Three Concepts: Information Lecture 6: MDL Principle (contd.)

Teemu Roos

Complex Systems Computation Group Department of Computer Science, University of Helsinki

Fall 2007



Lecture 6: MDL Principle (contd.)



A.N. Kolmogorov (left) introducing the structure function at a Bernoulli Society meeting, Tallinn, 1973

・ロト ・四ト ・モト・ モー

3

990

Outline

Kolmogorov Complexity Structure Function MDL Principle

- Kolmogorov Complexity
 - Definition
 - Basic Properties



< □ > < □ > < □</p>

Outline

Teemu Roos

Kolmogorov Complexity Structure Function MDL Principle

- Kolmogorov Complexity
 - Definition
 - Basic Properties

2 Structure Function

- Finite Set Models
- Structure Function
- Minimal Sufficient Statistic
- Ideal MDL

< □ > < 同 >

в

э

Sac

Outline

Kolmogorov Complexity Structure Function MDL Principle

Kolmogorov Complexity

- Definition
- Basic Properties

2 Structure Function

- Finite Set Models
- Structure Function
- Minimal Sufficient Statistic
- Ideal MDL

3 MDL Principle

- Definitions
- Universal Models
- Prediction & Model Selection



< 口 > < 同

Definition Basic Properties

Kolmogorov Complexity

We probably agree that the string

 $101010101010101010\dots 10$

is 'simple'.

Why?

Teemu Roos Three Concepts: Information

イロト イポト イヨト イヨト

1

SQR

Definition Basic Properties

Kolmogorov Complexity

We probably agree that the string

 $101010101010101010\dots 10$

is 'simple'.

Why?

(One) Solution: The string can be described briefly:

"10 repeated k times".

Definition Basic Properties

Kolmogorov Complexity

We probably agree that the string

 $101010101010101010\dots 10$

is 'simple'.

Why?

(One) Solution: The string can be described briefly:

"10 repeated k times".

Remark: 'Describe' should be understood as meaning "compute by an algorithm" (a formal procedure that halts).

< ロト < 同ト < ヨト < ヨト -

Definition Basic Properties

Kolmogorov Complexity

Let $U : \{0,1\}^* \to \{0,1\}^* \cup \emptyset$ be a computer that given a (binary) program $p \in \{0,1\}^*$ either produces a finite (binary) output $U(p) \in \{0,1\}^*$ or never halts. In the latter case, the output U(p) is said to be undefined (\emptyset).

Definition Basic Properties

Kolmogorov Complexity

Let $U : \{0,1\}^* \to \{0,1\}^* \cup \emptyset$ be a computer that given a (binary) program $p \in \{0,1\}^*$ either produces a finite (binary) output $U(p) \in \{0,1\}^*$ or never halts. In the latter case, the output U(p) is said to be undefined (\emptyset).

Kolmogorov Complexity

For a finite string $x \in \{0,1\}^*$, let $p^*(x)$ be the *shortest* program for which

$$U(p^*(x)) = x .$$

The **Kolmogorov complexity** of string x is defined as the length of $p^*(x)$:

$$K_U(x) = \min_{p: U(p)=x} |p| .$$

イロト イポト イヨト イヨト

3

Definition Basic Properties

Kolmogorov Complexity

We assume that the set of programs that halt forms a **prefix-free** set (like symbol codes).

イロト イポト イヨト イヨト

SQ (P

Definition Basic Properties

Kolmogorov Complexity

We assume that the set of programs that halt forms a **prefix-free** set (like symbol codes).

The advantage of prefix-free programs is that we can **concatenate** two programs, p and q to form the program pq so that the computer can separate the two programs.

Definition Basic Properties

Kolmogorov Complexity

Let U and V be two computers. If computer U is sufficiently 'rich', it can emulate computer V so that it outputs the same output as V for any program p.

イロト イポト イヨト イヨト

SQ (P

Definition Basic Properties

Kolmogorov Complexity

Let U and V be two computers. If computer U is sufficiently 'rich', it can emulate computer V so that it outputs the same output as V for any program p.

Universality

A computer U is said to be **universal**, if for *any* other computer V there is a 'translation' program $q \in \{0,1\}^*$ (which depends on V) such that for all programs p we have

$$U(qp) = V(p)$$
,

i.e., when given the concatenated program qp, computer U outputs the same string as computer V when given the program p.

Definition Basic Properties

Kolmogorov Complexity

For any universal computer U, and any other computer V, we have

$$K_U(x) \leq K_V(x) + C$$
,

where C is a constant independent of x.

イロト イポト イヨト イヨト

nar

Definition Basic Properties

Kolmogorov Complexity

For any *universal* computer U, and any other computer V, we have

$$K_U(x) \leq K_V(x) + C$$
,

where C is a constant independent of x.

Proof: Let q be a the translation program which translates programs of V into programs of U, and let $p_V^*(x)$ be the shortest program for which $V(p_V^*(x)) = x$. Then $U(qp_V^*(x)) = x$ so that

$${{K}_{U}(x)} \le \left| q p_{V}^{*}(x)
ight| = \left| p_{V}^{*}(x)
ight| + \left| q
ight| = {{K}_{V}(X)} + \left| q
ight| \; . \quad \Box$$

Definition Basic Properties

Kolmogorov Complexity

For any *universal* computer U, and any other computer V, we have

$$K_U(x) \leq K_V(x) + C \;\;,$$

where C is a constant independent of x.

Proof: Let q be a the translation program which translates programs of V into programs of U, and let $p_V^*(x)$ be the shortest program for which $V(p_V^*(x)) = x$. Then $U(qp_V^*(x)) = x$ so that

$${{K}_{U}(x)} \leq \left| {qp_{V}^{st }(x)}
ight| = \left| {p_{V}^{st }(x)}
ight| + \left| {q}
ight| = {{K}_{V}(X)} + \left| {q}
ight| \; . \quad \Box$$

Based on this property, it can be said that Kolmogorov complexity is the length of the **universally** shortest description of x.

< ロ > < 同 > < 三 > < 三 >

Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

◆ロ > ◆母 > ◆臣 > ◆臣 >

3

5900

Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

C (compiler + operating system + computer),

イロト イポト イヨト イヨト

1

SQR

Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),

イロト イポト イヨト イヨト

SQR

Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),
- your favorite programming language (compiler/interpreter + OS + computer),

< ロ > < 同 > < 三 > < 三 >

Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),
- your favorite programming language (compiler/interpreter + OS + computer),
- Universal Turing machine,



<ロト <同ト < ヨト < ヨト

Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),
- your favorite programming language (compiler/interpreter + OS + computer),
- Universal Turing machine,
- Universal recursive function,



<ロト <同ト < 三ト <

Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),
- your favorite programming language (compiler/interpreter + OS + computer),
- Universal Turing machine,
- Universal recursive function,
- 6 Lambda calculus,



Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),
- your favorite programming language (compiler/interpreter + OS + computer),
- Universal Turing machine,
- Universal recursive function,
- 6 Lambda calculus,
- Arithmetics,



Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),
- your favorite programming language (compiler/interpreter + OS + computer),
- Universal Turing machine,
- Universal recursive function,
- 6 Lambda calculus,
- Arithmetics,
- Game of Life





Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),
- your favorite programming language (compiler/interpreter + OS + computer),
- Universal Turing machine,
- Universal recursive function,
- 6 Lambda calculus,
- Arithmetics,
- Game of Life

9



Definition Basic Properties

Examples

Examples of (virtually) universal 'computers':

- C (compiler + operating system + computer),
- Java (compiler + operating system + computer),
- your favorite programming language (compiler/interpreter + OS + computer),
- Universal Turing machine,
- Universal recursive function,
- 6 Lambda calculus,
- Arithmetics,

9 ...

6 Game of Life





Definition Basic Properties

Invariance Theorem

From now on we restrict the choice of the computer U in K_U to *universal* computers.

< ロ > < 同 > < 三 > < 三 > .

1

SQR

Definition Basic Properties

Invariance Theorem

From now on we restrict the choice of the computer U in K_U to *universal* computers.

Invariance Theorem

Kolmogorov complexity is invariant (up to an additive constant) under a change of the universal computer. In other words, for any two universal computers, U and V, there is a constant C such that

 $|K_U(x) - K_V(x)| \leq C$ for all $x \in \{0,1\}^*$.

Definition Basic Properties

Invariance Theorem

From now on we restrict the choice of the computer U in K_U to *universal* computers.

Invariance Theorem

Kolmogorov complexity is invariant (up to an additive constant) under a change of the universal computer. In other words, for any two universal computers, U and V, there is a constant C such that

 $|K_U(x) - K_V(x)| \leq C$ for all $x \in \{0,1\}^*$.

Proof: Since U is universal, we have $K_U(x) \le K_V(x) + C_1$. Since V is universal, we have $K_V(x) \le K_U(x) + C_2$. The theorem follows by setting $C = \max\{C_1, C_2\}$.

Definition Basic Properties

Kolmogorov Complexity

Upper Bound 1

We have the following upper bound on $K_U(x)$:

$K_U(x) \leq 2|x| + C$

for some constant C which depends on the computer U but not on the string x.

Definition Basic Properties

Kolmogorov Complexity

Upper Bound 1

We have the following upper bound on $K_U(x)$:

 $K_U(x) \leq 2|x| + C$

for some constant C which depends on the computer U but not on the string x.

Proof: Let *q* be the program:

print every even bit that follows

until the next odd bit is $0: x_1 1 x_2 1 \dots x_n 0$.

The length of this program is 2|x| + C. Prefix-free.

Definition Basic Properties

Kolmogorov Complexity

Upper Bound 2

We have the following upper bound on $K_U(x)$:

$$K_U(x) \le |x| + 2\log_2|x| + C$$

for some constant C which depends on the computer U but not on the string x.

< ロ > < 同 > < 三 > < 三 >

Definition Basic Properties

Kolmogorov Complexity

Upper Bound 2

We have the following upper bound on $K_U(x)$:

$$K_U(x) \le |x| + 2\log_2|x| + C$$

for some constant C which depends on the computer U but not on the string x.

Proof: Let *q* be the program:

read integer n and print the following n bits:

 $n_1 1 n_2 1 \ldots n_{|n|} 0 x_1 x_2 \ldots x_n$

・ロト ・ 同ト ・ ヨト ・ ヨト

The length of n = |x| is at most $\lceil \log_2 |x| \rceil \le \log_2 |x| + 1$, so that the length of the program is at most $C' + 2 \log_2 |x| + 2 + |x|$.

Definition Basic Properties

Kolmogorov Complexity

Conditional Kolmogorov Complexity

The **conditional Kolmogorov complexity** is defined as the length of the shortest program to print *x* when *y* is given:

$$K_U(x \mid y) = \min_{p \colon U(\bar{y} \mid p) = x} |p| \;\;,$$

where \bar{y} is a 'self-delimiting' representation of y.

< ロ > < 同 > < 三 > < 三 >
Definition Basic Properties

Kolmogorov Complexity

Conditional Kolmogorov Complexity

The **conditional Kolmogorov complexity** is defined as the length of the shortest program to print *x* when *y* is given:

$$K_U(x \mid y) = \min_{p \colon U(\bar{y} \mid p) = x} |p| ,$$

where \bar{y} is a 'self-delimiting' representation of y.

Upper Bound 3

We have the following upper bound on $K_U(x \mid |x|)$:

 $K_U(x \mid |x|) \leq |x| + C$

for some constant C independent x.

 SQA

Definition Basic Properties

Examples

Let n = |x|.

Teemu Roos Three Concepts: Information

*ロト *部ト *注ト *注ト

-2

990

Definition Basic Properties

Examples

Let n = |x|. • $K_U(0101010101...01 | n) = C$. Program: print n/2 times 01.

イロト イポト イヨト イヨト

1

SQR

Definition Basic Properties

Examples

Let n = |x|.

- K_U(0101010101...01 | n) = C. Program: print n/2 times 01.
- $K_U(\pi_1 \pi_2 \dots \pi_n \mid n) = C.$ *Program:* print the first *n* bits of π .

< ロ > < 同 > < 三 > < 三 >

Definition Basic Properties

Examples

Let n = |x|.

- $K_U(0101010101...01 | n) = C.$ *Program:* print n/2 times 01.
- $K_U(\pi_1 \pi_2 \dots \pi_n \mid n) = C.$ *Program:* print the first *n* bits of π .
- $K_U(\text{English text} \mid n) \approx 1.3 \times n + C.$ *Program:* Huffman code. (Entropy of English is about 1.3 bits per symbol.)

イロト イポト イヨト イヨト

Definition Basic Properties

Examples

Let n = |x|.

- K_U(0101010101...01 | n) = C. Program: print n/2 times 01.
- $K_U(\pi_1 \pi_2 \dots \pi_n \mid n) = C.$ *Program:* print the first *n* bits of π .
- $K_U(\text{English text} \mid n) \approx 1.3 \times n + C.$ *Program:* Huffman code. (Entropy of English is about 1.3 bits per symbol.)
- $K_U(\text{fractal}) = C$. *Program:* print # of iterations until $z_{n+1} = z_n^2 + c > T$.

<ロト < 同ト < ヨト < ヨト -

Definition Basic Properties

Examples



Teemu Roos Three Concepts: Information

Definition Basic Properties

Martin-Löf Randomness

Examples (contd.):

イロト イポト イヨト イヨト

1

SQR

Definition Basic Properties

Martin-Löf Randomness

Examples (contd.):

K_U(x | n) ≈ n, for almost all x ∈ {0,1}ⁿ.
 Proof: Upper bound K_U(x | n) ≤ n + C. Lower bound by a counting argument: less than 2^{-k} of strings compressible by more than k bits (Lecture 1).

< ロ > < 同 > < 三 > < 三 >

Definition Basic Properties

Martin-Löf Randomness

Examples (contd.):

K_U(x | n) ≈ n, for almost all x ∈ {0,1}ⁿ.
 Proof: Upper bound K_U(x | n) ≤ n + C. Lower bound by a counting argument: less than 2^{-k} of strings compressible by more than k bits (Lecture 1).

Martin-Löf Randomness

String x is said to be Martin-Löf random iff $K_u(x \mid n) \ge n$.

<ロト <同ト < ヨト < ヨト

Definition Basic Properties

Martin-Löf Randomness

Examples (contd.):

K_U(x | n) ≈ n, for almost all x ∈ {0,1}ⁿ.
 Proof: Upper bound K_U(x | n) ≤ n + C. Lower bound by a counting argument: less than 2^{-k} of strings compressible by more than k bits (Lecture 1).

Martin-Löf Randomness

String x is said to be **Martin-Löf random** iff $K_u(x \mid n) \ge n$.

Consequence of point 5 above: An i.i.d. sequence of unbiased coin flips is with high probability Martin-Löf random.

イロト イポト イヨト イヨト

Definition Basic Properties

Universal Prediction

Since the set of valid (halting) programs is required to be **prefix-free** we can consider the probability distribution p_U^n :

$$p_U^n(x) = rac{2^{-\kappa_U(x|n)}}{C}$$
, where $C = \sum_{x \in \mathcal{X}^n} 2^{-\kappa_U(x|n)}$.

イロト イポト イヨト イヨト

SQR

Definition Basic Properties

Universal Prediction

Since the set of valid (halting) programs is required to be **prefix-free** we can consider the probability distribution p_{U}^{n} :

$$p_U^n(x) = rac{2^{-K_U(x|n)}}{C}$$
, where $C = \sum_{x \in \mathcal{X}^n} 2^{-K_U(x|n)}$

Universal Probability Distribution

The distribution p_U^n is universal in the sense that for any other computable distribution q, there is a constant C > 0 such that

$$p_U^n(x) \ge C q(x)$$
 for all $x \in \mathcal{X}^n$.

<ロト <同ト < ヨト < ヨト

Definition Basic Properties

Universal Prediction

Since the set of valid (halting) programs is required to be **prefix-free** we can consider the probability distribution p_{U}^{n} :

$$p_U^n(x) = rac{2^{-\kappa_U(x|n)}}{C}$$
, where $C = \sum_{x \in \mathcal{X}^n} 2^{-\kappa_U(x|n)}$

Universal Probability Distribution

The distribution p_U^n is universal in the sense that for any other computable distribution q, there is a constant C > 0 such that

$$p_U^n(x) \ge C q(x)$$
 for all $x \in \mathcal{X}^n$.

Proof idea: The universal computer *U* can imitate the Shannon-Fano prefix code with codelengths $\left[\log_2 \frac{1}{\sigma(v)}\right]$

$$\log \left| \log_2 \frac{1}{q(x)} \right|$$

Definition Basic Properties

Universal Prediction

The universal probability distribution p_U^n is a good predictor.

・ロト ・四ト ・ヨト ・ヨト

1

SQR

Definition Basic Properties

Universal Prediction

The universal probability distribution p_U^n is a good predictor.

This follows from the relationship between codelengths and probabilities (Kraft!):

 $K_U(x)$ is small $\Rightarrow p_U^n(x)$ is large

<ロト < 同ト < ヨト < ヨト -

Definition Basic Properties

Universal Prediction

The universal probability distribution p_U^n is a good predictor.

This follows from the relationship between codelengths and probabilities (Kraft!):

$$\mathcal{K}_U(x)$$
 is small $\Rightarrow p_U^n(x)$ is large
 $\Rightarrow \prod_{i=1}^n p_U^n(x_i \mid x_1, \dots, x_{i-1})$ is large

(日) (同) (三) (三)

SQR

Definition Basic Properties

Universal Prediction

The universal probability distribution p_U^n is a good predictor.

This follows from the relationship between codelengths and probabilities (Kraft!):

$$\begin{split} \mathcal{K}_U(x) \text{ is small } &\Rightarrow \ p_U^n(x) \text{ is large} \\ &\Rightarrow \ \prod_{i=1}^n p_U^n(x_i \mid x_1, \dots, x_{i-1}) \text{ is large} \\ &\Rightarrow \ p_U^n(x_i \mid x_1, \dots, x_{i-1}) \text{ is large for most } i \in \{1, \dots, n\}, \end{split}$$

where x_i denotes the *i*th bit in string x.

(日) (同) (三) (三)

Definition Basic Properties

Berry Paradox



The smallest integer that cannot be described in ten words?

Teemu Roos Three Concepts: Information

イロト イポト イヨト イヨト

SQR

Definition Basic Properties

Berry Paradox



The smallest integer that cannot be described in ten words?

Whatever this number is, we have just described (?) it in ten words.

イロト イポト イヨト イヨト

Definition Basic Properties

Berry Paradox



The smallest integer that cannot be described in ten words?

Whatever this number is, we have just described (?) it in ten words.

The smallest uninteresting number?

<ロト <同ト < 三ト <

1

Definition Basic Properties

Berry Paradox



The smallest integer that cannot be described in ten words?

Whatever this number is, we have just described (?) it in ten words.

The smallest uninteresting number?

Whatever this number is, it is quite interesting!

▲ 🗇 🕨 🔺

Definition Basic Properties

Non-computability

It is impossible to construct a general procedure (algorithm) to compute $K_U(x)$.

Non-Computability

Kolmogorov complexity K_U : $\{0,1\}^* \to \mathbb{N}$ is **non-computable**.

(日) (同) (三) (三)

Definition Basic Properties

Non-computability

It is impossible to construct a general procedure (algorithm) to compute $K_U(x)$.

Non-Computability

Kolmogorov complexity K_U : $\{0,1\}^* \to \mathbb{N}$ is **non-computable**.

Proof: Assume, by way of contradiction, that it would be possible to compute $K_U(x)$.

<ロト < 同ト < ヨト < ヨト -

Definition Basic Properties

Non-computability

It is impossible to construct a general procedure (algorithm) to compute $K_U(x)$.

Non-Computability

Kolmogorov complexity K_U : $\{0,1\}^* \to \mathbb{N}$ is **non-computable**.

Proof: Assume, by way of contradiction, that it would be possible to compute $K_U(x)$. Then for any M > 0, the program

```
print a string x for which K_U(x) > M.
```

would print a string with $K_U(x) > M$.

・ロト ・ 雪 ト ・ ヨ ト ・ ヨ ト

Definition Basic Properties

Non-computability

It is impossible to construct a general procedure (algorithm) to compute $K_U(x)$.

Non-Computability

Kolmogorov complexity K_U : $\{0,1\}^* \to \mathbb{N}$ is **non-computable**.

Proof: Assume, by way of contradiction, that it would be possible to compute $K_U(x)$. Then for any M > 0, the program

```
print a string x for which K_U(x) > M.
```

would print a string with $K_U(x) > M$. A contradiction follows by letting M be larger than the Kolmogorov complexity of this program. Hence, it cannot be possible to compute $K_U(x)$.

イロト イポト イヨト イヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

- 1 Kolmogorov Complexity
 - Definition
 - Basic Properties

2 Structure Function

- Finite Set Models
- Structure Function
- Minimal Sufficient Statistic
- Ideal MDL

3 MDL Principle

- Definitions
- Universal Models
- Prediction & Model Selection



< □ > < 同 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Finite Set Models

Each string $x \in \{0,1\}^n$ can be described in two parts:

① the regular features of x,

(日) (同) (三) (三)

1

SQR

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Finite Set Models

Each string $x \in \{0,1\}^n$ can be described in two parts:

- the regular features of x,
- the index of x in the set S of strings with those regular features.

< ロ > < 同 > < 三 > < 三 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Finite Set Models

Each string $x \in \{0,1\}^n$ can be described in two parts:

- the regular features of x,
- the index of x in the set S of strings with those regular features.

For set $S \subseteq \{0,1\}^n$, the length of such a two-part description is

$$K_U(S \mid n) + \log_2 |S| + C \geq K_U(x \mid n) \ ,$$

where $K_U(S \mid n)$ is the length of the shortest program to list the members of S (and then halt).

(日) (同) (三) (三)

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Finite Set Models

Each string $x \in \{0,1\}^n$ can be described in two parts:

- the regular features of x,
- the index of x in the set S of strings with those regular features.

For set $S \subseteq \{0,1\}^n$, the length of such a two-part description is

$$K_U(S \mid n) + \log_2 |S| + C \geq K_U(x \mid n) \ ,$$

where $K_U(S \mid n)$ is the length of the shortest program to list the members of S (and then halt).

We can consider *S* (regular features) a model.

・ロト ・ 同ト ・ ヨト ・ ヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Finite Set Models



3

990

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Example

For instance, if x is a sequence of biased coin flips, then with high probability the only regular feature is the number of 1s.

< ロ > < 同 > < 三 > < 三 >

SQ (P

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Example

For instance, if x is a sequence of biased coin flips, then with high probability the only regular feature is the number of 1s.

Let $S_k^n = \{x \in \{0,1\}^n : \sum x_i = k\}$. The size of this set is

$$|S_k^n| = \binom{n}{k} = \frac{n!}{k!(n-k)!} \approx 2^{nH(\frac{k}{n})} ,$$

where $H(\frac{k}{n})$ is the binary entropy with parameter $\frac{k}{n}$.

イロト イポト イヨト イヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Example

For instance, if x is a sequence of biased coin flips, then with high probability the only regular feature is the number of 1s.

Let $S_k^n = \{x \in \{0,1\}^n : \sum x_i = k\}$. The size of this set is

$$|S_k^n| = \binom{n}{k} = \frac{n!}{k!(n-k)!} \approx 2^{nH(\frac{k}{n})}$$

where $H(\frac{k}{n})$ is the binary entropy with parameter $\frac{k}{n}$.



<ロト <同ト < 三ト <

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

 $(X)_{E}^{0.5}$

<ロト <同ト < 三ト <

 $0.5 \\ Pr(X = 1)$

1.0

Example

For instance, if x is a sequence of biased coin flips, then with high probability the only regular feature is the number of 1s.

Let $S_k^n = \{x \in \{0,1\}^n : \sum x_i = k\}$. The size of this set is

$$|S_k^n| = \binom{n}{k} = \frac{n!}{k!(n-k)!} \approx 2^{nH(\frac{k}{n})}$$

where $H(\frac{k}{n})$ is the binary entropy with parameter $\frac{k}{n}$.

Thus, the two-part description has length

$$K_U(k) + nH(rac{k}{n}) + C \geq K_U(x \mid n)$$
.
Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

 $(X)_{E}^{0.5}$

 $0.5 \\ Pr(X = 1)$

1.0

Example

For instance, if x is a sequence of biased coin flips, then with high probability the only regular feature is the number of 1s.

Let $S_k^n = \{x \in \{0,1\}^n : \sum x_i = k\}$. The size of this set is

$$|S_k^n| = \binom{n}{k} = \frac{n!}{k!(n-k)!} \approx 2^{nH(\frac{k}{n})}$$

where $H(\frac{k}{n})$ is the binary entropy with parameter $\frac{k}{n}$.

Thus, the two-part description has length

$$K_U(k) + nH(\frac{k}{n}) + C \geq K_U(x \mid n)$$
.

By the **Asymptotic Equipartition Property** (Lecture 3), $nH(\frac{k}{n})$ is with high probability a lower bound (\approx) on $K_U(x \mid n)$.

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

 $(X)_{E}^{0.5}$

 $0.5 \\ Pr(X = 1)$

1.0

Example

For instance, if x is a sequence of biased coin flips, then with high probability the only regular feature is the number of 1s.

Let $S_k^n = \{x \in \{0,1\}^n : \sum x_i = k\}$. The size of this set is

$$|S_k^n| = \binom{n}{k} = \frac{n!}{k!(n-k)!} \approx 2^{nH(\frac{k}{n})}$$

where $H(\frac{k}{n})$ is the binary entropy with parameter $\frac{k}{n}$.

Thus, the two-part description has length

$$K_U(k) + nH(\frac{k}{n}) + C \approx K_U(x \mid n)$$
.

By the **Asymptotic Equipartition Property** (Lecture 3), $nH(\frac{k}{n})$ is with high probability a lower bound (\approx) on $K_U(x \mid n)$.

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

Kolmogorov Structure Function

The Kolmogorov structure function is defined as

$$h_{x}(\alpha) = \min_{\substack{S : x \in S \\ K_{U}(S|n) \leq \alpha}} \log_{2} |S| ,$$

i.e., the log-size of the smallest set containing x that can be described in α bits.

(日)

B b

MQ (P

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

Kolmogorov Structure Function

The Kolmogorov structure function is defined as

$$h_{x}(\alpha) = \min_{\substack{S: x \in S \\ K_{U}(S|n) \leq \alpha}} \log_{2} |S| ,$$

i.e., the log-size of the smallest set containing x that can be described in α bits.

The more bits we can use to describe S, the more regularities we can cover, which makes |S| smaller.

< ロ > < 同 > < 三 > < 三 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

Kolmogorov Structure Function

The Kolmogorov structure function is defined as

$$h_{x}(\alpha) = \min_{\substack{S: x \in S \\ K_{U}(S|n) \leq \alpha}} \log_{2} |S| ,$$

i.e., the log-size of the smallest set containing x that can be described in α bits.

The more bits we can use to describe S, the more regularities we can cover, which makes |S| smaller.

For all $\alpha > 0$, there is a two-part description of length $\alpha + h_x(\alpha)$.

イロト イボト イヨト イヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

Consider different values of α :

 $\underline{\alpha \approx 0}$:

We can only describe the whole set $\{0,1\}^n$, and not much else, so that $h_x(0) = \log_2 |\{0,1\}^n| = n$.

< ロ > < 同 > < 三 > < 三 >

SQA

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

Consider different values of α :

 $\underline{\alpha \approx 0}$:

We can only describe the whole set $\{0,1\}^n$, and not much else, so that $h_x(0) = \log_2 |\{0,1\}^n| = n$.

 $\underline{\alpha \approx K_U(x \mid n)}:$

We can use the singleton set $S = \{x\}$ since $\mathcal{K}_U(\{x\} \mid n) = \mathcal{K}_U(x \mid n) + C$. Thus, $h_x(\mathcal{K}_U(x)) = \log_2 |\{x\}| = 0$.

< ロ > < 同 > < 三 > < 三 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

Consider different values of α :

optimal:

$$h_x(\alpha) + \alpha \geq K_U(x \mid n)$$
.

イロト イポト イヨト イヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

Consider different values of α :

optimal:

$$h_x(\alpha) \geq K_U(x \mid n) - \alpha$$
.

イロト イポト イヨト イヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

The **slope** of the structure function $h_x(\alpha)$ is at least as steep as -1 (ignoring constants):

(日) (同) (三) (三)

SQA

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

The **slope** of the structure function $h_x(\alpha)$ is at least as steep as -1 (ignoring constants):

For k extra bits in α , we can reduce the set S in a fraction of $\frac{1}{2^k}$ by sorting S alphabetically, dividing in 2^k equal size parts, and encoding the index of the part including x.

<ロト <同ト < ヨト < ヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function

The **slope** of the structure function $h_x(\alpha)$ is at least as steep as -1 (ignoring constants):

For k extra bits in α , we can reduce the set S in a fraction of $\frac{1}{2^k}$ by sorting S alphabetically, dividing in 2^k equal size parts, and encoding the index of the part including x.

The constants related to instructions like "sort alphabetically and choose second half/third quarter/etc." cause **bumps** in the structure function.

イロト イポト イヨト イヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function



Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Sufficient Statistic

In statistics, a **sufficient statistic** is a function of the data which contains all the information relevant to a parameter. Examples:

< ロ > < 同 > < 三 > < 三 >

-

SQA

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Sufficient Statistic

In statistics, a **sufficient statistic** is a function of the data which contains all the information relevant to a parameter. Examples:

in coin flipping, the number of 1s is sufficient for the bias parameter.

< ロ > < 同 > < 三 > < 三 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Sufficient Statistic

In statistics, a **sufficient statistic** is a function of the data which contains all the information relevant to a parameter. Examples:

- in coin flipping, the number of 1s is sufficient for the bias parameter.
- in die tossing, the number of times each face is seen is sufficient for the parameters p₁,..., p₆.

< ロ > < 同 > < 三 > < 三 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Sufficient Statistic

In statistics, a **sufficient statistic** is a function of the data which contains all the information relevant to a parameter. Examples:

- in coin flipping, the number of 1s is sufficient for the bias parameter.
- in die tossing, the number of times each face is seen is sufficient for the parameters p₁,..., p₆.
- in a Gaussian density, the average $\frac{1}{n} \sum X_i$ is sufficient for the mean.

< ロ > < 同 > < 三 > < 三 >

SOR

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Sufficient Statistic

In statistics, a **sufficient statistic** is a function of the data which contains all the information relevant to a parameter. Examples:

- in coin flipping, the number of 1s is sufficient for the bias parameter.
- in die tossing, the number of times each face is seen is sufficient for the parameters p₁,..., p₆.
- in a Gaussian density, the average $\frac{1}{n} \sum X_i$ is sufficient for the mean.
- In all these, the order of the outcomes, for instance, is irrelevant.

(日) (同) (三) (三)

SOR

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Sufficient Statistic

In statistics, a **sufficient statistic** is a function of the data which contains all the information relevant to a parameter. Examples:

- in coin flipping, the number of 1s is sufficient for the bias parameter.
- in die tossing, the number of times each face is seen is sufficient for the parameters p₁,..., p₆.
- in a Gaussian density, the average $\frac{1}{n} \sum X_i$ is sufficient for the mean.
- In all these, the order of the outcomes, for instance, is irrelevant.

When estimating the parameter, it is *sufficient* to know the sufficient statistic.

< ロ > < 同 > < 三 > < 三 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Kolmogorov Sufficient Statistic

Kolmogorov Sufficient Statistic

A finite set S is a Kolmogorov sufficient statistic iff we have

$$K_U(S \mid n) + \log_2 |S| = K_U(x \mid n) + C$$
,

i.e., the two-part description using S is optimal.

< ロ > < 同 > < 三 > < 三 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Kolmogorov Sufficient Statistic

Kolmogorov Sufficient Statistic

A finite set S is a Kolmogorov sufficient statistic iff we have

$$K_U(S \mid n) + \log_2 |S| = K_U(x \mid n) + C ,$$

i.e., the two-part description using S is optimal.

A Kolmogorov sufficient statistic tells everything about the *structure* of the data x.

<ロト <同ト < ヨト < ヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Kolmogorov Sufficient Statistic

Kolmogorov Sufficient Statistic

A finite set S is a Kolmogorov sufficient statistic iff we have

$$K_U(S \mid n) + \log_2 |S| = K_U(x \mid n) + C ,$$

i.e., the two-part description using S is optimal.

A Kolmogorov sufficient statistic tells everything about the *structure* of the data x.

In coin flipping, the number of 1s is with high probability (also) a Kolmogorov sufficient statistic.

<ロト <同ト < ヨト < ヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function



Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Minimal Sufficient Statistic

If S is a Kolmogorov sufficient statistic, i.e., we have

$$K_U(S \mid n) + \log_2 |S| = K_U(x \mid n) \ ,$$

then for all α within the range $K_U(S \mid n) < \alpha \le K_U(x \mid n)$, there is another sufficient statistic with complexity α :

$$\alpha = K_U(S \mid n) + k : \quad K_U(S \mid n) + k + \log_2 \frac{|S|}{2^k} = K_U(x \mid n) .$$

< ロ > < 同 > < 三 > < 三 >

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Minimal Sufficient Statistic

If S is a Kolmogorov sufficient statistic, i.e., we have

$$K_U(S \mid n) + \log_2 |S| = K_U(x \mid n) \ ,$$

then for all α within the range $K_U(S \mid n) < \alpha \le K_U(x \mid n)$, there is another sufficient statistic with complexity α :

$$\alpha = K_U(S \mid n) + k : \quad K_U(S \mid n) + k + \log_2 \frac{|S|}{2^k} = K_U(x \mid n) .$$

Kolmogorov Minimal Sufficient Statistic

The least complex Kolmogorov sufficient statistic is called the **Kolmogorov minimal sufficient statistic**. It contains all the information about the *structure* of x but nothing more.

イロト イポト イラト イラト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Structure Function



Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Example: Mona Lisa





Figure 7.6. Mona Lisa.

Source: Cover & Thomas, 1991

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Example: Mona Lisa



(日)

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Ideal MDL

Ideal MDL

Given data x, choosing the Kolmogorov minimal sufficient statistic as the preferred model is called **"ideal MDL"**.

イロト イポト イヨト イヨト

SQA

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Ideal MDL

Ideal MDL

Given data x, choosing the Kolmogorov minimal sufficient statistic as the preferred model is called **"ideal MDL"**.

Extracts all regularity from data, and leaves out noise.

Teemu Roos Three Concepts: Information

イロト イポト イヨト イヨト

Finite Set Models Structure Function Minimal Sufficient Statistic Ideal MDL

Ideal MDL

Ideal MDL

Given data x, choosing the Kolmogorov minimal sufficient statistic as the preferred model is called **"ideal MDL"**.

Extracts all regularity from data, and leaves out noise.

- Complexity = Information + Noise
 - = Regularity + Randomness
 - = Algorithm + Compressed file

(日) (同) (三) (三)

 Outline
 Finite Set Models

 Kolmogorov Complexity
 Structure Function

 Structure Function
 Minimal Sufficient Statistic

 MDL Principle
 Ideal MDL

Ideal MDL

The finite set models can be replaced by (computable) probability distributions — distribution p is a *sufficient statistic* iff

$$K_U(p \mid n) + \log_2 \frac{1}{p(x)} = K_U(x \mid n)$$
,

i.e., two-part code optimal.

(日) (同) (三) (三)

 Outline
 Finite Set Models

 Kolmogorov Complexity
 Structure Function

 Structure Function
 Minimal Sufficient Statistic

 MDL Principle
 Ideal MDL

Ideal MDL

The finite set models can be replaced by (computable) probability distributions — distribution p is a *sufficient statistic* iff

$$K_U(p \mid n) + \log_2 \frac{1}{p(x)} = K_U(x \mid n)$$
,

i.e., two-part code optimal.

All things essentially unchanged:

- Finite set S can be replaced by a uniform distribution over S,
- Distribution *p* can be replaced by the typical set under *p*.

イロト イポト イヨト イヨト

Definitions Universal Models Prediction & Model Selection

- 1 Kolmogorov Complexity
 - Definition
 - Basic Properties
- 2 Structure Function
 - Finite Set Models
 - Structure Function
 - Minimal Sufficient Statistic
 - Ideal MDL

3 MDL Principle

- Definitions
- Universal Models
- Prediction & Model Selection



< A

Definitions Universal Models Prediction & Model Selection

Ideal vs. Practical MDL

There are two problematic issues with ideal MDL:

・ロト ・部 ト ・ヨト ・ヨト

3

SQR

Definitions Universal Models Prediction & Model Selection

Ideal vs. Practical MDL

There are two problematic issues with ideal MDL:

Uncomputability of Kolmogorov complexity.

(日) (同) (三) (三)

SQA
Definitions Universal Models Prediction & Model Selection

Ideal vs. Practical MDL

There are two problematic issues with ideal MDL:

- Uncomputability of Kolmogorov complexity.
- **2** Hidden constants in the definitions and theorems.
 - \Rightarrow Says nothing about individual strings.

Definitions Universal Models Prediction & Model Selection

Ideal vs. Practical MDL

There are two problematic issues with ideal MDL:

- Uncomputability of Kolmogorov complexity.
- ❷ Hidden constants in the definitions and theorems.
 ⇒ Says nothing about individual strings.

Practical MDL aims to solve these issues by:

Definitions Universal Models Prediction & Model Selection

Ideal vs. Practical MDL

There are two problematic issues with ideal MDL:

- Uncomputability of Kolmogorov complexity.
- e Hidden constants in the definitions and theorems.
 ⇒ Says nothing about individual strings.

Practical MDL aims to solve these issues by:

General Replace computer programs by probabilistic models.
 ⇒ Computable.

Definitions Universal Models Prediction & Model Selection

Ideal vs. Practical MDL

There are two problematic issues with ideal MDL:

- Uncomputability of Kolmogorov complexity.
- e Hidden constants in the definitions and theorems.
 ⇒ Says nothing about individual strings.

Practical MDL aims to solve these issues by:

- Q Replace computer programs by probabilistic models.
 ⇒ Computable.
- **②** Replace universal computer *U* by a **universal model**.
 ⇒ No hidden constants.

Definitions Universal Models Prediction & Model Selection

Definitions

We call a probability distribution $p : \mathcal{D} \rightarrow [0, 1]$ a **model**.

A model class $\mathcal{M} = \{p_{\theta} : \theta \in \Theta\}$ is a set of probability distributions (models).

<ロト < 同ト < ヨト < ヨト -

3

Definitions Universal Models Prediction & Model Selection

Definitions

We call a probability distribution $p : \mathcal{D} \rightarrow [0,1]$ a **model**.

A model class $\mathcal{M} = \{p_{\theta} : \theta \in \Theta\}$ is a set of probability distributions (models).

The model within \mathcal{M} that achieves the shortest codelength for data x is the **maximum likelihood (ML) model**:

$$\min_{ heta \in \Theta} \log_2 rac{1}{p_{ heta}(D)} = \log_2 rac{1}{p_{\hat{ heta}}(D)}$$

Definitions Universal Models Prediction & Model Selection

Definitions

We call a probability distribution $p : \mathcal{D} \rightarrow [0,1]$ a **model**.

A model class $\mathcal{M} = \{p_{\theta} : \theta \in \Theta\}$ is a set of probability distributions (models).

The model within \mathcal{M} that achieves the shortest codelength for data x is the **maximum likelihood (ML) model**:

$$\min_{\theta \in \Theta} \log_2 \frac{1}{p_{\theta}(D)} = \log_2 \frac{1}{p_{\hat{\theta}}(D)} \quad . \qquad \text{Depends on } D!$$

Definitions Universal Models Prediction & Model Selection

Definitions

We call a probability distribution $p : \mathcal{D} \rightarrow [0,1]$ a **model**.

A model class $\mathcal{M} = \{p_{\theta} : \theta \in \Theta\}$ is a set of probability distributions (models).

The model within \mathcal{M} that achieves the shortest codelength for data x is the **maximum likelihood (ML) model**:

$$\min_{\theta \in \Theta} \log_2 \frac{1}{p_{\theta}(D)} = \log_2 \frac{1}{p_{\hat{\theta}}(D)} \quad \text{Depends on } D!$$

For model q, the excess codelength or "**regret**" over the ML model in \mathcal{M} is given by

$$\log_2 rac{1}{q(D)} - \log_2 rac{1}{p_{\hat{ heta}}(D)}$$
 .

イロト イポト イヨト イヨト

SOR

Definitions Universal Models Prediction & Model Selection

Stochastic Complexity

A model (code) for which the regret grows slower than n is said to be a **universal model** (code) relative to model class M:

$$\lim_{n\to\infty}\frac{1}{n}\left[\log_2\frac{1}{q(D)}-\log_2\frac{1}{p_{\hat{\theta}}(D)}\right]=0 \ .$$

< ロ > < 同 > < 三 > < 三 >

MQ (P

Definitions Universal Models Prediction & Model Selection

Stochastic Complexity

A model (code) for which the regret grows slower than n is said to be a **universal model** (code) relative to model class M:

$$\lim_{n\to\infty} \frac{1}{n} \left[\log_2 \frac{1}{q(D)} - \log_2 \frac{1}{p_{\hat{\theta}}(D)} \right] = 0 \;\; .$$

The universal computer U is universal relative to all computers, while a universal model is universal relative to a model class \mathcal{M} .

Definitions Universal Models Prediction & Model Selection

Stochastic Complexity

A model (code) for which the regret grows slower than n is said to be a **universal model** (code) relative to model class M:

$$\lim_{n\to\infty}\frac{1}{n}\left[\log_2\frac{1}{q(D)}-\log_2\frac{1}{\rho_{\hat{\theta}}(D)}\right]=0 \ .$$

The universal computer U is universal relative to all computers, while a universal model is universal relative to a model class \mathcal{M} .

Stochastic Complexity

The **stochastic complexity** of data D relative to model class \mathcal{M} is the codelength achieved by a universal model.

Definitions Universal Models Prediction & Model Selection

Two-Part Universal Model

The two-part code (Lecture 5) consists of

- **(**) optimally quantized parameter values θ^q , and
- **2** encoding of the data under model p_{θ^q} .

< ロ > < 同 > < 三 > < 三 >

Definitions Universal Models Prediction & Model Selection

Two-Part Universal Model

The two-part code (Lecture 5) consists of

- **(**) optimally quantized parameter values θ^q , and
- 2 encoding of the data under model p_{θ^q} .

The total codelength is

$$\log_2 rac{1}{p_{ heta^q}(D)} + \ell(heta^q)$$

< ロ > < 同 > < 三 > < 三 >

Definitions Universal Models Prediction & Model Selection

Two-Part Universal Model

The two-part code (Lecture 5) consists of

- $\textbf{0} \text{ optimally quantized parameter values } \theta^q \text{, and } \\$
- **2** encoding of the data under model p_{θ^q} .

The total codelength is

$$\log_2 rac{1}{p_{ heta^q}(D)} + \ell(heta^q)$$

For 'smooth' parametric models, and optimal quantization, the codelength becomes (\approx)

$$\log_2 \frac{1}{p_{\hat{\theta}}(D)} + \frac{k}{2} \log_2 n \;\; ,$$

so that the regret is $\frac{k}{2}\log_2 n$. Since $\log_2 n$ grows slower than n, the two-part code is universal.

Definitions Universal Models Prediction & Model Selection

Mixture Universal Model

There are universal codes that are strictly better than the two-part code.

3

Definitions Universal Models Prediction & Model Selection

Mixture Universal Model

There are universal codes that are strictly better than the two-part code.

For instance, given a code for the parameters, let w be a distribution over the parameter space Θ (quantized if necessary) defined as

$$w(heta)=rac{2^{-\ell(heta)}}{c}\,\,,\,\,\,\,\,$$
 where $c=\sum_{ heta\in\Theta}2^{-\ell(heta)}.$

Definitions Universal Models Prediction & Model Selection

Mixture Universal Model

There are universal codes that are strictly better than the two-part code.

For instance, given a code for the parameters, let w be a distribution over the parameter space Θ (quantized if necessary) defined as

$$w(heta)=rac{2^{-\ell(heta)}}{c}\,\,,\,\,\,\,$$
 where $c=\sum_{ heta\in\Theta}2^{-\ell(heta)}.$

Let p^w be a **mixture distribution** over the data-sets $D \in \mathcal{D}$, defined as

$$p^w(D) = \sum_{\theta \in \Theta} p_{\theta}(D) w(\theta) \; ,$$

i.e., an "average" distribution, where each p is weighted by w.

イロト 人間ト イヨト イヨト

-

SOR

Definitions Universal Models Prediction & Model Selection

٠

3

SQR

Mixture Universal Model

The codelength of the **mixture model** p^{w} is given by

$$\begin{split} \log_2 \frac{1}{\sum_{\theta \in \Theta} p(D \mid \theta) \, w(\theta)} &\leq \log_2 \frac{1}{\max_{\theta \in \Theta} p(D \mid \theta) \, w(\theta)} \\ &= \log_2 \frac{1}{\max_{\theta \in \Theta} p(D \mid \theta)} + \log_2 \frac{c}{2^{-\ell(\theta)}} \end{split}$$

Definitions Universal Models Prediction & Model Selection

٠

1

SQR

< ロ > < 同 > < 回 > < 回 > :

Mixture Universal Model

The codelength of the **mixture model** p^w is given by

$$\begin{split} \log_2 \frac{1}{\sum_{\theta \in \Theta} p(D \mid \theta) \, w(\theta)} &\leq \log_2 \frac{1}{\max_{\theta \in \Theta} p(D \mid \theta) \, w(\theta)} \\ &= \log_2 \frac{1}{\max_{\theta \in \Theta} p(D \mid \theta)} + \log_2 \frac{c}{2^{-\ell(\theta)}} \end{split}$$

The right-hand side is equal to

$$\underbrace{ \underbrace{\log_2 \frac{1}{p_{\hat{\theta}}(D)} + \ell(\theta)}_{\text{two-part code}} \underbrace{- \underbrace{\log_2 \frac{1}{c}}_{\leq 0}}_{\leq 0} ,$$

Definitions Universal Models Prediction & Model Selection

.

3

SQR

Mixture Universal Model

The codelength of the **mixture model** p^w is given by

$$\begin{split} \log_2 \frac{1}{\sum_{\theta \in \Theta} p(D \mid \theta) \, w(\theta)} &\leq \log_2 \frac{1}{\max_{\theta \in \Theta} p(D \mid \theta) \, w(\theta)} \\ &= \log_2 \frac{1}{\max_{\theta \in \Theta} p(D \mid \theta)} + \log_2 \frac{c}{2^{-\ell(\theta)}} \end{split}$$

The right-hand side is equal to

$$\underbrace{ \underset{\text{two-part code}}{\log_2 \frac{1}{p_{\hat{\theta}}(D)} + \ell(\theta)} \underbrace{-\log_2 \frac{1}{c}}_{\leq 0} }_{\text{two-part code}} \ ,$$

The mixture code is always at least as good as the two-part code.

Definitions Universal Models Prediction & Model Selection

Normalized Maximum Likelihood

Consider the maximum likelihood model

$$p_{\hat{ heta}}(D) = \max_{ heta \in \Theta} p_{ heta}(D)$$
 .

It is the best probability assignment achievable under model $\mathcal{M}.$

-

Definitions Universal Models Prediction & Model Selection

Normalized Maximum Likelihood

Consider the maximum likelihood model

$$p_{\hat{ heta}}(D) = \max_{ heta \in \Theta} p_{ heta}(D)$$
 .

It is the best probability assignment achievable under model $\mathcal{M}.$

Unfortunately, it is not possible to use the ML model for coding because is not a probability distribution, i.e.,

$$C = \sum_{D \in \mathcal{D}} p_{\hat{\theta}}(D) > 1$$
,

unless $\hat{\theta}$ is constant wrt. D.

< ロ > < 同 > < 回 > < 回 > < 回 > <

Definitions Universal Models Prediction & Model Selection

Normalized Maximum Likelihood

Normalized Maximum Likelihood

The **normalized maximum likelihood (NML) model** is obtained by normalizing the ML model:

$$p_{\mathrm{nml}}(D) = rac{p_{\widehat{ heta}}(D)}{\mathcal{C}} \; , \; \; \; ext{where} \; \mathcal{C} = \sum_{D \in \mathcal{D}} p_{\widehat{ heta}}(D) \; .$$

(日) (同) (三) (三)

Definitions Universal Models Prediction & Model Selection

Normalized Maximum Likelihood

Normalized Maximum Likelihood

The **normalized maximum likelihood (NML) model** is obtained by normalizing the ML model:

$$p_{\mathrm{nml}}(D) = rac{p_{\hat{ heta}}(D)}{\mathcal{C}} \; , \; \; \; ext{where} \; \mathcal{C} = \sum_{D \in \mathcal{D}} p_{\hat{ heta}}(D) \; .$$

The regret of NML is given by

$$\log_2 \frac{1}{p_{\rm nml}(D)} - \log_2 \frac{1}{p_{\hat{\theta}}(D)} = \log_2 \frac{C}{p_{\hat{\theta}}(D)} - \log_2 \frac{1}{p_{\hat{\theta}}(D)} = \log_2 C \ ,$$

which is constant wrt. D.

Definitions Universal Models Prediction & Model Selection

Normalized Maximum Likelihood

Let q be any distribution other than p_{nml} . Then

• there must a data-set $D' \in \mathcal{D}$ for which we have

 $q(D') < p_{\mathrm{nml}}(D')$

-

Definitions Universal Models Prediction & Model Selection

Normalized Maximum Likelihood

Let q be any distribution other than $p_{\rm nml}$. Then

• there must a data-set $D' \in \mathcal{D}$ for which we have



<ロト <同ト < ヨト < ヨト

Definitions Universal Models Prediction & Model Selection

Normalized Maximum Likelihood

Let q be any distribution other than $p_{\rm nml}$. Then

• there must a data-set $D' \in \mathcal{D}$ for which we have



For D', the regret of q is greater than $\log_2 C$, the regret of p_{nml} .

イロト イポト イラト イラト

Definitions Universal Models Prediction & Model Selection

Normalized Maximum Likelihood

Let q be any distribution other than $p_{\rm nml}$. Then

• there must a data-set $D' \in \mathcal{D}$ for which we have



For D', the regret of q is greater than $\log_2 C$, the regret of p_{nml} .

Thus, the worst-case regret of q is greater than the (worst-case) regret of NML. \Rightarrow NML has the least possible **worst-case regret**.

イロト イポト イラト イラト

Definitions Universal Models Prediction & Model Selection

Universal Models

For 'smooth' parametric models, the regret of NML, $\log_2 C$, grows slower than *n*, so **NML is also a universal model**.

-

Definitions Universal Models Prediction & Model Selection

Universal Models

For 'smooth' parametric models, the regret of NML, $\log_2 C$, grows slower than *n*, so **NML is also a universal model**.

We have seen three kinds of universal models:

- two-part,
- e mixture,
- INML.

Definitions Universal Models Prediction & Model Selection

Universal Models

For 'smooth' parametric models, the regret of NML, $\log_2 C$, grows slower than *n*, so **NML is also a universal model**.

We have seen three kinds of universal models:

- two-part,
- Ø mixture,
- INML.

So what do we do with them?

Definitions Universal Models Prediction & Model Selection

Universal Models

For 'smooth' parametric models, the regret of NML, $\log_2 C$, grows slower than *n*, so **NML is also a universal model**.

We have seen three kinds of universal models:

- two-part,
- Ø mixture,
- INML.

So what do we do with them?

We can use them for (at least) three purposes:

compression,

Definitions Universal Models Prediction & Model Selection

Universal Models

For 'smooth' parametric models, the regret of NML, $\log_2 C$, grows slower than *n*, so **NML is also a universal model**.

We have seen three kinds of universal models:

- two-part,
- Ø mixture,
- INML.

So what do we do with them?

We can use them for (at least) three purposes:

- compression,
- 2 prediction,

Definitions Universal Models Prediction & Model Selection

Universal Models

For 'smooth' parametric models, the regret of NML, $\log_2 C$, grows slower than *n*, so **NML is also a universal model**.

We have seen three kinds of universal models:

- two-part,
- Ø mixture,
- INML.

So what do we do with them?

We can use them for (at least) three purposes:

- compression,
- 2 prediction,
- Image: model selection.

Definitions Universal Models Prediction & Model Selection

MDL Prediction

Prediction is done like in Kolmogorov complexity: universal probability distribution/universal model achieves

- good compression: $\ell(D)$ is small,
- good predictions: $p(D_i | D_1, ..., D_{i-1})$ is large for most $i \in \{1, ..., n\}$.

Definitions Universal Models Prediction & Model Selection

MDL Prediction

Prediction is done like in Kolmogorov complexity: universal probability distribution/universal model achieves

- good compression: $\ell(D)$ is small,
- good predictions: $p(D_i | D_1, ..., D_{i-1})$ is large for most $i \in \{1, ..., n\}$.

For instance, the mixture code gives a natural predictor which is equivalent to **Bayesian prediction**.

イロト イポト イヨト イヨト

SOR
Definitions Universal Models Prediction & Model Selection

MDL Prediction

Prediction is done like in Kolmogorov complexity: universal probability distribution/universal model achieves

- good compression: $\ell(D)$ is small,
- good predictions: $p(D_i | D_1, ..., D_{i-1})$ is large for most $i \in \{1, ..., n\}$.

For instance, the mixture code gives a natural predictor which is equivalent to **Bayesian prediction**.

The NML model gives predictions that are good relative to the best model in the model class, **no matter what happens**.

SOR

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

Recall (from Lecture 5) the multi-part codes used when multiple model classes, $\mathcal{M}_1, \mathcal{M}_2, \ldots$ are available:

3

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

Recall (from Lecture 5) the multi-part codes used when multiple model classes, $\mathcal{M}_1, \mathcal{M}_2, \ldots$ are available:

• Encoding of the model class index: $C_0(i), i \in \mathbb{N}$.

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

Recall (from Lecture 5) the multi-part codes used when multiple model classes, $\mathcal{M}_1, \mathcal{M}_2, \ldots$ are available:

- Encoding of the model class index: $C_0(i), i \in \mathbb{N}$.
- **2** Encoding of the parameter (vector): $C_i(\theta), \theta \in \Theta_i$.

イロト イポト イヨト イヨト

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

Recall (from Lecture 5) the multi-part codes used when multiple model classes, M_1, M_2, \ldots are available:

- Encoding of the model class index: $C_0(i), i \in \mathbb{N}$.
- **2** Encoding of the parameter (vector): $C_i(\theta), \theta \in \Theta_i$.
- Encoding of the data: $C_{\theta}(D), D \in \mathcal{D}$.

イロト イポト イヨト イヨト

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

Recall (from Lecture 5) the multi-part codes used when multiple model classes, M_1, M_2, \ldots are available:

- Encoding of the model class index: $C_0(i), i \in \mathbb{N}$.
- **2** Encoding of the parameter (vector): $C_i(\theta), \theta \in \Theta_i$.
- Solution Encoding of the data: $C_{\theta}(D), D \in \mathcal{D}$.

If we are interested in choosing a model class (and not the parameters), we can improve parts 2 & 3 by combining them into a better universal code than two-part:

イロト 人間ト イヨト イヨト

SOR

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

Recall (from Lecture 5) the multi-part codes used when multiple model classes, M_1, M_2, \ldots are available:

- Encoding of the model class index: $C_0(i), i \in \mathbb{N}$.
- **2** Encoding of the parameter (vector): $C_i(\theta), \theta \in \Theta_i$.
- Solution Encoding of the data: $C_{\theta}(D), D \in \mathcal{D}$.

If we are interested in choosing a model class (and not the parameters), we can improve parts 2 & 3 by combining them into a better universal code than two-part:

• Encoding of the model class index: $C_0(i), i \in \mathbb{N}$.

イロト 人間ト イヨト イヨト

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

Recall (from Lecture 5) the multi-part codes used when multiple model classes, M_1, M_2, \ldots are available:

- Encoding of the model class index: $C_0(i), i \in \mathbb{N}$.
- **2** Encoding of the parameter (vector): $C_i(\theta), \theta \in \Theta_i$.
- Encoding of the data: $C_{\theta}(D), D \in \mathcal{D}$.

If we are interested in choosing a model class (and not the parameters), we can improve parts 2 & 3 by combining them into a better universal code than two-part:

- Encoding of the model class index: $C_0(i), i \in \mathbb{N}$.
- **2** Encoding of the data: $C_{\mathcal{M}_i}(D)$, $D \in \mathcal{D}$, where $C_{\mathcal{M}_i}$ is a universal code (e.g., mixture, NML) based on model class \mathcal{M}_i .

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

The idea is the same as in the Kolmogorov minimal sufficient statistic (ideal MDL): **Extract all the structure from the data.**

イロト イポト イヨト イヨト

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

The idea is the same as in the Kolmogorov minimal sufficient statistic (ideal MDL): **Extract all the structure from the data.**

MDL Explanation of MDL

The success in extracting the structure from data can be measured by the codelength.

イロト イポト イヨト イヨト

Definitions Universal Models Prediction & Model Selection

MDL Model Selection

The idea is the same as in the Kolmogorov minimal sufficient statistic (ideal MDL): **Extract all the structure from the data.**

MDL Explanation of MDL

The success in extracting the structure from data can be measured by the codelength.

In practical MDL, we only find the structure that is 'visible' to the used model class(es). For instance, the Bernoulli (coin flipping) model only sees the number of 1s.

< ロ > < 同 > < 三 > < 三 >

Definitions Universal Models Prediction & Model Selection

MDL & Bayes

The MDL model selection criterion

minimize $\ell(\theta) + \ell_{\theta}(D)$ can be interpreted (via $p = 2^{-\ell}$) as maximize $p(\theta) \times p_{\theta}(D)$.

Teemu Roos Three Concepts: Information

3

Definitions Universal Models Prediction & Model Selection

MDL & Bayes

The MDL model selection criterion

minimize $\ell(heta) + \ell_{ heta}(D)$ can be interpreted (via $p = 2^{-\ell}$) as

maximize $p(heta) imes p_{ heta}(D)$.

In Bayesian probability, this is equivalent to **maximization of posterior probability**:

$$p(\theta \mid D) = rac{p(\theta) \, p(D \mid \theta)}{p(D)} \; ,$$

where the term p(D) (the marginal probability of D) is constant wrt. θ and doesn't affect model selection.

イロト 人間ト イヨト イヨト

SOR

Definitions Universal Models Prediction & Model Selection

MDL & Bayes

The MDL model selection criterion

minimize $\ell(heta) + \ell_{ heta}(D)$ can be interpreted (via $p = 2^{-\ell}$) as

maximize $p(heta) imes p_{ heta}(D)$.

In Bayesian probability, this is equivalent to **maximization of posterior probability**:

$$p(\theta \mid D) = rac{p(\theta) \, p(D \mid \theta)}{p(D)} \; ,$$

where the term p(D) (the marginal probability of D) is constant wrt. θ and doesn't affect model selection.

⇒ Three Conceps: Probability

・ロト ・ 同ト ・ ヨト ・ ヨト

Definitions Universal Models Prediction & Model Selection

Example: Denoising

Complexity	=	Information	+	Noise
	=	Regularity	+	Randomness
	=	Algorithm	+	Compressed file

イロト イロト イヨト イヨト

3

990

Definitions Universal Models Prediction & Model Selection

Example: Denoising

Complexity	=	Information	+	Noise
	=	Regularity	+	Randomness
	=	Algorithm	+	Compressed file

Denoising means the process of removing noise from a signal.

・ロト ・部 ト ・ヨト ・ヨト

1

Definitions Universal Models Prediction & Model Selection

Example: Denoising

Complexity	=	Information	+	Noise
	=	Regularity	+	Randomness
	=	Algorithm	+	Compressed file

Denoising means the process of removing noise from a signal.

The MDL principle gives a natural method for denoising since the very idea of MDL is to separate the total complexity of a signal into information and noise.

<ロト <同ト < ヨト < ヨト -

Definitions Universal Models Prediction & Model Selection

Example: Denoising

Complexity	=	Information	+	Noise
	=	Regularity	+	Randomness
	=	Algorithm	+	Compressed file

Denoising means the process of removing noise from a signal.

The MDL principle gives a natural method for denoising since the very idea of MDL is to separate the total complexity of a signal into information and noise.

First encode a smooth signal (information), and then the difference to the observed signal (noise).

SOR

Definitions Universal Models Prediction & Model Selection

Example: Denoising



イロト イポト イヨト イヨ

Definitions Universal Models Prediction & Model Selection

Example: Denoising



Definitions Universal Models Prediction & Model Selection

Example: Denoising



Definitions Universal Models Prediction & Model Selection

Example: Denoising



Definitions Universal Models Prediction & Model Selection

Example: Denoising



Teemu Roos Three Concepts: Information

Definitions Universal Models Prediction & Model Selection

Example: Denoising



Teemu Roos Three Concepts: Information

500

Definitions Universal Models Prediction & Model Selection

Last Slide

The End.

Teemu Roos Three Concepts: Information

<ロ> <部> < 部> < き> < き> <</p>

3

SQC