# Information-Theoretic Modeling
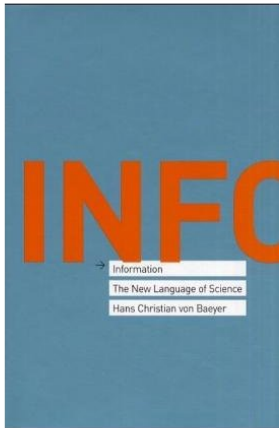
Teemu Roos

Department of Computer Science, University of Helsinki

Fall 2009

UNIVERSITY OF HELSINKI

"Whether on the internet, encoded in radio waves or coursing through wires, information is all around us. Our senses record it, our brains process it and our genes pass it on. But what exactly is information? Can it be analysed and measured? In this extraordinary book, Hans Christian von Baeyer illuminates a concept that could soon become as central to science as space, time mass or energy."

1. Administrative issues
   - Course details
   - Prerequisites
   - What do I need to do?
   - Grading

1. Administrative issues
   - Course details
   - Prerequisites
   - What do I need to do?
   - Grading

2. Overview of Contents
   - What is Information?
   - Why Information?
   - Information vs. Complexity
   - Information Theory

1. Administrative issues
   - Course details
   - Prerequisites
   - What do I need to do?
   - Grading

2. Overview of Contents
   - What is Information?
   - Why Information?
   - Information vs. Complexity
   - Information Theory

3. Compression
   - Dots and Dashes
   - Codes as Mappings
   - Data Compression
   - Information vs. Complexity (contd.)

Outline
**Administrative issues**
Overview of Contents
Compression

**Course details**
Prerequisites
What do I need to do?
Grading

## 582650, Information-Theoretic Modeling

- An advanced studies (laudatur) course.

Outline
**Administrative issues**
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582650, Information-Theoretic Modeling

- An advanced studies (laudatur) course.
- *Algorithms and Machine Learning* (was *Intelligent Systems*) sub-programme, optional.

Outline
**Administrative issues**
Overview of Contents
Compression

**Course details**
Prerequisites
What do I need to do?
Grading

# 582650, Information-Theoretic Modeling

- An advanced studies (laudatur) course.
- *Algorithms and Machine Learning* (was *Intelligent Systems*) sub-programme, optional.
- 4 credit units.

Outline
**Administrative issues**
Overview of Contents
Compression

**Course details**
Prerequisites
What do I need to do?
Grading

# 582650, Information-Theoretic Modeling



- An advanced studies (laudatur) course.
- *Algorithms and Machine Learning* (was *Intelligent Systems*) sub-programme, optional.
- 4 credit units.
- Lectures: 8.9.–16.10. Tue & Fri 10–12 in C222.

Outline
**Administrative issues**
Overview of Contents
Compression

**Course details**
Prerequisites
What do I need to do?
Grading

# 582650, Information-Theoretic Modeling

- An advanced studies (laudatur) course.

- *Algorithms and Machine Learning* (was *Intelligent Systems*) sub-programme, optional.

- 4 credit units.

- Lectures: 8.9.–16.10. Tue & Fri 10–12 in C222.

- Exercises: 15.9.–13.10. Tue (or Fri?) 12–14 in B119.

Outline
**Administrative issues**
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582650, Information-Theoretic Modeling

- An advanced studies (laudatur) course.

- *Algorithms and Machine Learning* (was *Intelligent Systems*) sub-programme, optional.

- 4 credit units.

- Lectures: 8.9.–16.10. Tue & Fri 10–12 in C222.

- Exercises: 15.9.–13.10. Tue (or Fri?) 12–14 in B119.

- Instructor: **Teemu Roos**, A322,
  `teemu.roos at cs.helsinki.fi`
  (Pls. make an appointment).

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582650, Information-Theoretic Modeling

- An advanced studies (laudatur) course.
- *Algorithms and Machine Learning* (was *Intelligent Systems*) sub-programme, optional.
- 4 credit units.
- Lectures: 8.9.–16.10. Tue & Fri 10–12 in C222.
- Exercises: 15.9.–13.10. Tue (or Fri?) 12–14 in B119.
- Instructor: **Teemu Roos**, A322,
  `teemu.roos at cs.helsinki.fi`
  (Pls. make an appointment).
- Course assistant: **Anupam Arohi**.
  `anupam.arohi at cs.helsinki.fi`

Outline
**Administrative issues**
Overview of Contents
Compression

**Course details**
Prerequisites
What do I need to do?
Grading

# 582650, Information-Theoretic Modeling

- An advanced studies (laudatur) course.
- *Algorithms and Machine Learning* (was *Intelligent Systems*) sub-programme, optional.
- 4 credit units.
- Lectures: 8.9.–16.10. Tue & Fri 10–12 in C222.
- Exercises: 15.9.–13.10. Tue (or Fri?) 12–14 in B119.
- Instructor: **Teemu Roos**, A322,
  teemu.roos at cs.helsinki.fi
  (Pls. make an appointment).
- Course assistant: **Anupam Arohi**.
  anupam.arohi at cs.helsinki.fi
- www.cs.helsinki.fi/group/cosco/Teaching/Information/2009/

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

## Resources

There is no required textbook on the course, but the following are recommended.

- **Highly recommended:** Cover & Thomas, *Elements of Information Theory*,
- MacKay, *Information Theory, Inference and Learning Algorithms*,
- Grünwald, *The Minimum Description Length Principle*,
- Solomon, *Data Compression: The Complete Reference*.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

## Resources

There is no required textbook on the course, but the following are recommended.

- Highly recommended: Cover & Thomas, *Elements of Information Theory*,
- MacKay, *Information Theory, Inference and Learning Algorithms*,
- Grünwald, *The Minimum Description Length Principle*,
- Solomon, *Data Compression: The Complete Reference*.

Copies of required material will be made available at the lectures, and afterwards in the *course folder* in room C127.

Outline
**Administrative issues**
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582651, Project in Information-Theoretic Modeling

There is also a related project:

- 2 credit units.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582651, Project in Information-Theoretic Modeling

There is also a related project:

- 2 credit units.
- Period II.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582651, Project in Information-Theoretic Modeling

There is also a related project:

- 2 credit units.
- Period II.
- This course is a prerequisite.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582651, Project in Information-Theoretic Modeling

There is also a related project:

- 2 credit units.
- Period II.
- This course is a prerequisite.
- Together they replace the old *Three Concepts: Information* course — both old and new version cannot be included in your degree.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

## 582651, Project in Information-Theoretic Modeling

There is also a related project:

- 2 credit units.
- Period II.
- This course is a prerequisite.
- Together they replace the old *Three Concepts: Information* course — both old and new version cannot be included in your degree.
- Groups of 2–3 persons.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582651, Project in Information-Theoretic Modeling

There is also a related project:

- 2 credit units.
- Period II.
- This course is a prerequisite.
- Together they replace the old *Three Concepts: Information* course — both old and new version cannot be included in your degree.
- Groups of 2–3 persons.
- Task: compress data.

Outline
**Administrative issues**
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582651, Project in Information-Theoretic Modeling

There is also a related project:

- 2 credit units.
- Period II.
- This course is a prerequisite.
- Together they replace the old *Three Concepts: Information* course — both old and new version cannot be included in your degree.
- Groups of 2–3 persons.
- Task: compress data.
- Best compressor wins! Intermediate results annonced periodically.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# 582651, Project in Information-Theoretic Modeling

There is also a related project:

- 2 credit units.
- Period II.
- This course is a prerequisite.
- Together they replace the old *Three Concepts: Information* course — both old and new version cannot be included in your degree.
- Groups of 2–3 persons.
- Task: compress data.
- Best compressor wins! Intermediate results annonced periodically.
- Programming + report.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

## Prerequisites

No formal prerequisites **but** you will need

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

## Prerequisites

No formal prerequisites **but** you will need

- Calculus: integrals, derivatives, convergence, ...

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

## Prerequisites

No formal prerequisites **but** you will need

- Calculus: integrals, derivatives, convergence, ...
- Probability theory: joint & conditional distributions, expectations, law of large numbers, ...

Outline
**Administrative issues**
Overview of Contents
Compression

Course details
**Prerequisites**
What do I need to do?
Grading

## Prerequisites

No formal prerequisites **but** you will need

- Calculus: integrals, derivatives, convergence, ...
- Probability theory: joint & conditional distributions, expectations, law of large numbers, ...
- Programming: language is up to you (but need to work in groups in project).

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# What do I need to do?

- Weekly exercises:

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

## What do I need to do?

- Weekly exercises:
  - Mathematical problems.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
**What do I need to do?**
Grading

## What do I need to do?

- Weekly exercises:
  - Mathematical problems.
  - Programming tasks.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# What do I need to do?

- Weekly exercises:
  - Mathematical problems.
  - Programming tasks.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# What do I need to do?

- Weekly exercises:
    - Mathematical problems.
    - Programming tasks.
- Final exam (date TBA).

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
**What do I need to do?**
Grading

## What do I need to do?

- Weekly exercises:
  - Mathematical problems.
  - Programming tasks.
- Final exam (date TBA).

You do *not* have to attend the classes, unless otherwise stated. However, we recommend that you do. (At least off-line.)

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# What do I need to do?

- Weekly exercises:
    - Mathematical problems.
    - Programming tasks.
- Final exam (date TBA).

You do *not* have to attend the classes, unless otherwise stated. However, we recommend that you do. (At least off-line.)

If you find that the course is not for you, please let us know *as soon as possible*. There are people willing to take your place.

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# Grading

The course grading is based on:

1. Exercises (40 %)

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# Grading

The course grading is based on:

1. Exercises (40 %)
2. Exam (60 %)

Outline
Administrative issues
Overview of Contents
Compression

Course details
Prerequisites
What do I need to do?
Grading

# Grading

The course grading is based on:

1. Exercises (40 %)
2. Exam (60 %)

Minimum 50 % of exercises have to be solved (or at least seriously attempted).

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

## What is Information?

- Etymology: *informare* = give form, 14th century.

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

## What is Information?

- Etymology: *informare* = give form, 14th century.
- *knowledge [...], intelligence, news, facts, data, [...], (as nucleotides in DNA or binary digits in a computer program) [...], a signal [...], a numerical quantity that measures the uncertainty in the outcome of an experiment to be performed.* (source: Merriam-Webster).

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# What is Information?

- Etymology: *informare* = give form, 14th century.
- *knowledge [...], intelligence, news, facts, data, [...], (as nucleotides in DNA or binary digits in a computer program) [...], a signal [...], a numerical quantity that measures the uncertainty in the outcome of an experiment to be performed.* (source: Merriam-Webster).
- Data < Information < Knowledge.

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# What is Information?

- Etymology: *informare* = give form, 14th century.
- *knowledge [...], intelligence, news, facts, data, [...], (as nucleotides in DNA or binary digits in a computer program) [...], a signal [...], a numerical quantity that measures the uncertainty in the outcome of an experiment to be performed.* (source: Merriam-Webster).
- Data < Information < Knowledge.
- Information technology.

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

## What is Information?

- Etymology: *informare* = give form, 14th century.

- *knowledge [...], intelligence, news, facts, data, [...], (as nucleotides in DNA or binary digits in a computer program) [...], a signal [...], a numerical quantity that measures the uncertainty in the outcome of an experiment to be performed.* (source: Merriam-Webster).

- Data < Information < Knowledge.

- Information technology.

- Physical information.

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# What is Information?

- Etymology: *informare* = give form, 14th century.
- *knowledge [...], intelligence, news, facts, data, [...], (as nucleotides in DNA or binary digits in a computer program) [...], a signal [...], a numerical quantity that measures the uncertainty in the outcome of an experiment to be performed.* (source: Merriam-Webster).
- Data < Information < Knowledge.
- Information technology.
- Physical information.
- This course: measuring *the amount* of information in data, and using such measures for automatically buiding *models*.

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# Why Information?

- The amount of information around us is exploding – internet!

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# Why Information?

- The amount of information around us is exploding – internet!
- Need to *store*, *transmit*, *and process* information efficiently.

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
**Why Information?**
Information vs. Complexity
Information Theory

# Why Information?

- The amount of information around us is exploding – internet!
- Need to *store*, *transmit*, *and process* information efficiently.
- Wish to *understand* more and more complex phenomena.

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# Why Information?

- The amount of information around us is exploding – internet!
- Need to *store*, *transmit*, *and process* information efficiently.
- Wish to *understand* more and more complex phenomena.
- Computer science: make things automatic (intelligent).

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# Information vs. Complexity

Is complexity the same as information?

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# Information vs. Complexity

Is complexity the same as information?

Is there a lot of *information* in a random string? **No.**

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
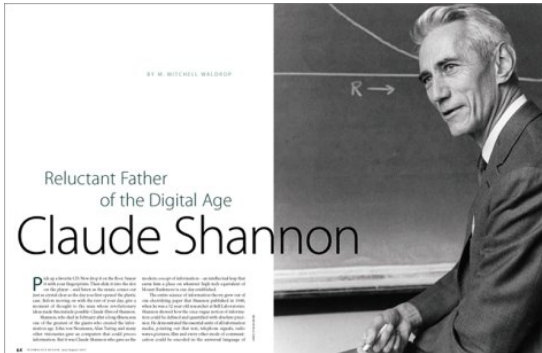Why Information?
Information vs. Complexity
Information Theory

# Information vs. Complexity

Is complexity the same as information?

Is there a lot of *information* in a random string? **No.**

$$
\begin{array}{rcccc}
Complexity & = & Information & + & Noise \\
& = & Regularity & + & Randomness \\
& = & Algorithm & + & Compressed\ file
\end{array}
$$

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
Why Information?
Information vs. Complexity
**Information Theory**

# Information Theory



Reluctant Father
of the Digital Age
Claude Shannon

"The real birth of modern information theory can be traced to the publication in 1948 of Claude Shannon's *"The Mathematical Theory of Communication"* in the Bell System Technical Journal. "
(Encyclopædia Britannica)

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
Why Information?
Information vs. Complexity
**Information Theory**

# Course Topics

Information Theory:

- entropy and information, bits,

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# Course Topics

Information Theory:

- entropy and information, bits,
- compression,

Outline
Administrative issues
Overview of Contents
Compression

What is Information?
Why Information?
Information vs. Complexity
Information Theory

# Course Topics

Information Theory:

- entropy and information, bits,

- compression,

- error correction.

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
Why Information?
Information vs. Complexity
**Information Theory**

# Course Topics

Information Theory:

- entropy and information, bits,

- compression,

- error correction.

Fundamental limits (mathematical and statistical) and practice (computer science).

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
Why Information?
Information vs. Complexity
**Information Theory**

# Course Topics

Information Theory:

- entropy and information, bits,
- compression,
- error correction.

Fundamental limits (mathematical and statistical) and practice (computer science).

Modeling:

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
Why Information?
Information vs. Complexity
**Information Theory**

# Course Topics

Information Theory:

- entropy and information, bits,

- compression,

- error correction.

Fundamental limits (mathematical and statistical) and practice (computer science).

Modeling:

- statistical models,

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
Why Information?
Information vs. Complexity
**Information Theory**

## Course Topics

Information Theory:

- entropy and information, bits,
- compression,
- error correction.

Fundamental limits (mathematical and statistical) and practice (computer science).

Modeling:

- statistical models,
- complexity (in data and models),

Outline
Administrative issues
**Overview of Contents**
Compression

What is Information?
Why Information?
Information vs. Complexity
**Information Theory**

# Course Topics

Information Theory:

- entropy and information, bits,
- compression,
- error correction.

Fundamental limits (mathematical and statistical) and practice (computer science).

Modeling:

- statistical models,
- complexity (in data and models),
- over-fitting, Occam's Razor, and MDL Principle.

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Coding Game

Form groups of 3–4 persons. Each group constructs a *code* for the
letters A–Z by using as *code-words* unique sequences of dots • and
dashes (—) like "•", "— •", " — • — —", etc.

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| A | _____ | G | _____ | M | _____ | S | _____ | Y | _____ |
| B | _____ | H | _____ | N | _____ | T | _____ | Z | _____ |
| C | _____ | I | _____ | O | _____ | U | _____ | | |
| D | _____ | J | _____ | P | _____ | V | _____ | | |
| E | _____ | K | _____ | Q | _____ | W | _____ | | |
| F | _____ | L | _____ | R | _____ | X | _____ | | |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Coding Game

Use your code to *encode* the message
"WHAT DOES THIS HAVE TO DO WITH INFORMATION".

Outline
Administrative issues
Overview of Contents
**Compression**

**Dots and Dashes**
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Coding Game

Use your code to *encode* the message
"WHAT DOES THIS HAVE TO DO WITH INFORMATION".

Now count how long the encoded message is using the rule:

- A dot •: 1 units.
- A dash —: 2 units.
- A space between words: 2 units.

Outline
Administrative issues
Overview of Contents
**Compression**

**Dots and Dashes**
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Coding Game

Use your code to *encode* the message
"WHAT DOES THIS HAVE TO DO WITH INFORMATION".

Now count how long the encoded message is using the rule:

- A dot •: 1 units.
- A dash —: 2 units.
- A space between words: 2 units.

• • • — — — • • •: $1 + 1 + 1 + 2 + 2 + 2 + 1 + 1 + 1 = 12$.

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Coding Game

Use your code to *encode* the message
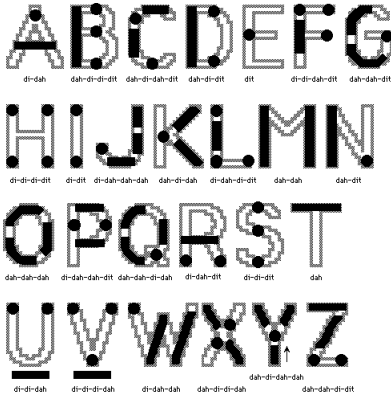"WHAT DOES THIS HAVE TO DO WITH INFORMATION".

Now count how long the encoded message is using the rule:

- A dot •: 1 units.
- A dash —: 2 units.
- A space between words: 2 units.

• • • — — — • • •: $1 + 1 + 1 + 2 + 2 + 2 + 1 + 1 + 1 = 12$.
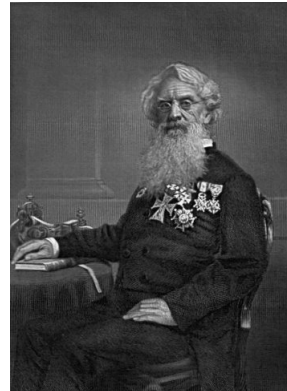
The *coding rate* of your code is the length of the encoded message divided by the length of the original message, including spaces (42).

Outline
Administrative issues
Overview of Contents
**Compression**

**Dots and Dashes**
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Coding Game



© 1989 A.G. Reinhold.



Samuel F.M. Morse (1791–1872)

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Coding Game

WHAT DOES THIS HAVE TO DO WITH INFORMATION

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Coding Game

WHAT DOES THIS HAVE TO DO WITH INFORMATION

```
.-- .... .- -   -.. --- . ...    - .... .. ...
.... .- ...- .    - ---   -.. ---   .-- .. - ....
.. -. ..-. --- .-. -- .- - .. --- -.
```

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Coding Game

WHAT DOES THIS HAVE TO DO WITH INFORMATION

```
.-- .... .- -   -.. --- . ...    - .... .. ...
.... .- ...- .    - ---   -.. ---   .-- .. - ....
.. -. ..-. --- .-. -- .- - .. --- -.
```

51 dots, 36 dashes, 7 spaces: $51 + 72 + 14 = 137$ units.

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Coding Game

WHAT DOES THIS HAVE TO DO WITH INFORMATION

```
.-- .... .- -   -.. --- . ...   - .... .. ...
.... .- ...- .   - ---   -.. ---   .-- .. - ....
.. -. ..-. --- .-. -- .- - .. --- -.
```

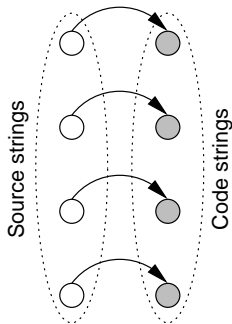51 dots, 36 dashes, 7 spaces: $51 + 72 + 14 = 137$ units.

**Morse code**

Coding rate: $\dfrac{137}{42} \approx 3.26$

Did you do better or worse? Why?

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Codes as Mappings

Lossless compression:
injective mapping

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Codes as Mappings



Lossless compression:
injective mapping

Lossy compression:
non-injective mapping

Source strings

Code strings

Source strings

Code strings

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Codes as Mappings

Lossless compression:
injective mapping

Lossy compression:
non-injective mapping



Only *lossless* codes are *uniquely decodable*.

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
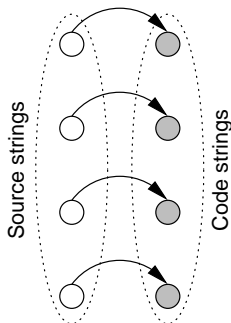**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Codes as Mappings



Lossless compression:
injective mapping

Lossy compression:
non-injective mapping

Source strings — Code strings

Source strings — Code strings

Only *lossless* codes are *uniquely decodable*.

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Examples

*general
purpose*

`gzip`

`bzip`

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Examples

| *general purpose* | `gzip` |
| | `bzip` |

| *image* | `png` |
| | `jpeg` |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Examples

| | |
|---|---|
| *general purpose* | `gzip` |
| | `bzip` |

| | |
|---|---|
| *image* | `png` |
| | `jpeg` |

| | |
|---|---|
| *music* | `mp3` |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Examples

*general purpose*    `gzip`

   `bzip`

*image*    `png`

   `jpeg`

*music*    `mp3`

*video*    `mpeg`

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Examples

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
**Codes as Mappings**
Data Compression
Information vs. Complexity (contd.)

# Examples



*compression ratio*

| | | | |
|---|---|---|---|
| *general purpose* | **gzip** | ~ 1 : 3 | *lossless* |
| | **bzip** | ~ 1 : 3.5 | |
| *image* | **png** | ~ 1 : 2.5 | |
| | **jpeg** | ~ 1 : 25 | *lossy* |
| *music* | **mp3** | ~ 1 : 12 | |
| *video* | **mpeg** | ~ 1 : 30 | |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# Compression

Is it always possible to compress data?

### Theorem

The proportion of binary strings compressible by more than $k$ bits is less than $2^{-k}$.

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# Compression

Is it always possible to compress data?

### Theorem

The proportion of binary strings compressible by more than $k$ bits is less than $2^{-k}$.

*Proof.* For all $n \geq 1$, the number of binary strings of length $n$ is $2^n$.

Outline — Dots and Dashes
Administrative issues — Codes as Mappings
Overview of Contents — **Data Compression**
**Compression** — Information vs. Complexity (contd.)

# Compression

Is it always possible to compress data?

### Theorem

The proportion of binary strings compressible by more than $k$ bits is less than $2^{-k}$.

*Proof.* For all $n \geq 1$, the number of binary strings of length $n$ is $2^n$. The number of binary code strings of length less than $n - k$ is
$$2^0 + 2^1 + 2^2 + \ldots + 2^{n-k-1}$$

Outline          Dots and Dashes
Administrative issues     Codes as Mappings
Overview of Contents    Data Compression
Compression        Information vs. Complexity (contd.)

# Compression

Is it always possible to compress data?

## Theorem

The proportion of binary strings compressible by more than $k$ bits is less than $2^{-k}$.

*Proof.* For all $n \geq 1$, the number of binary strings of length $n$ is $2^n$. The number of binary code strings of length less than $n - k$ is $2^0 + 2^1 + 2^2 + \ldots + 2^{n-k-1} = 2^{n-k} - 1$.

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Compression

Is it always possible to compress data?

### Theorem

The proportion of binary strings compressible by more than $k$ bits is less than $2^{-k}$.

*Proof.* For all $n \geq 1$, the number of binary strings of length $n$ is $2^n$. The number of binary code strings of length less than $n - k$ is $2^0 + 2^1 + 2^2 + \ldots + 2^{n-k-1} = 2^{n-k} - 1$. Thus the ratio is

$$\frac{2^{n-k} - 1}{2^n} < \frac{2^{n-k}}{2^n} = 2^{-k}.$$

$\square$

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Compression

Is it always possible to compress data?

## Theorem

The proportion of binary strings compressible by more than $k$ bits is less than $2^{-k}$.

*Proof.* For all $n \geq 1$, the number of binary strings of length $n$ is $2^n$. The number of binary code strings of length less than $n - k$ is $2^0 + 2^1 + 2^2 + \ldots + 2^{n-k-1} = 2^{n-k} - 1$. Thus the ratio is

$$\frac{2^{n-k} - 1}{2^n} < \frac{2^{n-k}}{2^n} = 2^{-k}.$$

☐

Less than 50 % of files are compressible by more than one bit.

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Compression

Is it always possible to compress data?

**Theorem**

The proportion of binary strings compressible by more than $k$ bits is less than $2^{-k}$.

*Proof.* For all $n \geq 1$, the number of binary strings of length $n$ is $2^n$. The number of binary code strings of length less than $n - k$ is $2^0 + 2^1 + 2^2 + \ldots + 2^{n-k-1} = 2^{n-k} - 1$. Thus the ratio is

$$\frac{2^{n-k} - 1}{2^n} < \frac{2^{n-k}}{2^n} = 2^{-k}.$$

$\square$

Less than 1 % of files are compressible by more than 7 bits.

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Compression

Is it always possible to compress data?

**Theorem**

The proportion of binary strings compressible by more than $k$ bits is less than $2^{-k}$.

*Proof.* For all $n \geq 1$, the number of binary strings of length $n$ is $2^n$. The number of binary code strings of length less than $n - k$ is $2^0 + 2^1 + 2^2 + \ldots + 2^{n-k-1} = 2^{n-k} - 1$. Thus the ratio is

$$\frac{2^{n-k} - 1}{2^n} < \frac{2^{n-k}}{2^n} = 2^{-k}.$$

$\square$

Less than 0.0000000000000000000000000000001 % of files are compressible by 100 bits.

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# How is it possible?

Why was the compression ratio greater than one in all the examples we saw?

What are those rare files that are compressible?

Why are the files we use in practice so often compressible?

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# Compression

```
echo <x> | gzip - | wc -c       # multiply by 8 for bits
```

| Source string, $x$ | | $\ell(C(x))$ | ratio |
|---|---|---|---|
| $aaa \ldots a$ | $(10000 \times a)$ | 368 | 27.2 : 1. |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# Compression

```
echo <x> | gzip - | wc -c       # multiply by 8 for bits
```

| Source string, $x$ | | $\ell(C(x))$ | ratio |
|---|---|---|---|
| $aaa \ldots a$ | $(10000 \times a)$ | 368 | 27.2 : 1. |
| $aabaabbbbabbbbb \ldots$ | (10000 random letters) | 13456 | 0.74 : 1 |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# Compression

```
echo <x> | gzip - | wc -c       # multiply by 8 for bits
```

| Source string, $x$ | | $\ell(C(x))$ | ratio |
|---|---|---|---|
| $aaa\ldots a$ | $(10000 \times a)$ | 368 | 27.2 : 1. |
| $aabaabbbbabbbbb\ldots$ | (10000 random letters) | 13456 | 0.74 : 1 |
| $abababab\ldots ab$ | $(5000 \times ab)$ | 368 | 27.2 : 1 |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# Compression

```
echo <x> | gzip - | wc -c          # multiply by 8 for bits
```

| Source string, $x$ | | $\ell(C(x))$ | ratio |
|---|---|---|---|
| $aaa\ldots a$ | $(10000 \times a)$ | 368 | 27.2 : 1. |
| $aabaabbbbabbbbb\ldots$ | (10000 random letters) | 13456 | 0.74 : 1 |
| $abababab\ldots ab$ | $(5000 \times ab)$ | 368 | 27.2 : 1 |
| $aaa\ldots abbb\ldots b$ | $(5000 \times a, 5000 \times b)$ | 376 | 26.6 : 1 |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# Compression

```
echo <x> | gzip - | wc -c        # multiply by 8 for bits
```

| Source string, $x$ | | $\ell(C(x))$ | ratio |
|---|---|---|---|
| $aaa \ldots a$ | $(10000 \times a)$ | 368 | 27.2 : 1. |
| $aabaabbbbabbbbb \ldots$ | (10000 random letters) | 13456 | 0.74 : 1 |
| $ababab \ldots ab$ | $(5000 \times ab)$ | 368 | 27.2 : 1 |
| $aaa \ldots abbb \ldots b$ | $(5000 \times a, 5000 \times b)$ | 376 | 26.6 : 1 |
| $abbaababba \ldots$ | $(1000 \times abbaababba)$ | 488 | 20.5 : 1 |

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

# Compression

```
echo <x> | gzip - | wc -c       # multiply by 8 for bits
```

| **Source string, $x$** | | $\ell(C(x))$ | **ratio** |
|---|---|---|---|
| *aaa . . . a* | (10000 × *a*) | 368 | 27.2 : 1. |
| *aabaabbbbabbbbb . . .* | (10000 random letters) | 13456 | 0.74 : 1 |
| *abababab . . . ab* | (5000 × *ab*) | 368 | 27.2 : 1 |
| *aaa . . . abbb . . . b* | (5000 × *a*, 5000 × *b*) | 376 | 26.6 : 1 |
| *abbaababba . . .* | (1000 × *abbaababba*) | 488 | 20.5 : 1 |

Strings following a rule are compressible?

Outline          Dots and Dashes
Administrative issues   Codes as Mappings
Overview of Contents   **Data Compression**
**Compression**      Information vs. Complexity (contd.)

# Compression

```
echo <x> | gzip - | wc -c        # multiply by 8 for bits
```

| **Source string, $x$** | | $\ell(C(x))$ | **ratio** |
|---|---|---|---|
| $aaa \dots a$ | $(10000 \times a)$ | 368 | 27.2 : 1. |
| $aabaabbbbabbbbb \dots$ | (10000 random letters) | 13456 | 0.74 : 1 |
| $abababab \dots ab$ | $(5000 \times ab)$ | 368 | 27.2 : 1 |
| $aaa \dots abbb \dots b$ | $(5000 \times a, 5000 \times b)$ | 376 | 26.6 : 1 |
| $abbaababba \dots$ | $(1000 \times abbaababba)$ | 488 | 20.5 : 1 |
| $aaabbabbabb \dots$ | $(\pi, 0\text{–}4 \mapsto a, 5\text{–}9 \mapsto b)$ | 13416 | 0.74 : 1 |

$\pi$ follows a rule but isn't compressed!

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

## Compression

```
echo <x> | gzip - | wc -c       # multiply by 8 for bits
```

| Source string, $x$ | | $\ell(C(x))$ | ratio |
|---|---|---|---|
| $aaa\ldots a$ | $(10000 \times a)$ | 368 | 27.2 : 1. |
| $aabaabbbbabbbbb\ldots$ | $(10000 \text{ random letters})$ | 13456 | 0.74 : 1 |
| $abababab\ldots ab$ | $(5000 \times ab)$ | 368 | 27.2 : 1 |
| $aaa\ldots abbb\ldots b$ | $(5000 \times a, 5000 \times b)$ | 376 | 26.6 : 1 |
| $abbaababba\ldots$ | $(1000 \times abbaababba)$ | 488 | 20.5 : 1 |
| $aaabbabbabb\ldots$ | $(\pi, 0\text{–}4 \mapsto a, 5\text{–}9 \mapsto b)$ | 13416 | 0.74 : 1 |

$\pi$ follows a rule but isn't compressed!

Maybe it's just gzip? It would be possible to create to *special program* to compress $\pi$ into a short file.

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
**Data Compression**
Information vs. Complexity (contd.)

## Compression

```
echo <x> | gzip - | wc -c        # multiply by 8 for bits
```

| **Source string, $x$** | | $\ell(C(x))$ | **ratio** |
|---|---|---|---|
| $aaa \ldots a$ | $(10000 \times a)$ | 368 | 27.2 : 1. |
| $aabaabbbbabbbbb \ldots$ | (10000 random letters) | 13456 | 0.74 : 1 |
| $abababab \ldots ab$ | $(5000 \times ab)$ | 368 | 27.2 : 1 |
| $aaa \ldots abbb \ldots b$ | $(5000 \times a, 5000 \times b)$ | 376 | 26.6 : 1 |
| $abbaababba \ldots$ | $(1000 \times abbaababba)$ | 488 | 20.5 : 1 |
| $aaabbabbabb \ldots$ | $(\pi, 0\text{–}4 \mapsto a, 5\text{–}9 \mapsto b)$ | 13416 | 0.74 : 1 |

$\pi$ follows a rule but isn't compressed!

Maybe it's just gzip? It would be possible to create to *special program* to compress $\pi$ into a short file.

But what does it mean to compress an *individual* string???

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Information

An individual string is "simple" (as opposed to "complex") if it can be compressed into a small file by a *prespecified* program.

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Information

An individual string is "simple" (as opposed to "complex") if it can be compressed into a small file by a *prespecified* program.

But which program? gzip is not good for images (or for $\pi$).

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Information

An individual string is "simple" (as opposed to "complex") if it can be compressed into a small file by a *prespecified* program.

But which program? `gzip` is not good for images (or for $\pi$).

We can use several compressors if we prefix the code string by an index of the used program.

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Information

An individual string is "simple" (as opposed to "complex") if it can be compressed into a small file by a *prespecified* program.

But which program? gzip is not good for images (or for $\pi$).

We can use several compressors if we prefix the code string by an index of the used program.

How about new compressors?

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Information

An individual string is "simple" (as opposed to "complex") if it can be compressed into a small file by a *prespecified* program.

But which program? gzip is not good for images (or for $\pi$).

We can use several compressors if we prefix the code string by an index of the used program.

How about new compressors? *Self-extracting files!*

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
**Information vs. Complexity (contd.)**

# Information

An individual string is "simple" (as opposed to "complex") if it can be compressed into a small file by a *prespecified* program.

But which program? `gzip` is not good for images (or for $\pi$).

We can use several compressors if we prefix the code string by an index of the used program.

How about new compressors? *Self-extracting files!*

Can it be made automatic? Find the shortest program to print $x$.

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Information

An individual string is "simple" (as opposed to "complex") if it can be compressed into a small file by a *prespecified* program.

But which program? `gzip` is not good for images (or for $\pi$).

We can use several compressors if we prefix the code string by an index of the used program.

How about new compressors? *Self-extracting files!*

Can it be made automatic? Find the shortest program to print $x$. **No.** *Kolmogorov complexity.*

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
**Information vs. Complexity (contd.)**

## Information

An individual string is "simple" (as opposed to "complex") if it can be compressed into a small file by a *prespecified* program.

But which program? `gzip` is not good for images (or for $\pi$).

We can use several compressors if we prefix the code string by an index of the used program.

How about new compressors? *Self-extracting files!*

Can it be made automatic? Find the shortest program to print $x$.
**No.** *Kolmogorov complexity.*

Project

Outline
Administrative issues
Overview of Contents
**Compression**

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Next lecture

On Friday:

- Brief excursion to *noisy* channel coding (error correction).

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Next lecture

On Friday:

- Brief excursion to *noisy* channel coding (error correction).

Next week:

- First exercises (will be posted on website tomorrow),

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

## Next lecture

On Friday:

- Brief excursion to *noisy* channel coding (error correction).

Next week:

- First exercises (will be posted on website tomorrow),
- Due Tuesday (or Friday?)

Outline
Administrative issues
Overview of Contents
Compression

Dots and Dashes
Codes as Mappings
Data Compression
Information vs. Complexity (contd.)

# Next lecture

On Friday:

- Brief excursion to *noisy* channel coding (error correction).

Next week:

- First exercises (will be posted on website tomorrow),
- Due Tuesday (or Friday?)
- Exercise sessions by Anupam.