



# Information in Spoken Language

## A quantitative approach

Rob van Son

Chair of Phonetic Sciences  
ACLC  
University of Amsterdam

LOT winterschool 2006



AMSTERDAM CENTER  
FOR LANGUAGE AND  
COMMUNICATION





## 1 Introduction to Information Theory

- Introduction
- Probability distributions
- Bayesian probabilities
- Information and probabilities
- Relative entropy
- Compression
- Markov Chains
- Maximum Entropy
- Bibliography



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Statistics is the bookkeeping of information

- Language is about communication
- Communication implies a message
- A message is only useful if it is “surprising” to some extent
- That is, the receiver must be uncertain about the content of the message
- Information and probability quantify uncertainty
- Information is the more fundamental concept



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Statistics is the bookkeeping of information

- Language is about communication
- Communication implies a message
- A message is only useful if it is “surprising” to some extent
- That is, the receiver must be uncertain about the content of the message
- Information and probability quantify uncertainty
- Information is the more fundamental concept



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Statistics is the bookkeeping of information

- Language is about communication
- Communication implies a message
- A message is only useful if it is “surprising” to some extent
- That is, the receiver must be uncertain about the content of the message
- Information and probability quantify uncertainty
- Information is the more fundamental concept



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Statistics is the bookkeeping of information

- Language is about communication
- Communication implies a message
- A message is only useful if it is “surprising” to some extent
- That is, the receiver must be uncertain about the content of the message
- Information and probability quantify uncertainty
- Information is the more fundamental concept



## Statistics is the bookkeeping of information

- Language is about communication
- Communication implies a message
- A message is only useful if it is “surprising” to some extent
- That is, the receiver must be uncertain about the content of the message
- Information and probability quantify uncertainty
- Information is the more fundamental concept



## Statistics is the bookkeeping of information

- Language is about communication
- Communication implies a message
- A message is only useful if it is “surprising” to some extent
- That is, the receiver must be uncertain about the content of the message
- Information and probability quantify uncertainty
- Information is the more fundamental concept





**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Probability is

- A measure of the frequency of outcomes
- A measure of chance given what is known
- A number between 0 and 1 (inclusive)
- A measure of our knowledge (or ignorance)
- Boring?

[Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas] [Schneider(1999)] [MacKay(2003)]



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Probability is

- A measure of the frequency of outcomes
- A measure of chance given what is known
- A number between 0 and 1 (inclusive)
- A measure of our knowledge (or ignorance)
- Boring?

[Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas] [Schneider(1999)] [MacKay(2003)]



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Probability is

- A measure of the frequency of outcomes
- A measure of chance given what is known
- A number between 0 and 1 (inclusive)
- A measure of our knowledge (or ignorance)
- Boring?

[Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas] [Schneider(1999)] [MacKay(2003)]



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Probability is

- A measure of the frequency of outcomes
- A measure of chance given what is known
- A number between 0 and 1 (inclusive)
- A measure of our knowledge (or ignorance)
- Boring?

[Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas] [Schneider(1999)] [MacKay(2003)]



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

## Probability is

- A measure of the frequency of outcomes
- A measure of chance given what is known
- A number between 0 and 1 (inclusive)
- A measure of our knowledge (or ignorance)
- Boring?

[Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas] [Schneider(1999)] [MacKay(2003)]

# Introduction: Axioms



Information in  
Speech

Introduction  
to Information  
Theory

## Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

Markov Chains

Maximum  
Entropy

Bibliography

Probability: if  $E_1, \dots, E_n$  are possible outcomes of an observation, then  $P(E_i)$  is the probability of outcome  $E_i$  iff

1  $0 \leq P(E_i) \leq 1$

2  $P(E_1 \vee \dots \vee E_i \vee \dots \vee E_n) = 1$

3 Additivity:  $P(E_1 \vee E_2) = P(E_1) + P(E_2)$   
where  $E_1$  and  $E_2$  are mutually exclusive.

4 Countable additivity:

$$P(\bigcup_{i=1}^n E_i) = \sum_{i=1}^n P(E_i) \text{ for } n = 1, 2, \dots, N$$

where  $E_1, E_2, \dots$  are mutually exclusive.



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

Markov Chains

Maximum  
Entropy

Bibliography

Probability: if  $E_1, \dots, E_n$  are possible outcomes of an observation, then  $P(E_i)$  is the probability of outcome  $E_i$  iff

1  $0 \leq P(E_i) \leq 1$

2  $P(E_1 \vee \dots \vee E_i \vee \dots \vee E_n) = 1$

3 Additivity:  $P(E_1 \vee E_2) = P(E_1) + P(E_2)$   
where  $E_1$  and  $E_2$  are mutually exclusive.

4 Countable additivity:

$$P(\bigcup_{i=1}^n E_i) = \sum_{i=1}^n P(E_i) \text{ for } n = 1, 2, \dots, N$$

where  $E_1, E_2, \dots$  are mutually exclusive.

# Introduction: Axioms



## Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

Markov Chains

Maximum  
Entropy

Bibliography

Probability: if  $E_1, \dots, E_n$  are possible outcomes of an observation, then  $P(E_i)$  is the probability of outcome  $E_i$  iff

1  $0 \leq P(E_i) \leq 1$

2  $P(E_1 \vee \dots \vee E_i \vee \dots \vee E_n) = 1$

3 **Additivity:**  $P(E_1 \vee E_2) = P(E_1) + P(E_2)$   
where  $E_1$  and  $E_2$  are mutually exclusive.

4 Countable additivity:

$$P(\bigcup_{i=1}^n E_i) = \sum_{i=1}^n P(E_i) \text{ for } n = 1, 2, \dots, N$$

where  $E_1, E_2, \dots$  are mutually exclusive.





## Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

Markov Chains

Maximum  
Entropy

Bibliography

Probability: if  $E_1, \dots, E_n$  are possible outcomes of an observation, then  $P(E_i)$  is the probability of outcome  $E_i$  iff

①  $0 \leq P(E_i) \leq 1$

②  $P(E_1 \vee \dots \vee E_i \vee \dots \vee E_n) = 1$

③ Additivity:  $P(E_1 \vee E_2) = P(E_1) + P(E_2)$   
where  $E_1$  and  $E_2$  are mutually exclusive.

④ Countable additivity:

$$P(\bigcup_{i=1}^n E_i) = \sum_{i=1}^n P(E_i) \text{ for } n = 1, 2, \dots, N$$

where  $E_1, E_2, \dots$  are mutually exclusive.

# Introduction: Conventions



Take three observables: Color, Flower, and Place

The following conventions will be used:

- 1  $P(C = \text{Red})$ : the probability of seeing a *Red flower*
- 2  $P(C = \text{Red}, F = \text{Rose})$ : the probability of seeing a *Red Rose*
- 3  $P(C = \text{Red} | F = \text{Rose})$ : the probability of seeing a *Red Flower*, given that the flower is a *Rose*
- 4  $P(C = \text{Red}, F = \text{Rose} | P = \text{Flower Shop}) \leq 1$
- 5  $P(\text{Red} \vee \text{Blue}) = P(\text{Red}) + P(\text{Blue})$   
Basic sum rule for probabilities
- 6  $P(C, F | P) = P(C | F, P) \cdot P(F | P) = P(F | C, P) \cdot P(C | P)$   
The basic product rule for probabilities

**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

Markov Chains

Maximum  
Entropy

Bibliography

# Introduction: Conventions



**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

Take three observables: Color, Flower, and Place

The following conventions will be used:

- 1  $P(C = \text{Red})$ : the probability of seeing a *Red flower*
- 2  $P(C = \text{Red}, F = \text{Rose})$ : the probability of seeing a *Red Rose*
- 3  $P(C = \text{Red} | F = \text{Rose})$ : the probability of seeing a *Red Flower*, given that the flower is a *Rose*
- 4  $P(C = \text{Red}, F = \text{Rose} | P = \text{Flower Shop}) \leq 1$
- 5  $P(\text{Red} \vee \text{Blue}) = P(\text{Red}) + P(\text{Blue})$   
Basic sum rule for probabilities
- 6  $P(C, F | P) = P(C | F, P) \cdot P(F | P) = P(F | C, P) \cdot P(C | P)$   
The basic product rule for probabilities

# Introduction: Conventions



Take three observables: Color, Flower, and Place

The following conventions will be used:

- 1  $P(C = \text{Red})$ : the probability of seeing a *Red flower*
- 2  $P(C = \text{Red}, F = \text{Rose})$ : the probability of seeing a *Red Rose*
- 3  $P(C = \text{Red} | F = \text{Rose})$ : the probability of seeing a *Red Flower*, given that the flower is a *Rose*
- 4  $P(C = \text{Red}, F = \text{Rose} | P = \text{Flower Shop}) \leq 1$
- 5  $P(\text{Red} \vee \text{Blue}) = P(\text{Red}) + P(\text{Blue})$   
Basic sum rule for probabilities
- 6  $P(C, F | P) = P(C | F, P) \cdot P(F | P) = P(F | C, P) \cdot P(C | P)$   
The basic product rule for probabilities

# Introduction: Conventions



Take three observables: Color, Flower, and Place

The following conventions will be used:

- 1  $P(C = \textit{Red})$ : the probability of seeing a *Red flower*
- 2  $P(C = \textit{Red}, F = \textit{Rose})$ : the probability of seeing a *Red Rose*
- 3  $P(C = \textit{Red} | F = \textit{Rose})$ : the probability of seeing a *Red Flower*, given that the flower is a *Rose*
- 4  $P(C = \textit{Red}, F = \textit{Rose} | P = \textit{Flower Shop}) \leq 1$
- 5  $P(\textit{Red} \vee \textit{Blue}) = P(\textit{Red}) + P(\textit{Blue})$   
Basic sum rule for probabilities
- 6  $P(C, F | P) = P(C | F, P) \cdot P(F | P) = P(F | C, P) \cdot P(C | P)$   
The basic product rule for probabilities

**Introduction**

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

# Introduction: Conventions



Take three observables: Color, Flower, and Place

The following conventions will be used:

- 1  $P(C = \textit{Red})$ : the probability of seeing a *Red flower*
- 2  $P(C = \textit{Red}, F = \textit{Rose})$ : the probability of seeing a *Red Rose*
- 3  $P(C = \textit{Red} | F = \textit{Rose})$ : the probability of seeing a *Red Flower*, given that the flower is a *Rose*
- 4  $P(C = \textit{Red}, F = \textit{Rose} | P = \textit{Flower Shop}) \leq 1$
- 5  $P(\textit{Red} \vee \textit{Blue}) = P(\textit{Red}) + P(\textit{Blue})$   
Basic sum rule for probabilities
- 6  $P(C, F | P) = P(C | F, P) \cdot P(F | P) = P(F | C, P) \cdot P(C | P)$   
The basic product rule for probabilities

# Introduction: Conventions



Take three observables: Color, Flower, and Place

The following conventions will be used:

- 1  $P(C = \text{Red})$ : the probability of seeing a *Red flower*
- 2  $P(C = \text{Red}, F = \text{Rose})$ : the probability of seeing a *Red Rose*
- 3  $P(C = \text{Red} | F = \text{Rose})$ : the probability of seeing a *Red Flower*, given that the flower is a *Rose*
- 4  $P(C = \text{Red}, F = \text{Rose} | P = \text{Flower Shop}) \leq 1$
- 5  $P(\text{Red} \vee \text{Blue}) = P(\text{Red}) + P(\text{Blue})$   
Basic sum rule for probabilities
- 6  $P(C, F | P) = P(C | F, P) \cdot P(F | P) = P(F | C, P) \cdot P(C | P)$   
The basic product rule for probabilities



## Useful distributions, called Probability Density Functions (pdf)

- Uniform distribution, discrete and uniform
- Poisson distribution
- Normal (Gaussian) distribution
- Zipf distribution
- Mean value,  $\mu$ , is called Expected value

$$\mu = E[x] = \int_{-\infty}^{+\infty} x \cdot P(x) dx$$

- Distribution width is called *Standard Deviation* which is defined as  $\sigma = \sqrt{E[(x - E(x))^2]}$





## Useful distributions, called Probability Density Functions (pdf)

- Uniform distribution, discrete and uniform
- Poisson distribution
- Normal (Gaussian) distribution
- Zipf distribution
- Mean value,  $\mu$ , is called Expected value

$$\mu = E[x] = \int_{-\infty}^{+\infty} x \cdot P(x) dx$$

- Distribution width is called *Standard Deviation* which is defined as  $\sigma = \sqrt{E[(x - E(x))^2]}$



## Useful distributions, called Probability Density Functions (pdf)

- Uniform distribution, discrete and uniform
- Poisson distribution
- Normal (Gaussian) distribution

- Zipf distribution

- Mean value,  $\mu$ , is called Expected value

$$\mu = E[x] = \int_{-\infty}^{+\infty} x \cdot P(x) dx$$

- Distribution width is called *Standard Deviation* which is defined as  $\sigma = \sqrt{E[(x - E(x))^2]}$



## Useful distributions, called Probability Density Functions (pdf)

- Uniform distribution, discrete and uniform
- Poisson distribution
- Normal (Gaussian) distribution
- Zipf distribution

- Mean value,  $\mu$ , is called Expected value

$$\mu = E[x] = \int_{-\infty}^{+\infty} x \cdot P(x) dx$$

- Distribution width is called *Standard Deviation* which is defined as  $\sigma = \sqrt{E[(x - E(x))^2]}$



## Useful distributions, called Probability Density Functions (pdf)

- Uniform distribution, discrete and uniform
- Poisson distribution
- Normal (Gaussian) distribution
- Zipf distribution
- Mean value,  $\mu$ , is called Expected value

$$\mu = E[x] = \int_{-\infty}^{+\infty} x \cdot P(x) dx$$

- Distribution width is called *Standard Deviation* which is defined as  $\sigma = \sqrt{E[(x - E(x))^2]}$



## Useful distributions, called Probability Density Functions (pdf)

- Uniform distribution, discrete and uniform
- Poisson distribution
- Normal (Gaussian) distribution
- Zipf distribution
- Mean value,  $\mu$ , is called Expected value

$$\mu = E[x] = \int_{-\infty}^{+\infty} x \cdot P(x) dx$$

- Distribution width is called *Standard Deviation* which is defined as  $\sigma = \sqrt{E[(x - E(x))^2]}$



# Probability distributions: Uniform Discrete

$N$ , equally probable and equally spaced values  $\{E_1, \dots, E_n\}$   
(possibly if  $N \rightarrow \infty$ )

- Each category,  $E_i$ , has the same probability
- $P(E_i) = 1/N$
- Example: Dice  $\{1, \dots, 6\}$  and coins  $\{Head, Tail\}$
- Most basic distribution
- Default if only the number of values is known
- Mean  $\mu = \frac{1}{N} \sum_{i=1}^N E_i = \frac{1}{2}(E_1 + E_N)$
- Variance  $\sigma^2 = \frac{1}{12}(E_N - E_1)^2$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28discrete%29,](http://en.wikipedia.org/wiki/Uniform_distribution_%28discrete%29)

<http://gwydir.demon.co.uk/jo/probability/diceinfo.htm>

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography



# Probability distributions: Uniform Discrete

$N$ , equally probable and equally spaced values  $\{E_1, \dots, E_n\}$   
(possibly if  $N \rightarrow \infty$ )

- Each category,  $E_i$ , has the same probability
- $P(E_i) = 1/N$
- Example: Dice  $\{1, \dots, 6\}$  and coins  $\{Head, Tail\}$
- Most basic distribution
- Default if only the number of values is known
- Mean  $\mu = \frac{1}{N} \sum_{i=1}^N E_i = \frac{1}{2}(E_1 + E_N)$
- Variance  $\sigma^2 = \frac{1}{12}(E_N - E_1)^2$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28discrete%29,](http://en.wikipedia.org/wiki/Uniform_distribution_%28discrete%29)

<http://gwydir.demon.co.uk/jo/probability/diceinfo.htm>

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography



# Probability distributions: Uniform Discrete

$N$ , equally probable and equally spaced values  $\{E_1, \dots, E_n\}$   
(possibly if  $N \rightarrow \infty$ )

- Each category,  $E_i$ , has the same probability
- $P(E_i) = 1/N$
- **Example: Dice  $\{1, \dots, 6\}$  and coins  $\{Head, Tail\}$**
- Most basic distribution
- Default if only the number of values is known
- Mean  $\mu = \frac{1}{N} \sum_{i=1}^N E_i = \frac{1}{2}(E_1 + E_N)$
- Variance  $\sigma^2 = \frac{1}{12}(E_N - E_1)^2$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28discrete%29,](http://en.wikipedia.org/wiki/Uniform_distribution_%28discrete%29)

<http://gwydir.demon.co.uk/jo/probability/diceinfo.htm>

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography





# Probability distributions: Uniform Discrete

$N$ , equally probable and equally spaced values  $\{E_1, \dots, E_n\}$   
(possibly if  $N \rightarrow \infty$ )

- Each category,  $E_i$ , has the same probability
- $P(E_i) = 1/N$
- Example: Dice  $\{1, \dots, 6\}$  and coins  $\{Head, Tail\}$
- **Most basic distribution**
- Default if only the number of values is known
- Mean  $\mu = \frac{1}{N} \sum_{i=1}^N E_i = \frac{1}{2}(E_1 + E_N)$
- Variance  $\sigma^2 = \frac{1}{12}(E_N - E_1)^2$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28discrete%29,](http://en.wikipedia.org/wiki/Uniform_distribution_%28discrete%29)

<http://gwydir.demon.co.uk/jo/probability/diceinfo.htm>

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography



# Probability distributions: Uniform Discrete

$N$ , equally probable and equally spaced values  $\{E_1, \dots, E_n\}$   
(possibly if  $N \rightarrow \infty$ )

- Each category,  $E_i$ , has the same probability
- $P(E_i) = 1/N$
- Example: Dice  $\{1, \dots, 6\}$  and coins  $\{Head, Tail\}$
- Most basic distribution
- **Default if only the number of values is known**

- Mean  $\mu = \frac{1}{N} \sum_{i=1}^N E_i = \frac{1}{2}(E_1 + E_N)$

- Variance  $\sigma^2 = \frac{1}{12}(E_N - E_1)^2$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28discrete%29,](http://en.wikipedia.org/wiki/Uniform_distribution_%28discrete%29)

<http://gwydir.demon.co.uk/jo/probability/diceinfo.htm>

Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography



# Probability distributions: Uniform Discrete

$N$ , equally probable and equally spaced values  $\{E_1, \dots, E_n\}$   
(possibly if  $N \rightarrow \infty$ )

- Each category,  $E_i$ , has the same probability
- $P(E_i) = 1/N$
- Example: Dice  $\{1, \dots, 6\}$  and coins  $\{Head, Tail\}$
- Most basic distribution
- Default if only the number of values is known
- Mean  $\mu = \frac{1}{N} \sum_{i=1}^N E_i = \frac{1}{2}(E_1 + E_N)$
- Variance  $\sigma^2 = \frac{1}{12}(E_N - E_1)^2$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28discrete%29,](http://en.wikipedia.org/wiki/Uniform_distribution_%28discrete%29)

<http://gwydir.demon.co.uk/jo/probability/diceinfo.htm>

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography



# Probability distributions: Uniform Discrete

$N$ , equally probable and equally spaced values  $\{E_1, \dots, E_n\}$   
(possibly if  $N \rightarrow \infty$ )

- Each category,  $E_i$ , has the same probability
- $P(E_i) = 1/N$
- Example: Dice  $\{1, \dots, 6\}$  and coins  $\{Head, Tail\}$
- Most basic distribution
- Default if only the number of values is known
- Mean  $\mu = \frac{1}{N} \sum_{i=1}^N E_i = \frac{1}{2}(E_1 + E_N)$
- Variance  $\sigma^2 = \frac{1}{12}(E_N - E_1)^2$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28discrete%29,](http://en.wikipedia.org/wiki/Uniform_distribution_%28discrete%29)

<http://gwydir.demon.co.uk/jo/probability/diceinfo.htm>

Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

# Probability distributions: Uniform Continuous



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

Equally probable values in interval  $[a, b]$

- Pdf  $f(x) = \frac{1}{b-a}$
- Most basic distribution (continuous case)
- Default if only the range is known
- Mean  $\mu = \frac{a+b}{2}$
- Variance  $\sigma^2 = \frac{(b-a)^2}{12}$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28continuous%29](http://en.wikipedia.org/wiki/Uniform_distribution_%28continuous%29)

# Probability distributions: Uniform Continuous



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

Equally probable values in interval  $[a, b]$

- Pdf  $f(x) = \frac{1}{b-a}$
- **Most basic distribution (continuous case)**
- Default if only the range is known
- Mean  $\mu = \frac{a+b}{2}$
- Variance  $\sigma^2 = \frac{(b-a)^2}{12}$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28continuous%29](http://en.wikipedia.org/wiki/Uniform_distribution_%28continuous%29)

# Probability distributions: Uniform Continuous



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

Equally probable values in interval  $[a, b]$

- Pdf  $f(x) = \frac{1}{b - a}$
- Most basic distribution (continuous case)
- **Default if only the range is known**
- Mean  $\mu = \frac{a + b}{2}$
- Variance  $\sigma^2 = \frac{(b - a)^2}{12}$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28continuous%29](http://en.wikipedia.org/wiki/Uniform_distribution_%28continuous%29)

# Probability distributions: Uniform Continuous



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

Equally probable values in interval  $[a, b]$

- Pdf  $f(x) = \frac{1}{b-a}$
- Most basic distribution (continuous case)
- Default if only the range is known
- Mean  $\mu = \frac{a+b}{2}$
- Variance  $\sigma^2 = \frac{(b-a)^2}{12}$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28continuous%29](http://en.wikipedia.org/wiki/Uniform_distribution_%28continuous%29)



# Probability distributions: Uniform Continuous



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

Equally probable values in interval  $[a, b]$

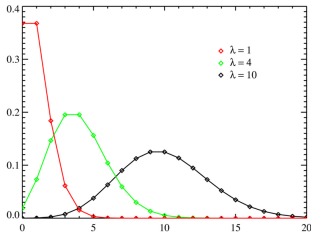
- Pdf  $f(x) = \frac{1}{b - a}$
- Most basic distribution (continuous case)
- Default if only the range is known
- Mean  $\mu = \frac{a + b}{2}$
- Variance  $\sigma^2 = \frac{(b - a)^2}{12}$

http:

[//en.wikipedia.org/wiki/Uniform\\_distribution\\_%28continuous%29](http://en.wikipedia.org/wiki/Uniform_distribution_%28continuous%29)



# Probability distributions: Poisson



$$Pdf(k; \lambda) = \frac{e^{-\lambda} \lambda^k}{k!}$$

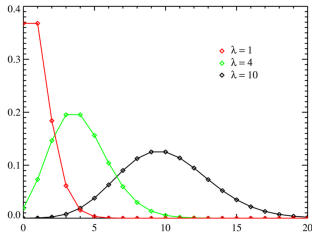
$k$ : count,  $\lambda$ : rate

## Rare events occurring with a fixed rate $\lambda$

- Mushrooms per meter of forest, typing errors per page, radio-active decay
- Average and variance are identical  $\mu = \sigma^2 = \lambda$
- Default if only an average is known



# Probability distributions: Poisson



$$Pdf(k; \lambda) = \frac{e^{-\lambda} \lambda^k}{k!}$$

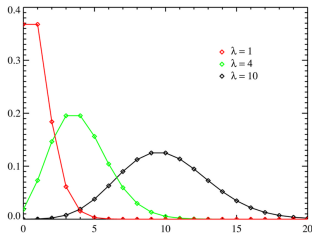
$k$ : count,  $\lambda$ : rate

## Rare events occurring with a fixed rate $\lambda$

- Mushrooms per meter of forest, typing errors per page, radio-active decay
- Average and variance are identical  $\mu = \sigma^2 = \lambda$
- Default if only an average is known



# Probability distributions: Poisson



$$Pdf(k; \lambda) = \frac{e^{-\lambda} \lambda^k}{k!}$$

$k$ : count,  $\lambda$ : rate

## Rare events occurring with a fixed rate $\lambda$

- Mushrooms per meter of forest, typing errors per page, radio-active decay
- Average and variance are identical  $\mu = \sigma^2 = \lambda$
- Default if only an average is known

# Probability distributions: Normal or Gaussian



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

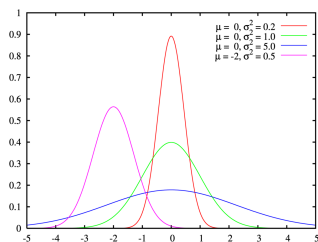
Compression

Markov Chains

Maximum

Entropy

Bibliography



$$Pdf(x; \mu, \sigma) = \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sigma\sqrt{2\pi}}$$

$x$ : observable

$\mu$ : Average

$\sigma^2$ : variance

## General measurements

- Many physical and physiological measurements, counting
- Default if both an average and a variance are known
- A sum of a large number of independent variables is approximately normal (under certain conditions)

[http://en.wikipedia.org/wiki/Normal\\_distribution](http://en.wikipedia.org/wiki/Normal_distribution)

# Probability distributions: Normal or Gaussian



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

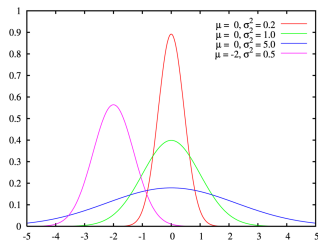
Compression

Markov Chains

Maximum

Entropy

Bibliography



$$Pdf(x; \mu, \sigma) = \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sigma\sqrt{2\pi}}$$

$x$ : observable

$\mu$ : Average

$\sigma^2$ : variance

## General measurements

- Many physical and physiological measurements, counting
- Default if both an average and a variance are known
- A sum of a large number of independent variables is approximately normal (under certain conditions)

[http://en.wikipedia.org/wiki/Normal\\_distribution](http://en.wikipedia.org/wiki/Normal_distribution)

# Probability distributions: Normal or Gaussian



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

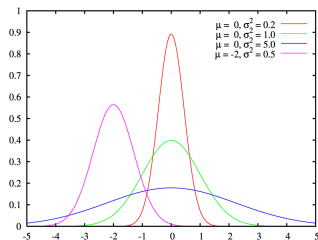
Compression

Markov Chains

Maximum

Entropy

Bibliography



$$Pdf(x; \mu, \sigma) = \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sigma\sqrt{2\pi}}$$

$x$ : observable

$\mu$ : Average

$\sigma^2$ : variance

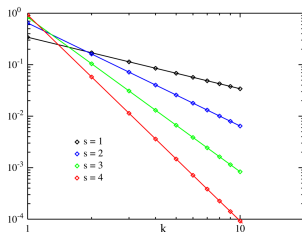
## General measurements

- Many physical and physiological measurements, counting
- Default if both an average and a variance are known
- A sum of a large number of independent variables is approximately normal (under certain conditions)

[http://en.wikipedia.org/wiki/Normal\\_distribution](http://en.wikipedia.org/wiki/Normal_distribution)



# Probability distributions: Zipf



$$Pdf(k; s, N) = \frac{1}{k^s} \frac{1}{\sum_{n=1}^N \frac{1}{n^s}}$$

$k$ : rank;  $s$ : exponent

$N$ : number of elements

note logarithmic scales

Product of frequency and rank is constant:  $f_i \approx C \cdot \frac{1}{r_i}$

- Word frequencies, city sizes, high incomes, earthquake sizes
- Default with power laws
- For word frequencies,  $s \approx 1$

[http://en.wikipedia.org/wiki/Zipf\\_distribution](http://en.wikipedia.org/wiki/Zipf_distribution)

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

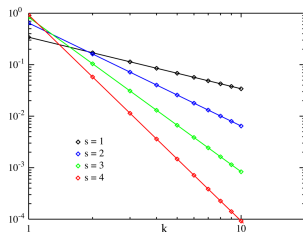
Entropy

Bibliography





# Probability distributions: Zipf



$$Pdf(k; s, N) = \frac{1}{k^s} \frac{1}{\sum_{n=1}^N \frac{1}{n^s}}$$

$k$ : rank;  $s$ : exponent

$N$ : number of elements

note logarithmic scales

Product of frequency and rank is constant:  $f_i \approx C \cdot \frac{1}{r_i}$

- Word frequencies, city sizes, high incomes, earthquake sizes
- Default with power laws
- For word frequencies,  $s \approx 1$

[http://en.wikipedia.org/wiki/Zipf\\_distribution](http://en.wikipedia.org/wiki/Zipf_distribution)

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

**Probability  
distributions**

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

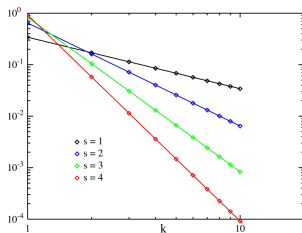
Maximum

Entropy

Bibliography



# Probability distributions: Zipf



$$Pdf(k; s, N) = \frac{1}{k^s} \frac{1}{\sum_{n=1}^N \frac{1}{n^s}}$$

$k$ : rank;  $s$ : exponent

$N$ : number of elements

note logarithmic scales

Product of frequency and rank is constant:  $f_i \approx C \cdot \frac{1}{r_i}$

- Word frequencies, city sizes, high incomes, earthquake sizes
- Default with power laws
- For word frequencies,  $s \approx 1$

[http://en.wikipedia.org/wiki/Zipf\\_distribution](http://en.wikipedia.org/wiki/Zipf_distribution)

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography



## Incorporating knowledge

- Probability is predicting outcomes from knowledge
- Explicitly formulate knowledge as probabilities
- Formalize the probability of a hypothesis
- Distinguish a priori (knowledge) and a posteriori (observations) probabilities
- Determine the information content of a single observation



## Incorporating knowledge

- Probability is predicting outcomes from knowledge
- Explicitly formulate knowledge as probabilities
- Formalize the probability of a hypothesis
- Distinguish a priori (knowledge) and a posteriori (observations) probabilities
- Determine the information content of a single observation



## Incorporating knowledge

- Probability is predicting outcomes from knowledge
- Explicitly formulate knowledge as probabilities
- Formalize the **probability of a hypothesis**
- Distinguish a priori (knowledge) and a posteriori (observations) probabilities
- Determine the information content of a single observation



## Incorporating knowledge

- Probability is predicting outcomes from knowledge
- Explicitly formulate knowledge as probabilities
- Formalize the *probability of a hypothesis*
- Distinguish **a priori** (knowledge) and **a posteriori** (observations) probabilities
- Determine the information content of a single observation



## Incorporating knowledge

- Probability is predicting outcomes from knowledge
- Explicitly formulate knowledge as probabilities
- Formalize the *probability of a hypothesis*
- Distinguish *a priori* (knowledge) and *a posteriori* (observations) probabilities
- Determine the information content of a single observation

# Bayesian probabilities

$$\begin{aligned}P(\text{Data}, \text{Hypothesis}) &= P(\text{Hypothesis}|\text{Data}) \cdot P(\text{Data}) \\ &= P(\text{Data}|\text{Hypothesis}) \cdot P(\text{Hypothesis}) \\ &\Leftrightarrow \\ P(\text{Hypothesis}|\text{Data}) &= \frac{P(\text{Data}|\text{Hypothesis}) \cdot P(\text{Hypothesis})}{P(\text{Data})}\end{aligned}$$

Express  $P(\text{Hypothesis}|\text{Data})$ :

- As a function of the measurements
- And the a priori probability of the hypothesis
- Normalized by the a priori probability of the data
- The normalization probability can often be ignored, as it will be identical for all hypotheses





# Bayesian probabilities



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

**Bayesian  
probabilities**

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

$$\begin{aligned}P(Data, Hypothesis) &= P(Hypothesis|Data) \cdot P(Data) \\ &= P(Data|Hypothesis) \cdot P(Hypothesis) \\ &\Leftrightarrow \\ P(Hypothesis|Data) &= \frac{P(Data|Hypothesis) \cdot P(Hypothesis)}{P(Data)}\end{aligned}$$

Express  $P(Hypothesis|Data)$ :

- As a function of the measurements
- And the **a priori** probability of the hypothesis
- Normalized by the a priori probability of the data
- The normalization probability can often be ignored, as it will be identical for all hypotheses

# Bayesian probabilities



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

**Bayesian  
probabilities**

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

$$\begin{aligned}P(Data, Hypothesis) &= P(Hypothesis|Data) \cdot P(Data) \\ &= P(Data|Hypothesis) \cdot P(Hypothesis) \\ &\Leftrightarrow \\ P(Hypothesis|Data) &= \frac{P(Data|Hypothesis) \cdot P(Hypothesis)}{P(Data)}\end{aligned}$$

Express  $P(Hypothesis|Data)$ :

- As a function of the measurements
- And the *a priori* probability of the hypothesis
- Normalized by the **a priori** probability of the data
- The normalization probability can often be ignored, as it will be identical for all hypotheses

# Bayesian probabilities



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

**Bayesian  
probabilities**

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

$$\begin{aligned}P(\text{Data}, \text{Hypothesis}) &= P(\text{Hypothesis}|\text{Data}) \cdot P(\text{Data}) \\ &= P(\text{Data}|\text{Hypothesis}) \cdot P(\text{Hypothesis}) \\ &\Leftrightarrow \\ P(\text{Hypothesis}|\text{Data}) &= \frac{P(\text{Data}|\text{Hypothesis}) \cdot P(\text{Hypothesis})}{P(\text{Data})}\end{aligned}$$

Express  $P(\text{Hypothesis}|\text{Data})$ :

- As a function of the measurements
- And the *a priori* probability of the hypothesis
- Normalized by the *a priori* probability of the data
- The normalization probability can often be ignored, as it will be identical for all hypotheses



# Bayesian probabilities: Toy example

Where has Watson most likely been: *Market, Garden, Meadow, Park?*

- Watson carries a **Yellow Buttercup**
- He divides his walks equally along these “*places*” (uniform prior)
- Which is most likely, obtaining a **Yellow Buttercup** in a *Market*, a *Garden*, a *Meadow*, or a *Park*?
- In formula:

$$\operatorname{argmax}_p P(p|Y, B) = \operatorname{argmax}_p P(Y, B|p) \cdot P(p)$$

where  $p \in \{\text{Market, Garden, Meadow, Park}\}$

Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

**Bayesian  
probabilities**

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography



# Bayesian probabilities: Toy example

Where has Watson most likely been: *Market, Garden, Meadow, Park?*

- Watson carries a **Yellow Buttercup**
- He divides his walks equally along these “*places*” (uniform prior)
- Which is most likely, obtaining a **Yellow Buttercup** in a *Market*, a *Garden*, a *Meadow*, or a *Park*?
- In formula:

$$\operatorname{argmax}_p P(p|Y, B) = \operatorname{argmax}_p P(Y, B|p) \cdot P(p)$$

where  $p \in \{\text{Market, Garden, Meadow, Park}\}$

Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

**Bayesian  
probabilities**

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography



# Bayesian probabilities: Toy example

Where has Watson most likely been: *Market*, *Garden*, *Meadow*, *Park*?

- Watson carries a **Yellow Buttercup**
- He divides his walks equally along these “*places*” (uniform prior)
- Which is most likely, obtaining a **Yellow Buttercup** in a *Market*, a *Garden*, a *Meadow*, or a *Park*?
- In formula:

$$\operatorname{argmax}_p P(p|Y, B) = \operatorname{argmax}_p P(Y, B|p) \cdot P(p)$$

where  $p \in \{\textit{Market}, \textit{Garden}, \textit{Meadow}, \textit{Park}\}$

Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

**Bayesian  
probabilities**

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

# Bayesian probabilities: Toy example



Where has Watson most likely been: *Market*, *Garden*, *Meadow*, *Park*?

- Watson carries a **Yellow Buttercup**
- He divides his walks equally along these “*places*” (uniform prior)
- Which is most likely, obtaining a **Yellow Buttercup** in a *Market*, a *Garden*, a *Meadow*, or a *Park*?
- In formula:

$$\operatorname{argmax}_p P(p|Y, B) = \operatorname{argmax}_p P(Y, B|p) \cdot P(p)$$

where  $p \in \{\textit{Market}, \textit{Garden}, \textit{Meadow}, \textit{Park}\}$

# Information and probabilities: Surprise!



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions  
Bayesian  
probabilities

**Information and  
probabilities**

Relative entropy  
Compression  
Markov Chains  
Maximum  
Entropy  
Bibliography

Information is a quantification of *surprise*

- Information depends on probability  $p_i$
- A more surprising observation, ie, a lower  $p_i$ , carries more information
- Information should be additive, two CD's can carry twice the information of one CD
- Define information in observation  $O_i$  with probability  $p_i$  as  $h(p_i) = -\log_2 p_i$



# Information and probabilities: Surprise!



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions  
Bayesian  
probabilities

**Information and  
probabilities**

Relative entropy  
Compression  
Markov Chains  
Maximum  
Entropy  
Bibliography

Information is a quantification of *surprise*

- Information depends on probability  $p_i$
- A more surprising observation, ie, a lower  $p_i$ , carries more information
- Information should be additive, two CD's can carry twice the information of one CD
- Define information in observation  $O_i$  with probability  $p_i$  as  $h(p_i) = -\log_2 p_i$

# Information and probabilities: Surprise!



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

**Information and  
probabilities**

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

Information is a quantification of *surprise*

- Information depends on probability  $p_i$
- A more surprising observation, ie, a lower  $p_i$ , carries more information
- Information should be additive, two CD's can carry twice the information of one CD
- Define information in observation  $O_i$  with probability  $p_i$  as  $h(p_i) = -\log_2 p_i$

# Information and probabilities: Surprise!



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

Information is a quantification of *surprise*

- Information depends on probability  $p_i$
- A more surprising observation, ie, a lower  $p_i$ , carries more information
- Information should be additive, two CD's can carry twice the information of one CD
- Define information in observation  $O_i$  with probability  $p_i$  as
$$h(p_i) = -\log_2 p_i$$

# Information and probabilities: Uncertainty



The uncertainty is the average information content and is called *Entropy*,  $H(p_1, p_2, \dots, p_n)$ . *Entropy* should be:

- **Independent** of the labeling, ie, numbering, of  $p_i$
- Decomposable, splitting a category in two gives:  
$$H'(p'_1, p''_1, \dots) = H(p_1, \dots) + p_1 \cdot H\left(\frac{p'_1}{p_1}, \frac{p''_1}{p_1}\right)$$
- Continuous, a small change in the probabilities should result in a small change in *entropy*
- Monotonic, for a uniform distribution of  $n$  items, entropy increases monotonically with the number of categories  
 $n \geq 1$
- $\Rightarrow H(p_1, p_2, \dots, p_n) = - \sum_{i=1}^n p_i \log_2(p_i)$

See chapter 1 of [Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas]

# Information and probabilities: Uncertainty



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions  
Bayesian  
probabilities

**Information and  
probabilities**

Relative entropy  
Compression  
Markov Chains  
Maximum  
Entropy  
Bibliography

The uncertainty is the average information content and is called *Entropy*,  $H(p_1, p_2, \dots, p_n)$ . *Entropy* should be:

- *Independent* of the labeling, ie, numbering, of  $p_i$
- **Decomposable**, splitting a category in two gives:

$$H'(p'_1, p''_1, \dots) = H(p_1, \dots) + p_1 \cdot H\left(\frac{p'_1}{p_1}, \frac{p''_1}{p_1}\right)$$

- Continuous, a small change in the probabilities should result in a small change in *entropy*
- Monotonic, for a uniform distribution of  $n$  items, entropy increases monotonically with the number of categories  $n \geq 1$
- $\Rightarrow H(p_1, p_2, \dots, p_n) = - \sum_{i=1}^n p_i \log_2(p_i)$

See chapter 1 of [Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas]

# Information and probabilities: Uncertainty



The uncertainty is the average information content and is called *Entropy*,  $H(p_1, p_2, \dots, p_n)$ . *Entropy* should be:

- *Independent* of the labeling, ie, numbering, of  $p_i$
- *Decomposable*, splitting a category in two gives:  
$$H'(p'_1, p''_1, \dots) = H(p_1, \dots) + p_1 \cdot H\left(\frac{p'_1}{p_1}, \frac{p''_1}{p_1}\right)$$
- **Continuous**, a small change in the probabilities should result in a small change in *entropy*
- Monotonic, for a uniform distribution of  $n$  items, entropy increases monotonically with the number of categories  
 $n \geq 1$
- $\Rightarrow H(p_1, p_2, \dots, p_n) = -\sum_{i=1}^n p_i \log_2(p_i)$

See chapter 1 of [Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas]

# Information and probabilities: Uncertainty



The uncertainty is the average information content and is called *Entropy*,  $H(p_1, p_2, \dots, p_n)$ . *Entropy* should be:

- *Independent* of the labeling, ie, numbering, of  $p_i$
- *Decomposable*, splitting a category in two gives:  
$$H'(p'_1, p''_1, \dots) = H(p_1, \dots) + p_1 \cdot H\left(\frac{p'_1}{p_1}, \frac{p''_1}{p_1}\right)$$
- *Continuous*, a small change in the probabilities should result in a small change in *entropy*
- **Monotonic**, for a uniform distribution of  $n$  items, entropy increases monotonically with the number of categories  
 $n \geq 1$

- $\Rightarrow H(p_1, p_2, \dots, p_n) = - \sum_{i=1}^n p_i \log_2(p_i)$

See chapter 1 of [Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas]

# Information and probabilities: Uncertainty



The uncertainty is the average information content and is called *Entropy*,  $H(p_1, p_2, \dots, p_n)$ . *Entropy* should be:

- *Independent* of the labeling, ie, numbering, of  $p_i$
- *Decomposable*, splitting a category in two gives:  
$$H'(p'_1, p''_1, \dots) = H(p_1, \dots) + p_1 \cdot H\left(\frac{p'_1}{p_1}, \frac{p''_1}{p_1}\right)$$
- *Continuous*, a small change in the probabilities should result in a small change in *entropy*
- *Monotonic*, for a uniform distribution of  $n$  items, entropy increases monotonically with the number of categories  
 $n \geq 1$
- $\Rightarrow H(p_1, p_2, \dots, p_n) = - \sum_{i=1}^n p_i \log_2(p_i)$

See chapter 1 of [Bavaud et al.(2005)Bavaud, Chappelier, and Kohlas]





## Probability distributions have *entropies*: Examples

- Discrete Uniform distribution:  $H\left(\frac{1}{N}\right) = \log_2(N)$

- Continuous Uniform distribution  $[a, b]$ :

$$H\left(\frac{1}{b-a}\right) = \log_2(b-a)$$

- Poisson distribution:

$$H(k; \lambda) = \lambda \left[ \frac{1}{\ln(2)} - \log_2(\lambda) \right] + e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k \log_2(k!)}{k!}$$

- Normal (Gaussian) distribution:

$$H(x; \mu, \sigma) = \log_2(\sigma \sqrt{2\pi e})$$

- Zipf distribution ( $C_{N,s} = \sum_{k=1}^N \frac{1}{k^s}$ ):

$$H(k; s, N) = \frac{s}{C_{N,s}} \sum_{k=1}^N \frac{\log_2(k)}{k^s} + \log_2(C_{N,s})$$



## Probability distributions have *entropies*: Examples

- Discrete Uniform distribution:  $H\left(\frac{1}{N}\right) = \log_2(N)$
- Continuous Uniform distribution  $[a, b]$ :

$$H\left(\frac{1}{b-a}\right) = \log_2(b-a)$$

- Poisson distribution:

$$H(k; \lambda) = \lambda \left[ \frac{1}{\ln(2)} - \log_2(\lambda) \right] + e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k \log_2(k!)}{k!}$$

- Normal (Gaussian) distribution:

$$H(x; \mu, \sigma) = \log_2(\sigma \sqrt{2\pi e})$$

- Zipf distribution ( $C_{N,s} = \sum_{k=1}^N \frac{1}{k^s}$ ):

$$H(k; s, N) = \frac{s}{C_{N,s}} \sum_{k=1}^N \frac{\log_2(k)}{k^s} + \log_2(C_{N,s})$$



## Probability distributions have *entropies*: Examples

- Discrete Uniform distribution:  $H\left(\frac{1}{N}\right) = \log_2(N)$
- Continuous Uniform distribution  $[a, b]$ :

$$H\left(\frac{1}{b-a}\right) = \log_2(b-a)$$

- Poisson distribution:

$$H(k; \lambda) = \lambda \left[ \frac{1}{\ln(2)} - \log_2(\lambda) \right] + e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k \log_2(k!)}{k!}$$

- Normal (Gaussian) distribution:

$$H(x; \mu, \sigma) = \log_2(\sigma \sqrt{2\pi e})$$

- Zipf distribution ( $C_{N,s} = \sum_{k=1}^N \frac{1}{k^s}$ ):

$$H(k; s, N) = \frac{s}{C_{N,s}} \sum_{k=1}^N \frac{\log_2(k)}{k^s} + \log_2(C_{N,s})$$



## Probability distributions have *entropies*: Examples

- Discrete Uniform distribution:  $H\left(\frac{1}{N}\right) = \log_2(N)$
- Continuous Uniform distribution  $[a, b]$ :

$$H\left(\frac{1}{b-a}\right) = \log_2(b-a)$$

- Poisson distribution:

$$H(k; \lambda) = \lambda \left[ \frac{1}{\ln(2)} - \log_2(\lambda) \right] + e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k \log_2(k!)}{k!}$$

- Normal (Gaussian) distribution:

$$H(x; \mu, \sigma) = \log_2(\sigma \sqrt{2\pi e})$$

- Zipf distribution ( $C_{N,s} = \sum_{k=1}^N \frac{1}{k^s}$ ):

$$H(k; s, N) = \frac{s}{C_{N,s}} \sum_{k=1}^N \frac{\log_2(k)}{k^s} + \log_2(C_{N,s})$$



## Probability distributions have *entropies*: Examples

- Discrete Uniform distribution:  $H\left(\frac{1}{N}\right) = \log_2(N)$

- Continuous Uniform distribution  $[a, b]$ :

$$H\left(\frac{1}{b-a}\right) = \log_2(b-a)$$

- Poisson distribution:

$$H(k; \lambda) = \lambda \left[ \frac{1}{\ln(2)} - \log_2(\lambda) \right] + e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k \log_2(k!)}{k!}$$

- Normal (Gaussian) distribution:

$$H(x; \mu, \sigma) = \log_2(\sigma \sqrt{2\pi e})$$

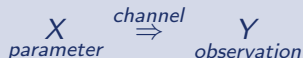
- Zipf distribution ( $C_{N,s} = \sum_{k=1}^N \frac{1}{k^s}$ ):

$$H(k; s, N) = \frac{s}{C_{N,s}} \sum_{k=1}^N \frac{\log_2(k)}{k^s} + \log_2(C_{N,s})$$

# Information and probabilities: Measuring



Information is the reduction of uncertainty

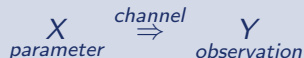


- Entropy in  $X$  before the observation:  $H(X)$
- Entropy after the observation of value of  $Y$ :  $H(X|Y)$
- Average information gained through observing  $Y$ :  
 $I(X|Y) = H(X) - H(X|Y)$
- If there is *no* uncertainty left after observing  $Y$ , ie,  
 $H(X|Y) = 0$ :  $I(X|Y) = H(X)$
- If  $X$  and  $Y$  are independent, ie,  $H(X|Y) = H(X)$ , then  
 $I(X|Y) = 0$
- Always,  $H(X|Y) \leq H(X) \Rightarrow I(X|Y) \leq H(X)$
- It is common to use  $H(\cdot)$  as a synonym of  $I(\cdot)$

# Information and probabilities: Measuring



Information is the reduction of uncertainty



- Entropy in  $X$  before the observation:  $H(X)$
- Entropy after the observation of value of  $Y$ :  $H(X|Y)$
- Average information gained through observing  $Y$ :  
 $I(X|Y) = H(X) - H(X|Y)$
- If there is *no* uncertainty left after observing  $Y$ , ie,  
 $H(X|Y) = 0$ :  $I(X|Y) = H(X)$
- If  $X$  and  $Y$  are independent, ie,  $H(X|Y) = H(X)$ , then  
 $I(X|Y) = 0$
- Always,  $H(X|Y) \leq H(X) \Rightarrow I(X|Y) \leq H(X)$
- It is common to use  $H(\cdot)$  as a synonym of  $I(\cdot)$

# Information and probabilities: Measuring



Information is the reduction of uncertainty



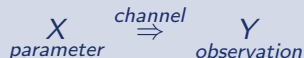
- Entropy in  $X$  before the observation:  $H(X)$
- Entropy after the observation of value of  $Y$ :  $H(X|Y)$
- Average information gained through observing  $Y$ :  
 $I(X|Y) = H(X) - H(X|Y)$
- If there is *no* uncertainty left after observing  $Y$ , ie,  
 $H(X|Y) = 0$ :  $I(X|Y) = H(X)$
- If  $X$  and  $Y$  are independent, ie,  $H(X|Y) = H(X)$ , then  
 $I(X|Y) = 0$
- Always,  $H(X|Y) \leq H(X) \Rightarrow I(X|Y) \leq H(X)$
- It is common to use  $H(\cdot)$  as a synonym of  $I(\cdot)$



# Information and probabilities: Measuring



Information is the reduction of uncertainty



- Entropy in  $X$  before the observation:  $H(X)$
- Entropy after the observation of value of  $Y$ :  $H(X|Y)$
- Average information gained through observing  $Y$ :  
 $I(X|Y) = H(X) - H(X|Y)$
- If there is *no* uncertainty left after observing  $Y$ , ie,  
 $H(X|Y) = 0$ :  $I(X|Y) = H(X)$
- If  $X$  and  $Y$  are independent, ie,  $H(X|Y) = H(X)$ , then  
 $I(X|Y) = 0$
- Always,  $H(X|Y) \leq H(X) \Rightarrow I(X|Y) \leq H(X)$
- It is common to use  $H(\cdot)$  as a synonym of  $I(\cdot)$

# Information and probabilities: Measuring



Information is the reduction of uncertainty

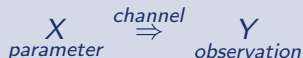


- Entropy in  $X$  before the observation:  $H(X)$
- Entropy after the observation of value of  $Y$ :  $H(X|Y)$
- Average information gained through observing  $Y$ :  
 $I(X|Y) = H(X) - H(X|Y)$
- If there is *no* uncertainty left after observing  $Y$ , ie,  
 $H(X|Y) = 0$ :  $I(X|Y) = H(X)$
- If  $X$  and  $Y$  are independent, ie,  $H(X|Y) = H(X)$ , then  
 $I(X|Y) = 0$
- Always,  $H(X|Y) \leq H(X) \Rightarrow I(X|Y) \leq H(X)$
- It is common to use  $H(\cdot)$  as a synonym of  $I(\cdot)$

# Information and probabilities: Measuring



Information is the reduction of uncertainty



- Entropy in  $X$  before the observation:  $H(X)$
- Entropy after the observation of value of  $Y$ :  $H(X|Y)$
- Average information gained through observing  $Y$ :  
 $I(X|Y) = H(X) - H(X|Y)$
- If there is *no* uncertainty left after observing  $Y$ , ie,  
 $H(X|Y) = 0$ :  $I(X|Y) = H(X)$
- If  $X$  and  $Y$  are independent, ie,  $H(X|Y) = H(X)$ , then  
 $I(X|Y) = 0$
- Always,  $H(X|Y) \leq H(X) \Rightarrow I(X|Y) \leq H(X)$
- It is common to use  $H(\cdot)$  as a synonym of  $I(\cdot)$

# Information and probabilities: Measuring



Information is the reduction of uncertainty



- Entropy in  $X$  before the observation:  $H(X)$
- Entropy after the observation of value of  $Y$ :  $H(X|Y)$
- Average information gained through observing  $Y$ :  
 $I(X|Y) = H(X) - H(X|Y)$
- If there is *no* uncertainty left after observing  $Y$ , ie,  
 $H(X|Y) = 0$ :  $I(X|Y) = H(X)$
- If  $X$  and  $Y$  are independent, ie,  $H(X|Y) = H(X)$ , then  
 $I(X|Y) = 0$
- Always,  $H(X|Y) \leq H(X) \Rightarrow I(X|Y) \leq H(X)$
- It is common to use  $H(\cdot)$  as a synonym of  $I(\cdot)$



Relative entropy:  $KL(p : q) = H(p, q) - H(p)$

$$KL(p : q) = \sum_i p_i \log_2 \frac{p_i}{q_i} \quad \vee \quad \int_{-\infty}^{\infty} p(x) \log_2 \frac{p(x)}{q(x)} dx$$

*discontinuous*                      *continuous*

$H(p, q) = \sum_i p_i \log_2 q_i$ : Cross Entropy

### Kullback-Leibler distance

- A non-symmetric divergence:  $KL(p : q) \neq KL(q : p)$
- Measures “distance” between prob. distributions
- Information gain between Prior and Posterior distribution
- Example: Word distributions as a distance between document types



Relative entropy:  $KL(p : q) = H(p, q) - H(p)$

$$KL(p : q) = \sum_i p_i \log_2 \frac{p_i}{q_i} \quad \vee \quad \int_{-\infty}^{\infty} p(x) \log_2 \frac{p(x)}{q(x)} dx$$

*discontinuous*                      *continuous*

$H(p, q) = \sum_i p_i \log_2 q_i$ : Cross Entropy

### Kullback-Leibler distance

- A non-symmetric divergence:  $KL(p : q) \neq KL(q : p)$
- Measures “distance” between prob. distributions
- Information gain between Prior and Posterior distribution
- Example: Word distributions as a distance between document types



Relative entropy:  $KL(p : q) = H(p, q) - H(p)$

$$KL(p : q) = \sum_i p_i \log_2 \frac{p_i}{q_i} \quad \vee \quad \int_{-\infty}^{\infty} p(x) \log_2 \frac{p(x)}{q(x)} dx$$

*discontinuous*                      *continuous*

$H(p, q) = \sum_i p_i \log_2 q_i$ : Cross Entropy

### Kullback-Leibler distance

- A non-symmetric divergence:  $KL(p : q) \neq KL(q : p)$
- Measures “distance” between prob. distributions
- Information gain between **Prior** and **Posterior** distribution
- Example: Word distributions as a distance between document types



Relative entropy:  $KL(p : q) = H(p, q) - H(p)$

$$KL(p : q) = \sum_i p_i \log_2 \frac{p_i}{q_i} \quad \vee \quad \int_{-\infty}^{\infty} p(x) \log_2 \frac{p(x)}{q(x)} dx$$

*discontinuous*                      *continuous*

$H(p, q) = \sum_i p_i \log_2 q_i$ : Cross Entropy

### Kullback-Leibler distance

- A non-symmetric divergence:  $KL(p : q) \neq KL(q : p)$
- Measures “distance” between prob. distributions
- Information gain between *Prior* and *Posterior* distribution
- Example: Word distributions as a distance between document types



# Compression: Minimum size



Entropy,  $H(A)$  can be understood as the minimal number of bits needed to fully *specify*  $A$  given a known production process

- In an **unknown process**,  $K(A)$  replaces  $H(A)$  as the information content
- $K(A)$ : Minimum number of bits to reconstruct  $A$
- $K(A)$  is the theoretical lower limit of compression size  $C(A)$
- Practical (lossless) compression packages,  $C(A)$ , eg, ZIP, GZIP, BZIP2 etc. never reach this limit

$K(A)$  is called the Kolmogorov complexity [Vitanyi(2005)][Chater and Vitanyi(2001)]

# Compression: Minimum size



Entropy,  $H(A)$  can be understood as the minimal number of bits needed to fully *specify*  $A$  given a *known production process*

- In an *unknown process*,  $K(A)$  replaces  $H(A)$  as the information content
- $K(A)$ : Minimum number of bits to **reconstruct**  $A$
- $K(A)$  is the theoretical lower limit of compression size  $C(A)$
- Practical (lossless) compression packages,  $C(A)$ , eg, ZIP, GZIP, BZIP2 etc. never reach this limit

$K(A)$  is called the Kolmogorov complexity [Vitanyi(2005)][Chater and Vitanyi(2001)]

# Compression: Minimum size



Entropy,  $H(A)$  can be understood as the minimal number of bits needed to fully *specify*  $A$  given a *known production process*

- In an *unknown process*,  $K(A)$  replaces  $H(A)$  as the information content
- $K(A)$ : Minimum number of bits to *reconstruct*  $A$
- $K(A)$  is the theoretical lower limit of **compression size**  $C(A)$
- Practical (lossless) compression packages,  $C(A)$ , eg, ZIP, GZIP, BZIP2 etc. never reach this limit

$K(A)$  is called the Kolmogorov complexity [Vitanyi(2005)][Chater and Vitanyi(2001)]

# Compression: Minimum size



Entropy,  $H(A)$  can be understood as the minimal number of bits needed to fully *specify*  $A$  given a *known production process*

- In an *unknown process*,  $K(A)$  replaces  $H(A)$  as the information content
- $K(A)$ : Minimum number of bits to *reconstruct*  $A$
- $K(A)$  is the theoretical lower limit of *compression size*  $C(A)$
- Practical (lossless) compression packages,  $C(A)$ , eg, ZIP, GZIP, BZIP2 etc. never reach this limit

$K(A)$  is called the Kolmogorov complexity [Vitanyi(2005)][Chater and Vitanyi(2001)]

# Compression: Similarity metric



$$\begin{aligned} NCD(A, B) &= \frac{\min\{C(A|B), C(B|A)\}}{\max\{C(A), C(B)\}} \\ &= \frac{C(AB) - \min\{C(A), C(B)\}}{\max\{C(A), C(B)\}} \end{aligned}$$

NCD: Normalized Compression Distance

## Similarity by compression

- Always  $C(AB) \leq C(A) + C(B)$  (+constant)
- Estimate entropy by suitable “long range” compression
- $K(\text{text}) \leq C(\text{text})$  in bits

<http://www.complearn.org/ncd.html>

[Chen et al.(2004)Chen, Li, Ma, and Vitányi][Vitányi(2005)][Chater and Vitányi(2001)]

# Compression: Similarity metric



$$\begin{aligned} NCD(A, B) &= \frac{\min\{C(A|B), C(B|A)\}}{\max\{C(A), C(B)\}} \\ &= \frac{C(AB) - \min\{C(A), C(B)\}}{\max\{C(A), C(B)\}} \end{aligned}$$

NCD: Normalized Compression Distance

## Similarity by compression

- Always  $C(AB) \leq C(A) + C(B)$  (+constant)
- Estimate entropy by suitable “long range” compression
- $K(\text{text}) \leq C(\text{text})$  in bits

<http://www.complearn.org/ncd.html>

[Chen et al.(2004)Chen, Li, Ma, and Vitányi][Vitányi(2005)][Chater and Vitányi(2001)]

# Compression: Similarity metric



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

**Compression**

Markov Chains

Maximum

Entropy

Bibliography

$$\begin{aligned} NCD(A, B) &= \frac{\min\{C(A|B), C(B|A)\}}{\max\{C(A), C(B)\}} \\ &= \frac{C(AB) - \min\{C(A), C(B)\}}{\max\{C(A), C(B)\}} \end{aligned}$$

NCD: Normalized Compression Distance

## Similarity by compression

- Always  $C(AB) \leq C(A) + C(B)$  (+constant)
- Estimate entropy by suitable “long range” compression
- $K(\text{text}) \leq C(\text{text})$  in bits

<http://www.complearn.org/ncd.html>

[Chen et al.(2004)Chen, Li, Ma, and Vitányi][Vitányi(2005)][Chater and Vitányi(2001)]

# Markov Chains



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

**Markov Chains**

Maximum  
Entropy

Bibliography

## Words and letters never follow each other at random

- The simplest language “model” predicts the next word based on the previous word

- Markov chain: 
$$P(w_{i+1}|w_i) = \frac{P(w_{i+1}, w_i)}{P(w_i)}$$

- Can be extended to more words
- Large amounts of text are needed to determine  $P(w_{i+1}, w_i)$  reliably

- Example Markov text:

*Step which one could go be grabbed. People to Do that my the former Netscape brand's fortunes that means indent command to The user visible displays a.*

[http://en.wikipedia.org/wiki/Markov\\_chain](http://en.wikipedia.org/wiki/Markov_chain)

Generate texts: <http://www.jwz.org/dadadodo/>





## Words and letters never follow each other at random

- The simplest language “model” predicts the next word based on the previous word

- Markov chain: 
$$P(w_{i+1}|w_i) = \frac{P(w_{i+1}, w_i)}{P(w_i)}$$

- Can be extended to more words
- Large amounts of text are needed to determine  $P(w_{i+1}, w_i)$  reliably

- Example Markov text:

*Step which one could go be grabbed. People to Do that my the former Netscape brand's fortunes that means indent command to The user visible displays a.*

[http://en.wikipedia.org/wiki/Markov\\_chain](http://en.wikipedia.org/wiki/Markov_chain)

Generate texts: <http://www.jwz.org/dadadodo/>



## Words and letters never follow each other at random

- The simplest language “model” predicts the next word based on the previous word

- Markov chain: 
$$P(w_{i+1}|w_i) = \frac{P(w_{i+1}, w_i)}{P(w_i)}$$

- Can be extended to more words
- Large amounts of text are needed to determine  $P(w_{i+1}, w_i)$  reliably

- Example Markov text:

*Step which one could go be grabbed. People to Do that my the former Netscape brand's fortunes that means indent command to The user visible displays a.*

[http://en.wikipedia.org/wiki/Markov\\_chain](http://en.wikipedia.org/wiki/Markov_chain)

Generate texts: <http://www.jwz.org/dadadodo/>



## Words and letters never follow each other at random

- The simplest language “model” predicts the next word based on the previous word

- Markov chain: 
$$P(w_{i+1}|w_i) = \frac{P(w_{i+1}, w_i)}{P(w_i)}$$

- Can be extended to more words
- Large amounts of text are needed to determine  $P(w_{i+1}, w_i)$  reliably

- Example Markov text:

*Step which one could go be grabbed. People to Do that my the former Netscape brand's fortunes that means indent command to The user visible displays a.*

[http://en.wikipedia.org/wiki/Markov\\_chain](http://en.wikipedia.org/wiki/Markov_chain)

Generate texts: <http://www.jwz.org/dadadodo/>



## Words and letters never follow each other at random

- The simplest language “model” predicts the next word based on the previous word

- Markov chain: 
$$P(w_{i+1}|w_i) = \frac{P(w_{i+1}, w_i)}{P(w_i)}$$

- Can be extended to more words
- Large amounts of text are needed to determine  $P(w_{i+1}, w_i)$  reliably

- Example Markov text:

*Step which one could go be grabbed. People to Do that my the former Netscape brand's fortunes that means indent command to The user visible displays a.*

[http://en.wikipedia.org/wiki/Markov\\_chain](http://en.wikipedia.org/wiki/Markov_chain)

Generate texts: <http://www.jwz.org/dadadodo/>

# Markov Chains: Language models



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

**Markov Chains**

Maximum

Entropy

Bibliography

With Markov chains, or N-grams, the probability of a sequence can be calculated

- What is the probability of encountering a sentence  $(w_1, \dots, w_n)$ ?
- A human style language model is not known
- Use *N-gram* Markov chains
- $P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$  (exact)
- $P(w_1, \dots, w_n) \approx \prod_{i=1}^n P(w_i | w_{i-N+1}, \dots, w_{i-1})$  (*N-gram* approximation)

# Markov Chains: Language models



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

**Markov Chains**

Maximum  
Entropy

Bibliography

With Markov chains, or N-grams, the probability of a sequence can be calculated

- What is the probability of encountering a sentence  $(w_1, \dots, w_n)$ ?
- A human style language model is not known
- Use *N-gram* Markov chains
- $P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$  (exact)
- $P(w_1, \dots, w_n) \approx \prod_{i=1}^n P(w_i | w_{i-N+1}, \dots, w_{i-1})$  (*N-gram* approximation)

# Markov Chains: Language models



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

**Markov Chains**

Maximum  
Entropy

Bibliography

With Markov chains, or N-grams, the probability of a sequence can be calculated

- What is the probability of encountering a sentence  $(w_1, \dots, w_n)$ ?
- A human style language model is not known
- Use *N-gram* Markov chains
- $P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$  (exact)
- $P(w_1, \dots, w_n) \approx \prod_{i=1}^n P(w_i | w_{i-N+1}, \dots, w_{i-1})$  (*N-gram* approximation)

# Markov Chains: Language models



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

**Markov Chains**

Maximum

Entropy

Bibliography

With Markov chains, or N-grams, the probability of a sequence can be calculated

- What is the probability of encountering a sentence  $(w_1, \dots, w_n)$ ?
- A human style language model is not known
- Use *N-gram* Markov chains
- $P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$  (exact)
- $P(w_1, \dots, w_n) \approx \prod_{i=1}^n P(w_i | w_{i-N+1}, \dots, w_{i-1})$  (*N-gram* approximation)



# Markov Chains: Language models



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

**Markov Chains**

Maximum  
Entropy  
Bibliography

With Markov chains, or  $N$ -grams, the probability of a sequence can be calculated

- What is the probability of encountering a sentence  $(w_1, \dots, w_n)$ ?
- A human style language model is not known
- Use  $N$ -gram Markov chains
- $P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$  (exact)
- $P(w_1, \dots, w_n) \approx \prod_{i=1}^n P(w_i | w_{i-N+1}, \dots, w_{i-1})$  ( $N$ -gram approximation)



# Markov Chains: Perplexity

Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

**Markov Chains**

Maximum  
Entropy

Bibliography

$$P_X(\text{Model}) = 2^{H_X(w_i|W_{1\dots i-1})}$$
$$H_X(w_i|W_{1\dots i-1}) = - \sum_{\{W\}} P_{\text{observed}}(w_i|\dots) \log P_{\text{model}}(w_i|\dots)$$

$H_X(\cdot)$ : Cross Entropy

Perplexity: “average” number of choices for the next word

- Matches observed with modelled word order
- A better language model has a lower perplexity
- For an  $N$ -gram Markov chain the perplexity is well defined
- Using the model entropy i.e. the cross entropy estimates the quality of the model on the training corpus

# Markov Chains: Perplexity



$$P_X(\text{Model}) = 2^{H_X(w_i|W_{1\dots i-1})}$$
$$H_X(w_i|W_{1\dots i-1}) = - \sum_{\{W\}} P_{\text{observed}}(w_i|\dots) \log P_{\text{model}}(w_i|\dots)$$

$H_X(\cdot)$ : Cross Entropy

Perplexity: “average” number of choices for the next word

- Matches observed with modelled word order
- A better language model has a lower perplexity
- For an  $N$ -gram Markov chain the perplexity is well defined
- Using the model entropy i.e. the cross entropy estimates the quality of the model on the training corpus



# Markov Chains: Perplexity

Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

**Markov Chains**

Maximum  
Entropy

Bibliography

$$P_X(\text{Model}) = 2^{H_X(w_i|W_{1\dots i-1})}$$
$$H_X(w_i|W_{1\dots i-1}) = - \sum_{\{W\}} P_{\text{observed}}(w_i|\dots) \log P_{\text{model}}(w_i|\dots)$$

$H_X(\cdot)$ : Cross Entropy

Perplexity: “average” number of choices for the next word

- Matches observed with modelled word order
- A better language model has a lower perplexity
- For an  $N$ -gram Markov chain the perplexity is well defined
- Using the model entropy i.e. the cross entropy estimates the quality of the model on the training corpus

# Markov Chains: Perplexity



$$P_X(\text{Model}) = 2^{H_X(w_i|W_{1\dots i-1})}$$
$$H_X(w_i|W_{1\dots i-1}) = - \sum_{\{W\}} P_{\text{observed}}(w_i|\dots) \log P_{\text{model}}(w_i|\dots)$$

$H_X(\cdot)$ : Cross Entropy

Perplexity: “average” number of choices for the next word

- Matches observed with modelled word order
- A better language model has a lower perplexity
- For an  $N$ -gram Markov chain the perplexity is well defined
- Using the model entropy i.e. the cross entropy estimates the quality of the model on the training corpus



# Maximum Entropy

Information in  
Speech

$$\begin{aligned} \text{find } p^* &= \underset{p \in \mathcal{C}}{\operatorname{argmax}} H(p) \\ &= \underset{p \in \mathcal{C}}{\operatorname{argmax}} \left( - \sum_{x,y} \tilde{p}(x) p(y|x) \log p(y|x) \right) \end{aligned}$$

Which model,  $p^*$ , fits my data best and by what criterium?

- Quantify all constraints (knowledge) and determine the set of possible distributions  $p \in \mathcal{C}$
- Determine the average entropy,  $H(y|x)$ , over the observed (measured) probabilities  $\tilde{p}(x)$
- The best distribution,  $p^*$  has the highest entropy

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

Markov Chains

**Maximum  
Entropy**

Bibliography

[Berger()][Berger(1996)] [Berger et al.(1996)Berger, della Pietra, and della Pietra] [Maxent()]

[Roni Rosenfeld(1996)]



# Maximum Entropy

$$\begin{aligned} \text{find } p^* &= \underset{p \in \mathcal{C}}{\operatorname{argmax}} H(p) \\ &= \underset{p \in \mathcal{C}}{\operatorname{argmax}} \left( - \sum_{x,y} \tilde{p}(x) p(y|x) \log p(y|x) \right) \end{aligned}$$

Which model,  $p^*$ , fits my data best and by what criterium?

- Quantify all constraints (knowledge) and determine the set of possible distributions  $p \in \mathcal{C}$
- Determine the average entropy,  $H(y|x)$ , over the observed (measured) probabilities  $\tilde{p}(x)$
- The best distribution,  $p^*$  has the highest entropy

[Berger()][Berger(1996)] [Berger et al.(1996)Berger, della Pietra, and della Pietra] [Maxent()]

[Roni Rosenfeld(1996)]

# Maximum Entropy



$$\begin{aligned} \text{find } p^* &= \underset{p \in \mathcal{C}}{\operatorname{argmax}} H(p) \\ &= \underset{p \in \mathcal{C}}{\operatorname{argmax}} \left( - \sum_{x,y} \tilde{p}(x) p(y|x) \log p(y|x) \right) \end{aligned}$$

Which model,  $p^*$ , fits my data best and by what criterium?

- Quantify all constraints (knowledge) and determine the set of possible distributions  $p \in \mathcal{C}$
- Determine the average entropy,  $H(y|x)$ , over the observed (measured) probabilities  $\tilde{p}(x)$
- The **best** distribution,  $p^*$  has the highest entropy





# Maximum Entropy: Kangaroo example

$\frac{1}{3}$  of all kangaroos have **blue** eyes and  $\frac{1}{3}$  are left handed

blue eyed	Left true	Handed false	tot
true	$x$	$\frac{1}{3} - x$	$\frac{1}{3}$
false	$\frac{1}{3} - x$	$\frac{1}{3} + x$	$\frac{2}{3}$
tot	$\frac{1}{3}$	$\frac{2}{3}$	$1$

How many are both **blue** eyed and left handed?

- All  $0 \leq x \leq \frac{1}{3}$  are possible
- $H(x = \frac{1}{9}) \approx 1.84$  has maximum entropy
- $x = \frac{1}{9}$  is the only solution with uncorrelated eye color and handedness

# Bibliography I



Abelard.

cause, chance and Bayesian statistics a briefing document.  
Web.

URL <http://www.abelard.org/briefings/bayes.htm>.



F. Bavaud, J.-C. Chappelier, and J. Kohlas.

*An Introduction to Information Theory and Applications.*

UniFr course, 3 September 2005.

URL <http://diuf.unifr.ch/tcs/courses/it04-05/script/information-theory.pdf>.



Adam L. Berger.

MaxEnt and Exponential Models.

Website.

URL <http://www.cs.cmu.edu/~aberger/maxent.html>.



Adam L. Berger.

A Brief Maxent Tutorial.

Web, 5 July 1996.

URL <http://www.cs.cmu.edu/afs/cs/user/aberger/www/html/tutorial/tutorial.html>.



Adam L. Berger, Stephen A. della Pietra, and Vincent J. della Pietra.

A maximum entropy approach to natural language processing.

*Computational Linguistics*, 1:22-1, March 1996.

URL <http://www.cs.cmu.edu/afs/cs/user/aberger/www/ps/compling.ps>.

# Bibliography II



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

**Bibliography**



Tom Carter.

An introduction to information theory and entropy.

Web presentation, June 2005.

URL <http://astarte.csustan.edu/~tom/SFI-CSSS/2005/info-1ec.pdf>.



Nick Chater and Paul Vitányi.

The generalized universal law of generalization.

*Arxiv.org Computer Science*, 29 January 2001.

URL <http://arxiv.org/pdf/cs.CV/0101036>.



Xin Chen, Xin Li, Bin Ma, and Paul M. B. Vitányi.

The Similarity Metric.

*IEEE TRANSACTIONS ON INFORMATION THEORY*, 50(12):3250–32364, December 2004.

URL <http://homepages.cwi.nl/~paulv/papers/similarity.pdf>.



Minh N. Do.

Fast Approximation of Kullback-Leibler Distance for Dependence Trees and Hidden Markov Models.

*IEEE Signal Processing Letters*, 10:115–118, April 2003.

URL [http://www.ifp.uiuc.edu/~minhdo/publications/KLD\\_HMM.pdf](http://www.ifp.uiuc.edu/~minhdo/publications/KLD_HMM.pdf).



Yaniv Dover.

A short account of a connection of power laws to the information entropy.

*PHYSICA A*, 334:591–597, 2004.

URL <http://arxiv.org/pdf/cond-mat/0309383>.

# Bibliography III



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

**Bibliography**



FSF.

GNU General Public License.

Web, June 1991.

URL <http://www.gnu.org/licenses/gpl.html>.



E. T. Jaynes.

Probability Theory as Logic.

In P.F. Fougère, editor, *Maximum-Entropy and Bayesian Methods*, page 1. Kluwer, Dordrecht, 1990.

URL <http://bayes.wustl.edu/etj/articles/prob.as.logic.pdf>.



E. T. Jaynes.

*Probability Theory: The logic of Science.*

1995.

URL <http://bayes.wustl.edu/etj/prob/book.pdf>.

Web publication.



Kenji Kawamura and Naomichi Hatano.

Universality of zipf's law.

Preprint archive, 2002.

URL <http://arxiv.org/pdf/cond-mat/0203455>.



David J. C. MacKay.

*Information Theory, Inference, and Learning Algorithms.*

Cambridge University Press, September 2003.

URL <http://www.inference.phy.cam.ac.uk/mackay/itila/>.

# Bibliography IV



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions  
Bayesian  
probabilities  
Information and  
probabilities  
Relative entropy  
Compression  
Markov Chains  
Maximum  
Entropy

**Bibliography**



David J.C. MacKay.

**A Short Course in Information Theory.**

Web, January 1995.

URL <http://www.inference.phy.cam.ac.uk/mackay/info-theory/course.html>.



Maxent.

**Maximum Entropy Modeling.**

Web.

URL <http://homepages.inf.ed.ac.uk/s0450736/maxent.html>.

Mainly links.



Roni Rosenfeld.

***Adaptive Statistical Language Modeling: A Maximum Entropy Approach.***

PhD thesis, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, 19 April 1994.

URL <http://www.cs.cmu.edu/~roni/papers/me-thesis-TR-94-138.pdf>.



Roni Rosenfeld.

**A maximum entropy approach to adaptive statistical language modelling.**

*Computer, Speech and Language*, 10:187–228, 1996.

URL <http://www.cs.cmu.edu/~roni/me-csl-revised.ps>.



Tom Schneider.

**Information Theory Primer, With an Appendix on Logarithms.**

Web, 9 February 1999.

URL <http://www.lecb.ncifcrf.gov/~toms/paper/primer/>.  
updated 2005.

# Bibliography V



Information in  
Speech

Introduction  
to Information  
Theory

Introduction  
Probability  
distributions  
Bayesian  
probabilities  
Information and  
probabilities  
Relative entropy  
Compression  
Markov Chains  
Maximum  
Entropy

**Bibliography**



**D. J. Strom.**

**Introduction to Bayesian Statistics.**

Web.

URL <http://bidug.pnl.gov/presentations/PEP/>.  
Slides.



**Nick Szabo.**

**Introduction to Algorithmic Information Theory.**

Web, 1996.

URL <http://szabo.best.vwh.net/kolmogorov.html>.



**H. Thornburg.**

**Introduction to Bayesian Statistics.**

Web, 2003.

URL <http://ccrma.stanford.edu/~jos/bayes/bayes.html>.



**Paul Vitanyi.**

**Universal Similarity.**

In *Proceedings of ITW2005*, 29 August 2005.

URL <http://arxiv.org/pdf/cs.IR/0504089>.  
arXiv:cs.IR/0504089.



Copyright ©2005,2006 R.J.J.H. van Son, GNU General Public License [FSF(1991)]

*This program is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2 of the License, or (at your option) any later version.*

*This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.*

*You should have received a copy of the GNU General Public License along with this program; if not, write to the Free Software Foundation, Inc., 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301, USA.*



Version 2, June 1991

Copyright © 1989, 1991 Free Software Foundation, Inc.

51 Franklin Street, Fifth Floor, Boston, MA 02110-1301, USA

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

## Preamble

The licenses for most software are designed to take away your freedom to share and change it. By contrast, the GNU General Public License is intended to guarantee your freedom to share and change free software—to make sure the software is free for all its users. This General Public License applies to most of the Free Software Foundation's software and to any other program whose authors commit to using it. (Some other Free Software Foundation software is covered by the GNU Library General Public License instead.) You can apply it to your programs, too.

When we speak of free software, we are referring to freedom, not price. Our General Public Licenses are designed to make sure that you have the freedom to distribute copies of free software (and charge for this service if you wish), that you receive source code or can get it if you want it, that you can change the software or use pieces of it in new free programs; and that you know you can do these things.

To protect your rights, we need to make restrictions that forbid anyone to deny you these rights or to ask you to surrender the rights. These restrictions translate to certain responsibilities for you if you distribute copies of the software, or if you modify it.



# Bibliography II



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum  
Entropy

Bibliography

For example, if you distribute copies of such a program, whether gratis or for a fee, you must give the recipients all the rights that you have. You must make sure that they, too, receive or can get the source code. And you must show them these terms so they know their rights.

We protect your rights with two steps: (1) copyright the software, and (2) offer you this license which gives you legal permission to copy, distribute and/or modify the software.

Also, for each author's protection and ours, we want to make certain that everyone understands that there is no warranty for this free software. If the software is modified by someone else and passed on, we want its recipients to know that what they have is not the original, so that any problems introduced by others will not reflect on the original authors' reputations.

Finally, any free program is threatened constantly by software patents. We wish to avoid the danger that redistributors of a free program will individually obtain patent licenses, in effect making the program proprietary. To prevent this, we have made it clear that any patent must be licensed for everyone's free use or not licensed at all.

The precise terms and conditions for copying, distribution and modification follow.

## TERMS AND CONDITIONS FOR COPYING, DISTRIBUTION AND MODIFICATION

# Bibliography III



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

- 0 This License applies to any program or other work which contains a notice placed by the copyright holder saying it may be distributed under the terms of this General Public License. The “Program”, below, refers to any such program or work, and a “work based on the Program” means either the Program or any derivative work under copyright law: that is to say, a work containing the Program or a portion of it, either verbatim or with modifications and/or translated into another language. (Hereinafter, translation is included without limitation in the term “modification”.) Each licensee is addressed as “you”.

Activities other than copying, distribution and modification are not covered by this License; they are outside its scope. The act of running the Program is not restricted, and the output from the Program is covered only if its contents constitute a work based on the Program (independent of having been made by running the Program). Whether that is true depends on what the Program does.

- 1 You may copy and distribute verbatim copies of the Program's source code as you receive it, in any medium, provided that you conspicuously and appropriately publish on each copy an appropriate copyright notice and disclaimer of warranty; keep intact all the notices that refer to this License and to the absence of any warranty; and give any other recipients of the Program a copy of this License along with the Program.

You may charge a fee for the physical act of transferring a copy, and you may at your option offer warranty protection in exchange for a fee.

- 2 You may modify your copy or copies of the Program or any portion of it, thus forming a work based on the Program, and copy and distribute such modifications or work under the terms of Section 1 above, provided that you also meet all of these conditions:

- 1 You must cause the modified files to carry prominent notices stating that you changed the files and the date of any change.
- 2 You must cause any work that you distribute or publish, that in whole or in part contains or is derived from the Program or any part thereof, to be licensed as a whole at no charge to all third parties under the terms of this License.

# Bibliography IV



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy  
Compression

Markov Chains

Maximum  
Entropy

Bibliography

- 3 If the modified program normally reads commands interactively when run, you must cause it, when started running for such interactive use in the most ordinary way, to print or display an announcement including an appropriate copyright notice and a notice that there is no warranty (or else, saying that you provide a warranty) and that users may redistribute the program under these conditions, and telling the user how to view a copy of this License. (Exception: if the Program itself is interactive but does not normally print such an announcement, your work based on the Program is not required to print an announcement.)

These requirements apply to the modified work as a whole. If identifiable sections of that work are not derived from the Program, and can be reasonably considered independent and separate works in themselves, then this License, and its terms, do not apply to those sections when you distribute them as separate works. But when you distribute the same sections as part of a whole which is a work based on the Program, the distribution of the whole must be on the terms of this License, whose permissions for other licensees extend to the entire whole, and thus to each and every part regardless of who wrote it.

Thus, it is not the intent of this section to claim rights or contest your rights to work written entirely by you; rather, the intent is to exercise the right to control the distribution of derivative or collective works based on the Program.

In addition, mere aggregation of another work not based on the Program with the Program (or with a work based on the Program) on a volume of a storage or distribution medium does not bring the other work under the scope of this License.

- 3 You may copy and distribute the Program (or a work based on it, under Section 2) in object code or executable form under the terms of Sections 1 and 2 above provided that you also do one of the following:
  - 1 Accompany it with the complete corresponding machine-readable source code, which must be distributed under the terms of Sections 1 and 2 above on a medium customarily used for software interchange; or,

# Bibliography V



- 2 Accompany it with a written offer, valid for at least three years, to give any third party, for a charge no more than your cost of physically performing source distribution, a complete machine-readable copy of the corresponding source code, to be distributed under the terms of Sections 1 and 2 above on a medium customarily used for software interchange; or,
- 3 Accompany it with the information you received as to the offer to distribute corresponding source code. (This alternative is allowed only for noncommercial distribution and only if you received the program in object code or executable form with such an offer, in accord with Subsection b above.)

The source code for a work means the preferred form of the work for making modifications to it. For an executable work, complete source code means all the source code for all modules it contains, plus any associated interface definition files, plus the scripts used to control compilation and installation of the executable. However, as a special exception, the source code distributed need not include anything that is normally distributed (in either source or binary form) with the major components (compiler, kernel, and so on) of the operating system on which the executable runs, unless that component itself accompanies the executable.

If distribution of executable or object code is made by offering access to copy from a designated place, then offering equivalent access to copy the source code from the same place counts as distribution of the source code, even though third parties are not compelled to copy the source along with the object code.

- 4 You may not copy, modify, sublicense, or distribute the Program except as expressly provided under this License. Any attempt otherwise to copy, modify, sublicense or distribute the Program is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.

# Bibliography VI



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum

Entropy

Bibliography

- 5 You are not required to accept this License, since you have not signed it. However, nothing else grants you permission to modify or distribute the Program or its derivative works. These actions are prohibited by law if you do not accept this License. Therefore, by modifying or distributing the Program (or any work based on the Program), you indicate your acceptance of this License to do so, and all its terms and conditions for copying, distributing or modifying the Program or works based on it.
- 6 Each time you redistribute the Program (or any work based on the Program), the recipient automatically receives a license from the original licensor to copy, distribute or modify the Program subject to these terms and conditions. You may not impose any further restrictions on the recipients' exercise of the rights granted herein. You are not responsible for enforcing compliance by third parties to this License.
- 7 If, as a consequence of a court judgment or allegation of patent infringement or for any other reason (not limited to patent issues), conditions are imposed on you (whether by court order, agreement or otherwise) that contradict the conditions of this License, they do not excuse you from the conditions of this License. If you cannot distribute so as to satisfy simultaneously your obligations under this License and any other pertinent obligations, then as a consequence you may not distribute the Program at all. For example, if a patent license would not permit royalty-free redistribution of the Program by all those who receive copies directly or indirectly through you, then the only way you could satisfy both it and this License would be to refrain entirely from distribution of the Program. If any portion of this section is held invalid or unenforceable under any particular circumstance, the balance of the section is intended to apply and the section as a whole is intended to apply in other circumstances.

It is not the purpose of this section to induce you to infringe any patents or other property right claims or to contest validity of any such claims; this section has the sole purpose of protecting the integrity of the free software distribution system, which is implemented by public license practices. Many people have made generous contributions to the wide range of software distributed through that system in reliance on consistent application of that system; it is up to the author/donor to

# Bibliography VII



decide if he or she is willing to distribute software through any other system and a licensee cannot impose that choice.

This section is intended to make thoroughly clear what is believed to be a consequence of the rest of this License.

- 8 If the distribution and/or use of the Program is restricted in certain countries either by patents or by copyrighted interfaces, the original copyright holder who places the Program under this License may add an explicit geographical distribution limitation excluding those countries, so that distribution is permitted only in or among countries not thus excluded. In such case, this License incorporates the limitation as if written in the body of this License.
- 9 The Free Software Foundation may publish revised and/or new versions of the General Public License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns.  
Each version is given a distinguishing version number. If the Program specifies a version number of this License which applies to it and "any later version", you have the option of following the terms and conditions either of that version or of any later version published by the Free Software Foundation. If the Program does not specify a version number of this License, you may choose any version ever published by the Free Software Foundation.
- 10 If you wish to incorporate parts of the Program into other free programs whose distribution conditions are different, write to the author to ask for permission. For software which is copyrighted by the Free Software Foundation, write to the Free Software Foundation; we sometimes make exceptions for this. Our decision will be guided by the two goals of preserving the free status of all derivatives of our free software and of promoting the sharing and reuse of software generally.

NO WARRANTY



- 11 BECAUSE THE PROGRAM IS LICENSED FREE OF CHARGE, THERE IS NO WARRANTY FOR THE PROGRAM, TO THE EXTENT PERMITTED BY APPLICABLE LAW. EXCEPT WHEN OTHERWISE STATED IN WRITING THE COPYRIGHT HOLDERS AND/OR OTHER PARTIES PROVIDE THE PROGRAM "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE ENTIRE RISK AS TO THE QUALITY AND PERFORMANCE OF THE PROGRAM IS WITH YOU. SHOULD THE PROGRAM PROVE DEFECTIVE, YOU ASSUME THE COST OF ALL NECESSARY SERVICING, REPAIR OR CORRECTION.
- 12 IN NO EVENT UNLESS REQUIRED BY APPLICABLE LAW OR AGREED TO IN WRITING WILL ANY COPYRIGHT HOLDER, OR ANY OTHER PARTY WHO MAY MODIFY AND/OR REDISTRIBUTE THE PROGRAM AS PERMITTED ABOVE, BE LIABLE TO YOU FOR DAMAGES, INCLUDING ANY GENERAL, SPECIAL, INCIDENTAL OR CONSEQUENTIAL DAMAGES ARISING OUT OF THE USE OR INABILITY TO USE THE PROGRAM (INCLUDING BUT NOT LIMITED TO LOSS OF DATA OR DATA BEING RENDERED INACCURATE OR LOSSES SUSTAINED BY YOU OR THIRD PARTIES OR A FAILURE OF THE PROGRAM TO OPERATE WITH ANY OTHER PROGRAMS), EVEN IF SUCH HOLDER OR OTHER PARTY HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

END OF TERMS AND CONDITIONS



## Appendix: How to Apply These Terms to Your New Programs

If you develop a new program, and you want it to be of the greatest possible use to the public, the best way to achieve this is to make it free software which everyone can redistribute and change under these terms.

To do so, attach the following notices to the program. It is safest to attach them to the start of each source file to most effectively convey the exclusion of warranty; and each file should have at least the "copyright" line and a pointer to where the full notice is found.

*one line to give the program's name and a brief idea of what it does.*

*Copyright (C) yyyy name of author*

*This program is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2 of the License, or (at your option) any later version.*

*This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.*

*You should have received a copy of the GNU General Public License along with this program; if not, write to the Free Software Foundation, Inc., 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301, USA.*

Also add information on how to contact you by electronic and paper mail.

If the program is interactive, make it output a short notice like this when it starts in an interactive mode:

*Gnomovision version 69, Copyright (C) yyyy name of author*

*Gnomovision comes with ABSOLUTELY NO WARRANTY; for details type 'show w'.*

*This is free software, and you are welcome to redistribute it under certain conditions; type 'show c' for details.*



# Bibliography X



Information in  
Speech

Introduction  
to Information  
Theory

Introduction

Probability  
distributions

Bayesian  
probabilities

Information and  
probabilities

Relative entropy

Compression

Markov Chains

Maximum  
Entropy

**Bibliography**

The hypothetical commands `show w` and `show c` should show the appropriate parts of the General Public License. Of course, the commands you use may be called something other than `show w` and `show c`; they could even be mouse-clicks or menu items—whatever suits your program. You should also get your employer (if you work as a programmer) or your school, if any, to sign a “copyright disclaimer” for the program, if necessary. Here is a sample; alter the names:

*Yoyodyne, Inc., hereby disclaims all copyright interest in the program  
'Gnomovision' (which makes passes at compilers) written by James Hacker.  
signature of Ty Coon, 1 April 1989  
Ty Coon, President of Vice*

This General Public License does not permit incorporating your program into proprietary programs. If your program is a subroutine library, you may consider it more useful to permit linking proprietary applications with the library. If this is what you want to do, use the GNU Library General Public License instead of this License.