:
$\Sigma = \{e, f, g, h\}$   $P(e) = .1$   $P(f) = .7$   $P(g) = .15$   $P(h) = .05$
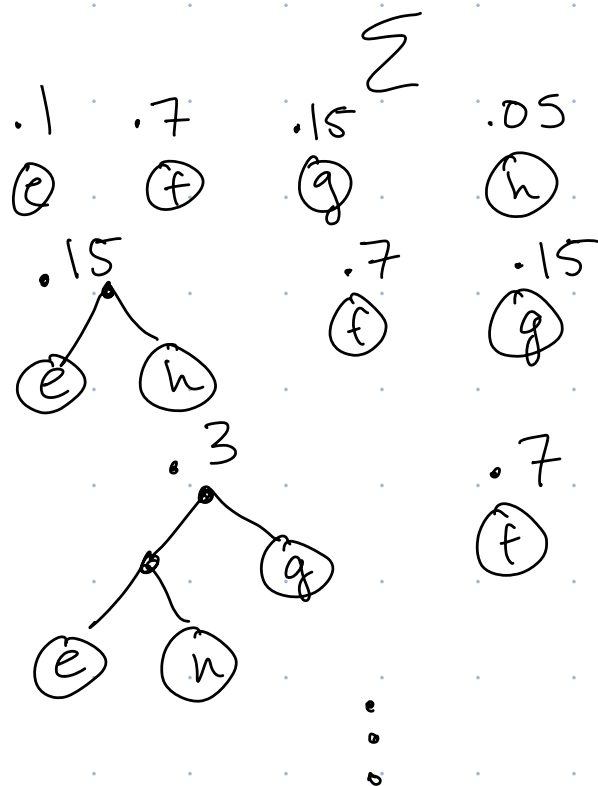then

ex:

$\Sigma = \{e, f, g, h\}$    $P(e) = .1$    $P(f) = .7$    $P(g) = .15$    $P(h) = .05$
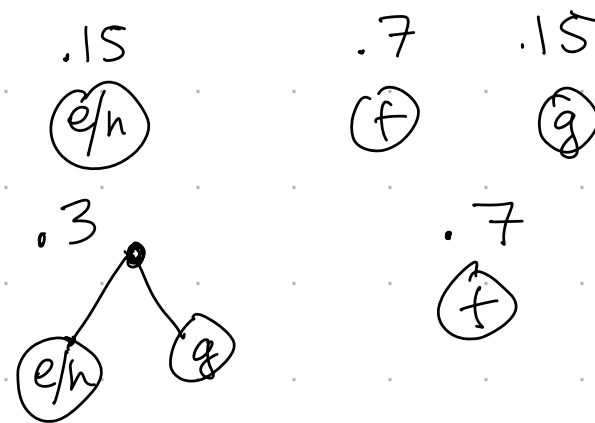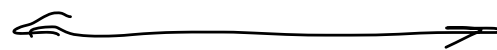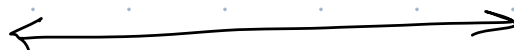
$\Sigma^- = \{e/h, g, f\}$    $P(e/h) = .15$    $P(f) = .7$    $P(g) = .15$
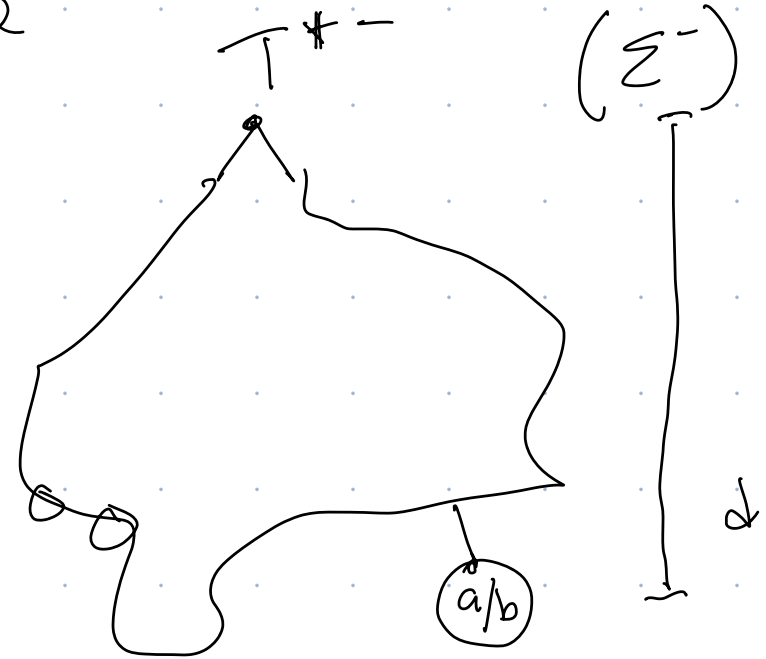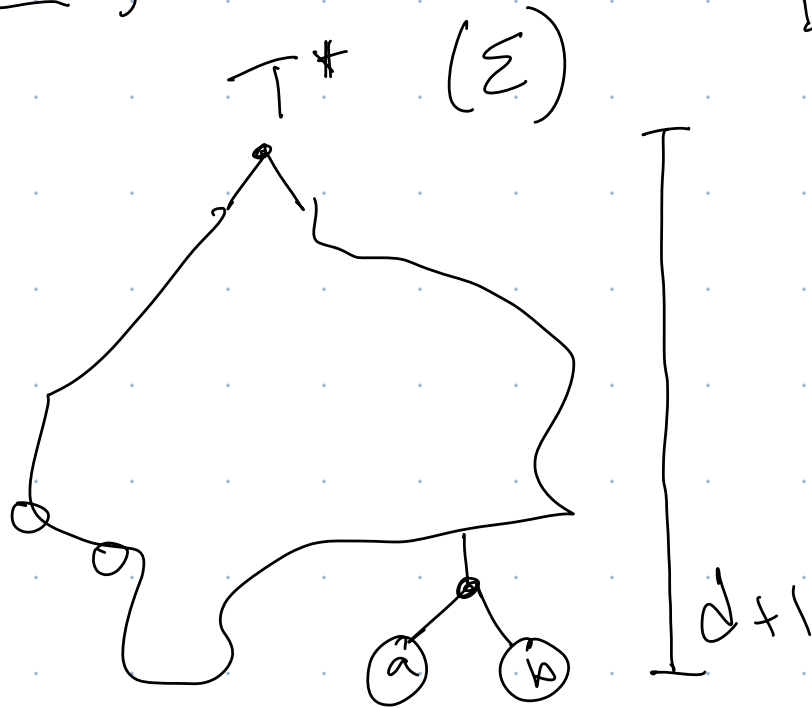


$\Sigma$

Huffman

$\Sigma^-$

In general:



By inductive assumption $T^-$ is optimal because it has K characters and was produced by Huffman's.

Lemma: There is an optimal tree for $\Sigma$ where $a, b$ are siblings.

Suppose for contradiction that $T$ is not optimal. Let $T^* \neq T$ be optimal with $a, b$ siblings (can do this by Lemma).



$T^*$ $(\Sigma)$

Define

$T^{*-}$ $(\Sigma^-)$

$d+1$

$d$

$a/b$

$P(a), P(b), d$

Then

$$L(T^*) = \sum_{i \neq a, b \in \Sigma} P(i) d(i) + \boxed{\phantom{XXXXXXXXXXX}}$$

$$L(T^{*-}) = \sum_{i \neq a/b} P(i) d(i) + \boxed{\phantom{XXXXXXXXXXX}}$$

So
$$L(T^*) - L(T^{* -}) = \boxed{\phantom{xxxxxxxxxxxxxxxxxxxxxxx}}$$

Similarly
$$L(T) - L(T^-) = \boxed{\phantom{xxxxxxxxxxxxxxxxxxxx}}$$

Thus
$$L(T^*) - L(T^{* -}) = L(T) - L(T^-)$$

Rearranging: $L(T) - L(T^*) = \boxed{\phantom{xxxxxxxxxxxxxxxx}}$

This is a contradiction because

$$\boxed{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}$$

<u>Lemma</u>: There is an optimal tree for $\Sigma$ with $a, b$ (characters with smallest p-values) siblings.