

# Entropy & coding

Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL  
*of* ENGINEERING

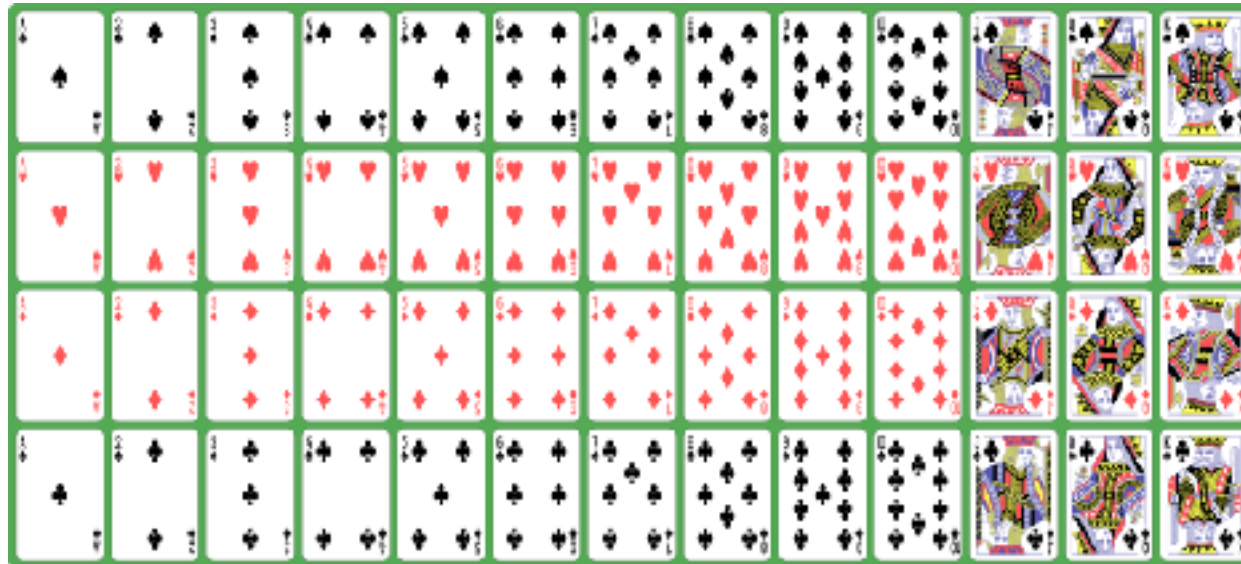
Department of Computer Science



Please sign guestbook ([www.langmead-lab.org/teaching-materials](http://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](mailto:ben.langmead@gmail.com)).

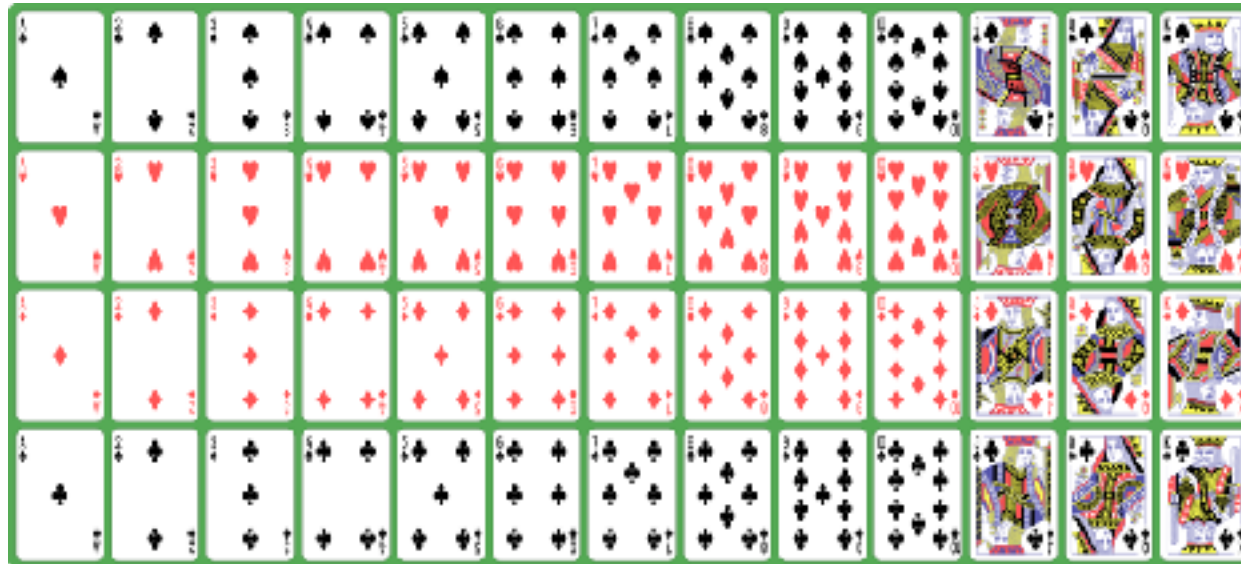
# Entropy & coding

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



# Entropy & coding

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



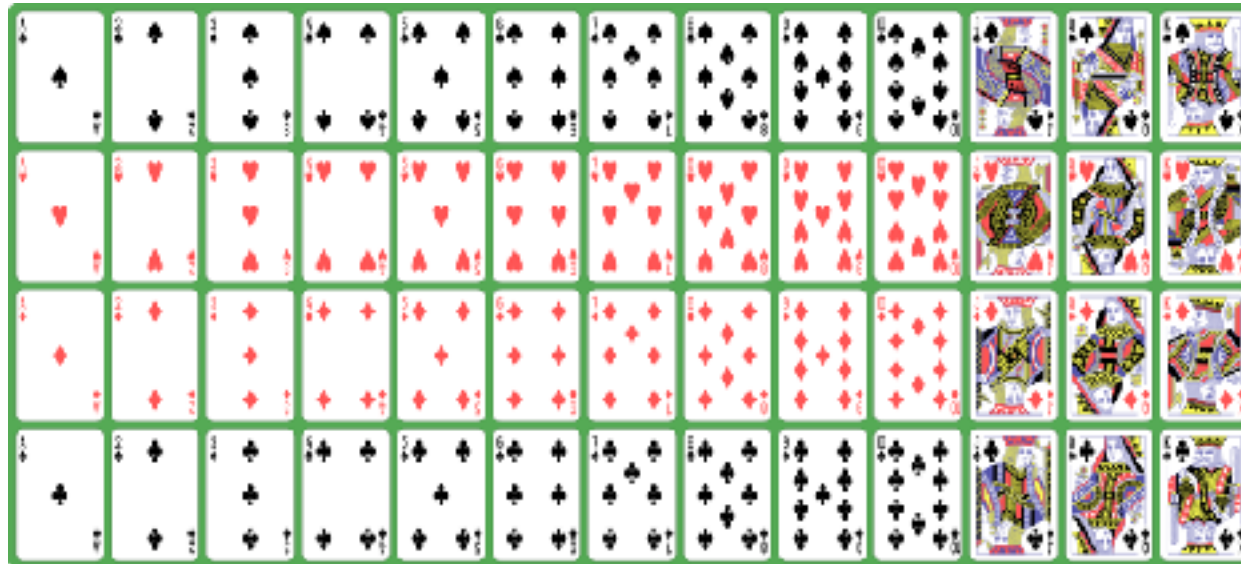
Say, code = rank + (13 \* suit)

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

♠ = 0, ♥ = 1, ...

# Entropy & coding

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



A ♠	
2 ♠	
3 ♠	
4 ♠	
⋮	
10 ♣	
J ♣	
Q ♣	
K ♣	

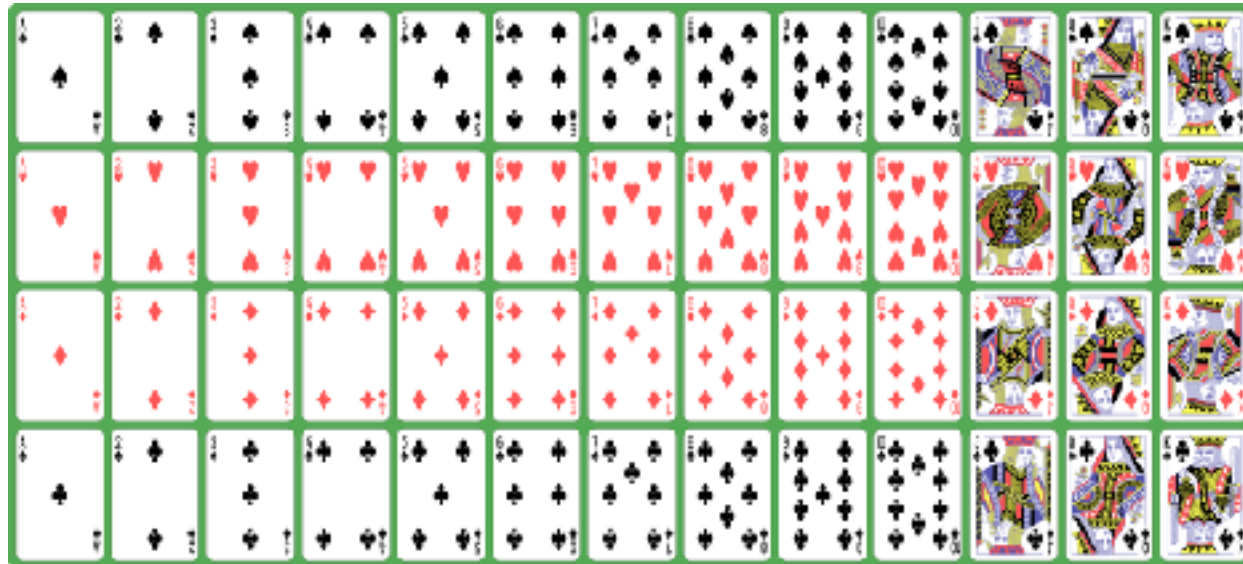
Say, code = rank + (13 \* suit)

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

♠ = 0, ♥ = 1, ...

# Entropy & coding

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



Say, code = rank + (13 \* 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  $U$ ?

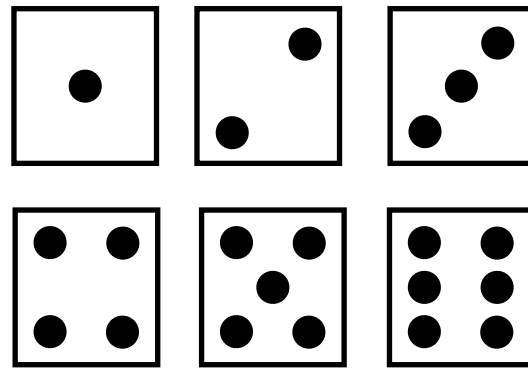
$$H_{wc}(U) = \log_2 |U|$$

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?



(37 or 38 slots)

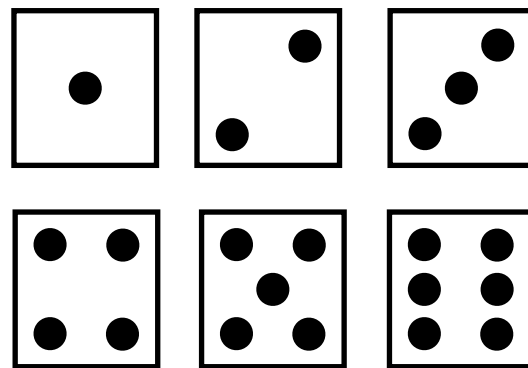


# Entropy

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

# Entropy

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

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

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

# 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

# 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	— • —	1	• — — — —
L	• — • •	2	• • — — —
M	— —	3	• • • — —
N	— •	4	• • • • —
O	— — —	5	• • • • •
P	• — — •	6	— • • • •
Q	— — • —	7	— — • • •
R	• — •	8	— — — • •
S	• • •	9	— — — — •
T	—	0	— — — — —

# Entropy

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

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

# Entropy

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

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

This is *Shannon entropy*



# Entropy

$$X = \left\{ \begin{array}{l} \text{Lincoln} : 0.5, \\ \text{Lincoln} : 0.5 \end{array} \right\}$$

$$H(X) =$$

# Entropy

$$X = \left\{ \begin{array}{l} \text{Lincoln} : 0.5, \\ \text{Washington} : 0.5 \end{array} \right\}$$

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

# Entropy

$$X = \left\{ \begin{array}{l} \text{Lincoln} : 0.9, \\ \text{Lincoln} : 0.1 \end{array} \right\}$$

$$H(X) =$$

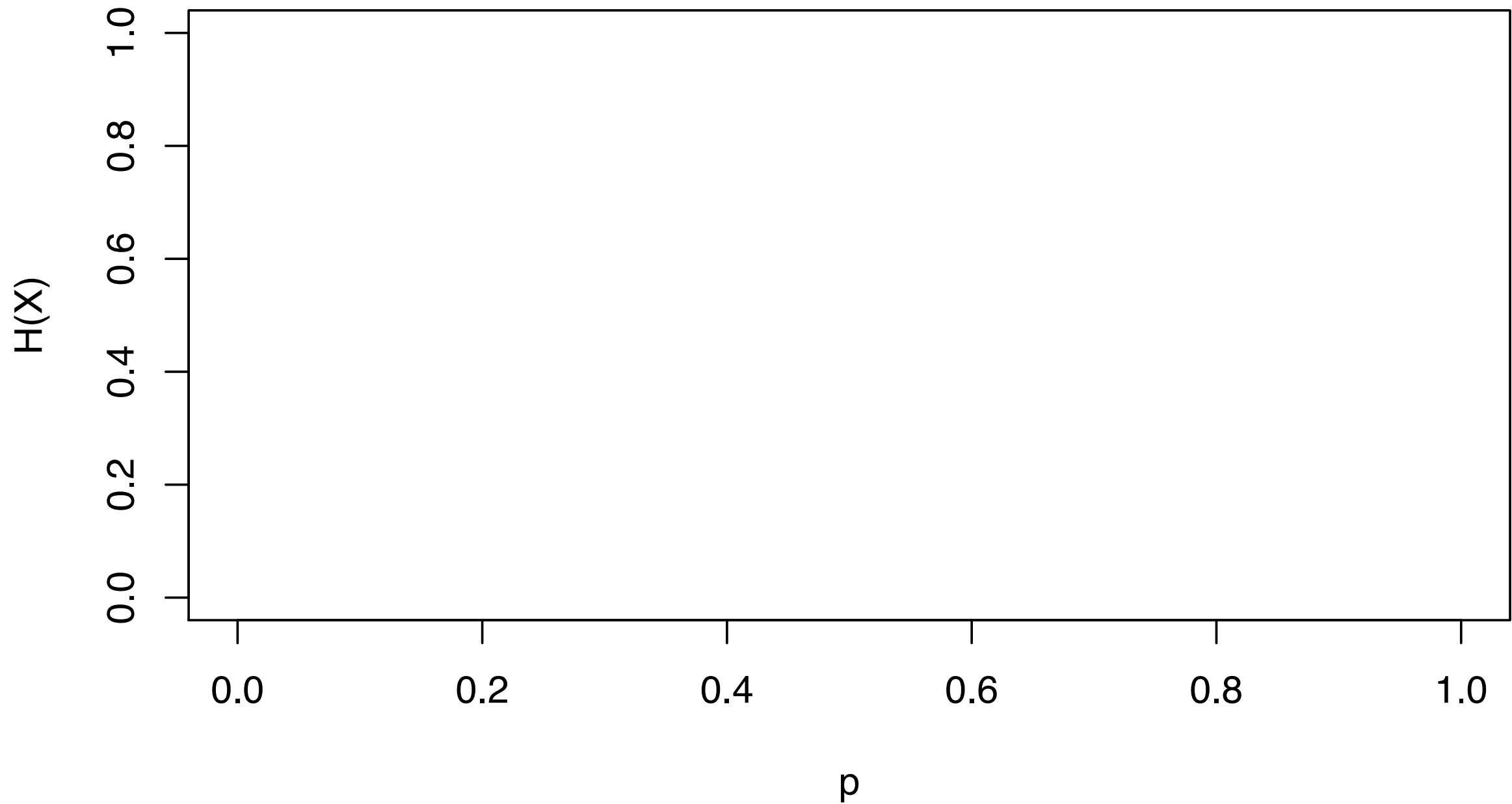
# Entropy

$$X = \left\{ \begin{array}{l} \text{Lincoln Penny} : 0.9, \\ \text{Reverse Penny} : 0.1 \end{array} \right\}$$

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

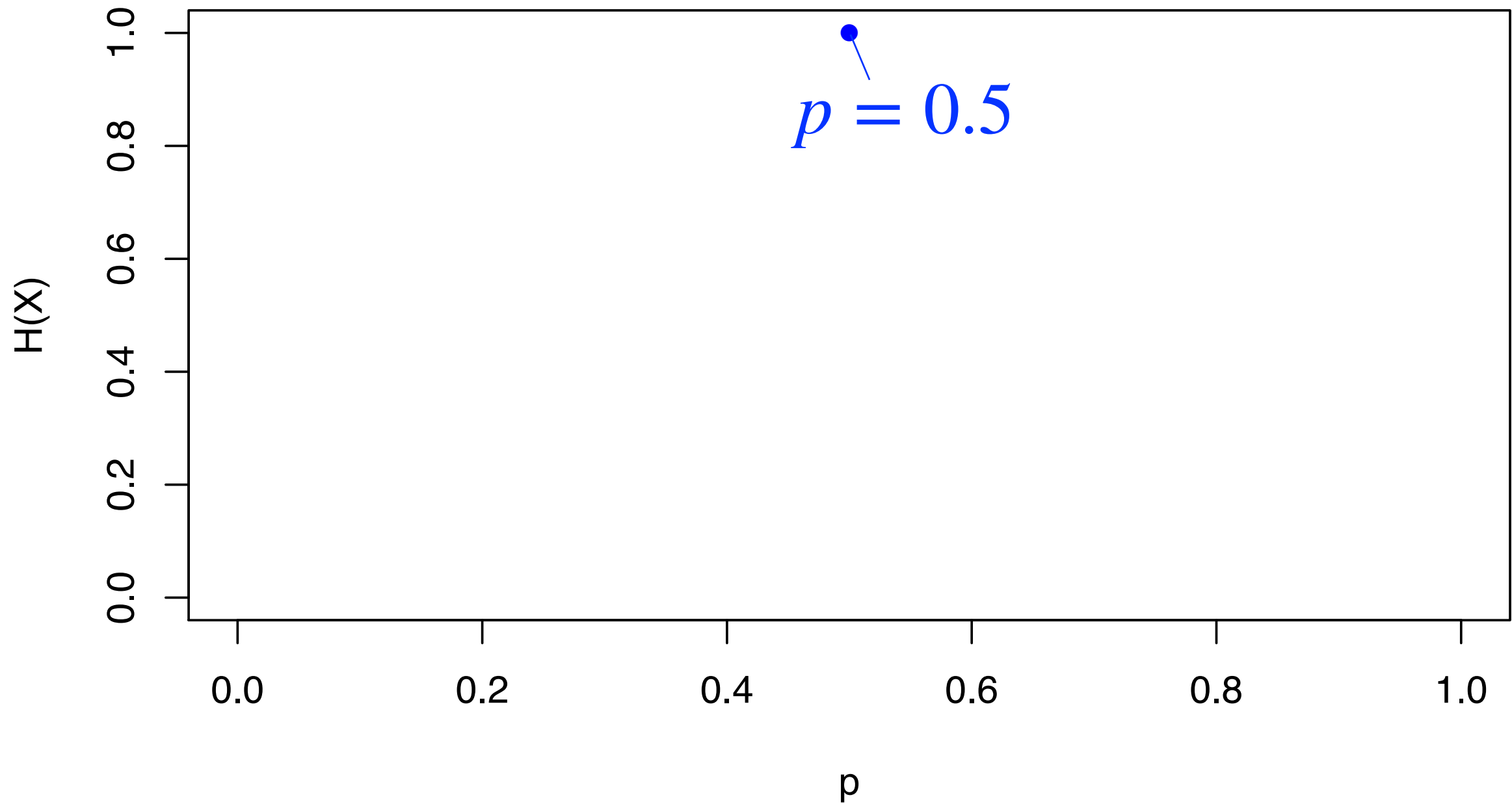
# Entropy

$$X = \left\{ \text{Lincoln} : p, \text{Lincoln Memorial} : 1 - p \right\}$$



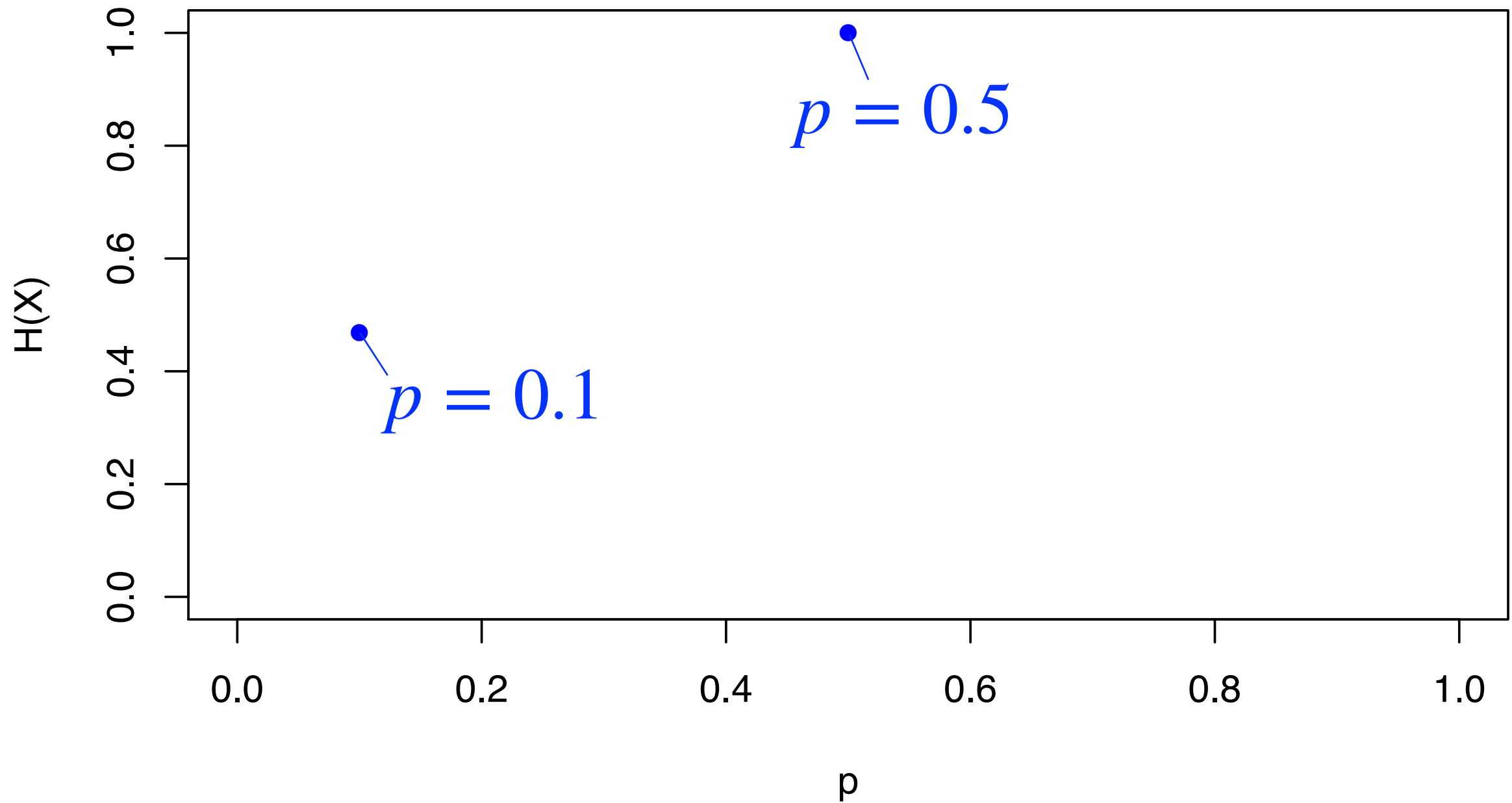
# Entropy

$$X = \left\{ \text{Lincoln} : p, \text{Lincoln Memorial} : 1 - p \right\}$$



# Entropy

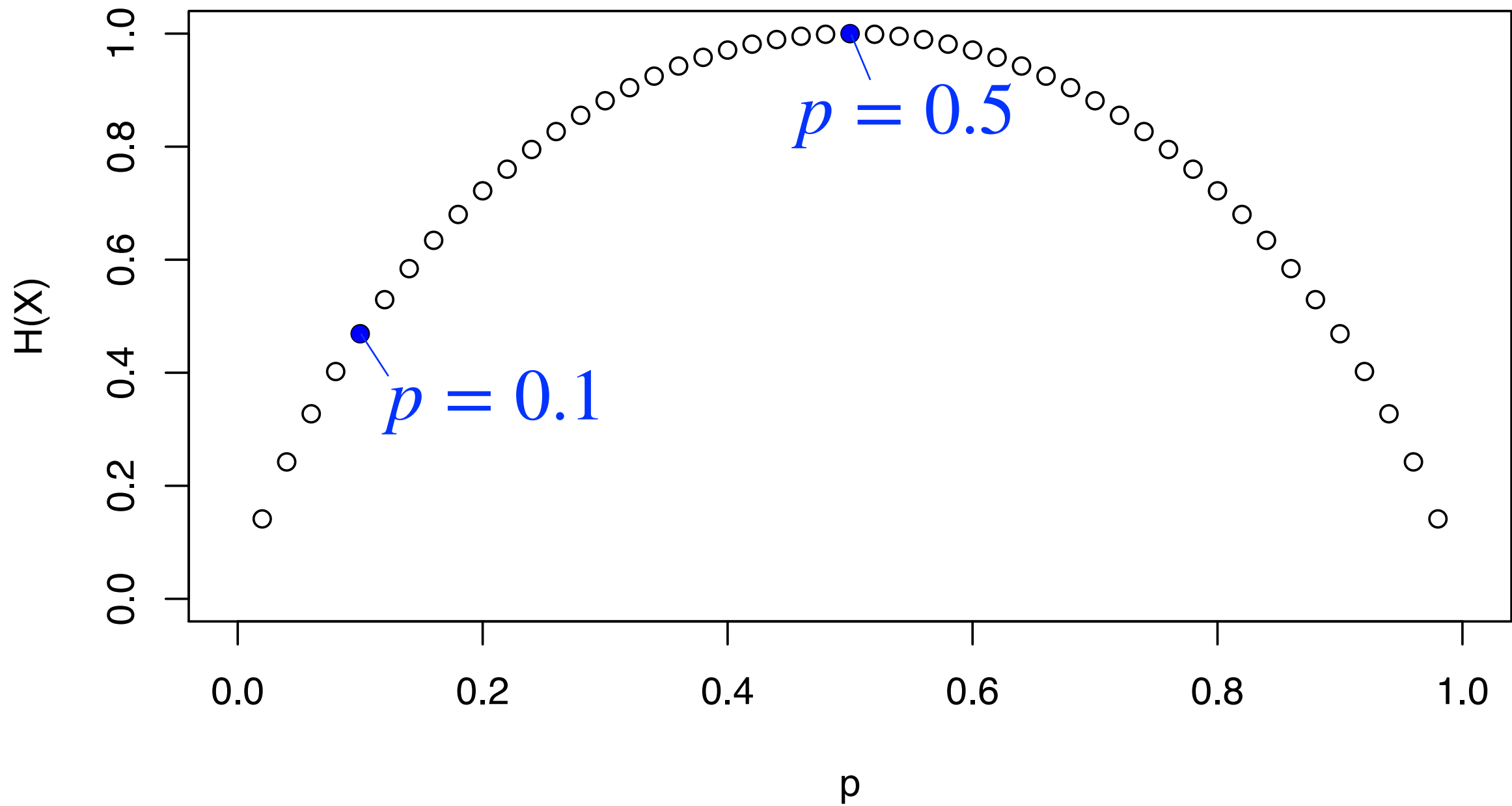
$$X = \left\{ \text{penny} : p, \text{nickel} : 1 - p \right\}$$





# Entropy

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



# Entropy

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

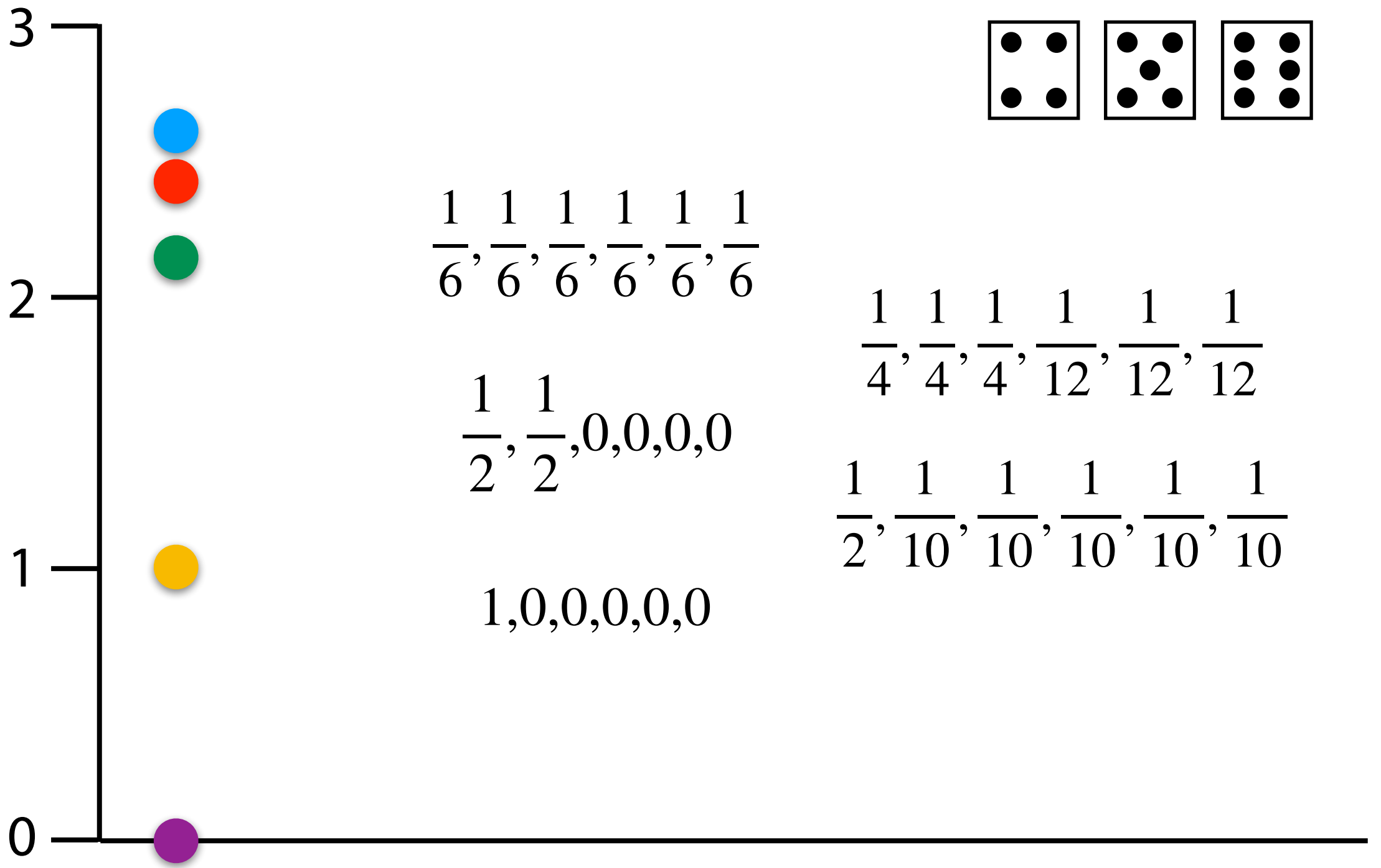
$$H(X) =$$

# Entropy

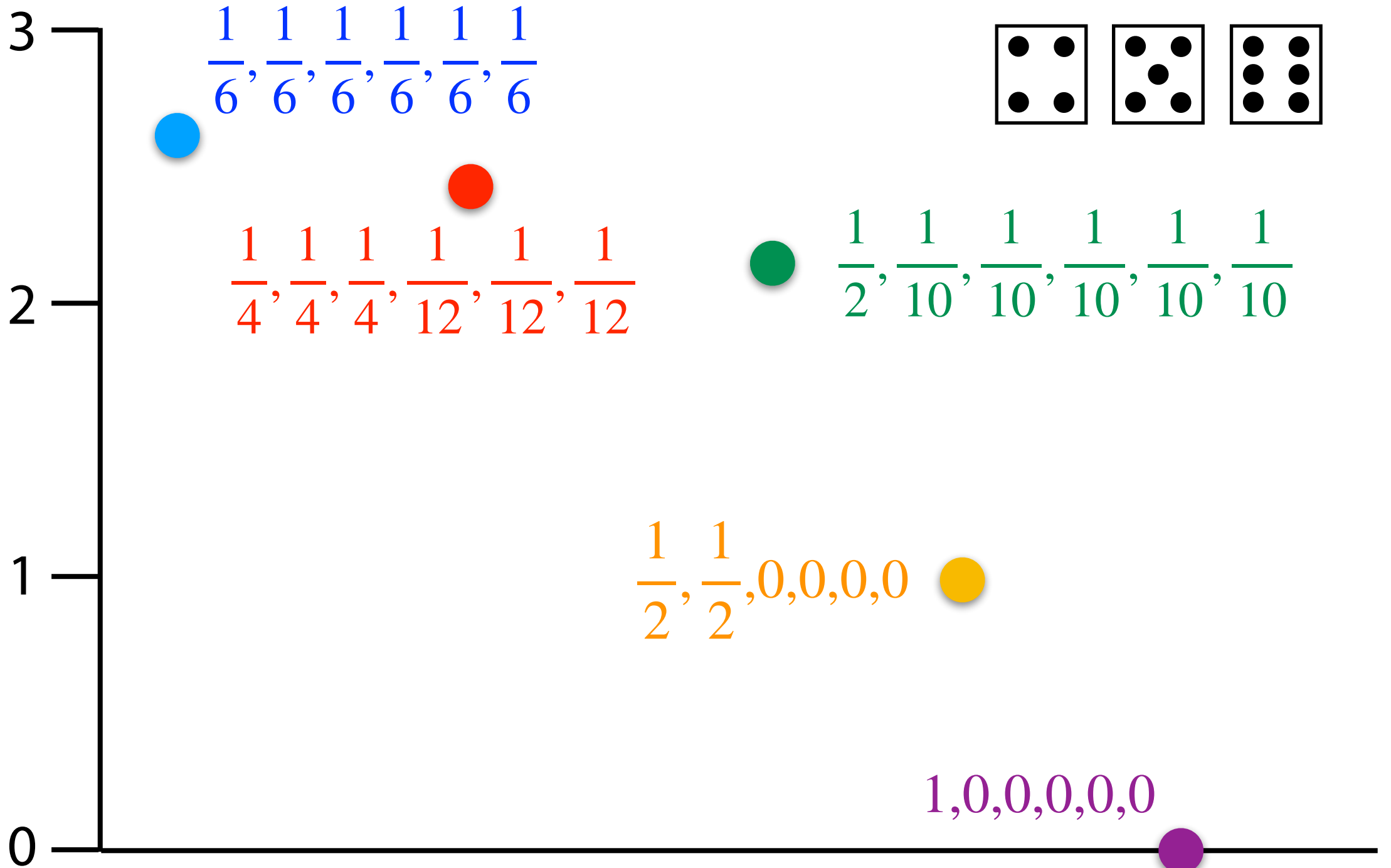
$$X = \left\{ \begin{array}{c} \square \cdot \\ \square \cdot \cdot \\ \square \cdot \cdot \cdot \\ \square \cdot \cdot \cdot \cdot \\ \square \cdot \cdot \cdot \cdot \cdot \\ \square \cdot \cdot \cdot \cdot \cdot \cdot \end{array} : \frac{1}{6} \text{ each} \right\}$$

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

# Entropy



# Entropy



# Entropy

When outcomes are equally probable:

$$H(X) =$$

# 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

$$H_0(B) =$$

$$m = \# \text{ 1s in } B$$

$$n = |B|$$

# 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

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

# 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	— — — — —



# Codes

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

Shannon entropy equation hints at codes of

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

# 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

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

## Proposal 1

$$C(a) = 0$$

$$C(c) = 10$$

$$C(g) = 110$$

$$C(t) = 111$$

a a g c  
↓

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

# Codes

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

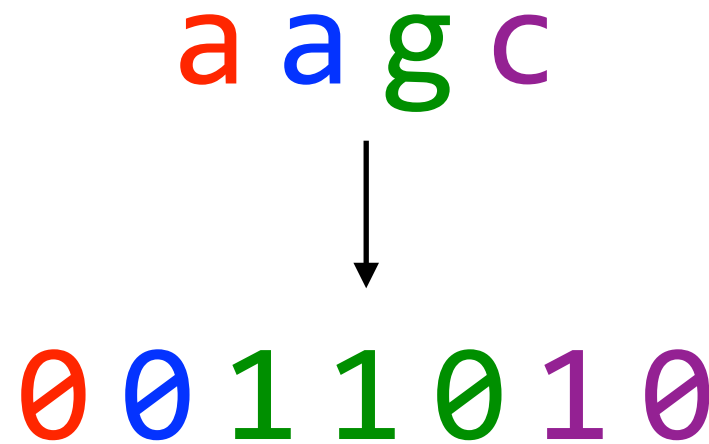
## Proposal 1

$$C(a) = 0$$

$$C(c) = 10$$

$$C(g) = 110$$

$$C(t) = 111$$



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

# 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**



# 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

# Codes

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

## Proposal 2

$$C(a) = 0$$

$$C(c) = 1$$

$$C(g) = 01$$

$$C(t) = 10$$

a a g c  
↓

Again,  $C$  is injective

# Codes

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

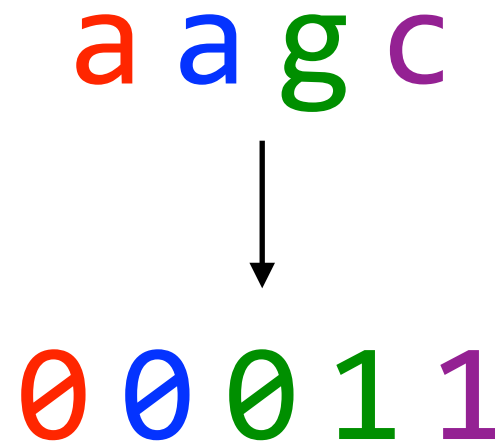
## Proposal 2

$$C(a) = 0$$

$$C(c) = 1$$

$$C(g) = 01$$

$$C(t) = 10$$



Again,  $C$  is injective



# 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  
          ↑ ?

# 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 ?

# Codes

Let  $C'$  be the code extended to *sequences*

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

$$C(a) = 0$$

$$C'(a) = 0$$

$$C(c) = 10$$

$$C'(ag) = 0110$$

$$C(g) = 110$$

$$C'(tt) = 111111$$

$$C(t) = 111$$

$$C'(aaaac) = 000010$$

$C'$  should be injective, giving *unambiguous* code

( $C$  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:

	<b>A</b>
$C(a)$	$= 1$
$C(c)$	$= 10$
$C(g)$	$= 00$

1 1 0 0 1 0

Either we have

ac...

aag...

# Codes

Now we decode:

	<b>A</b>
$C(a)$	$= 1$
$C(c)$	$= 10$
$C(g)$	$= 00$

1 1 0 0 1 0

Either we have

acg...

aag...



# Codes

Now we decode:

	<b>A</b>
$C(a) =$	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:

	<b>A</b>
$C(a) =$	1
$C(c) =$	10
$C(g) =$	00

**1 1 0 0 1 0**  
└──────────┘

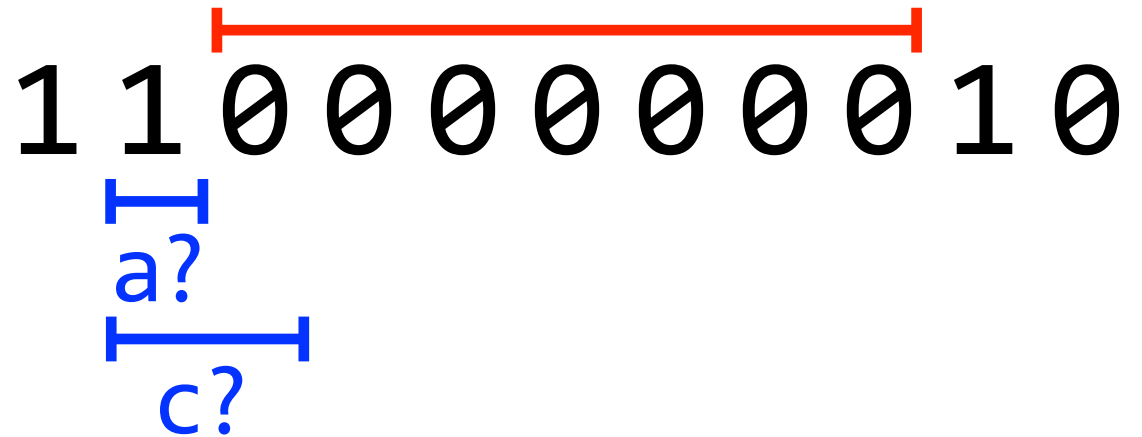
Now we're sure we have:

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**c**gggc

# Codes

Now we decode:

	<b>B</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>B</b>
$C(a)$	$= 1$
$C(c)$	$= 01$
$C(g)$	$= 00$


**1 1** 0 0 0 1  


Definitely **aa**

# Codes

Now we decode:

	<b>B</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>B</b>
$C(a)$	$= 1$
$C(c)$	$= 01$
$C(g)$	$= 00$

**1 1 0 0 0 1**  


Definitely aag

# Codes

Now we decode:

	<b>B</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>B</b>
$C(a)$	$= 1$
$C(c)$	$= 01$
$C(g)$	$= 00$

**1 1 0 0 0 1**

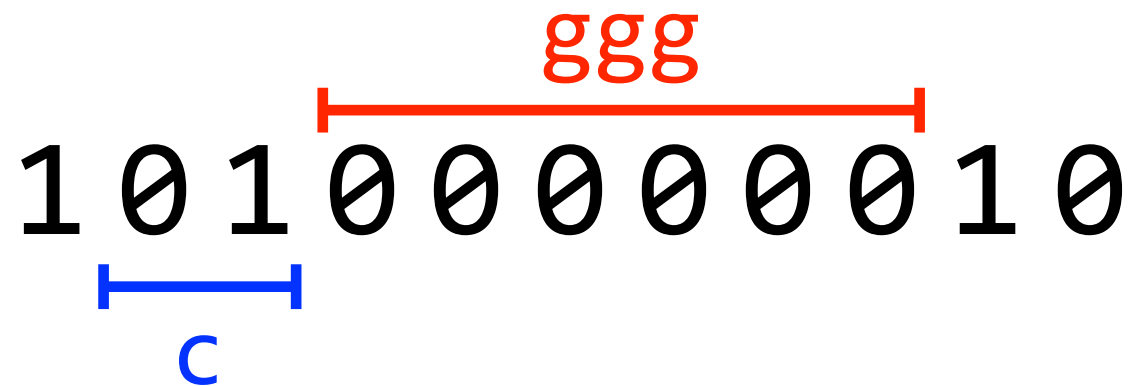


Definitely aagc

# Codes

No problems with decoding efficiency here.

	<b>B</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 }

# Huffman

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

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

a<sup>5</sup>

b<sup>2</sup>

c<sup>1</sup>

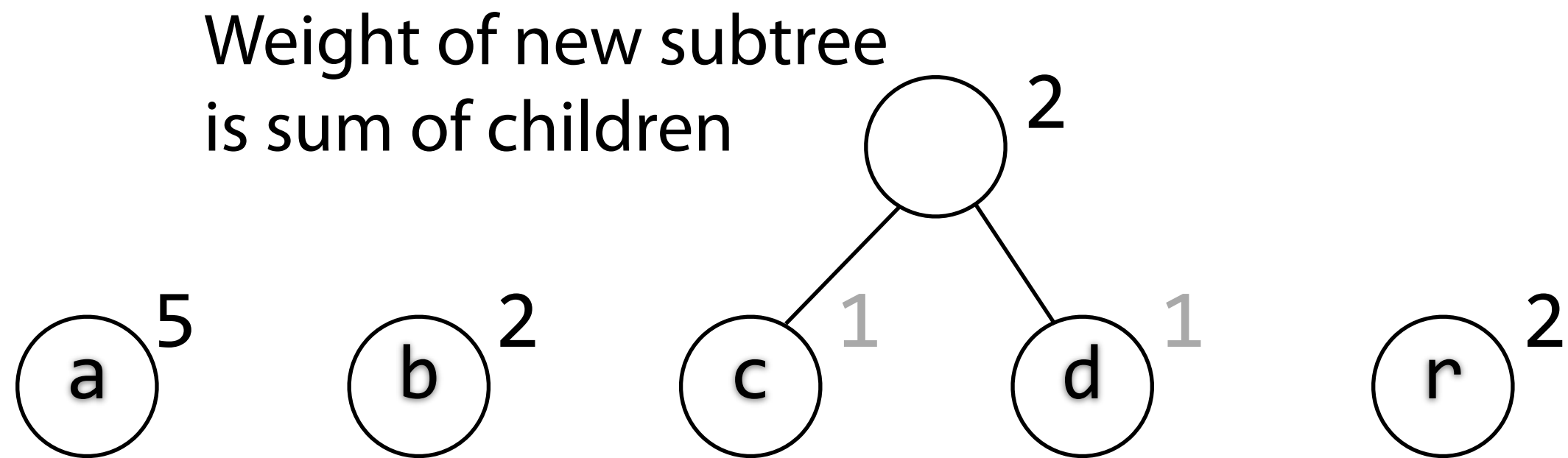
d<sup>1</sup>

r<sup>2</sup>

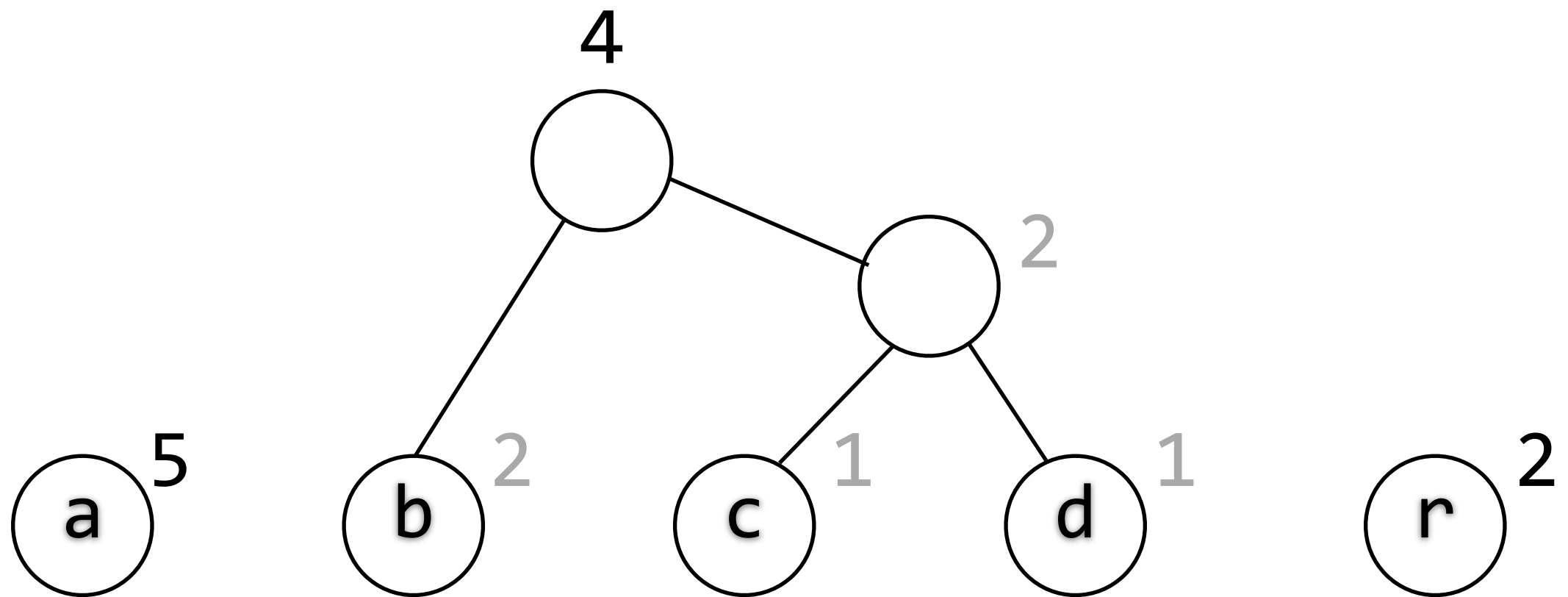
# Huffman

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

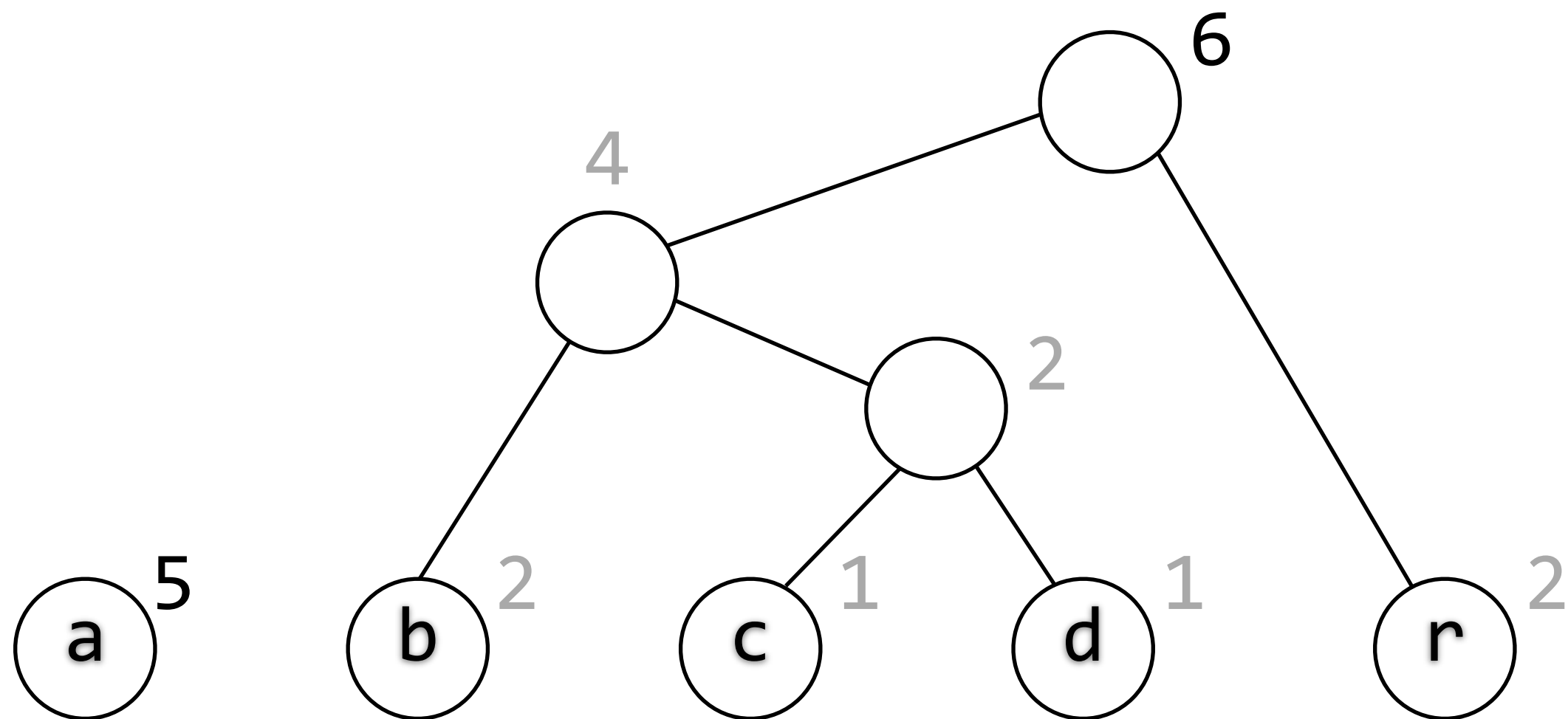
In each round, **join** the 2 subtrees with lowest total weight



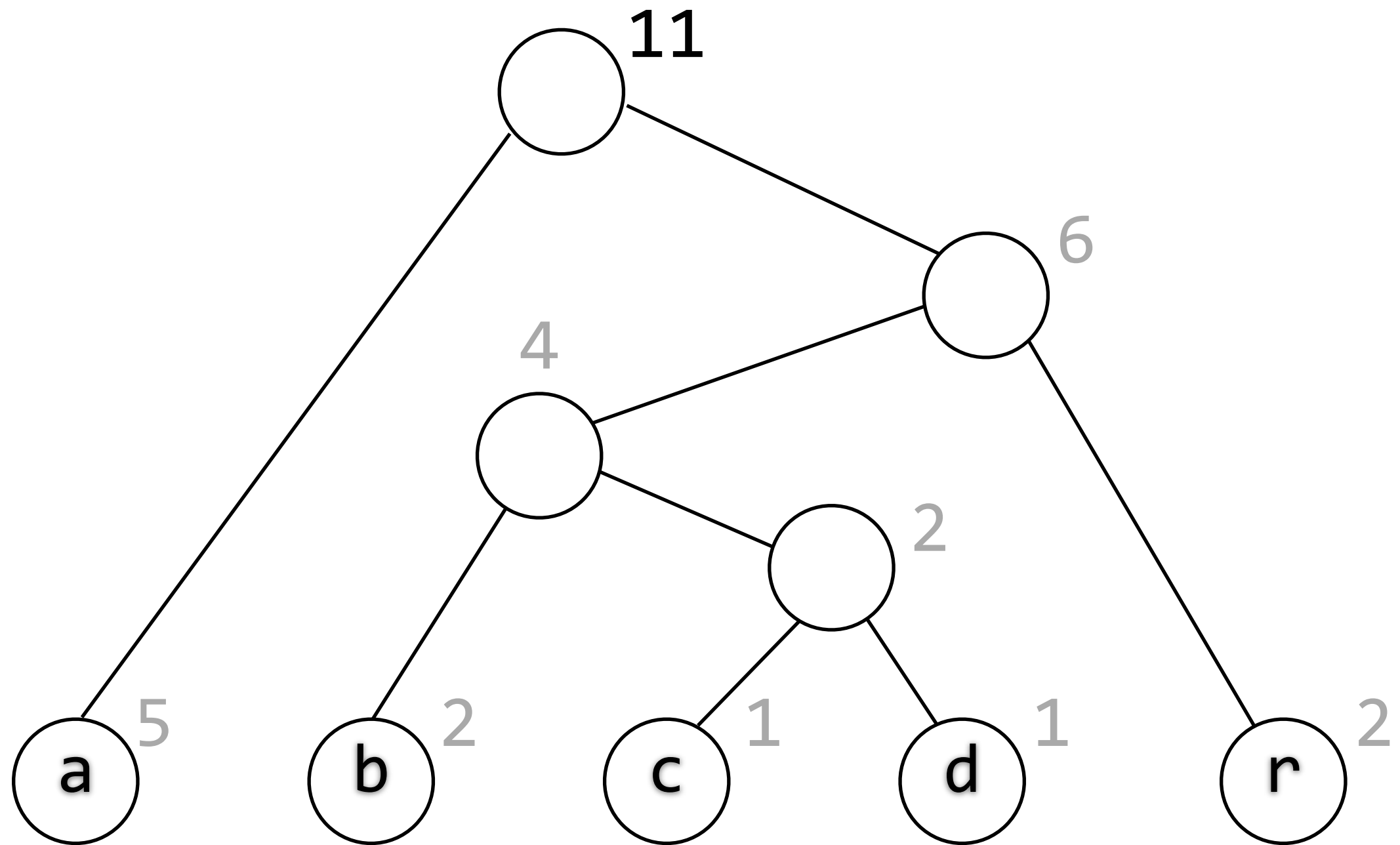
# Huffman



# Huffman



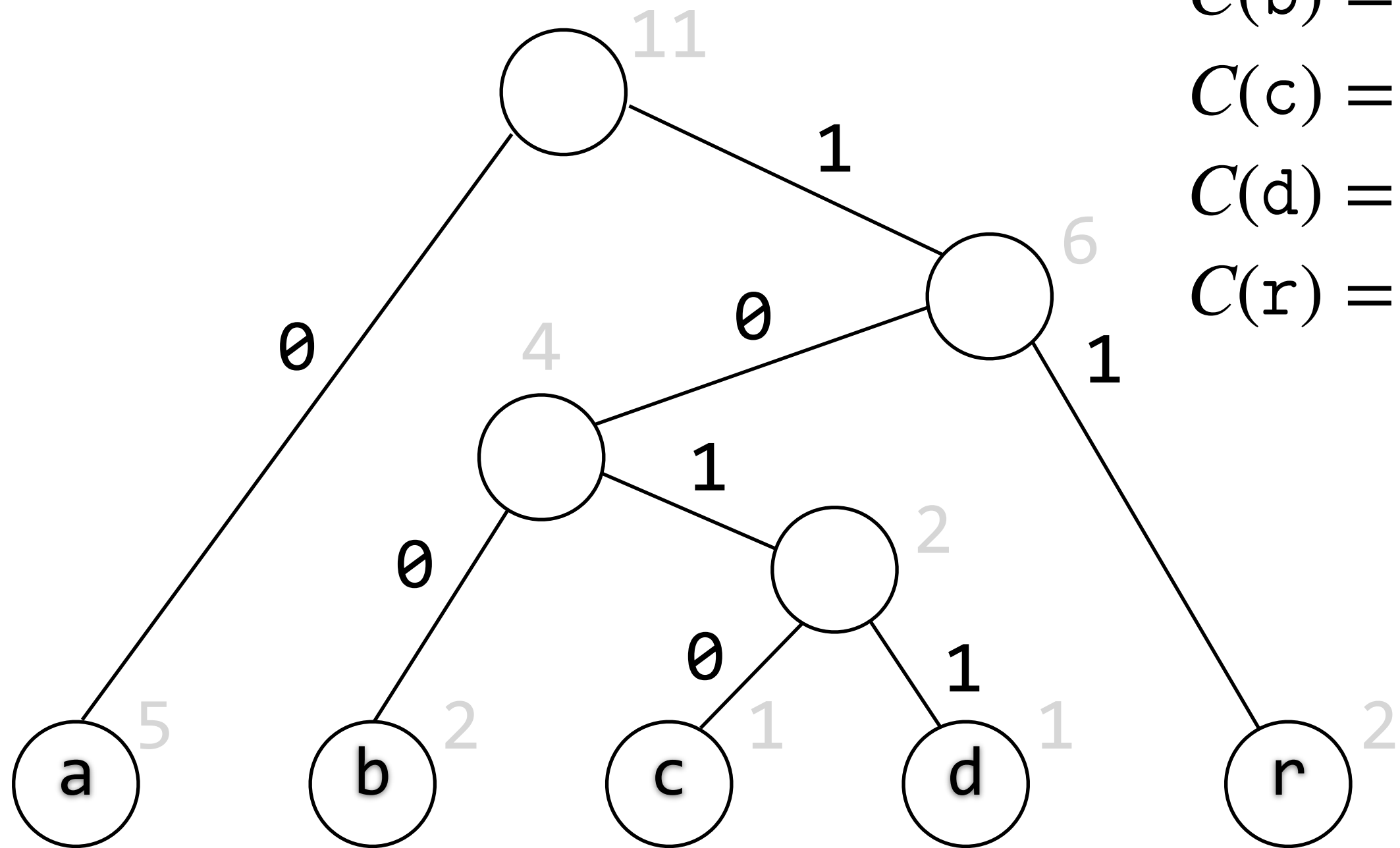
# Huffman



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" in that each code is at most 1 bit longer than optimal

In other words, for an input string  $S$ :

$$|C'(S)| \leq$$

# Huffman

Huffman codes are "optimal" in that each code is at most 1 bit longer than optimal

In other words, for an input string  $S$ :

$$|C'(S)| \leq |S| (H_0(S) + 1)$$