

**CSE 312**

# **Foundations of Computing II**

**Lecture 16: Information Theory and Data Compression**



**Stefano Tessaro**

tessaro@cs.washington.edu

# Announcements

- Office hours: I am available 1-3pm.
- Please make sure to read the instructions for the midterm.
- Practice midterm solutions posted in the afternoon.

# Today

How much can we compress data?

How much information is really contained in data?

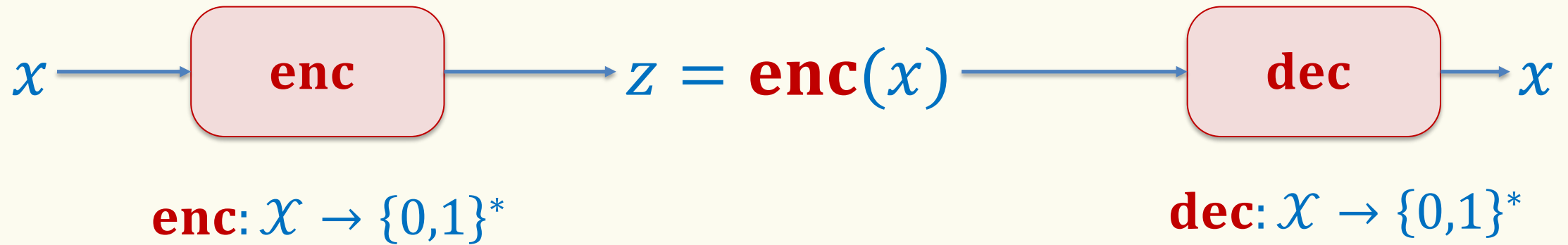
Central topic in **information theory**, a discipline based on probability which has been extremely useful across electrical engineering, computer science, statistics, physics, ...

Claude Shannon, "A Mathematical Theory of Communication", 1948

<http://www.math.harvard.edu/~ctm/home/text/others/shannon/entropy/entropy.pdf>



# Encoding Scheme



**Decodability.** For all values  $x \in \mathcal{X}$ :  $\mathbf{dec}(\mathbf{enc}(x)) = x$

**Goal:** Encoding should “compress”

[We will formalize this using the language of probability theory]

# Encoding – Example

Say we need to encode a word from the set  $\mathcal{X} = \{\text{hello, world, cse312}\}$

hello → 0  
world → 1  
cse312 → 11  
**enc**

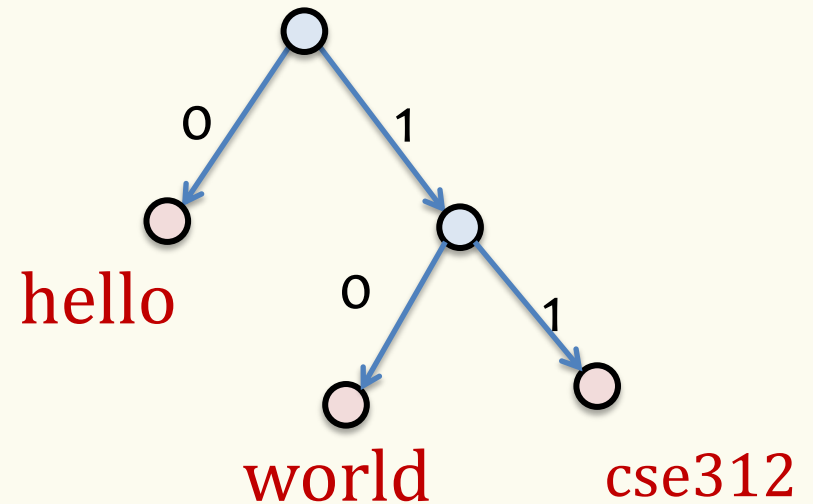
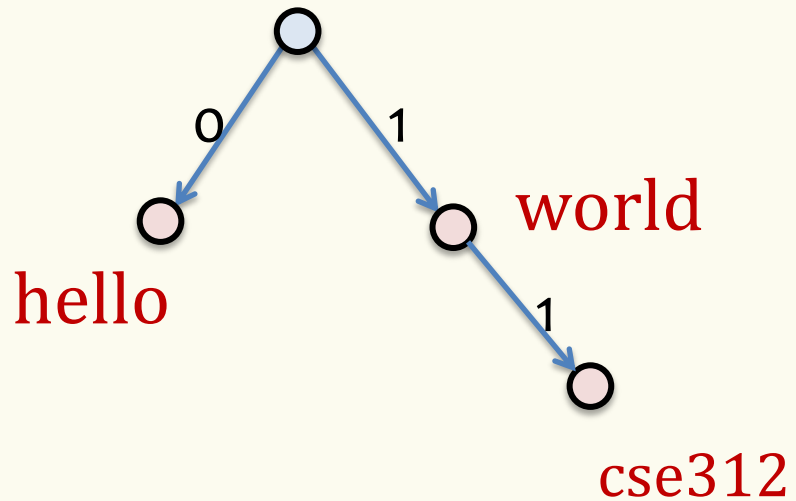
hello → 0  
world → 10  
cse312 → 11  
**enc**

hello → 0  
world → 11  
cse312 → 100000000  
**enc**

# Better Visualization – Trees

hello → 0  
world → 1  
cse312 → 11

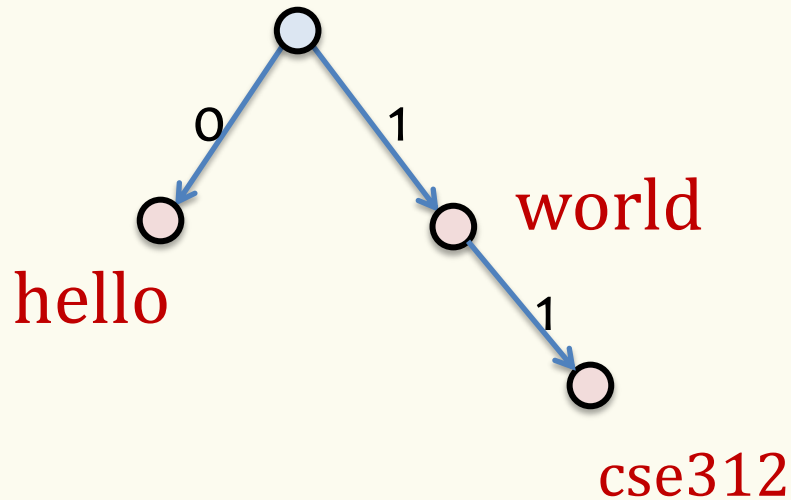
hello → 0  
world → 10  
cse312 → 11



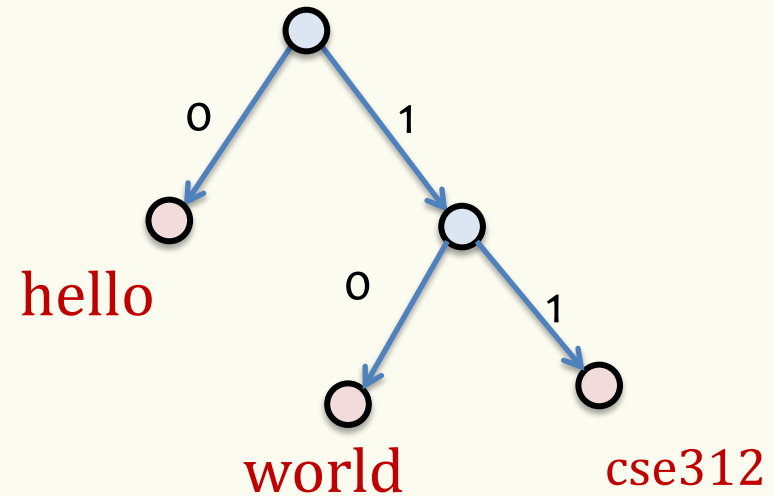
# Focus – Prefix-free codes

A code is **prefix-free** if no encoding is a **prefix** of another one.

i.e. every encoding is a leaf



**Not prefix-free!**  
1 is a prefix of 11



**Prefix-free!!**

# Random Variables – Arbitrary Values

We will consider random variables  $X: \Omega \rightarrow \mathcal{X}$  taking values from a (finite) set  $\mathcal{X}$ . [We refer to these as a “random variable over the alphabet  $\mathcal{X}$ .”]

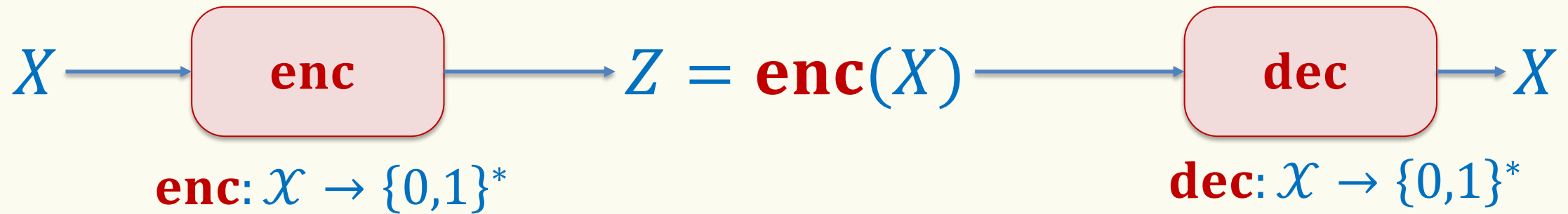
**Example:**  $\mathcal{X} = \{\text{hello, world, cse312}\}$

$$\mathbb{P}_X(\text{hello}) = \frac{1}{2} \quad \mathbb{P}_X(\text{world}) = \frac{1}{4} \quad \mathbb{P}_X(\text{cse312}) = \frac{1}{4}$$



# The Data Compression Problem

Data = random variable  $X$  over alphabet  $\mathcal{X}$



Two goals:

1. **Decodability.** For all values  $x \in \mathcal{X}$ :  $\mathbf{dec}(\mathbf{enc}(x)) = x$
2. **Minimal length.** The length  $|Z|$  of  $Z$  should be as small as possible

**More formally:** minimize  $\mathbb{E}(|Z|)$

# Expected Length – Example

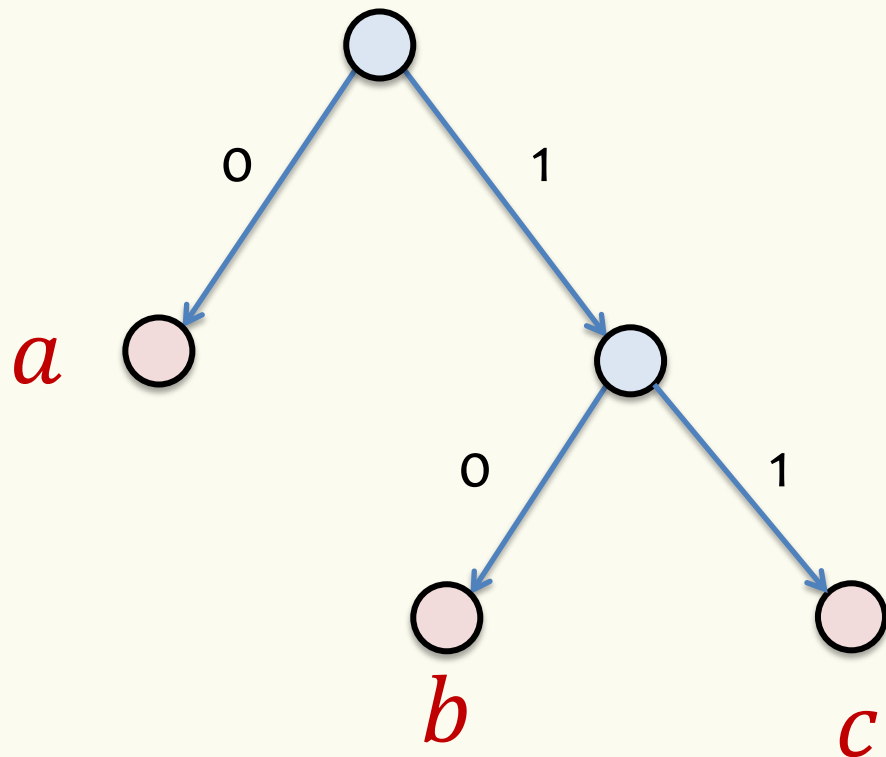
$$\mathcal{X} = \{a, b, c\}$$

$$\mathbb{P}_X(a) = \frac{1}{2} \quad \mathbb{P}_X(b) = \frac{1}{4} \quad \mathbb{P}_X(c) = \frac{1}{4}$$

$$\mathbb{P}_Z(0) = \frac{1}{2}$$

$$\mathbb{P}_Z(10) = \frac{1}{4}$$

$$\mathbb{P}_Z(11) = \frac{1}{4}$$



$$\mathbb{E}(|Z|) = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{4} \cdot 2 = \frac{3}{2}$$

# Expected Length – Example

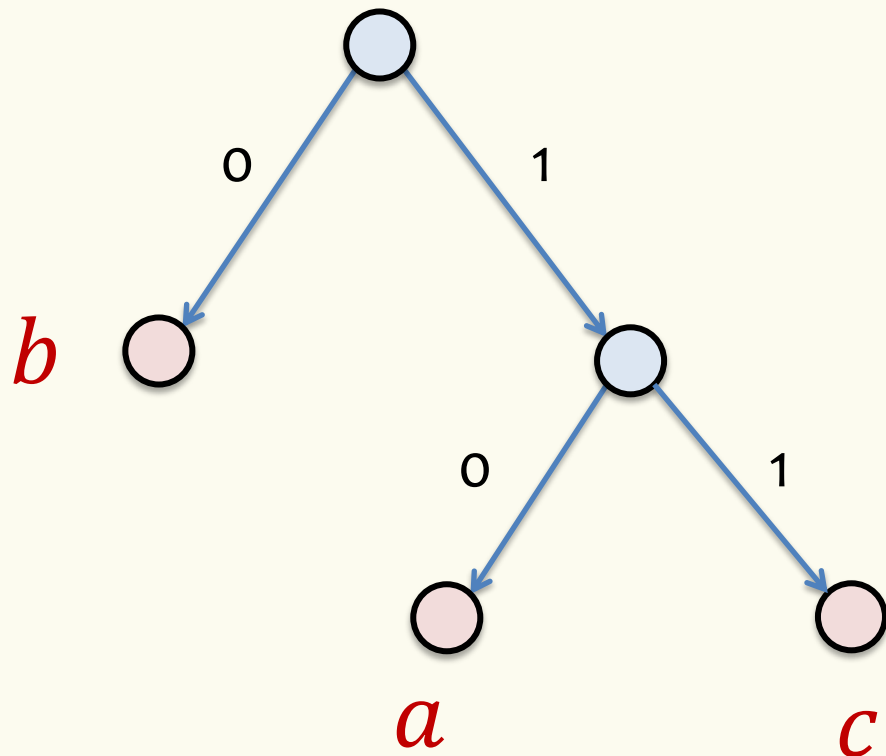
$$\mathcal{X} = \{a, b, c\}$$

$$\mathbb{P}_X(a) = \frac{1}{2} \quad \mathbb{P}_X(b) = \frac{1}{4} \quad \mathbb{P}_X(c) = \frac{1}{4}$$

$$\mathbb{P}_Z(0) = \frac{1}{4}$$

$$\mathbb{P}_Z(10) = \frac{1}{2}$$

$$\mathbb{P}_Z(11) = \frac{1}{4}$$



$$\mathbb{E}(|Z|) = \frac{1}{4} \cdot 1 + \frac{1}{2} \cdot 2 + \frac{1}{4} \cdot 2 = \frac{7}{4}$$

## What is the shortest encoding?

**Problem.** Given a random variable  $X$ , find optimal (**enc, dec**), i.e.,  $\mathbb{E}(|\mathbf{enc}(X)|)$  is as small as possible.

Next: There is an inherent limit on how short the encoding can be (in expectation).

# Random Variables – Arbitrary Values

Assume you are given a random variable  $X$  with the following PMF:

$x$	$a$	$b$	$c$	$d$
$\mathbb{P}_X(x)$	$\frac{15}{16}$	$\frac{1}{32}$	$\frac{1}{64}$	$\frac{1}{64}$

You learn  $X = a$ ; surprised?

$$s(a) = \log_2 16/15 \approx 0.09$$

You learn  $X = d$ ; surprised?

$$s(d) = 6$$

**Definition.** The **surprise** of outcome  $x$  is  $s(x) = \log_2 \left( \frac{1}{\mathbb{P}_X(x)} \right)$

# Entropy = Expected Surprise

**Definition.** The **entropy** of a discrete RV  $X$  over alphabet  $\mathcal{X}$  is

$$\mathbb{H}(X) = \mathbb{E}(s(X)) = \sum_{x \in \mathcal{X}} \mathbb{P}_X(x) \cdot \log_2 \left( \frac{1}{\mathbb{P}_X(x)} \right)$$

**Weird convention:**  $0 \log_2 1/0 = 0$

Intuitively: Captures how surprising outcome of random variable is.

# Entropy = Expected Surprise

**Definition** The entropy of a discrete RV  $X$  over alphabet  $\mathcal{X}$  is

$$\mathbb{H}(X) = \mathbb{E}(s(X)) = \sum_{x \in \mathcal{X}} \mathbb{P}_X(x) \cdot \log_2 \left( \frac{1}{\mathbb{P}_X(x)} \right)$$

$x$	$a$	$b$	$c$	$d$
$\mathbb{P}_X(x)$	$\frac{15}{16}$	$\frac{1}{32}$	$\frac{1}{64}$	$\frac{1}{64}$

$$\begin{aligned} \mathbb{H}(X) &= \frac{15}{16} \cdot \log_2 \frac{16}{15} + \frac{1}{32} \cdot 5 + \frac{1}{64} \cdot 6 + \frac{1}{64} \cdot 6 \\ &= \frac{15}{16} \log_2 \frac{16}{15} + \frac{11}{32} \approx 0.431 \dots \end{aligned}$$

# Entropy = Expected Surprise

**Definition.** The **entropy** of a discrete RV  $X$  over alphabet  $\mathcal{X}$  is

$$\mathbb{H}(X) = \mathbb{E}(s(X)) = \sum_{x \in \mathcal{X}} \mathbb{P}_X(x) \cdot \log_2 \left( \frac{1}{\mathbb{P}_X(x)} \right)$$

$x$	$a$	$b$	$c$	$d$
$\mathbb{P}_X(x)$	1	0	0	0

$$\mathbb{H}(X) = 1 \cdot 0 + 3 \cdot 0 \log_2 \frac{1}{0} = 0$$

$x$	$a$	$b$	$c$	$d$
$\mathbb{P}_X(x)$	1/4	1/4	1/4	1/4

$$\mathbb{H}(X) = 4 \cdot \frac{1}{4} \log_2(4) = 2$$



# Entropy = Expected Surprise

**Definition** The **entropy** of a discrete RV  $X$  over alphabet  $\mathcal{X}$  is

$$\mathbb{H}(X) = \mathbb{E}(s(X)) = \sum_{x \in \mathcal{X}} \mathbb{P}_X(x) \cdot \log_2 \left( \frac{1}{\mathbb{P}_X(x)} \right)$$

**Proposition.**  $0 \leq \mathbb{H}(X) \leq \log_2 |\mathcal{X}|$

Uniform distribution

Takes one value with prob 1

# Shannon's Source Coding Theorem

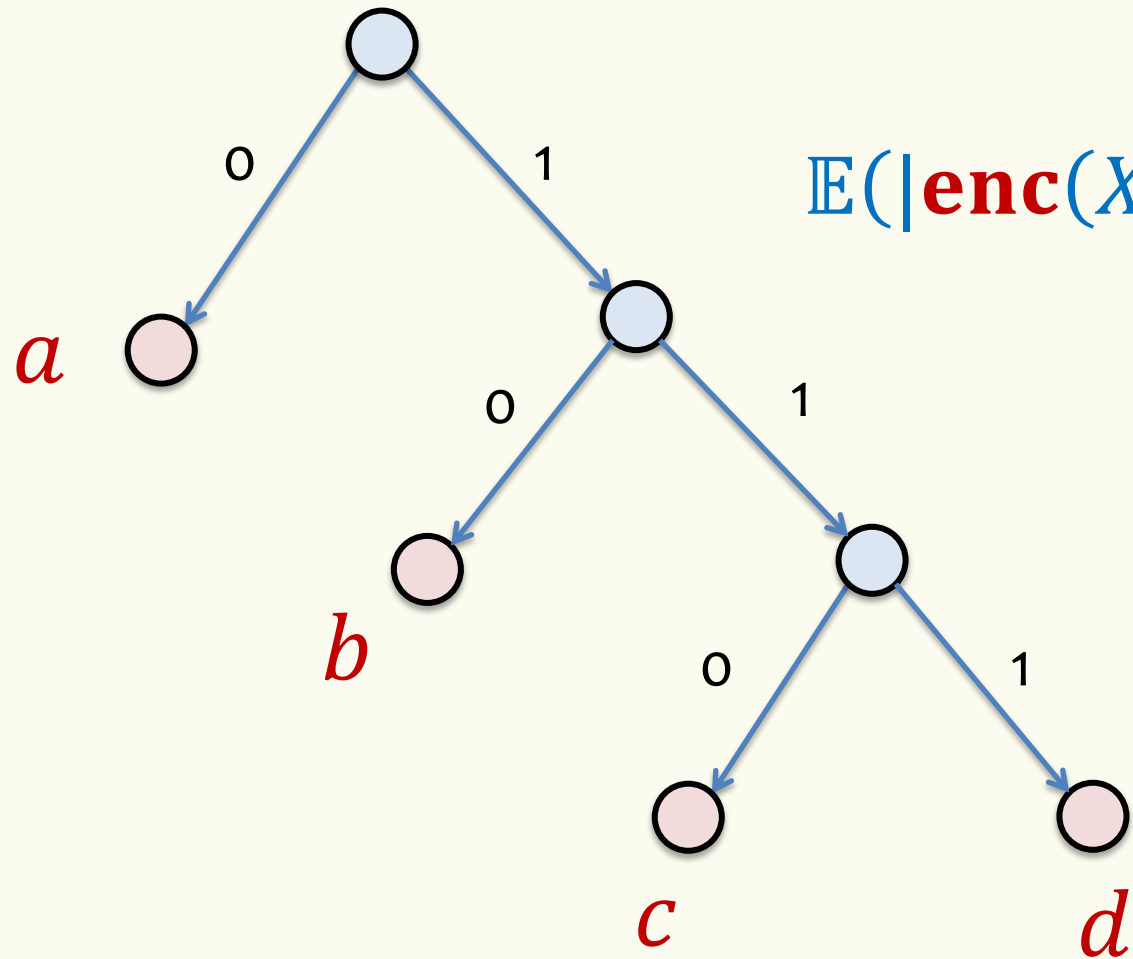
**Theorem. (Source Coding Theorem)** Let  $(\mathbf{enc}, \mathbf{dec})$  be an optimal prefix-free encoding scheme for a RV  $X$ , then

$$\mathbb{H}(X) \leq \mathbb{E}(|\mathbf{enc}(X)|) \leq \mathbb{H}(X) + 1$$

- We cannot compress beyond the entropy
  - Corollary: "uniform" data cannot be compressed
- We can get within one bit of it.
- Example of optimal code: Huffman Code (CSE 143?)
- Result can be extended to uniquely decodable codes. (E.g., suffix free)

# Example

$x$	$a$	$b$	$c$	$d$
$\mathbb{P}_X(x)$	$\frac{15}{16}$	$\frac{1}{32}$	$\frac{1}{64}$	$\frac{1}{64}$



$$\begin{aligned}\mathbb{E}(|\mathbf{enc}(X)|) &= \frac{15}{16} \cdot 1 + \frac{1}{32} \cdot 2 + 2 \cdot \frac{1}{64} \cdot 3 \\ &= \frac{15}{16} + \frac{10}{64} = \frac{70}{64} \leq \mathbb{H}(X) + 1\end{aligned}$$

# Data Compression in the Real World

**Main issue:** we do not know the distribution of  $X$

- Universal compression: Lempel/Ziv/Welch
  - See <http://web.mit.edu/6.02/www/f2011/handouts/3.pdf>
  - Used in GIF, UNIX compress.
  - General idea: Assume data is sequence of symbols generated from a random process to be “estimated”.
- Whole area of computer science dedicated to the topic.
- This is lossless compression, very different from “lossy compression” used in images, videos, audio etc.
  - Assumes humans can be “fooled” with some loss of data