**Digital Communications**
**Lectures 2, 3, and 4**

**A Measure of Information, Source Compression**
**(Chapters 2 and 3)**

Christian Häger
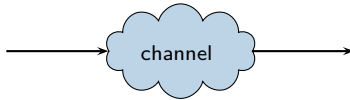Slides prepared by Alexandre Graell i Amat

November 1 and 6, 2023

**CHALMERS**

digital data, voice
movie, audio

satellite link,
fiber, telephone line,
hard-disk drive, DVD

100110101

100011101

### In These Lectures...

- How do we measure information
- How can we mathematically describe an information source
- Data compression

# Discrete Information Source

## Discrete Information Source

The information source generates a sequence of symbols that take values on a finite alphabet $\mathcal{X}$.

# Discrete Information Source

## Discrete Information Source

The information source generates a sequence of symbols that take values on a finite alphabet $\mathcal{X}$.

## A Mathematical Model

The output of an information source is a sequence of random variables $X^{(1)} X^{(2)} X^{(3)} \ldots$, at times $t = 1, 2, 3, \ldots$,

# Discrete Information Source

## Discrete Information Source

The information source generates a sequence of symbols that take values on a finite alphabet $\mathcal{X}$.

## A Mathematical Model

The output of an information source is a sequence of random variables $X^{(1)} X^{(2)} X^{(3)} \ldots$, at times $t = 1, 2, 3, \ldots,$

- $X^{(i)} \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, $|\mathcal{X}| = M$

# Discrete Information Source

## Discrete Information Source

The information source generates a sequence of symbols that take values on a finite alphabet $\mathcal{X}$.

## A Mathematical Model

The output of an information source is a sequence of random variables $X^{(1)} X^{(2)} X^{(3)} \ldots$, at times $t = 1, 2, 3, \ldots$,

- $X^{(i)} \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, $|\mathcal{X}| = M$
- Probability mass function (PMF) $P_X(x)$; $\Pr(X = x_i) = P_X(x_i) = p_i$,

$$p_i \geq 0 \ \forall i, \quad \text{and} \quad \sum_{x_i \in \mathcal{X}} p_i = 1.$$

# Discrete Information Source

## Discrete Information Source

The information source generates a sequence of symbols that take values on a finite alphabet $\mathcal{X}$.

## A Mathematical Model

The output of an information source is a sequence of random variables $X^{(1)} X^{(2)} X^{(3)} \ldots$, at times $t = 1, 2, 3, \ldots$,

- $X^{(i)} \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, $|\mathcal{X}| = M$
- Probability mass function (PMF) $P_X(x)$; $\Pr(X = x_i) = P_X(x_i) = p_i$,

$$p_i \geq 0 \ \forall i, \quad \text{and} \quad \sum_{x_i \in \mathcal{X}} p_i = 1.$$

The source is completely described by $\mathcal{X}$ and $P_X(x)$.

# Discrete Information Source

## Discrete Information Source

The information source generates a sequence of symbols that take values on a finite alphabet $\mathcal{X}$.

## A Mathematical Model

The output of an information source is a sequence of random variables $X^{(1)} X^{(2)} X^{(3)} \ldots$, at times $t = 1, 2, 3, \ldots$,

- $X^{(i)} \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, $|\mathcal{X}| = M$
- Probability mass function (PMF) $P_X(x)$; $\Pr(X = x_i) = P_X(x_i) = p_i$,

$$p_i \geq 0 \ \forall i, \quad \text{and} \quad \sum_{x_i \in \mathcal{X}} p_i = 1.$$

The source is completely described by $\mathcal{X}$ and $P_X(x)$.

## In This Course...

Discrete memoryless sources, i.e., sources where the random variables $X^{(i)}$ are independent of each other.

# A Measure of Information

- How can we measure the information content of a <span style="color:red">particular outcome</span> $X = x_i$ from the source?

- How much information does a <span style="color:red">source</span> contain in <span style="color:red">average</span>?

# A Measure of Information

> Example: Consider the following two claims (about Göteborg's weather)
>
> - This weekend it will rain
> - This weekend there will be more than $30°$C
>
> Which sentence conveys more information?

# A Measure of Information

**Example:** Consider the following two claims (about Göteborg's weather)

- This weekend it will rain
- This weekend there will be more than $30°$C

Which sentence conveys more information?

# A Measure of Information

> **Example: Consider the following two claims (about Göteborg's weather)**
>
> - This weekend it will rain
> - This weekend there will be more than $30°$C
>
> Which sentence conveys more information?

> **Discrete Information Source**
>
> $\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$
> Which output conveys more information, $x_1$ or $x_M$?

# A Measure of Information

**Example: Consider the following two claims (about Göteborg's weather)**

- This weekend it will rain
- This weekend there will be more than $30°$C

Which sentence conveys more information?

**Discrete Information Source**

$\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$
Which output conveys more information, $x_1$ or $x_M$?

# A Measure of Information

**Example: Consider the following two claims (about Göteborg's weather)**

- This weekend it will rain
- This weekend there will be more than $30°$C

Which sentence conveys more information?

**Discrete Information Source**

$\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$
Which output conveys more information, $x_1$ or $x_M$?

The information associated to a particular message $x$ is related to its probability: The less probable the message is, the more information it conveys!

# A Measure of Information

A good measure of information $i(x)$, associated to a symbol $x$ should:

P: Be a decreasing function of the probability of the symbol, i.e.,

$$i(x_i) > i(x_j) \quad \text{if} \quad p_i < p_j$$

# A Measure of Information

Example: Consider the following three claims

- This weekend it will rain
- This weekend there will be more than $30°$C
- This weekend there will be more that $30.2°$C

Does claim 3 convey more/less information than claim 2?

# A Measure of Information

Example: Consider the following three claims

- This weekend it will rain
- This weekend there will be more than $30°$C
- This weekend there will be more that $30.2°$C

Does claim 3 convey more/less information than claim 2? A little bit more, but roughly the same!

# A Measure of Information

Example: Consider the following three claims

- This weekend it will rain
- This weekend there will be more than $30°$C
- This weekend there will be more that $30.2°$C

Does claim 3 convey more/less information than claim 2? A little bit more, but roughly the same!

Discrete Information Source

$\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$, and $p_i = 0.25$ and $p_j = 0.25001$

# A Measure of Information

Example: Consider the following three claims

- This weekend it will rain
- This weekend there will be more than $30°$C
- This weekend there will be more that $30.2°$C

Does claim 3 convey more/less information than claim 2? A little bit more, but roughly the same!

### Discrete Information Source

$\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$, and $p_i = 0.25$ and $p_j = 0.25001$

We would expect that $\mathsf{i}(x_i) \gtrapprox \mathsf{i}(x_j)$

# A Measure of Information

**Example: Consider the following three claims**

- This weekend it will rain
- This weekend there will be more than $30°$C
- This weekend there will be more that $30.2°$C

Does claim 3 convey more/less information than claim 2? A little bit more, but roughly the same!

**Discrete Information Source**

$\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$, and $p_i = 0.25$ and $p_j = 0.25001$

We would expect that $\mathsf{i}(x_i) \gtreqless \mathsf{i}(x_j)$

- A slight change in the probability of a symbol should only cause a slight change in the information it conveys

# A Measure of Information

Example: Consider the following three claims

- This weekend it will rain
- This weekend there will be more than $30°$C
- This weekend there will be more that $30.2°$C

Does claim 3 convey more/less information than claim 2? A little bit more, but roughly the same!

## Discrete Information Source

$\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$, and $p_i = 0.25$ and $p_j = 0.25001$

We would expect that $\mathsf{i}(x_i) \gtrapprox \mathsf{i}(x_j)$

- A slight change in the probability of a symbol should only cause a slight change in the information it conveys $\longrightarrow$ A measure of information should be a continuous function of the probability

# A Measure of Information

- This weekend it will rain
- This weekend there will be more than $30°$C
- This weekend there will be more that $30.2°$C

Does claim 3 convey more/less information than claim 2? A little bit more, but roughly the same!

## Discrete Information Source

$\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$, and $p_i = 0.25$ and $p_j = 0.25001$
We would expect that $\mathsf{i}(x_i) \gtrapprox \mathsf{i}(x_j)$

- A slight change in the probability of a symbol should only cause a slight change in the information it conveys $\longrightarrow$ A measure of information should be a continuous function of the probability
- Two equiprobable symbols carry the same amount of information

# A Measure of Information

**Example: Consider the following three claims**

- This weekend it will rain
- This weekend there will be more than $30°$C
- This weekend there will be more that $30.2°$C

Does claim 3 convey more/less information than claim 2? A little bit more, but roughly the same!

**Discrete Information Source**

$\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$, with $p_1 \geq p_2 \ldots \geq p_M$, and $p_i = 0.25$ and $p_j = 0.25001$

We would expect that $\mathsf{i}(x_i) \gtrapprox\lessapprox \mathsf{i}(x_j)$

- A slight change in the probability of a symbol should only cause a slight change in the information it conveys $\longrightarrow$ A measure of information should be a continuous function of the probability
- Two equiprobable symbols carry the same amount of information $\longrightarrow$ A measure of information should depend only on the probability of occurrence

# A Measure of Information

A good measure of information $i(x)$ associated to a symbol $x$ should:

P1: Not depend on the symbol itself, but only on its probability,

$$i(x_i) = i(x_j) \quad \text{if} \quad p_i = p_j$$

P2: Be a continuous, decreasing function of the probability of the symbol,

$$i(x_i) > i(x_j) \quad \text{if} \quad p_i < p_j$$

# A Measure of Information

# A Measure of Information

## Example: Rolling a dice

Information conveyed by the outcome of rolling a dice once: $i_1$
How much information is conveyed by the outcome of rolling a dice twice?

# A Measure of Information

## Example: Rolling a dice

Information conveyed by the outcome of rolling a dice once: $i_1$
How much information is conveyed by the outcome of rolling a dice twice?
$i_2 = 2i_1$

# A Measure of Information

## Example: Rolling a dice

Information conveyed by the outcome of rolling a dice once: $i_1$
How much information is conveyed by the outcome of rolling a dice twice?
$i_2 = 2i_1$

For two independent random variables $X$ and $Y$, the information we gain when we learn $X$ and $Y$ should equal the sum of the information gained if we learn $X$ and $Y$ separately

# A Measure of Information

> **Example: Rolling a dice**
>
> Information conveyed by the outcome of rolling a dice once: $i_1$
> How much information is conveyed by the outcome of rolling a dice twice?
> $i_2 = 2i_1$

For two independent random variables $X$ and $Y$, the information we gain when we learn $X$ and $Y$ should equal the sum of the information gained if we learn $X$ and $Y$ separately $\longrightarrow$ The information should be additive

# A Measure of Information

A good measure of information $i(x)$ associated to a symbol $x$ should satisfy the following three properties:

# A Measure of Information

A good measure of information $i(x)$ associated to a symbol $x$ should satisfy the following three properties:

P1: Not depend on the symbol itself, but only on its probability,

$$i(x_i) = i(x_j) \quad \text{if} \quad p_i = p_j$$

# A Measure of Information

A good measure of information $i(x)$ associated to a symbol $x$ should satisfy the following three properties:

P1: Not depend on the symbol itself, but only on its probability,

$$i(x_i) = i(x_j) \quad \text{if} \quad p_i = p_j$$

P2: Be a continuous, decreasing function of the probability of the symbol,

$$i(x_i) > i(x_j) \quad \text{if} \quad p_i < p_j$$

# A Measure of Information

A good measure of information $i(x)$ associated to a symbol $x$ should satisfy the following three properties:

P1: Not depend on the symbol itself, but only on its probability,

$$i(x_i) = i(x_j) \quad \text{if} \quad p_i = p_j$$

P2: Be a continuous, decreasing function of the probability of the symbol,

$$i(x_i) > i(x_j) \quad \text{if} \quad p_i < p_j$$

P3: For two independent symbols $x_i$ and $x_j$,

$$i(x_i, x_j) = i(x_i) + i(x_j)$$

# Shannon's Information Measure

> **Definition (Shannon's Information Content)**
>
> The Shannon information content of the outcome $X = x_i$ is
>
> $$\mathrm{i}(x_i) = \log_a \frac{1}{p_i}.$$
>
> The base $a$ of the logarithm determines the unit of information. If $a = 2$, the unit of information is the bit. If the base is $e$, then the information unit is called nat.

# Shannon's Information Measure

## Definition (Shannon's Information Content)

The Shannon information content of the outcome $X = x_i$ is

$$\mathrm{i}(x_i) = \log_a \frac{1}{p_i}.$$

The base $a$ of the logarithm determines the unit of information. If $a = 2$, the unit of information is the bit. If the base is $e$, then the information unit is called nat.

## Properties

# Shannon's Information Measure

## Definition (Shannon's Information Content)

The Shannon information content of the outcome $X = x_i$ is

$$i(x_i) = \log_a \frac{1}{p_i}.$$

The base $a$ of the logarithm determines the unit of information. If $a = 2$, the unit of information is the bit. If the base is $e$, then the information unit is called nat.

## Properties

1. ✓

# Shannon's Information Measure

---

**Definition (Shannon's Information Content)**

The Shannon information content of the outcome $X = x_i$ is

$$\mathrm{i}(x_i) = \log_a \frac{1}{p_i}.$$

The base $a$ of the logarithm determines the unit of information. If $a = 2$, the unit of information is the bit. If the base is $e$, then the information unit is called nat.

---

**Properties**

1. ✓
2. ✓

---

# Shannon's Information Measure

## Definition (Shannon's Information Content)

The Shannon information content of the outcome $X = x_i$ is

$$\mathrm{i}(x_i) = \log_a \frac{1}{p_i}.$$

The base $a$ of the logarithm determines the unit of information. If $a = 2$, the unit of information is the bit. If the base is $e$, then the information unit is called nat.

## Properties

1. ✓
2. ✓
3. For two independent symbols $x$ and $y$,

$$\mathrm{i}(x, y) = \log \frac{1}{P(x,y)} =$$

# Shannon's Information Measure

## Definition (Shannon's Information Content)

The Shannon information content of the outcome $X = x_i$ is

$$\mathsf{i}(x_i) = \log_a \frac{1}{p_i}.$$

The base $a$ of the logarithm determines the unit of information. If $a = 2$, the unit of information is the bit. If the base is $e$, then the information unit is called nat.

## Properties

1. ✓
2. ✓
3. For two independent symbols $x$ and $y$,

$$\mathsf{i}(x, y) = \log \frac{1}{P(x, y)} = \log \frac{1}{P(x)P(y)}$$

# Shannon's Information Measure

## Definition (Shannon's Information Content)

The Shannon information content of the outcome $X = x_i$ is

$$\mathrm{i}(x_i) = \log_a \frac{1}{p_i}.$$

The base $a$ of the logarithm determines the unit of information. If $a = 2$, the unit of information is the bit. If the base is $e$, then the information unit is called nat.

## Properties

1. ✓
2. ✓
3. For two independent symbols $x$ and $y$,

$$\mathrm{i}(x,y) = \log \frac{1}{P(x,y)} = \log \frac{1}{P(x)P(y)} = \log \frac{1}{P(x)} + \log \frac{1}{P(y)}$$

# Shannon's Information Measure

## Definition (Shannon's Information Content)

The Shannon information content of the outcome $X = x_i$ is

$$\text{i}(x_i) = \log_a \frac{1}{p_i}.$$

The base $a$ of the logarithm determines the unit of information. If $a = 2$, the unit of information is the bit. If the base is $e$, then the information unit is called nat.

## Properties

1. ✓
2. ✓
3. For two independent symbols $x$ and $y$,

$$\text{i}(x,y) = \log \frac{1}{P(x,y)} = \log \frac{1}{P(x)P(y)} = \log \frac{1}{P(x)} + \log \frac{1}{P(y)}$$
$$= \text{i}(x) + \text{i}(y).$$

# Information Content of a Source

## Average Information Content of a Source

A complete characterization of the source can then be obtained by defining the average information content,

$$\mathsf{H}(X) = \sum_{i=1}^{M} p_i \mathsf{i}(x_i) = \sum_{i=1}^{M} p_i \log \frac{1}{p_i}.$$

# Average Information Content of a Source

A complete characterization of the source can then be obtained by defining the average information content,

$$\mathsf{H}(X) = \sum_{i=1}^{M} p_i \mathsf{i}(x_i) = \sum_{i=1}^{M} p_i \log \frac{1}{p_i}.$$

---

**Definition (Entropy)**

The entropy of a random variable (random symbol) $X$ that takes values on the alphabet $\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$ with probabilities $p_1, p_2, \ldots, p_M$ is defined as

$$\mathsf{H}(X) = \sum_{x \in \mathcal{X}} P(x) \log \frac{1}{P(x)} = \sum_{i=1}^{M} p_i \log \frac{1}{p_i}.$$

---

# Average Information Content of a Source

A complete characterization of the source can then be obtained by defining the average information content,

$$\mathsf{H}(X) = \sum_{i=1}^{M} p_i \mathsf{i}(x_i) = \sum_{i=1}^{M} p_i \log \frac{1}{p_i}.$$

> **Definition (Entropy)**
>
> The entropy of a random variable (random symbol) $X$ that takes values on the alphabet $\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$ with probabilities $p_1, p_2, \ldots, p_M$ is defined as
>
> $$\mathsf{H}(X) = \sum_{x \in \mathcal{X}} P(x) \log \frac{1}{P(x)} = \sum_{i=1}^{M} p_i \log \frac{1}{p_i}.$$

The entropy measures the average information content (or *uncertainty*) of $X$.

# Average Information Content of a Source

> **Definition (The Binary Entropy Function)**
>
> Let $X$ be a binary source with two possible values $\mathcal{X} = \{x_1, x_2\}$ such that $P(x_1) = p$ and $P(x_2) = 1 - p$. Then
>
> $$\mathsf{H}(X) = \mathsf{H_b}(p),$$
>
> where $\mathsf{H_b}(p)$ is called the binary entropy function,
>
> $$\mathsf{H_b}(p) \triangleq p \log \frac{1}{p} + (1 - p) \log \frac{1}{1 - p}.$$
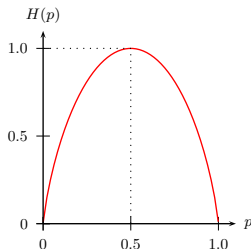
# Average Information Content of a Source



Figure: The binary entropy function as a function of the probability $p$.

**Definition (The Binary Entropy Function)**

Let $X$ be a binary source with two possible values $\mathcal{X} = \{x_1, x_2\}$ such that $P(x_1) = p$ and $P(x_2) = 1 - p$. Then

$$H(X) = \mathsf{H_b}(p),$$

where $\mathsf{H_b}(p)$ is called the binary entropy function,

$$\mathsf{H_b}(p) \triangleq p \log \frac{1}{p} + (1 - p) \log \frac{1}{1 - p}.$$

# Average Information Content of a Source

### Lemma

*The entropy of a random variable $X$ that takes values on the alphabet $\mathcal{X} = \{x_1, \ldots, x_M\}$, $|\mathcal{X}| = M$, is bounded by*

$$0 \leq \mathsf{H}(X) \leq \log M \quad \textit{bits},$$

# Average Information Content of a Source

### Lemma

*The entropy of a random variable $X$ that takes values on the alphabet $\mathcal{X} = \{x_1, \ldots, x_M\}$, $|\mathcal{X}| = M$, is bounded by*

$$0 \leq \mathsf{H}(X) \leq \log M \quad \text{bits,}$$

*where*

$$\mathsf{H}(X) = 0 \qquad \text{if and only if } p_i = 1 \text{ for some } i,$$
$$\mathsf{H}(X) = \log M \qquad \text{if and only if } p_i = \frac{1}{M} \ \ \forall i.$$

$\mathsf{H}(X) \geq 0$

$\mathsf{H}(X) \leq \log M$     (hint: use $\ln x \leq x - 1$)

# Data Compression

# Data Compression

Wh_ c_n y_u r_co_er this q_est_on?

# Data Compression

## Wh_ c_n y_u r_co_er this q_est_on?

Most information sources contain <span style="color:red">redundancy</span>!

# Data Compression

## Wh_ c_n y_u r_co_er this q_est_on?

Most information sources contain <span style="color:red">redundancy</span>!

### Data (Source) Compression

Find an <span style="color:red">efficient representation</span> of the output of the source which results in <span style="color:magenta">zero or little redundancy</span>

# Wh_ c_n y_u r_co_er this q_est_on?

Most information sources contain redundancy!

## Data (Source) Compression

Find an efficient representation of the output of the source which results in zero or little redundancy $\longrightarrow$ Encode (map) the symbols (or messages) at the output of the source to a sequence of symbols, called codeword, that is shorter in average.

# Wh_ c_n y_u r_co_er this q_est_on?

Most information sources contain redundancy!

## Data (Source) Compression

Find an efficient representation of the output of the source which results in zero or little redundancy $\longrightarrow$ Encode (map) the symbols (or messages) at the output of the source to a sequence of symbols, called codeword, that is shorter in average.

## Source Encoder

- Can encode source symbols separately or groups of symbols.

# Wh_ c_n y_u r_co_er this q_est_on?

Most information sources contain redundancy!

## Data (Source) Compression

Find an efficient representation of the output of the source which results in zero or little redundancy ⟶ Encode (map) the symbols (or messages) at the output of the source to a sequence of symbols, called codeword, that is shorter in average.

## Source Encoder

- Can encode source symbols separately or groups of symbols.
- Typically requires knowledge of the statistics of the source

# Wh_ c_n y_u r_co_er this q_est_on?

Most information sources contain redundancy!

## Data (Source) Compression

Find an efficient representation of the output of the source which results in zero or little redundancy $\longrightarrow$ Encode (map) the symbols (or messages) at the output of the source to a sequence of symbols, called codeword, that is shorter in average.

## Source Encoder

- Can encode source symbols separately or groups of symbols.
- Typically requires knowledge of the statistics of the source
- We can achieve compression on average by assigning shorter codewords to more probable symbols and longer codewords to less likely ones: Variable-length encoder

# Data Compression

## Example: The Morse Code

The Morse code assigns the most frequent letters to shorter codewords, and less frequent letters to longer codewords:

- "e" $\longrightarrow$ "."
- "q" $\longrightarrow$ "$- - \cdot -$"

# Data Compression

- How much can we compress such that we lose no information?

# Data Compression

- How much can we compress such that we lose no information?
- Are there efficient methods to assign codewords to source symbols?

# Data Compression

- How much can we compress such that we lose no information?
- Are there efficient methods to assign codewords to source symbols?
- How can we make sure that the source code is easy to decode?

# Data Compression

| symbol | probability | $\mathcal{C}_1$ | $\mathcal{C}_2$ | $\mathcal{C}_3$ | $\mathcal{C}_4$ | $\mathcal{C}_5$ | $\mathcal{C}_6$ |
|--------|-------------|------|------|------|------|------|------|
| $a$ | $1/2$ | 000 | 00 | 00 | 0 | 0 | 0 |
| $b$ | $1/4$ | 001 | 00 | 01 | 01 | 01 | 10 |
| $c$ | $1/8$ | 101 | 10 | 10 | 001 | 011 | 110 |
| $d$ | $1/8$ | 111 | 11 | 11 | 111 | 111 | 111 |

Example: Alphabet $\{a, b, c, d\}$ with probabilities $\{1/2, 1/4, 1/8, 1/8\}$

We would like to design a good (binary) source code to compress the language produced by this alphabet.

# Data Compression

| symbol | probability | $\mathcal{C}_1$ | $\mathcal{C}_2$ | $\mathcal{C}_3$ | $\mathcal{C}_4$ | $\mathcal{C}_5$ | $\mathcal{C}_6$ |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| $a$ | $1/2$ | 000 | 00 | 00 | 0 | 0 | 0 |
| $b$ | $1/4$ | 001 | 00 | 01 | 01 | 01 | 10 |
| $c$ | $1/8$ | 101 | 10 | 10 | 001 | 011 | 110 |
| $d$ | $1/8$ | 111 | 11 | 11 | 111 | 111 | 111 |

Example: Alphabet $\{a, b, c, d\}$ with probabilities $\{1/2, 1/4, 1/8, 1/8\}$

We would like to design a good (binary) source code to compress the language produced by this alphabet.

$\mathcal{C}_1$ Not efficient

# Data Compression

| symbol | probability | $\mathcal{C}_1$ | $\mathcal{C}_2$ | $\mathcal{C}_3$ | $\mathcal{C}_4$ | $\mathcal{C}_5$ | $\mathcal{C}_6$ |
|:------:|:-----------:|:---:|:--:|:--:|:---:|:---:|:---:|
| $a$ | $1/2$ | 000 | 00 | 00 | 0 | 0 | 0 |
| $b$ | $1/4$ | 001 | 00 | 01 | 01 | 01 | 10 |
| $c$ | $1/8$ | 101 | 10 | 10 | 001 | 011 | 110 |
| $d$ | $1/8$ | 111 | 11 | 11 | 111 | 111 | 111 |

Example: Alphabet $\{a, b, c, d\}$ with probabilities $\{1/2, 1/4, 1/8, 1/8\}$

We would like to design a good (binary) source code to compress the language produced by this alphabet.

$\mathcal{C}_1$ Not efficient

$\mathcal{C}_2$ Not uniquely decodable

# Data Compression

| symbol | probability | $\mathcal{C}_1$ | $\mathcal{C}_2$ | $\mathcal{C}_3$ | $\mathcal{C}_4$ | $\mathcal{C}_5$ | $\mathcal{C}_6$ |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| $a$ | 1/2 | 000 | 00 | 00 | 0 | 0 | 0 |
| $b$ | 1/4 | 001 | 00 | 01 | 01 | 01 | 10 |
| $c$ | 1/8 | 101 | 10 | 10 | 001 | 011 | 110 |
| $d$ | 1/8 | 111 | 11 | 11 | 111 | 111 | 111 |

Example: Alphabet $\{a, b, c, d\}$ with probabilities $\{1/2, 1/4, 1/8, 1/8\}$

We would like to design a good (binary) source code to compress the language produced by this alphabet.

$\mathcal{C}_1$ Not efficient

$\mathcal{C}_2$ Not uniquely decodable

$\mathcal{C}_3$ Uniquely decodable, not efficient

# Data Compression

| symbol | probability | $\mathcal{C}_1$ | $\mathcal{C}_2$ | $\mathcal{C}_3$ | $\mathcal{C}_4$ | $\mathcal{C}_5$ | $\mathcal{C}_6$ |
|:------:|:-----------:|:---:|:---:|:---:|:---:|:---:|:---:|
| $a$ | $1/2$ | 000 | 00 | 00 | 0 | 0 | 0 |
| $b$ | $1/4$ | 001 | 00 | 01 | 01 | 01 | 10 |
| $c$ | $1/8$ | 101 | 10 | 10 | 001 | 011 | 110 |
| $d$ | $1/8$ | 111 | 11 | 11 | 111 | 111 | 111 |

Example: Alphabet $\{a, b, c, d\}$ with probabilities $\{1/2, 1/4, 1/8, 1/8\}$

We would like to design a good (binary) source code to compress the language produced by this alphabet.

$\mathcal{C}_1$ Not efficient

$\mathcal{C}_2$ Not uniquely decodable

$\mathcal{C}_3$ Uniquely decodable, not efficient

$\mathcal{C}_4$ Not uniquely decodable

# Data Compression

| symbol | probability | $\mathcal{C}_1$ | $\mathcal{C}_2$ | $\mathcal{C}_3$ | $\mathcal{C}_4$ | $\mathcal{C}_5$ | $\mathcal{C}_6$ |
|--------|-------------|------|------|------|------|------|------|
| $a$ | $1/2$ | 000 | 00 | 00 | 0 | 0 | 0 |
| $b$ | $1/4$ | 001 | 00 | 01 | 01 | 01 | 10 |
| $c$ | $1/8$ | 101 | 10 | 10 | 001 | 011 | 110 |
| $d$ | $1/8$ | 111 | 11 | 11 | 111 | 111 | 111 |

Example: Alphabet $\{a, b, c, d\}$ with probabilities $\{1/2, 1/4, 1/8, 1/8\}$

We would like to design a good (binary) source code to compress the language produced by this alphabet.

$\mathcal{C}_1$ Not efficient

$\mathcal{C}_2$ Not uniquely decodable

$\mathcal{C}_3$ Uniquely decodable, not efficient

$\mathcal{C}_4$ Not uniquely decodable

$\mathcal{C}_5$ Uniquely decodable, not instantaneous

# Data Compression

| symbol | probability | $\mathcal{C}_1$ | $\mathcal{C}_2$ | $\mathcal{C}_3$ | $\mathcal{C}_4$ | $\mathcal{C}_5$ | $\mathcal{C}_6$ |
|--------|-------------|------|------|------|------|------|------|
| $a$ | $1/2$ | 000 | 00 | 00 | 0 | 0 | 0 |
| $b$ | $1/4$ | 001 | 00 | 01 | 01 | 01 | 10 |
| $c$ | $1/8$ | 101 | 10 | 10 | 001 | 011 | 110 |
| $d$ | $1/8$ | 111 | 11 | 11 | 111 | 111 | 111 |

Example: Alphabet $\{a, b, c, d\}$ with probabilities $\{1/2, 1/4, 1/8, 1/8\}$

We would like to design a good (binary) source code to compress the language produced by this alphabet.

$\mathcal{C}_1$ Not efficient

$\mathcal{C}_2$ Not uniquely decodable

$\mathcal{C}_3$ Uniquely decodable, not efficient

$\mathcal{C}_4$ Not uniquely decodable

$\mathcal{C}_5$ Uniquely decodable, not instantaneous

$\mathcal{C}_6$ Uniquely decodable, instantaneous (prefix-free)

# Data Compression

A source code should satisfy the following properties:

P1 The code must be uniquely decodable, i.e., the symbols generated by the source must be uniquely retrieved from the encoded string of bits.

# Data Compression

A source code should satisfy the following properties:

P1 The code must be uniquely decodable, i.e., the symbols generated by the source must be uniquely retrieved from the encoded string of bits.

P2 The code should be easy to decode.

# Data Compression

A source code should satisfy the following properties:

P1 The code must be uniquely decodable, i.e., the symbols generated by the source must be uniquely retrieved from the encoded string of bits.

P2 The code should be easy to decode.

P3 The code should compress the source as much as possible.

# Data Compression

### Definition

Given a set $\Sigma$, $\Sigma^+$ denotes the set of all strings over $\Sigma$ of any (nonzero) finite length.

# Data Compression

### Definition

Given a set $\Sigma$, $\Sigma^+$ denotes the set of all strings over $\Sigma$ of any (nonzero) finite length.

### Example

If $\Sigma = \{a, b, c\}$ then $a$, $ab$, $aac$ and $bbac$ are strings over $\Sigma$ of lengths one, two, three and four respectively.

# Data Compression

**Definition (Binary Source Code)**

Consider a random variable $X$ which takes values on $\mathcal{X} = \{x_1, \ldots, x_M\}$.

- A binary encoder $\mathcal{E}$ is a function $\mathcal{E} : \mathcal{X} \to \{0,1\}^+$ that maps each source symbol $x_i \in \mathcal{X}$ to a binary sequence $\boldsymbol{c}_i \in \{0,1\}^+$ (codeword). The number of bits in $\boldsymbol{c}_i$ is called its length and is denoted by $\ell_i$.

# Data Compression

## Definition (Binary Source Code)

Consider a random variable $X$ which takes values on $\mathcal{X} = \{x_1, \ldots, x_M\}$.

- A binary encoder $\mathcal{E}$ is a function $\mathcal{E} : \mathcal{X} \to \{0, 1\}^+$ that maps each source symbol $x_i \in \mathcal{X}$ to a binary sequence $c_i \in \{0, 1\}^+$ (codeword). The number of bits in $c_i$ is called its length and is denoted by $\ell_i$.

- The binary source code $\mathcal{C}$ is the set of all codewords, i.e., $\mathcal{C} = \{c_1, \ldots, c_M\}$.

# Data Compression

## Definition (Binary Source Code)

Consider a random variable $X$ which takes values on $\mathcal{X} = \{x_1, \ldots, x_M\}$.

- A binary encoder $\mathcal{E}$ is a function $\mathcal{E} : \mathcal{X} \to \{0,1\}^+$ that maps each source symbol $x_i \in \mathcal{X}$ to a binary sequence $\boldsymbol{c}_i \in \{0,1\}^+$ (codeword). The number of bits in $\boldsymbol{c}_i$ is called its length and is denoted by $\ell_i$.
- The binary source code $\mathcal{C}$ is the set of all codewords, i.e., $\mathcal{C} = \{\boldsymbol{c}_1, \ldots, \boldsymbol{c}_M\}$.

## Example: Code $\mathcal{C}_6$

| $x_i$ | | $\boldsymbol{c}_i$ | $\ell_i$ |
|-------|-----|------|-----|
| $a$ | $\to$ | 0 | 1 |
| $b$ | $\to$ | 10 | 2 |
| $c$ | $\to$ | 110 | 3 |
| $d$ | $\to$ | 111 | 3 |

# Data Compression

## Definition (Extended Code)

The extended code $\mathcal{C}^+$ is the set of binary sequences resulting from the mapping $\mathcal{E}^+ : \mathcal{X}^+ \to \{0,1\}^+$ from a sequence of symbols in the alphabet $\mathcal{X}$ to a binary sequence obtained by concatenating the corresponding codewords, such that, for a sequence of symbols of length $n$, at times $t = 1, \ldots, n$,

$$\boldsymbol{c}^+(x^{(1)} x^{(2)} \ldots x^{(n)}) = \boldsymbol{c}(x^{(1)}) \boldsymbol{c}(x^{(2)}) \ldots \boldsymbol{c}(x^{(n)}),$$

# Data Compression

## Definition (Extended Code)

The extended code $\mathcal{C}^+$ is the set of binary sequences resulting from the mapping $\mathcal{E}^+ : \mathcal{X}^+ \to \{0,1\}^+$ from a sequence of symbols in the alphabet $\mathcal{X}$ to a binary sequence obtained by concatenating the corresponding codewords, such that, for a sequence of symbols of length $n$, at times $t = 1, \ldots, n$,

$$\boldsymbol{c}^+(x^{(1)}x^{(2)} \ldots x^{(n)}) = \boldsymbol{c}(x^{(1)})\boldsymbol{c}(x^{(2)}) \ldots \boldsymbol{c}(x^{(n)}),$$

## Example: Extended Code $\mathcal{C}_6^+$

$$
\begin{array}{ccc}
aab & \to & 0010 \\
bca & \to & 101100
\end{array}
$$

# Uniquely Decodable Code

P1 The code must be uniquely decodable, i.e., the symbols generated by the source must be uniquely retrieved from the encoded string of bits.

# Uniquely Decodable Code

**P1** The code must be <span style="color:red">uniquely decodable</span>, i.e., the symbols generated by the source must be uniquely retrieved from the encoded string of bits.

---

**Definition (Uniquely Decodable Code)**

A code $\mathcal{C}$ is uniquely decodable if

$$\forall \boldsymbol{x}, \boldsymbol{y} \in \mathcal{X}^{+}, \boldsymbol{x} \neq \boldsymbol{y} \implies \boldsymbol{c}^{+}(\boldsymbol{x}) \neq \boldsymbol{c}^{+}(\boldsymbol{y}).$$

# Uniquely Decodable Code

**P1** The code must be uniquely decodable, i.e., the symbols generated by the source must be uniquely retrieved from the encoded string of bits.

---

### Definition (Uniquely Decodable Code)

A code $\mathcal{C}$ is uniquely decodable if

$$\forall \boldsymbol{x}, \boldsymbol{y} \in \mathcal{X}^+, \boldsymbol{x} \neq \boldsymbol{y} \implies \boldsymbol{c}^+(\boldsymbol{x}) \neq \boldsymbol{c}^+(\boldsymbol{y}).$$

---

### Example

| symbol | $\mathcal{C}_4$ | $\mathcal{C}_5$ | $\mathcal{C}_6$ |
|:------:|:---:|:---:|:---:|
| $a$ | 0 | 0 | 0 |
| $b$ | 01 | 01 | 10 |
| $c$ | 001 | 011 | 110 |
| $d$ | 111 | 111 | 111 |

$\mathcal{C}_4$ is not uniquely decodable: $ab \rightarrow 001$ and $c \rightarrow 001$.

# Prefix-Free Code

P2 The code should be easy to decode

# Prefix-Free Code

P2 The code should be easy to decode (decode each codeword immediately at the arrival of the last bit of that codeword)

# Prefix-Free Code

P2 The code should be easy to decode (decode each codeword immediately at the arrival of the last bit of that codeword) $\longrightarrow$ no codeword can be a prefix of another codeword.

# Prefix-Free Code

P2 The code should be easy to decode (decode each codeword immediately at the arrival of the last bit of that codeword) $\longrightarrow$ no codeword can be a prefix of another codeword.

### Definition (Prefix-Free Code)

A code is called a prefix-free code if no codeword is the prefix of any other codeword.

# Prefix-Free Code

P2 The code should be easy to decode (decode each codeword immediately at the arrival of the last bit of that codeword) $\longrightarrow$ no codeword can be a prefix of another codeword.

## Definition (Prefix-Free Code)

A code is called a prefix-free code if no codeword is the prefix of any other codeword.

- We can decode codewords instantaneously, i.e., as soon as they are received

# Prefix-Free Code

P2  The code should be easy to decode (decode each codeword immediately at the arrival of the last bit of that codeword) $\longrightarrow$ no codeword can be a prefix of another codeword.

### Definition (Prefix-Free Code)

A code is called a prefix-free code if no codeword is the prefix of any other codeword.

- We can decode codewords instantaneously, i.e., as soon as they are received $\longrightarrow$ instantaneous code

# Prefix-Free Code

P2 The code should be easy to decode (decode each codeword immediately at the arrival of the last bit of that codeword) $\longrightarrow$ no codeword can be a prefix of another codeword.

---

### Definition (Prefix-Free Code)

A code is called a prefix-free code if no codeword is the prefix of any other codeword.

---

- We can decode codewords instantaneously, i.e., as soon as they are received $\longrightarrow$ instantaneous code
- Prefix-free code $\Longrightarrow$ Uniquely decodable

# Binary Code Tree

# Binary Code Tree



- Each node has two descendants

# Binary Code Tree



- Each node has two descendants
- Each node in the tree represents the binary string corresponding to the branches from the root to the node

# Binary Code Tree



- Each node has two descendants
- Each node in the tree represents the binary string corresponding to the branches from the root to the node
- Every codeword can be represented by a particular path through the tree

# Binary Code Tree



For prefix-free codes...

# Binary Code Tree



For prefix-free codes...

- Every codeword corresponds to a leaf

# Binary Code Tree



## For prefix-free codes...

- Every codeword corresponds to a leaf
- Intermediate nodes correspond to prefixes of some codewords

# Binary Code Tree



| symbol | $\mathcal{C}_6$ |
|:------:|:---------------:|
| $a$ | 0 |
| $b$ | 10 |
| $c$ | 110 |
| $d$ | 111 |

**For prefix-free codes...**

- Every codeword corresponds to a leaf
- Intermediate nodes correspond to prefixes of some codewords

# Kraft's Inequality

## Theorem (Kraft's Inequality)

*There exists a binary prefix-free code of cardinality $M$ and codeword lengths $\ell_1, \ell_2, \ldots, \ell_M$ if and only if*

$$\sum_{i=1}^{M} 2^{-\ell_i} \leq 1.$$

# Kraft's Inequality: proof

(Necessity: Prefix-free $\implies \sum_{i=1}^{M} 2^{-\ell_i} \leq 1$)

# Kraft's Inequality: proof

(Necessity: Prefix-free $\implies \sum_{i=1}^{M} 2^{-\ell_i} \le 1$)

# Kraft's Inequality: proof

(Sufficiency: $\sum_{i=1}^{M} 2^{-\ell_i} \leq 1 \implies$ Prefix-free)

# Kraft's Inequality: proof

(Sufficiency: $\sum_{i=1}^{M} 2^{-\ell_i} \leq 1 \implies$ Prefix-free)

## Prefix-Free Codes

Very nice properties: Uniquely decodable and easy to decode.

# Kraft-McMillan's Inequality

## Prefix-Free Codes

Very nice properties: Uniquely decodable and easy to decode.

However...is there a non prefix-free uniquely-decodable code that achieves lower average codeword length?

# Kraft-McMillan's Inequality

## Prefix-Free Codes

Very nice properties: Uniquely decodable and easy to decode.

However...is there a non prefix-free uniquely-decodable code that achieves lower average codeword length?

## Theorem (Kraft-McMillan's Inequality)

*A uniquely decodable code with codeword lengths $\ell_1, \ldots, \ell_M$ exists if and only if*

$$\sum_{i=1}^{M} 2^{-\ell_i} \leq 1.$$

# Kraft-McMillan's Inequality

## Prefix-Free Codes

Very nice properties: Uniquely decodable and easy to decode.

However...is there a non prefix-free uniquely-decodable code that achieves lower average codeword length?

## Theorem (Kraft-McMillan's Inequality)

*A uniquely decodable code with codeword lengths $\ell_1, \ldots, \ell_M$ exists if and only if*

$$\sum_{i=1}^{M} 2^{-\ell_i} \leq 1.$$

$\implies$ We can restrict ourselves to prefix-free codes with no loss in performance!

# Efficient Codes

P3 The code should compress the source as much as possible.

# Efficient Codes

P3 The code should compress the source as much as possible.

Our goal will be to minimize the average number of code bits per source symbol (the average codeword length),

# Efficient Codes

**P3** The code should compress the source as much as possible.

Our goal will be to minimize the average number of code bits per source symbol (the average codeword length),

$$\bar{L} = \sum_{i=1}^{M} p_i \ell_i.$$

# Efficient Codes

**P3** The code should compress the source as much as possible.

Our goal will be to minimize the average number of code bits per source symbol (the average codeword length),

$$\bar{L} = \sum_{i=1}^{M} p_i \ell_i.$$

## Example (Code $\mathcal{C}_6$)

Code $\mathcal{C}_6$ has average codeword length

# Efficient Codes

**P3** The code should compress the source as much as possible.

Our goal will be to minimize the average number of code bits per source symbol (the average codeword length),

$$\bar{L} = \sum_{i=1}^{M} p_i \ell_i.$$

> ## Example (Code $\mathcal{C}_6$)
>
> Code $\mathcal{C}_6$ has average codeword length
>
> $$\bar{L} = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + 2 \cdot \frac{1}{8} \cdot 3 = 1.75 \text{ bits}.$$

# The Limits of Source Coding

- What is the best achievable compression?

---

**Theorem**

*The average codeword length of a uniquely decodable code is lower bounded by*

---

# The Limits of Source Coding

- What is the best achievable compression?

---

**Theorem**

*The average codeword length of a uniquely decodable code is lower bounded by*

$$\bar{L} \geq \mathsf{H}(X).$$

---

*Proof*: (hint: use $\ln x \leq x - 1$)

# The Limits of Source Coding

# The Limits of Source Coding

- How close to the entropy can we expect to compress?

# The Limits of Source Coding

- How close to the entropy can we expect to compress?

---

**Theorem (Source Coding Theorem for a Single Random Symbol)**

*Let $X$ be a random variable generated by a discrete memoryless source with entropy $\mathsf{H}(X)$. There exists a prefix-free code $\mathcal{C}$ with average codeword length satisfying*

$$\mathsf{H}(X) \leq \bar{L}(\mathcal{C}) < \mathsf{H}(X) + 1.$$

---

# The Limits of Source Coding

*Proof:*

# Optimal Source Coding: Huffman Coding

- Given a probability distribution $P_X$, how can we design an optimal prefix-free code? (i.e., with minimum $\bar{L}$)

# Optimal Source Coding: Huffman Coding

- Given a probability distribution $P_X$, how can we design an optimal prefix-free code? (i.e., with minimum $\bar{L}$)

## The Huffman Coding Algorithm

# Optimal Source Coding: Huffman Coding

- Given a probability distribution $P_X$, how can we design an optimal prefix-free code? (i.e., with minimum $\bar{L}$)

## The Huffman Coding Algorithm

1. Order the $M$ source symbols in decreasing order of their probabilities

# Optimal Source Coding: Huffman Coding

- Given a probability distribution $P_X$, how can we design an optimal prefix-free code? (i.e., with minimum $\bar{L}$)

## The Huffman Coding Algorithm

1. Order the $M$ source symbols in decreasing order of their probabilities
2. Create $M$ leaves, one for each symbol

# Optimal Source Coding: Huffman Coding

- Given a probability distribution $P_X$, how can we design an <span style="color:red">optimal prefix-free code</span>? (i.e., with minimum $\bar{L}$)

## The Huffman Coding Algorithm

1. Order the $M$ source symbols in decreasing order of their probabilities
2. Create $M$ leaves, one for each symbol
3. Combine the two leaves corresponding to the two least probable symbols into a new node. Assign to this node the sum of the probabilities of the two original nodes

# Optimal Source Coding: Huffman Coding

- Given a probability distribution $P_X$, how can we design an **optimal prefix-free code**? (i.e., with minimum $\bar{L}$)

## The Huffman Coding Algorithm

1. Order the $M$ source symbols in decreasing order of their probabilities
2. Create $M$ leaves, one for each symbol
3. Combine the two leaves corresponding to the two least probable symbols into a new node. Assign to this node the sum of the probabilities of the two original nodes
4. Repeat Step 2 until only one node is free, and root it

# Optimal Source Coding: Huffman Coding

- Given a probability distribution $P_X$, how can we design an optimal prefix-free code? (i.e., with minimum $\bar{L}$)

### The Huffman Coding Algorithm

1. Order the $M$ source symbols in decreasing order of their probabilities
2. Create $M$ leaves, one for each symbol
3. Combine the two leaves corresponding to the two least probable symbols into a new node. Assign to this node the sum of the probabilities of the two original nodes
4. Repeat Step 2 until only one node is free, and root it
5. Starting from the root of the obtained tree, assign the binary symbols $0$ and $1$ to each pair of branches that arise from each node. The codeword for each symbol is read as the binary sequence from the root to the leaf associated to the symbol

# Optimal Source Coding: Huffman Coding

| $p$ | $x$ |
|---|---|
| 0.5 | $x_1$ ● |
| 0.15 | $x_2$ ● |
| 0.15 | $x_3$ ● |
| 0.10 | $x_4$ ● |
| 0.05 | $x_5$ ● |
| 0.05 | $x_6$ ● |

### Example (Huffman Coding)

Coding of $\mathcal{X} = \{x_1, x_2, x_3, x_4, x_5, x_6\}$, $P_X = \{0.5, 0.15, 0.15, 0.10, 0.05, 0.05\}$.
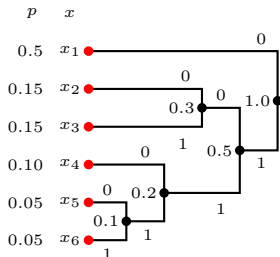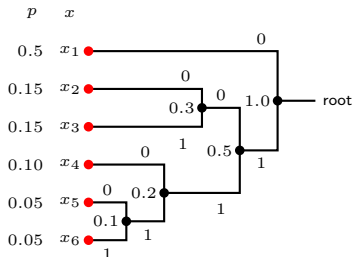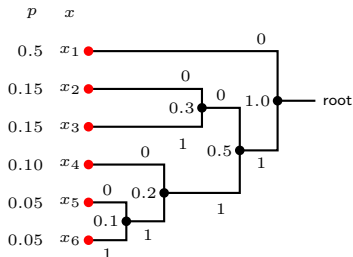$\mathsf{H}(X) = 2.08548$ bits.

# Optimal Source Coding: Huffman Coding

### Example (Huffman Coding)

Coding of $\mathcal{X} = \{x_1, x_2, x_3, x_4, x_5, x_6\}$, $P_X = \{0.5, 0.15, 0.15, 0.10, 0.05, 0.05\}$.
$\mathsf{H}(X) = 2.08548$ bits.

# Optimal Source Coding: Huffman Coding

# Optimal Source Coding: Huffman Coding



## Example (Huffman Coding)

Coding of $\mathcal{X} = \{x_1, x_2, x_3, x_4, x_5, x_6\}$, $P_X = \{0.5, 0.15, 0.15, 0.10, 0.05, 0.05\}$. $\mathsf{H}(X) = 2.08548$ bits.
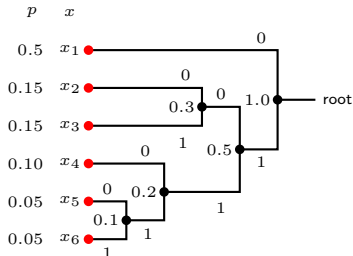
# Optimal Source Coding: Huffman Coding



## Example (Huffman Coding)

Coding of $\mathcal{X} = \{x_1, x_2, x_3, x_4, x_5, x_6\}$, $P_X = \{0.5, 0.15, 0.15, 0.10, 0.05, 0.05\}$.
$\mathsf{H}(X) = 2.08548$ bits.

# Optimal Source Coding: Huffman Coding



### Example (Huffman Coding)

Coding of $\mathcal{X} = \{x_1, x_2, x_3, x_4, x_5, x_6\}$, $P_X = \{0.5, 0.15, 0.15, 0.10, 0.05, 0.05\}$.
$\mathsf{H}(X) = 2.08548$ bits.

# Optimal Source Coding: Huffman Coding



## Example (Huffman Coding)

Coding of $\mathcal{X} = \{x_1, x_2, x_3, x_4, x_5, x_6\}$, $P_X = \{0.5, 0.15, 0.15, 0.10, 0.05, 0.05\}$.
$\mathsf{H}(X) = 2.08548$ bits.

# Optimal Source Coding: Huffman Coding



| symbol | codeword |
|--------|----------|
| $x_1$ | 0 |
| $x_2$ | 100 |
| $x_3$ | 101 |
| $x_4$ | 110 |
| $x_5$ | 1110 |
| $x_6$ | 1111 |

### Example (Huffman Coding)

Coding of $\mathcal{X} = \{x_1, x_2, x_3, x_4, x_5, x_6\}$, $P_X = \{0.5, 0.15, 0.15, 0.10, 0.05, 0.05\}$.
$\mathsf{H}(X) = 2.08548$ bits.

# Optimal Source Coding: Huffman Coding



| symbol | codeword |
|--------|----------|
| $x_1$ | 0 |
| $x_2$ | 100 |
| $x_3$ | 101 |
| $x_4$ | 110 |
| $x_5$ | 1110 |
| $x_6$ | 1111 |

## Example (Huffman Coding)

Coding of $\mathcal{X} = \{x_1, x_2, x_3, x_4, x_5, x_6\}$, $P_X = \{0.5, 0.15, 0.15, 0.10, 0.05, 0.05\}$.
$\mathsf{H}(X) = 2.08548$ bits. $\bar{L} = 2.1$ bits.

## What We Have Achieved up to Now...

Source compression by encoding each output of the source independently using an optimal prefix-free (Huffman) code.

# Encoding Blocks of Symbols

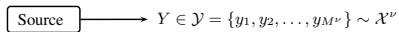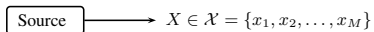**What We Have Achieved up to Now...**

Source compression by encoding each output of the source <span style="color:red">independently</span> using an optimal prefix-free (Huffman) code.

- Is there a more efficient approach?

# Encoding Blocks of Symbols

## What We Have Achieved up to Now...

Source compression by encoding each output of the source <span style="color:red">independently</span> using an optimal prefix-free (Huffman) code.
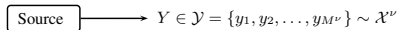
- Is there a more efficient approach?

## Encoding Blocks of Symbols

We can do better than encoding each symbol separately!

# Encoding Blocks of Symbols

## What We Have Achieved up to Now...

Source compression by encoding each output of the source independently using an optimal prefix-free (Huffman) code.
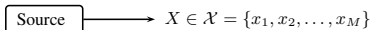
- Is there a more efficient approach?

## Encoding Blocks of Symbols

We can do better than encoding each symbol separately! $\longrightarrow$ Encoding blocks of symbols.

# Encoding Blocks of Symbols

Source $\longrightarrow X \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\}$

Source $\longrightarrow Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$

# Encoding Blocks of Symbols

Source $\longrightarrow X \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\}$     Source $\longrightarrow Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$
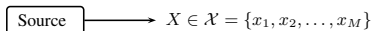
1. Combine $\nu$ consecutive symbols at the output of the source in a new symbol (message),

$$(x^{(1)}, x^{(2)}, \ldots, x^{(\nu)}),$$

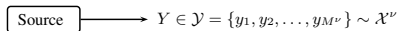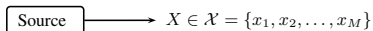where $x^{(j)} \in \mathcal{X}$.

# Encoding Blocks of Symbols

$$\boxed{\text{Source}} \longrightarrow X \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\} \qquad \boxed{\text{Source}} \longrightarrow Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$$

1. Combine $\nu$ consecutive symbols at the output of the source in a new symbol (message),

$$(x^{(1)}, x^{(2)}, \ldots, x^{(\nu)}),$$

where $x^{(j)} \in \mathcal{X}$.

2. Encode this message using a Huffman code.

# Encoding Blocks of Symbols

Source ⟶ $X \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\}$

Source ⟶ $Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$

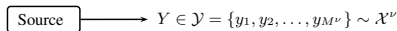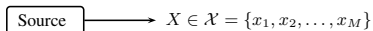1. Combine $\nu$ consecutive symbols at the output of the source in a new symbol (message),

$$(x^{(1)}, x^{(2)}, \ldots, x^{(\nu)}),$$

where $x^{(j)} \in \mathcal{X}$.

2. Encode this message using a Huffman code.

We design a code for an equivalent source $Y$ with alphabet $\mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\}$, $|\mathcal{Y}| = M^\nu$, consisting of all possible tuples $(x^{(1)}, x^{(2)}, \ldots, x^{(\nu)})$ of length $\nu$.

# Encoding Blocks of Symbols

Source $\longrightarrow$ $X \in \mathcal{X} = \{x_1, x_2, \ldots, x_M\}$

Source $\longrightarrow$ $Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$

1. Combine $\nu$ consecutive symbols at the output of the source in a new symbol (message),

$$(x^{(1)}, x^{(2)}, \ldots, x^{(\nu)}),$$

where $x^{(j)} \in \mathcal{X}$.
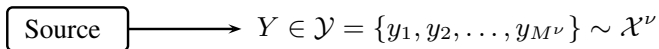
2. Encode this message using a Huffman code.

We design a code for an equivalent source $Y$ with alphabet $\mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\}$, $|\mathcal{Y}| = M^\nu$, consisting of all possible tuples $(x^{(1)}, x^{(2)}, \ldots, x^{(\nu)})$ of length $\nu$.

For a memoryless source, the probability of a message $y_i$ is

$$p_i = \prod_{j=1}^{\nu} p_{i,j},$$

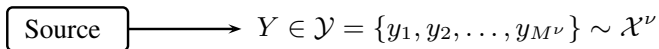where $p_{i,j}$ is the probability of the $j$th symbol of $y_i$.

# Encoding Blocks of Symbols

$$\boxed{\text{Source}} \longrightarrow Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$$

**From the Source Coding Theorem for a Single Random Symbol...**

We can build a code for $Y$ with average length $\bar{L}_\nu$,
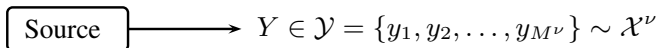
# Encoding Blocks of Symbols

$$\boxed{\text{Source}} \longrightarrow Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$$

From the Source Coding Theorem for a Single Random Symbol...

We can build a code for $Y$ with average length $\bar{L}_\nu$,

$$\mathsf{H}(Y) \leq \bar{L}_\nu < \mathsf{H}(Y) + 1.$$

# Encoding Blocks of Symbols

$$\boxed{\text{Source}} \longrightarrow Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$$
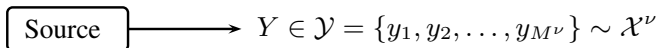
**From the Source Coding Theorem for a Single Random Symbol...**

We can build a code for $Y$ with average length $\bar{L}_\nu$,

$$\mathsf{H}(Y) \leq \bar{L}_\nu < \mathsf{H}(Y) + 1.$$

How do we compare $\bar{L}_\nu$ for different values of $\nu$?

# Encoding Blocks of Symbols

$$\boxed{\text{Source}} \longrightarrow Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$$

**From the Source Coding Theorem for a Single Random Symbol...**

We can build a code for $Y$ with average length $\bar{L}_\nu$,

$$\mathsf{H}(Y) \leq \bar{L}_\nu < \mathsf{H}(Y) + 1.$$

**How do we compare $\bar{L}_\nu$ for different values of $\nu$?**

The average codeword length necessary to describe one source symbol, $\bar{L}_\nu/\nu$!

# Encoding Blocks of Symbols



$$\boxed{\text{Source}} \longrightarrow Y \in \mathcal{Y} = \{y_1, y_2, \ldots, y_{M^\nu}\} \sim \mathcal{X}^\nu$$

# The Source Coding Theorem

### Theorem (Source Coding Theorem)

*Let $X$ be a random variable generated by a discrete memoryless source with entropy $\mathsf{H}(X)$. There exists a prefix-free code $\mathcal{C}$ that encodes messages of length $\nu$ symbols with average codeword length per source symbol $\frac{\bar{L}_\nu}{\nu}$ satisfying*

$$\mathsf{H}(X) \leq \frac{\bar{L}_\nu}{\nu} < \mathsf{H}(X) + \frac{1}{\nu}.$$

# The Source Coding Theorem

**Theorem (Source Coding Theorem)**

*Let $X$ be a random variable generated by a discrete memoryless source with entropy $\mathsf{H}(X)$. There exists a prefix-free code $\mathcal{C}$ that encodes messages of length $\nu$ symbols with average codeword length per source symbol $\frac{\bar{L}_\nu}{\nu}$ satisfying*

$$\mathsf{H}(X) \leq \frac{\bar{L}_\nu}{\nu} < \mathsf{H}(X) + \frac{1}{\nu}.$$

Choosing $\nu$ large enough we can approach the ultimate limit of compression, $\mathsf{H}(X)$, arbitrarily closely using Huffman codes!

# The Source Coding Theorem

> **Theorem (Source Coding Theorem)**
>
> *Let $X$ be a random variable generated by a discrete memoryless source with entropy $\mathsf{H}(X)$. There exists a prefix-free code $\mathcal{C}$ that encodes messages of length $\nu$ symbols with average codeword length per source symbol $\frac{\bar{L}_\nu}{\nu}$ satisfying*
>
> $$\mathsf{H}(X) \leq \frac{\bar{L}_\nu}{\nu} < \mathsf{H}(X) + \frac{1}{\nu}.$$

Choosing $\nu$ large enough we can approach the ultimate limit of compression, $\mathsf{H}(X)$, arbitrarily closely using Huffman codes!

However...encoding large blocks of symbols requires long codes $\longrightarrow$ difficult to construct, delay.

# The Source Coding Theorem

## Efficiency of a Source Code

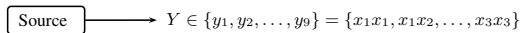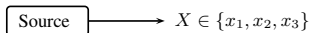$$\eta = \frac{\nu \mathsf{H}(X)}{\bar{L}_\nu}, \quad 0 \leq \eta \leq 1.$$

# The Source Coding Theorem

Source → $X \in \{x_1, x_2, x_3\}$

Source → $Y \in \{y_1, y_2, \ldots, y_9\} = \{x_1 x_1, x_1 x_2, \ldots, x_3 x_3\}$

Example: Source $X$, $\mathcal{X} = \{x_1, x_2, x_3\}$, with probabilities $\{0.45, 0.35, 0.20\}$. $\mathsf{H}(X) = 1.513$ bits.

# The Source Coding Theorem

Source ⟶ $X \in \{x_1, x_2, x_3\}$

Source ⟶ $Y \in \{y_1, y_2, \ldots, y_9\} = \{x_1x_1, x_1x_2, \ldots, x_3x_3\}$

Example: Source $X$, $\mathcal{X} = \{x_1, x_2, x_3\}$, with probabilities $\{0.45, 0.35, 0.20\}$. $\mathsf{H}(X) = 1.513$ bits.

$\mathcal{C}_1$: Encode symbols separately using a Huffman code. $\bar{L} = 1.55$ bits, $\eta = 0.976$.

# The Source Coding Theorem

| Source | $\longrightarrow$ | $X \in \{x_1, x_2, x_3\}$ |

| Source | $\longrightarrow$ | $Y \in \{y_1, y_2, \ldots, y_9\} = \{x_1x_1, x_1x_2, \ldots, x_3x_3\}$ |

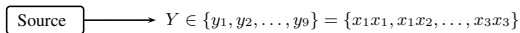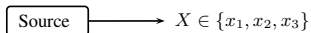Example: Source $X$, $\mathcal{X} = \{x_1, x_2, x_3\}$, with probabilities $\{0.45, 0.35, 0.20\}$. $\mathsf{H}(X) = 1.513$ bits.

$\mathcal{C}_1$: Encode symbols separately using a Huffman code. $\bar{L} = 1.55$ bits, $\eta = 0.976$.

$\mathcal{C}_2$: Encode pairs of symbols, i.e., we have equivalent source $Y = (X^1, X^2)$, with $\mathcal{Y} = \{y_1, y_2, \ldots, y_9\} = \{x_1x_1, x_1x_2, x_1x_3, x_2x_1, x_2x_2, x_2x_3, x_3x_1, x_3x_2, x_3x_3\}$, $|\mathcal{Y}| = |\mathcal{X}|^2 = 9$, and encode messages $Y$ using a Huffman code. $\frac{\bar{L}_\nu}{\nu} = 1.53375$ bits, $\eta = 0.989$.

# The Source Coding Theorem

Table: Huffman code, encoding symbols separately, $\bar{L} = 1.55$, $\eta = 0.976$

| symbol | probability | codeword |
|--------|-------------|----------|
| $x_1$  | 0.45        | 1        |
| $x_2$  | 0.35        | 00       |
| $x_3$  | 0.20        | 01       |

Table: Huffman code, encoding blocks of two symbols, $\bar{L}_\nu = 3.0675$, $\eta = 0.989$

| symbol   | probability | codeword |
|----------|-------------|----------|
| $x_1x_1$ | 0.2025      | 10       |
| $x_1x_2$ | 0.1575      | 001      |
| $x_2x_1$ | 0.1575      | 010      |
| $x_2x_2$ | 0.1225      | 011      |
| $x_1x_3$ | 0.09        | 111      |
| $x_3x_1$ | 0.09        | 0000     |
| $x_2x_3$ | 0.07        | 0001     |
| $x_3x_2$ | 0.07        | 1100     |
| $x_3x_3$ | 0.04        | 1101     |

# Further Discussion

## Source Compression using Huffman Codes

- Requires knowledge of the source statistics.

# Further Discussion

## Source Compression using Huffman Codes

- Requires knowledge of the source statistics.
- Lossless: we can always recover the transmitted sequence of data symbols from the coded sequence with no loss of information.

# Further Discussion

## Source Compression using Huffman Codes

- Requires knowledge of the source statistics.
- Lossless: we can always recover the transmitted sequence of data symbols from the coded sequence with no loss of information.

## Universal Source Compression

A source compression that achieves the ultimate limit regardless of the source: The Lempel-Ziv compression algorithm (gzip,GIF), proven asymptotically to compress down to the entropy of the source.

# Further Discussion

## Source Compression using Huffman Codes

- Requires knowledge of the source statistics.
- Lossless: we can always recover the transmitted sequence of data symbols from the coded sequence with no loss of information.

## Universal Source Compression

A source compression that achieves the ultimate limit regardless of the source: The Lempel-Ziv compression algorithm (gzip,GIF), proven asymptotically to compress down to the entropy of the source.

## Lossy Source Compression

- Compression beyond the entropy of the source $\longrightarrow$ we cannot recover the original data identically after decompression.

# Further Discussion

## Source Compression using Huffman Codes

- Requires knowledge of the source statistics.
- Lossless: we can always recover the transmitted sequence of data symbols from the coded sequence with no loss of information.

## Universal Source Compression

A source compression that achieves the ultimate limit regardless of the source: The Lempel-Ziv compression algorithm (gzip,GIF), proven asymptotically to compress down to the entropy of the source.

## Lossy Source Compression

- Compression beyond the entropy of the source $\longrightarrow$ we cannot recover the original data identically after decompression.

# Further Discussion

## Source Compression using Huffman Codes

- Requires knowledge of the source statistics.
- Lossless: we can always recover the transmitted sequence of data symbols from the coded sequence with no loss of information.

## Universal Source Compression

A source compression that achieves the ultimate limit regardless of the source: The Lempel-Ziv compression algorithm (gzip,GIF), proven asymptotically to compress down to the entropy of the source.

## Lossy Source Compression

- Compression beyond the entropy of the source $\longrightarrow$ we cannot recover the original data identically after decompression.
- Typically closely related to the human perception.

# Further Discussion

## Source Compression using Huffman Codes

- Requires knowledge of the source statistics.
- Lossless: we can always recover the transmitted sequence of data symbols from the coded sequence with no loss of information.

## Universal Source Compression

A source compression that achieves the ultimate limit regardless of the source: The Lempel-Ziv compression algorithm (gzip,GIF), proven asymptotically to compress down to the entropy of the source.

## Lossy Source Compression

- Compression beyond the entropy of the source $\longrightarrow$ we cannot recover the original data identically after decompression.
- Typically closely related to the human perception.
- Used to compress sound, video or image (jpeg compression).