

Entropy & coding

Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

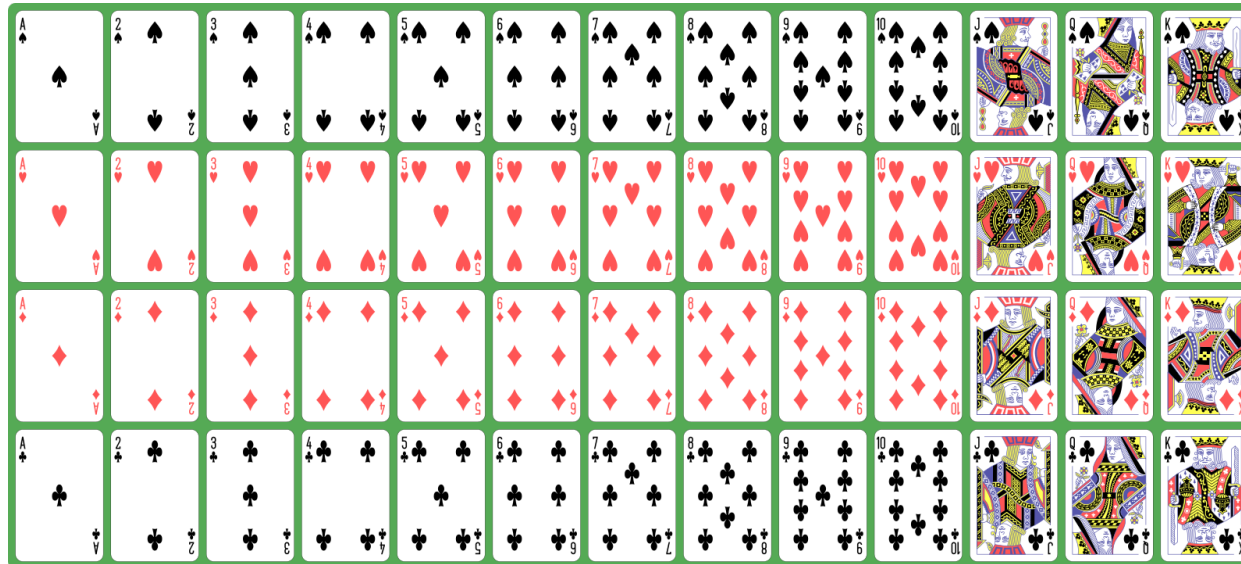
Department of Computer Science



Please sign guestbook (www.langmead-lab.org/teaching-materials) to tell me briefly how you are using the slides. For original Keynote files, email me (ben.langmead@gmail.com).

Entropy & coding

Let's identify items with **codes**, made of bits



Say, $\text{code} = \text{rank} + (13 * \text{suit})$

Where Ace = 0, Jack = 10, ...

♠ = 0, ♥ = 1, ...

A ♠	0
2 ♠	1
3 ♠	10
4 ♠	11

⋮

10 ♣	110000
J ♣	110001
Q ♣	110010
K ♣	110011

Entropy & coding

How many bits are required to encode items from universe U ?

$$H_{wc}(U) = \log_2 |U|$$

If codes must have **same** length, length must be $\geq \log_2(|U|)$, best choice is $\lceil \log_2(|U|) \rceil$

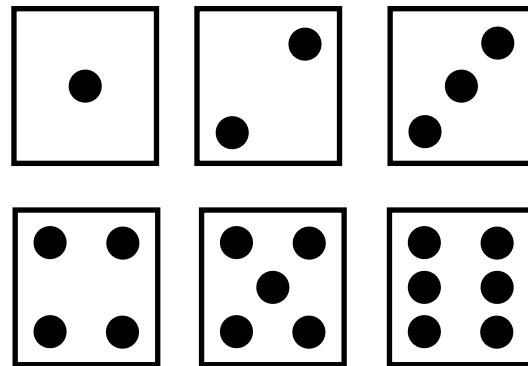
If codes can have **various** lengths, *longest* code must be $\geq \log_2(|U|)$

Entropy

How many bits required to identify an item from this set?



1 bit



3 bits



(37 or 38 slots)

6 bits

Entropy

$$H_{wc}(U) = \log_2 |U|$$

This is ***worst-case entropy***

If $|U| = 2^n$, then $H_{wc}(U) = n$

If $U = \{\text{length-}n \text{ strings from } \Sigma = \{1, \dots, \sigma\}\}$,
then $H_{wc}(U) = \log_2 \sigma^n = n \log_2 \sigma$

Entropy

If codes can vary in length,
we can use shorter codes
for more frequent events

Seeking to minimize
average (or *expected*)
code length $\bar{\ell}$

$$\bar{\ell} = \sum_{u \in U} \text{Pr}(u) \cdot \ell(u)$$

$\ell(u)$ = length of code for u

International Morse Code

A	• —	U	• • —
B	— • • •	V	• • • —
C	— • — •	W	• — —
D	— • •	X	— • • —
E	•	Y	— • — —
F	• • — •	Z	— — • •
G	— — •		
H	• • • •		
I	• •		
J	• — — —		
K	— • —		
L	• — • •		
M	— —		
N	— •		
O	— — —		
P	• — — •		
Q	— — • —		
R	• — •		
S	• • •		
T	—		
		1	• — — —
		2	• • — —
		3	• • • —
		4	• • • •
		5	• • • • •
		6	— • • • •
		7	— — • • •
		8	— — — • •
		9	— — — — •
		0	— — — — —

Entropy

Instead of items $u \in U$, let's think of a discrete r.v. X and its sample space Ω & probability function Pr

$$\begin{aligned} H(X) &= \sum_{s \in \Omega} \text{Pr}(s) \cdot \log_2 \frac{1}{\text{Pr}(s)} \\ &= - \sum_{s \in \Omega} \text{Pr}(s) \cdot \log_2 \text{Pr}(s) \end{aligned}$$

This is *Shannon entropy*

Entropy

$$X = \left\{ \text{Lincoln Penny} : 0.5, \text{ Reverse Penny} : 0.5 \right\}$$

$$\begin{aligned} H(X) &= 0.5 \cdot \log_2 \frac{1}{0.5} + 0.5 \cdot \log_2 \frac{1}{0.5} \\ &= 0.5 \cdot 1 + 0.5 \cdot 1 \\ &= 1 \end{aligned}$$

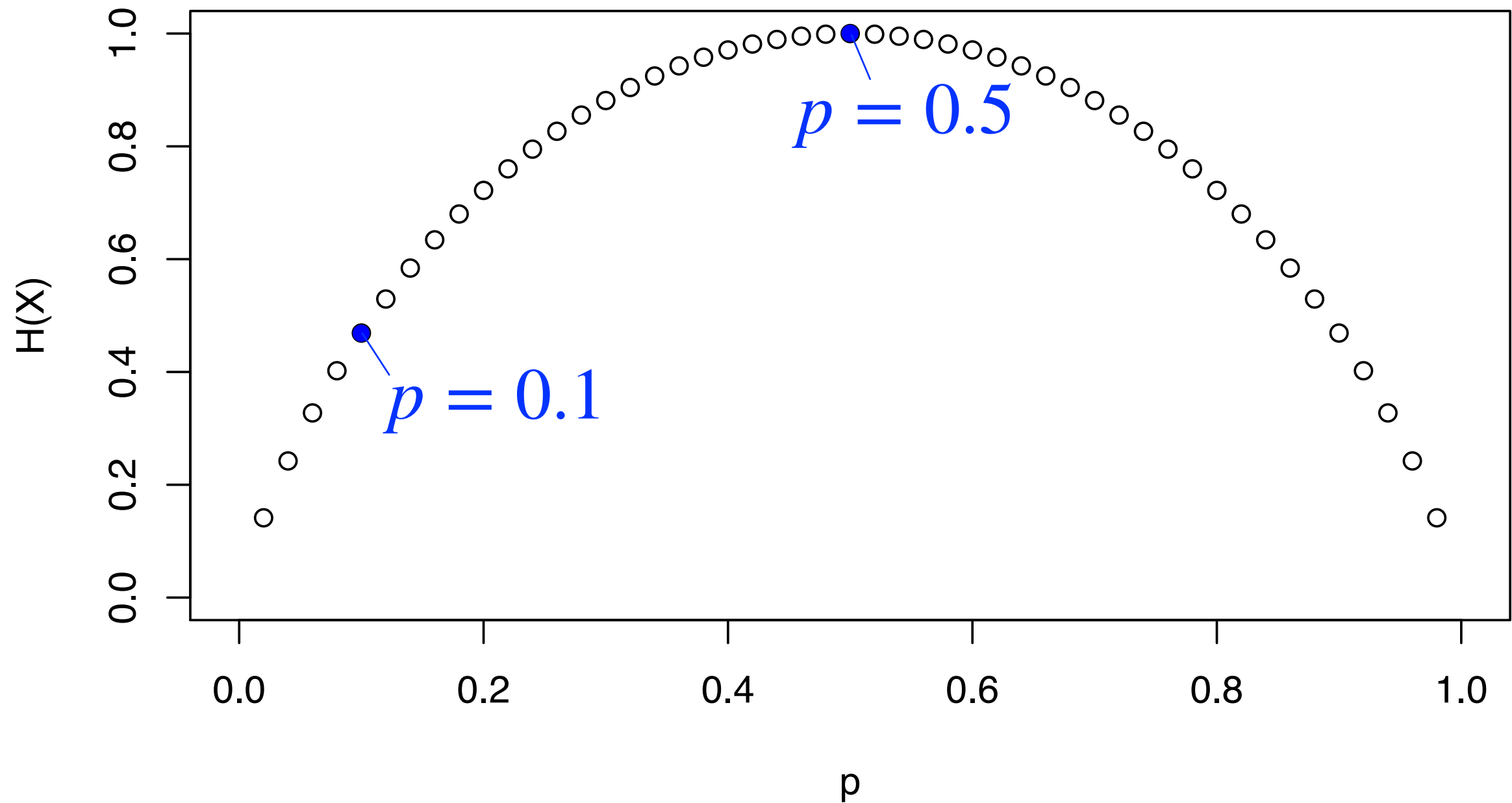
Entropy

$$X = \left\{ \text{Lincoln} : 0.9, \text{Lincoln Memorial} : 0.1 \right\}$$

$$\begin{aligned} H(X) &= 0.9 \cdot \log_2 \frac{1}{0.9} + 0.1 \cdot \log_2 \frac{1}{0.1} \\ &= 0.9 \cdot 0.15 + 0.1 \cdot 3.32 \\ &= 0.47 \end{aligned}$$

Entropy

$$X = \left\{ \text{🇺🇸} : p, \text{🇺🇸} : 1 - p \right\}$$

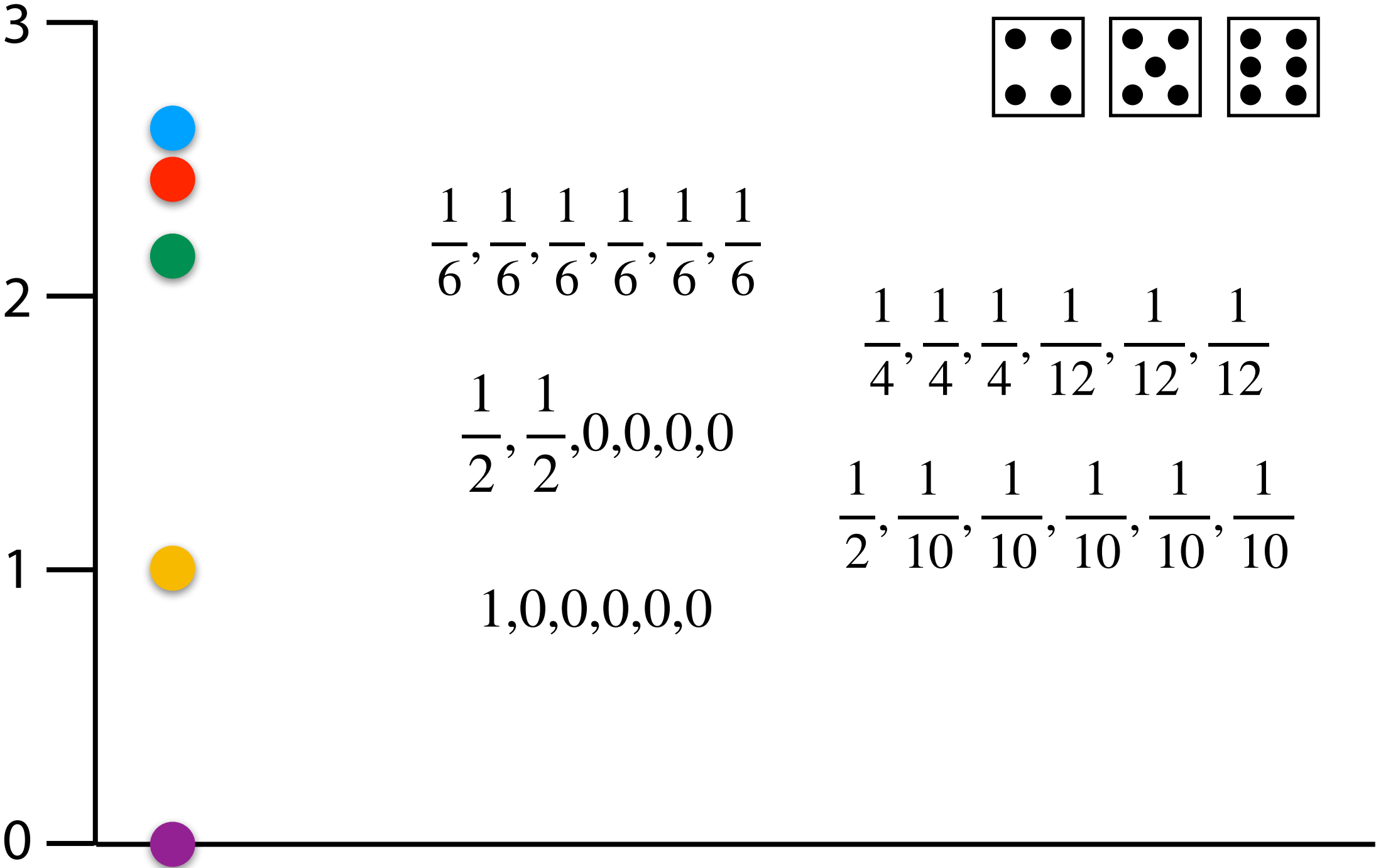


Entropy

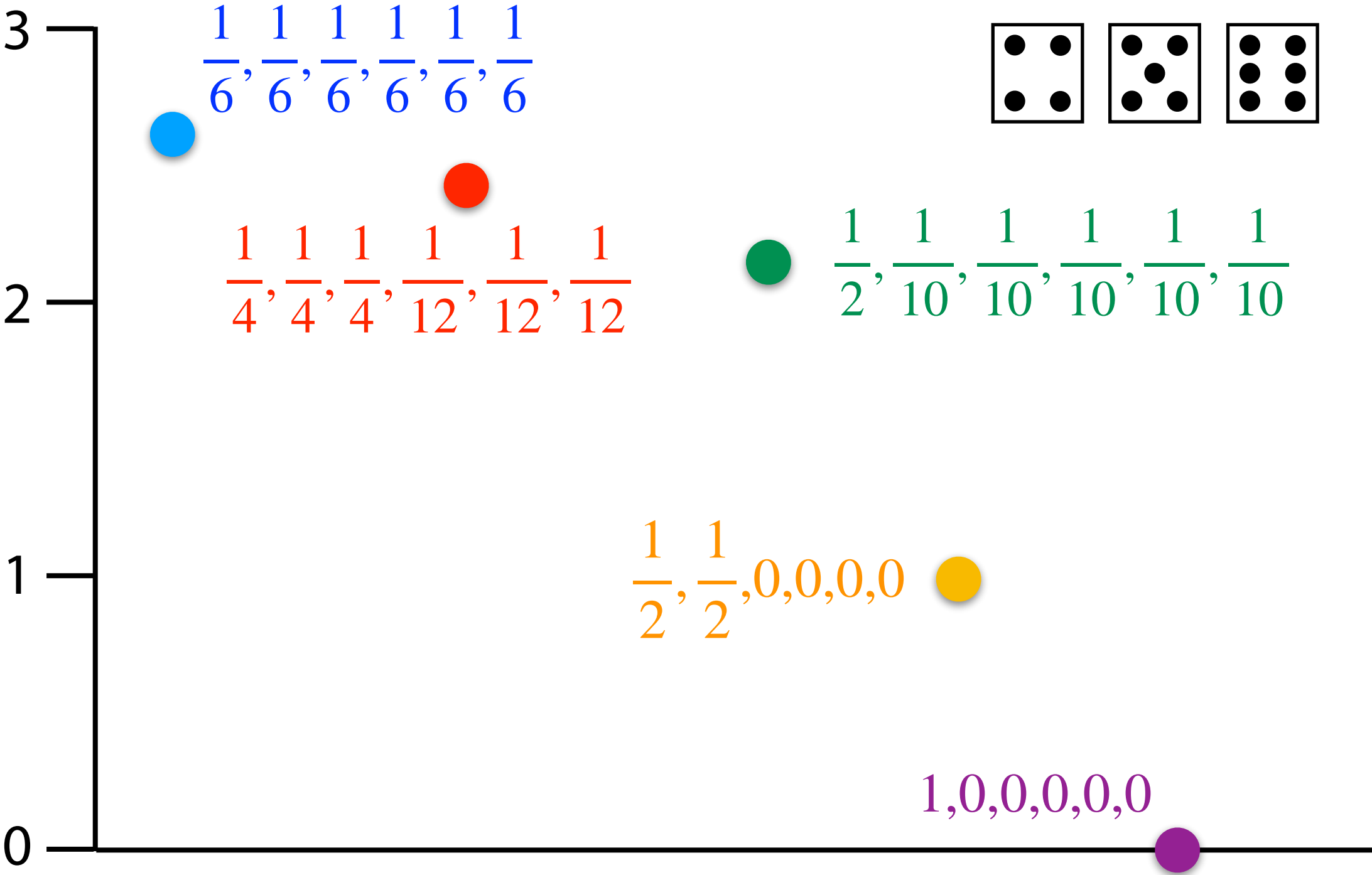
$$X = \left\{ \begin{array}{c} \square \cdot \quad \square \cdot \quad \square \cdot \quad \square \cdot \quad \square \cdot \quad \square \cdot \\ \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \end{array} : \frac{1}{6} \text{ each} \right\}$$

$$\begin{aligned} H(X) &= \sum_{i=1}^6 \frac{1}{6} \log_2 6 \\ &= \log_2 6 = 2.58 \end{aligned}$$

Entropy



Entropy



Entropy

When outcomes are equally probable:

$$\begin{aligned} H(X) &= \sum_{s \in \Omega_X} \Pr(s) \cdot \log_2 \frac{1}{\Pr(s)} \\ &= \sum_{s \in \Omega_X} \frac{1}{|\Omega_X|} \cdot \log_2 |\Omega_X| \\ &= \log_2 |\Omega_X| \end{aligned}$$

Matching the definition of worst-case entropy

Entropy

Shannon entropy $H(X)$ is a function of a random variable

The r.v. models a data ***source***; e.g. a person speaking, or letters of a DNA string

Assumes a ***memoryless*** source; each item is an i.i.d. draw

Entropy

So far we've seen

Worst-case entropy $H_{wc}(U)$ is a function of a **set**

Shannon entropy $H(X)$, a function of a **random variable**

When outcomes are equiprobable, $H(X) = H_{wc}(\Omega_X)$

Entropy

Say we have a memoryless binary source and an ***example string*** B it emitted

We can count B 's 0s & 1s to "train" a model

$$H_0(B) = H \left(X \sim \text{Bern} \left(\frac{m}{n} \right) \right) \quad \begin{array}{l} m = \# \text{ 1s in } B \\ n = |B| \end{array}$$
$$= \frac{m}{n} \log_2 \frac{n}{m} + \frac{n-m}{n} \log_2 \frac{n}{n-m}$$

H_0 is the ***empirical zero order entropy***

Entropy

So:

Worst-case entropy $H_{wc}(U)$ is a function of a **set**

Shannon entropy $H(X)$, a function of a **random variable**

Empirical zero order entropy $H_0(B)$ of a **sequence** B is the Shannon entropy of a memoryless source "trained" to B

Codes

A good code will:

Minimize average code length (approach H_0)

Give ***unambiguous*** mappings for encoding & decoding

Allow efficient encoding & decoding

International Morse Code

A	• —	U	• • —
B	— • • •	V	• • • —
C	— • — •	W	• — —
D	— • •	X	— • • —
E	•	Y	— • — —
F	• • — •	Z	— — • •
G	— — •		
H	• • • •		
I	• •		
J	• — — —		
K	— • —		
L	• — • •		
M	— —		
N	— •		
O	— — —		
P	• — — •		
Q	— — • —		
R	• — •		
S	• • •		
T	—		
		1	• — — —
		2	• • — —
		3	• • • —
		4	• • • •
		5	• • • • •
		6	— • • • •
		7	— — • • •
		8	— — — • •
		9	— — — — •
		0	— — — — —

(Note: letters are separated by a pause of duration equal to three dots; words separated by 7-dot pause.)

Codes

$$H(X) = \sum_{s \in \Omega} \text{Pr}(s) \cdot \log_2 \frac{1}{\text{Pr}(s)}$$

Shannon entropy equation hints at codes of

$$\text{length } \log_2 \frac{1}{\text{Pr}(s)}$$

Codes

Say we have a source emitting ***symbols*** from ***alphabet*** $\Sigma = \{a, c, g, t\}$

Source is ***memoryless***, modeled by r.v.:

$$X = \left\{ a : \frac{1}{2}, c : \frac{1}{4}, g : \frac{1}{8}, t : \frac{1}{8} \right\}$$

C is a function mapping symbols to binary code sequences. $C : \Sigma \rightarrow \{0,1\}^*$

What kind of C do we want?

Codes

$$X = \{ a : \frac{1}{2}, c : \frac{1}{4}, g : \frac{1}{8}, t : \frac{1}{8} \}$$

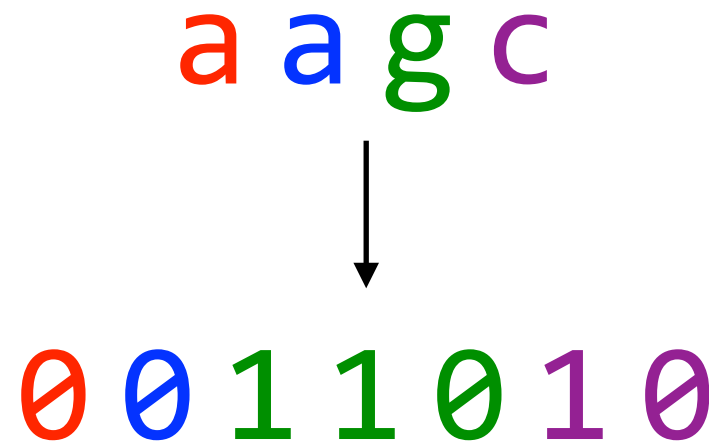
Proposal 1

$$C(a) = 0$$

$$C(c) = 10$$

$$C(g) = 110$$

$$C(t) = 111$$



Each codeword is
unique; i.e. C is injective

Example courtesy of Mathematicalmonk videos on information theory
<https://youtu.be/9MCxXJn7TPU>

Codes

Can we go **recover original string** from code?


Proposal 1

$$C(a) = 0$$

$$C(c) = 10$$

$$C(g) = 110$$

$$C(t) = 111$$


1 1 1 0 0 1 0

Example courtesy of Mathematicalmonk videos on information theory

<https://youtu.be/9MCxXJn7TPU>

Codes

Can we go **recover original string** from code?

Proposal 1

$$C(a) = 0$$

$$C(c) = 10$$

$$C(g) = 110$$

$$C(t) = 111$$

t a a c
 ↑ yes
1 1 1 0 0 1 0

Example courtesy of Mathematicalmonk videos on information theory

<https://youtu.be/9MCxXJn7TPU>

Codes

$$X = \{ a : \frac{1}{2}, c : \frac{1}{4}, g : \frac{1}{8}, t : \frac{1}{8} \}$$

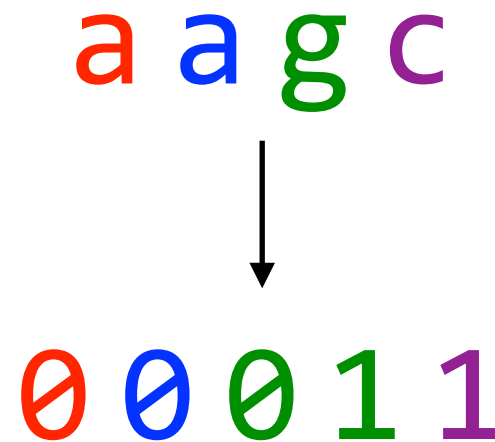
Proposal 2

$$C(a) = 0$$

$$C(c) = 1$$

$$C(g) = 01$$

$$C(t) = 10$$



Again, C is injective

Example courtesy of Mathematicalmonk videos on information theory

<https://youtu.be/9MCxXJn7TPU>

Codes

Can we go **recover original string** from code?

Proposal 2

$$C(a) = 0$$

$$C(c) = 1$$

$$C(g) = 01$$

$$C(t) = 10$$

0 0 0 1 1
 ↑ ?

Example courtesy of Mathematicalmonk videos on information theory

<https://youtu.be/9MCxXJn7TPU>

Codes

Can we go **recover original string** from code?

Proposal 2

$$C(a) = 0$$

$$C(c) = 1$$

$$C(g) = 01$$

$$C(t) = 10$$

no
↑?
0 0 0 1 1
—
a a a c ?
a a g ?

Example courtesy of Mathematicalmonk videos on information theory

<https://youtu.be/9MCxXJn7TPU>

Codes

Let C' be the code extended to *sequences*

$$C' : \Sigma^* \rightarrow \{0,1\}^*$$

$$C(a) = 0$$

$$C'(a) = 0$$

$$C(c) = 10$$

$$C'(ag) = 0110$$

$$C(g) = 110$$

$$C'(tt) = 111111$$

$$C(t) = 111$$

$$C'(aaaac) = 000010$$

Goal is for C' to be injective

(C being injective is not enough)

Codes

Consider two codes, both unambiguous

A

$$C(a) = 1$$

$$C(c) = 10$$

$$C(g) = 00$$

a a g c



1 1 0 0 1 0

B

$$C(a) = 1$$

$$C(c) = 01$$

$$C(g) = 00$$

a a g c



1 1 0 0 0 1

Codes

Now we decode:

A

$$C(a) = 1$$

$$C(c) = 10$$

$$C(g) = 00$$

1 1 0 0 1 0
1

Considering first **1**, can't yet tell if it's an a or part of a c

Codes

Now we decode:

A

$C(a) = 1$

$C(c) = 10$

$C(g) = 00$

1 1 0 0 1 0

Now sure that first 1 is a.

Not sure about second 1.

Codes

Now we decode:

A

$$C(a) = 1$$

$$C(c) = 10$$

$$C(g) = 00$$

1 1 0 0 1 0

Either we have

ac...

aag...

Codes

Now we decode:

A

$$C(a) = 1$$

$$C(c) = 10$$

$$C(g) = 00$$

1 1 0 0 1 0

Either we have

acg...

aag...

Codes

Now we decode:

A

$$C(a) = 1$$

$$C(c) = 10$$

$$C(g) = 00$$

1 1 0 0 1 0

Now we're sure we have:

aag...

But could still be aaga...

or aagc...

Codes

Now we decode:

A

$C(a) = 1$

$C(c) = 10$

$C(g) = 00$

1 1 0 0 1 0

Now we're sure we have:

aagc

Codes

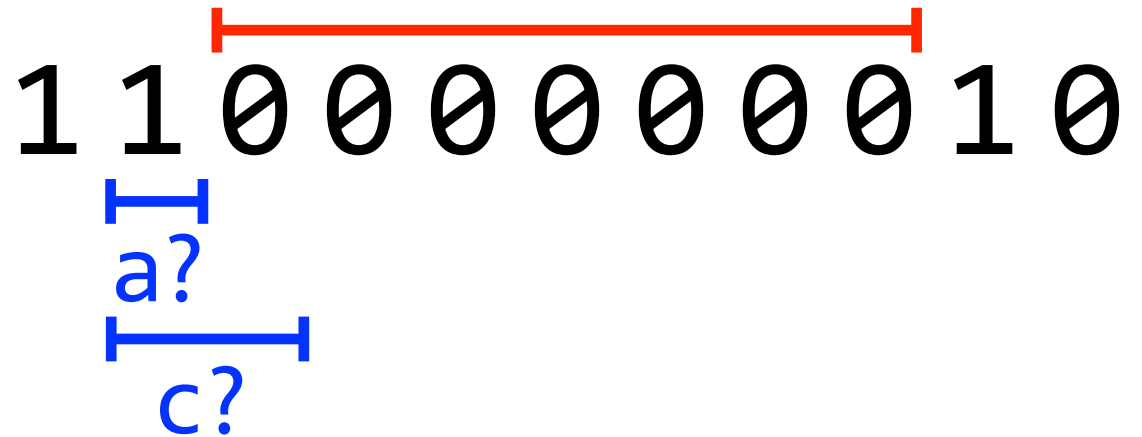
Consider an example with a longer run of 0s:

A

$$C(a) = 1$$

$$C(c) = 10$$

$$C(g) = 00$$



Can't distinguish **a** from **c** until we see whether **run of 0s** is odd or even

Since it's odd, must be a c: a**cgggc**

Codes


Now we decode:

B

$$C(a) = 1$$

$$C(c) = 01$$

$$C(g) = 00$$

1 1 0 0 0 1


Considering first **1**, we're immediately sure it's an a

Codes

Now we decode:

B

$$C(a) = 1$$

$$C(c) = 01$$

$$C(g) = 00$$

1 1 0 0 0 1

Definitely **aa**

Codes

Now we decode:

B

$C(a) = 1$

$C(c) = 01$

$C(g) = 00$

1 1 0 0 0 1

Could be aac or aag

Codes

Now we decode:

B

$$C(a) = 1$$

$$C(c) = 01$$

$$C(g) = 00$$

1 1 0 0 0 1

Definitely aag

Codes

Now we decode:

B

$$C(a) = 1$$

$$C(c) = 01$$

$$C(g) = 00$$

1 1 0 0 0 1

Could be aagc or aagg

Codes

Now we decode:

B

$$C(a) = 1$$

$$C(c) = 01$$

$$C(g) = 00$$

1 1 0 0 0 1

Definitely aagc

Codes

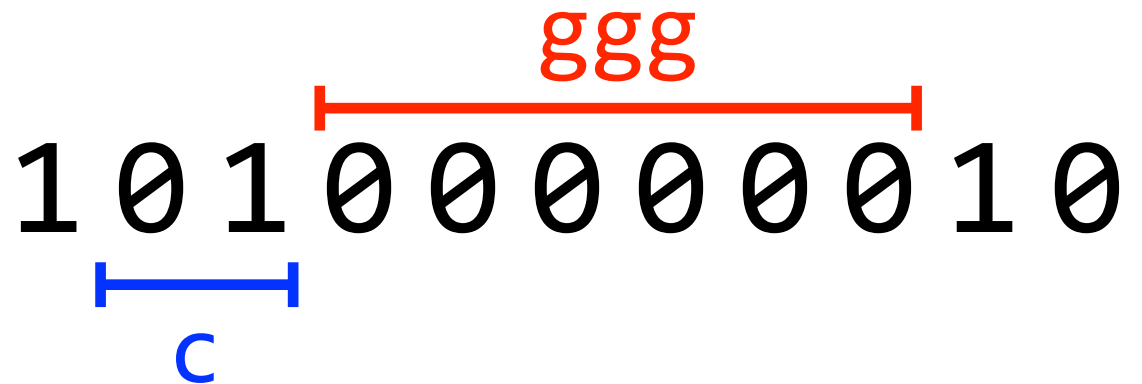
No problems with decoding efficiency here.

B

$$C(a) = 1$$

$$C(c) = 01$$

$$C(g) = 00$$



Code is ***prefix-free***; no code is a prefix of another.
Also called a ***prefix code*** for short.

AKA ***instantaneous***

Huffman

Say we start with a string: abracadabra

Can compile symbols and their frequencies:

$$\{ a : 5, b : 2, c : 1, d : 1, r : 2 \}$$

Or equivalently, a r.v.:

$$X = \{ a : \frac{5}{11}, b : \frac{2}{11}, c : \frac{1}{11}, d : \frac{1}{11}, r : \frac{2}{11} \}$$

Huffman

{ a : 5, b : 2, c : 1, d : 1 r : 2 }

In each round, ***join*** the 2 subtrees with lowest total weight

a⁵

b²

c¹

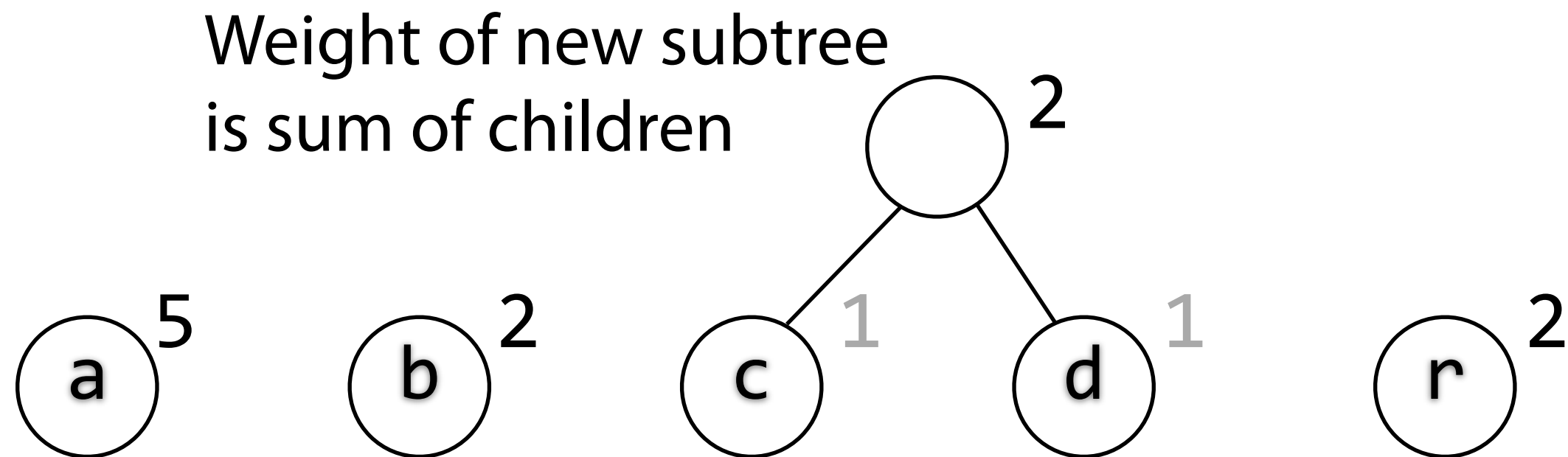
d¹

r²

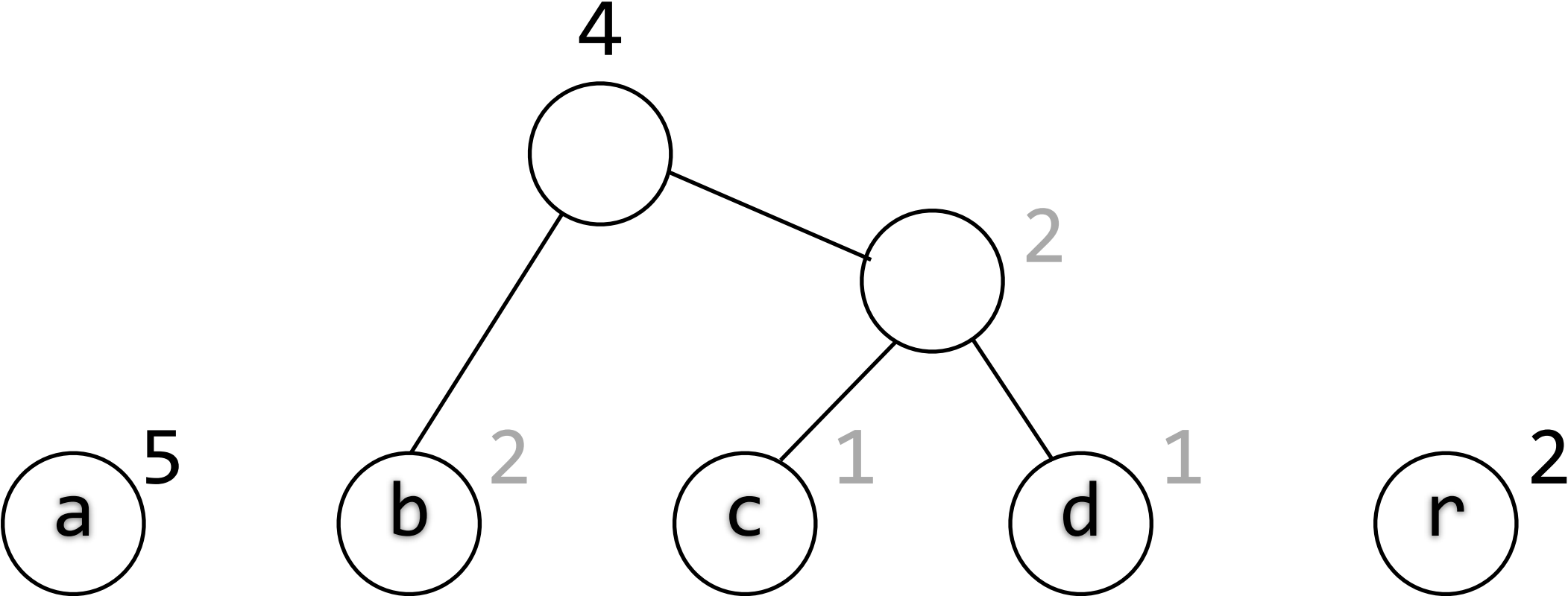
Huffman

{ a : 5, b : 2, c : 1, d : 1 r : 2 }

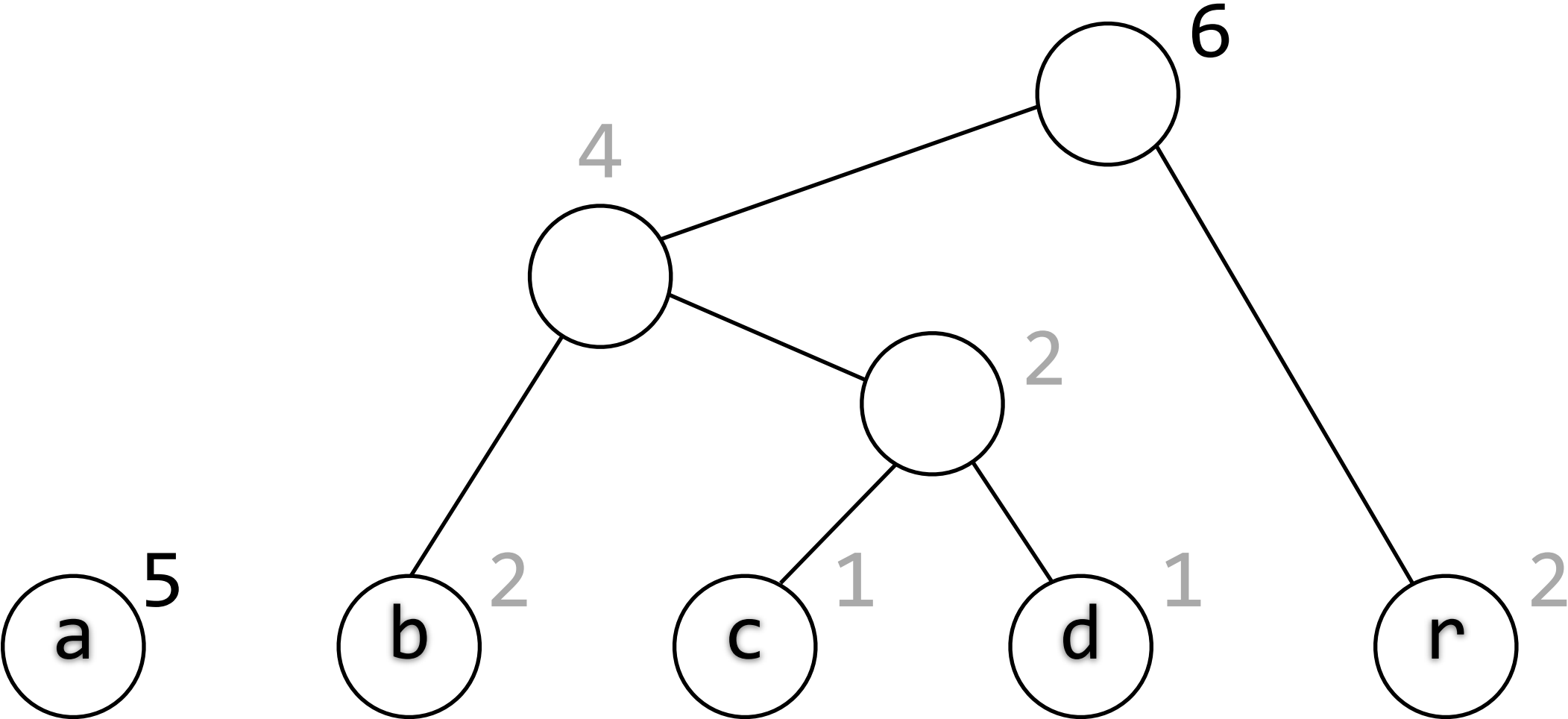
In each round, **join** the 2 subtrees with lowest total weight



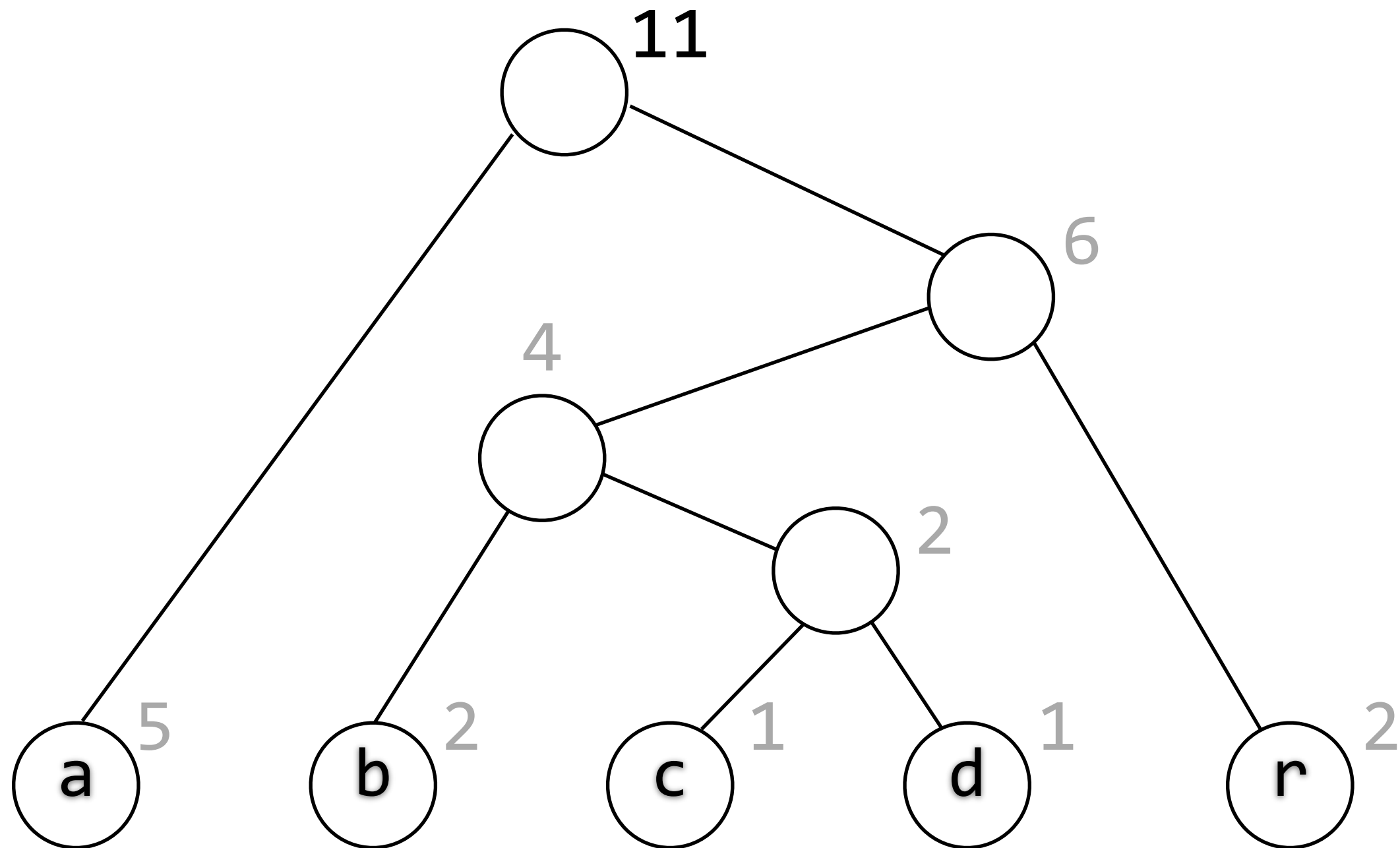
Huffman



Huffman

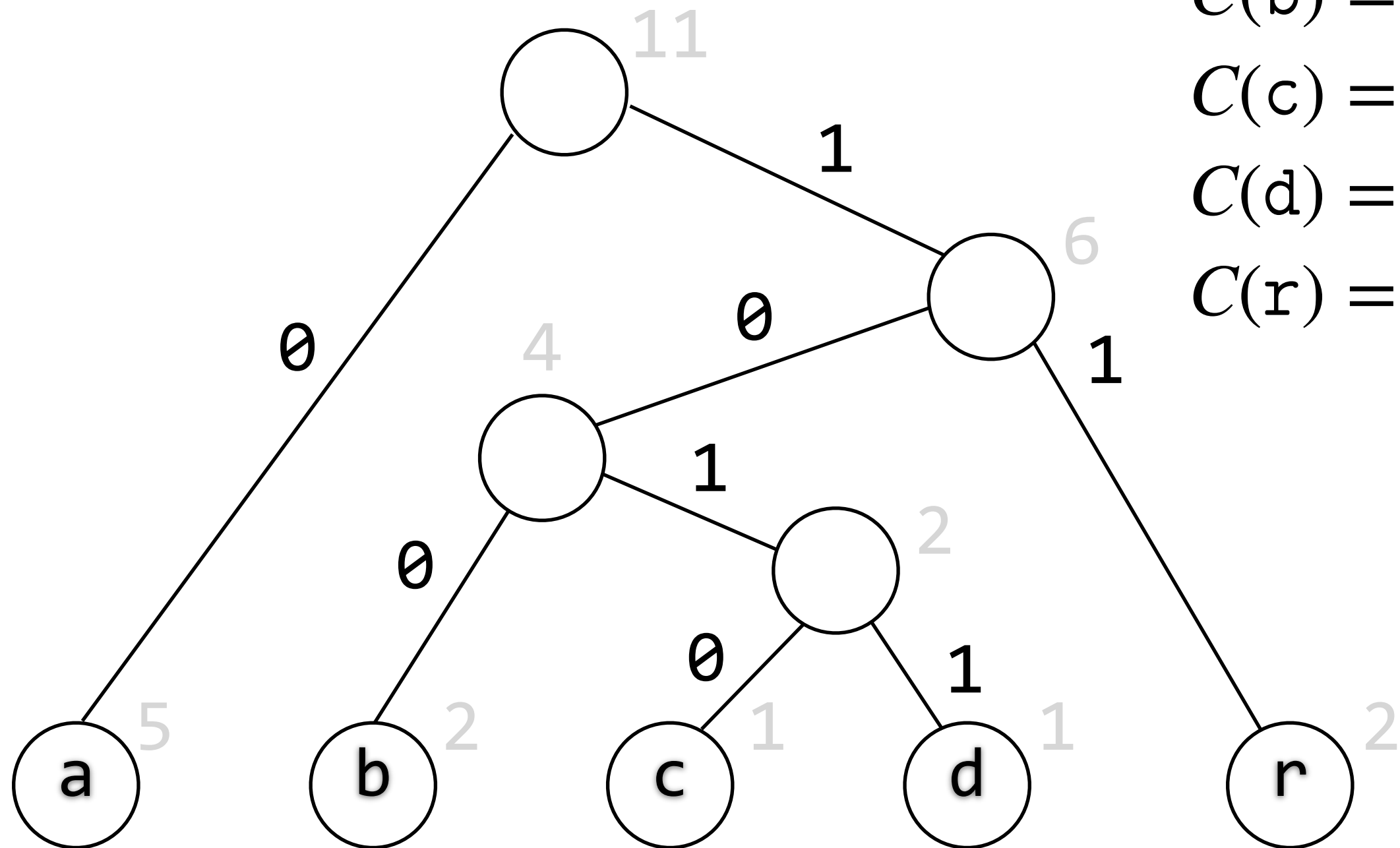


Huffman



This is the tree but what is the code?

Huffman



$$C(a) = 0$$

$$C(b) = 100$$

$$C(c) = 1010$$

$$C(d) = 1011$$

$$C(r) = 11$$

Label edges with 0/1 according to left/right child of parent

Codes equal root-to-leaf concatenations of 0/1's

Huffman

Huffman codes are "optimal," wasting at most 1 bit per symbol

In other words, if c is the number of bits in the Huffman code for an input string S of length n

$$c \leq n(H_0(S) + 1) \text{ bits}$$