Introduction to Information-Theoretic Modeling

Teemu Roos

Helsinki Institute for Information Technology HIIT, Department of Computer Science, University of Helsinki

BCSMIF, Maresias, April 11-12, 2011





Jorma Rissanen (left) receiving the IEEE Information Theory Society Best Paper Award from Claude Shannon in 1986.

IEEE Golden Jubilee Award for Technological Innovation (for the invention of arithmetic coding) 1998; IEEE Richard W. Hamming Medal (for fundamental contribution to information theory, statistical inference, control theory, and the theory of complexity) 1993; Kolmogorov Medal 2006; **IBM Outstanding Innovation** Award (for work in statistical inference, information theory, and the theory of complexity) 1988: IEEE Claude E. Shannon Award 2009: ...













- 2 MDL Principle
- Oliversal Source Coding





- 2 MDL Principle
- Oniversal Source Coding
- MDL Principle (contd.)



Occam's Razor MDL Principle (contd.)

House Razor

Occam's Razor House

Razor

MDL Principle

MDL Principle (contd.)



Occam's Razor Universal Source Coding MDL Principle (contd.)

House Razor

House



Teemu Roos Introduction to Information-Theoretic Modeling

House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

イロト イポト イヨト イヨト

Э

House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

No single disease causes all of these.

◆ロ > ◆母 > ◆臣 > ◆臣 >

3

DQ P

House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

No single disease causes all of these.

Each symptom can be caused by some (possibly different) disease...

イロト イポト イヨト イヨト

House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

No single disease causes all of these.

Each symptom can be caused by some (possibly different) disease...



House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

No single disease causes all of these.

Each symptom can be caused by *some* (possibly different) disease...

- pneumonia,
- appendicitis,

House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

No single disease causes all of these.

Each symptom can be caused by *some* (possibly different) disease...



- appendicitis,
- food poisoning,

House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

No single disease causes all of these.

Each symptom can be caused by *some* (possibly different) disease...

- pneumonia,
- appendicitis,
- food poisoning,
- hemorrhage,

House Razor

House

Brandon has

- cough,
- evere abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

- pneumonia,
- appendicitis,
- food poisoning,
- hemorrhage,
- o meningitis.

No single disease causes all of these.

Each symptom can be caused by some (possibly different) disease...

House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

- pneumonia,
 appendicitis,
- food poisoning,
- hemorrhage,
- o meningitis.

No single disease causes all of these.

Each symptom can be caused by some (possibly different) disease...

Dr. House explains the symptoms with two simple causes:

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

House Razor

House

Brandon has

- cough,
- severe abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

- common cold,
 appendicitis,
 food poisoning,
- hemorrhage,
- common cold.

No single disease causes all of these.

Each symptom can be caused by some (possibly different) disease...

Dr. House explains the symptoms with two simple causes:

common cold, causing the cough and fever,

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

House Razor

House

Brandon has

- cough,
- evere abdominal pain,
- Inausea,
- Iow blood pressure,
- 6 fever.

- common cold,
 gout medicine,
 gout medicine,
 gout medicine,
- ommon cold.

No single disease causes all of these.

Each symptom can be caused by some (possibly different) disease...

Dr. House explains the symptoms with two simple causes:

- common cold, causing the cough and fever,
- In pharmacy error: cough medicine replaced by gout medicine.

House Razor

William of Ockham (c. 1288–1348)



→ E + < E +</p>

House Razor

Occam's Razor

Occam's Razor

Entities should not be multiplied beyond necessity.

イロト イポト イヨト イヨト

nar

House Razor

Occam's Razor

Occam's Razor

Entities should not be multiplied beyond necessity.

Isaac Newton: "We are to admit no more causes of natural things than such as are both true and sufficient to explain their appearances."

House Razor

Occam's Razor

Occam's Razor

Entities should not be multiplied beyond necessity.

Isaac Newton: "We are to admit no more causes of natural things than such as are both true and sufficient to explain their appearances."

Diagnostic parsimony: Find the fewest possible causes that explain the symptoms.

◆ロ > ◆母 > ◆臣 > ◆臣 >

House Razor

Occam's Razor

Occam's Razor

Entities should not be multiplied beyond necessity.

Isaac Newton: "We are to admit no more causes of natural things than such as are both true and sufficient to explain their appearances."

Diagnostic parsimony: Find the fewest possible causes that explain the symptoms.

(Hickam's dictum: "Patients can have as many diseases as they damn well please.")

House Razor

Guessing Game



◆ロ > ◆母 > ◆臣 > ◆臣 >

E

House Razor

Guessing Game



◆ロ > ◆母 > ◆臣 > ◆臣 >

E

House Razor

Guessing Game



◆ロ > ◆母 > ◆臣 > ◆臣 >

E

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

・ロト ・回ト ・モト ・モト

E

House Razor

Guessing Game



・ロト ・回ト ・モト ・モト

E

House Razor

Guessing Game



・ロト ・回ト ・モト ・モト

E

House Razor

Guessing Game



・ロト ・回ト ・モト ・モト

E

House Razor

Guessing Game



・ロト ・回ト ・モト ・モト

E

House Razor

Guessing Game



・ロト ・回ト ・モト ・モト

E

House Razor

Guessing Game



・ロト ・回ト ・モト ・モト

E

House Razor

Guessing Game



・ロト ・四ト ・ヨト ・ヨト

E

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

・ロト ・四ト ・ヨト ・ヨト

E
House Razor

Guessing Game



・ロト ・四ト ・ヨト ・ヨト

E

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

・ロト ・四ト ・ヨト ・ヨト

E

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

・ロト ・日 ・ ・ 日 ・ ・ 日 ・

E

House Razor

Guessing Game



▲ロト ▲圖ト ▲国ト ▲国ト

E

House Razor

Guessing Game



・ロト ・四ト ・ヨト ・ヨト

1

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

◆ロト ◆部 ト ◆注 ト ◆注 ト

1

House Razor

Guessing Game



◆ロト ◆部 ト ◆注 ト ◆注 ト

1

House Razor

Guessing Game



◆ロト ◆部 ト ◆注 ト ◆注 ト

1

House Razor

Guessing Game



◆ロト ◆部 ト ◆注 ト ◆注 ト

1

House Razor

Guessing Game



<ロ> <同> <同> <同> < 同> < 同> < 同> <

1

House Razor

Guessing Game



▲ロト ▲圖ト ▲国ト ▲国ト

1

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

<ロ> <同> <同> <同> < 同> < 同> < 同> <

1

House Razor

Guessing Game



▲ロト ▲圖ト ▲国ト ▲国ト

1

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

<ロト < 回 > < 回 > < 回 > < 回 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

House Razor

Guessing Game



<ロト < 回 > < 回 > < 回 > < 回 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

House Razor

Guessing Game



<ロト < 回 > < 回 > < 回 > < 回 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

House Razor

Guessing Game



<ロト < 回 > < 回 > < 回 > < 回 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

<ロト < 回 > < 回 > < 回 > < 回 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

House Razor

Guessing Game



Teemu Roos Introduction to Information-Theoretic Modeling

House Razor

Guessing Game



Introduction to Information-Theoretic Modeling

House Razor

Guessing Game



Teemu Roos

Introduction to Information-Theoretic Modeling

House Razor

Guessing Game



House Razor

Guessing Game



Teemu Roos

Introduction to Information-Theoretic Modeling

House Razor

Guessing Game



Teemu Roos

Introduction to Information-Theoretic Modeling

Rules & Exceptions Probabilistic Models

Occam's Razor

- 2 MDL Principle
 Rules & Exceptions
 Probabilistic Models
- 3 Universal Source Coding
- MDL Principle (contd.)



Rules & Exceptions Probabilistic Models

MDL Principle

Minimum Description Length (MDL) Principle (2-part)

Choose the hypothesis which minimizes the sum of

- the codelength of the hypothesis, and
- the codelength of the data with the help of the hypothesis.

イロト イポト イヨト イヨト

Rules & Exceptions Probabilistic Models

MDL Principle

Minimum Description Length (MDL) Principle (2-part)

Choose the hypothesis which minimizes the sum of

- the codelength of the hypothesis, and
- the codelength of the data with the help of the hypothesis.

 $\min_{h\in\mathcal{H}}(\ell(h)+\ell(D;h))$

イロト イポト イヨト イヨト

Rules & Exceptions Probabilistic Models

MDL Principle

Minimum Description Length (MDL) Principle (2-part)

Choose the hypothesis which minimizes the sum of

- the codelength of the hypothesis, and
- the codelength of the data with the help of the hypothesis.

 $\min_{h\in\mathcal{H}}(\ell(h)+\ell(D\,;\,h))$

How to encode data with the help of a hypothesis?

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.

・ロト ・回ト ・ヨト ・ヨト

1

DQ P

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

- 4 同 ト 4 三 ト

MQ (P

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

For image of size n = 625, there are 2^n different images, and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

different groups of k exceptions.

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

For image of size n = 625, there are 2^n different images, and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

different groups of k exceptions.

$$k = 1: \binom{n}{1} = 625 \ll 2^{625} \approx 1.4 \times 10^{188}.$$

Codelength $\log_2(n+1) + \log_2\binom{n}{k} \approx 19$ vs. $\log_2 2^{625} = 625$

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

For image of size n = 625, there are 2^n different images, and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

different groups of
$$k$$
 exceptions.
 $k = 2: \binom{n}{2} = 195\,000 \ll 2^{625} \approx 1.4 \times 10^{188}.$
Codelength $\log_2(n+1) + \log_2\binom{n}{k} \approx 27 \text{ vs.} \log_2 2^{625} = 625$

Teemu Roos

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

For image of size n = 625, there are 2^n different images, and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

different groups of
$$k$$
 exceptions.
 $k = 3: \binom{n}{3} = 40\,495\,000 \ll 2^{625} \approx 1.4 \times 10^{188}.$
Codelength $\log_2(n+1) + \log_2\binom{n}{k} \approx 35$ vs. $\log_2 2^{625} = 625$

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

For image of size n = 625, there are 2^n different images, and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

different groups of k exceptions. $k = 10: \binom{n}{10} = 2\,331\,354\,000\,000\,000\,000\,000 \ll 2^{625}.$ Codelength $\log_2(n+1) + \log_2\binom{n}{k} \approx 80$ vs. $\log_2 2^{625} = 625$

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

For image of size n = 625, there are 2^n different images, and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

different groups of
$$k$$
 exceptions.
 $k = 100 : \binom{n}{100} \approx 9.5 \times 10^{117} \ll 2^{625} \approx 1.4 \times 10^{188}.$
Codelength $\log_2(n+1) + \log_2\binom{n}{k} \approx 401$ vs. $\log_2 2^{625} = 625$

Teemu Roos
Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

For image of size n = 625, there are 2^n different images, and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

different groups of k exceptions.

$$k = 300: \binom{n}{300} \approx 2.7 \times 10^{186} < 2^{625} \approx 1.4 \times 10^{188}.$$
Codelength log₂(n+1) + log₂ $\binom{n}{k} \approx 629$ vs. log₂ $2^{625} = 625$

Teemu Roos

Rules & Exceptions Probabilistic Models

Encoding Data: Rules & Exceptions

Idea 1: Hypothesis = rule; encode exceptions.



Black box of size $25 \times 25 = 625$, white dots at $(x_1, y_1), (x_2, y_2), (x_3, y_3)$.

For image of size n = 625, there are 2^n different images, and

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

different groups of
$$k$$
 exceptions.
 $k = 372: \binom{n}{372} \approx 5.1 \times 10^{181} \ll 2^{625} \approx 1.4 \times 10^{188}.$
Codelength $\log_2(n+1) + \log_2\binom{n}{k} \approx 613$ vs. $\log_2 2^{625} = 625$

Teemu Roos

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 2: Hypothesis = probability distribution.

<ロト < 同ト < ヨト < ヨト -

3

SQA

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 2: Hypothesis = probability distribution.



▲ 同 ▶ ▲ 国 ▶ ▲ 国 ▶

MQ (P

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 2: Hypothesis = probability distribution.

Important Observation

Probability distributions are codes are probability distributions!

イロト イポト イヨト イヨト

nar

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 2: Hypothesis = probability distribution.

Important Observation

Probability distributions are codes are probability distributions!

The code-length of the data is given by

$$\ell(D; h) = \log_2 \frac{1}{p_h(D)}$$

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 2: Hypothesis = probability distribution.

Important Observation

Probability distributions are codes are probability distributions!

The code-length of the data is given by

$$\ell(D; h) = \log_2 \frac{1}{p_h(D)}$$

Remember to encode distribution too: $\ell(h) > 0$.

Rules & Exceptions Probabilistic Models

MDL & Bayes

The MDL model selection criterion

minimize $\ell(h) + \ell(D; h)$

can be interpreted (via $p = 2^{-\ell}$) as

maximize $p(h) \times p_h(D)$.

◆ロ > ◆母 > ◆臣 > ◆臣 >

3

DQ P

Rules & Exceptions Probabilistic Models

MDL & Bayes

The MDL model selection criterion

minimize $\ell(h) + \ell(D; h)$

can be interpreted (via $p=2^{-\ell}$) as

maximize $p(h) \times p_h(D)$.

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

$$p(h \mid D) = \frac{p(h) p(D \mid h)}{p(D)} ,$$

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

化口下 化晶下 化原下化原下

Rules & Exceptions Probabilistic Models

MDL & Bayes

ENVIRONMETRICS Environmetrics 2001; 12: 559–568 (DOI: 10.1002/env.482)

Model selection: Full Bayesian approach

Carlos Alberto de Bragança Pereira*,† and Julio Michael Stern‡

BIOINFO and IME-USP - University of Sao Paulo, Brazil

(日) (同) (三) (三)

1

DQ P

Rules & Exceptions Probabilistic Models

MDL & Bayes

ENVIRONMETRICS Environmetrics 2001; 12: 559–568 (DOI: 10.1002/env.482)

We can use the FBST as a model selection criterion, testing the hypothesis of some of its parameters being null, and using the following version of the 'Ockham razor: Do not include in the model a new parameter unless there is strong evidence it is not null.'

The FBST selection criterion has an intrinsic regularization mechanism, under some general circumstances discussed later.

Carlos Alberto de Bragança Pereira*,† and Julio Michael Stern[‡]

BIOINFO and IME-USP - University of Sao Paulo, Brazil

(日) (同) (三) (三)

Rules & Exceptions Probabilistic Models

MDL & Bayes

ENVIRONMETRICS Environmetrics 2001; 12: 559–568 (DOI: 10.1002/env.482)

We can use the FBST as a model selection criterion, testing the hypothesis of some of its parameters being null, and using the following version of the 'Ockham razor: Do not include in the model a new parameter unless there is strong evidence it is not null.'

The FBST selection criterion has an intrinsic regularization mechanism, under some general circumstances discussed later.

Carlos Alberto de Bragança Pereira*,† and Julio Michael Stern[‡]

BIOINFO and IME-USP - University of Sao Paulo, Brazil

"Do not include in the model a new parameter unless there is strong evidence it is not null."

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 3: Hypothesis = set of probability distributions = model class.

・ロト ・部 ト ・ヨト ・ヨト

3

DQ P

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 3: Hypothesis = set of probability distributions = model class.



Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 3: Hypothesis = set of probability distributions = model class.

Universal Coding

A **universal code** achieves almost as short a code-length as the code based on the **best distribution** in the model class.

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 3: Hypothesis = set of probability distributions = model class.

Universal Coding

A **universal code** achieves almost as short a code-length as the code based on the **best distribution** in the model class.

Different types of universal codes:

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 3: Hypothesis = set of probability distributions = model class.

Universal Coding

A **universal code** achieves almost as short a code-length as the code based on the **best distribution** in the model class.

Different types of universal codes:

two-part code,

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 3: Hypothesis = set of probability distributions = model class.

Universal Coding

A **universal code** achieves almost as short a code-length as the code based on the **best distribution** in the model class.

Different types of universal codes:

- two-part code,
- Inixture code,

Rules & Exceptions Probabilistic Models

Encoding Data: Probabilistic Models

Idea 3: Hypothesis = set of probability distributions = model class.

Universal Coding

A **universal code** achieves almost as short a code-length as the code based on the **best distribution** in the model class.

Different types of universal codes:

- two-part code,
- 2 mixture code,
- onrmalized maximum likelihood (NML) code.

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Occam's Razor

- 2 MDL Principle
- Oniversal Source Coding
 - Two-Part Codes
 - Mixture Codes
 - Normalized Maximum Likelihood
 - Universal Prediction

4 MDL Principle (contd.)

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal models

The typical situation might be as follows:

Teemu Roos Introduction to Information-Theoretic Modeling

(日) (同) (三) (三)

3

DQ P

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal models

The typical situation might be as follows:

• We know (think) that the source symbols are generated by a Bernoulli model with parameter $p \in [0, 1]$.

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal models

The typical situation might be as follows:

- We know (think) that the source symbols are generated by a Bernoulli model with parameter $p \in [0, 1]$.
- **2** However, we do not know p in advance.

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal models

The typical situation might be as follows:

- We know (think) that the source symbols are generated by a Bernoulli model with parameter p ∈ [0, 1].
- **2** However, we do not know p in advance.
- We'd like to encode data at rate H(p).

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Two-Part Codes

Let $\mathcal{M} = \{p_{\theta} : \theta \in \Theta\}$ be a parametric probabilistic model class, i.e., a set of distributions p_{θ} indexed by parameter θ .

(日) (同) (三) (三)

MQ (P

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Two-Part Codes

Let $\mathcal{M} = \{p_{\theta} : \theta \in \Theta\}$ be a parametric probabilistic model class, i.e., a set of distributions p_{θ} indexed by parameter θ .

For any distribution p_{θ} , the Shannon code-lengths satisfy

$$\ell_ heta(D) = \log_2 rac{1}{
ho_ heta(D)}$$
 .

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Two-Part Codes

Let $\mathcal{M} = \{p_{\theta} : \theta \in \Theta\}$ be a parametric probabilistic model class, i.e., a set of distributions p_{θ} indexed by parameter θ .

For any distribution p_{θ} , the Shannon code-lengths satisfy

$$\ell_{ heta}(D) = \log_2 rac{1}{p_{ heta}(D)}$$

Using parameter value θ , the total code-length becomes

$$\ell_1(heta) + \log_2 rac{1}{p_ heta(D)}$$

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

◆ロ > ◆母 > ◆臣 > ◆臣 >

1

DQ P

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

Solution: Quantization. Choose a discrete subset of points, $\theta^{(1)}, \theta^{(2)}, \ldots$, and use only them.

(日) (同) (三) (三)

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

Solution: Quantization. Choose a discrete subset of points, $\theta^{(1)}, \theta^{(2)}, \ldots$, and use only them.



Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

Solution: Quantization. Choose a discrete subset of points, $\theta^{(1)}, \theta^{(2)}, \ldots$, and use only them.



Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

Solution: Quantization. Choose a discrete subset of points, $\theta^{(1)}, \theta^{(2)}, \ldots$, and use only them.



Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

Solution: Quantization. Choose a discrete subset of points, $\theta^{(1)}, \theta^{(2)}, \ldots$, and use only them.



Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

Solution: Quantization. Choose a discrete subset of points, $\theta^{(1)}, \theta^{(2)}, \ldots$, and use only them.



イロト イポト イヨト イヨト

MQ (P

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

Solution: Quantization. Choose a discrete subset of points, $\theta^{(1)}, \theta^{(2)}, \ldots$, and use only them.



If the points are sufficiently *dense* (in a code-length sense) then the code-length for data is still almost as short as $\min_{\theta \in \Theta} \ell_{\theta}(D)$.

・ロト ・ 一 ト ・ 日 ト ・ 日 ト

MQ (P

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Continuous Parameters

What if the parameters are continuous (like polynomial coefficients)?

Solution: Quantization. Choose a discrete subset of points, $\theta^{(1)}, \theta^{(2)}, \ldots$, and use only them.



If the points are sufficiently *dense* (in a code-length sense) then the code-length for data is still almost as short as $\min_{\theta \in \Theta} \ell_{\theta}(D)$.

・ロト ・ 一 ト ・ 日 ト ・ 日 ト

DQ P
Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Mixture Universal Model

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

(日) (同) (三) (三)

3

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

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

$$w(heta)=2^{-\ell(heta)}$$
 .

(4月) (4日) (4日)

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

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

$$w(heta) = 2^{-\ell(heta)}$$

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_{θ} is weighted by $w(\theta)$.

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood

Consider again 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} .

・ロト ・四ト ・ヨト ・ヨト

1

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood

Consider again 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.

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood

Normalized Maximum Likelihood

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

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

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood

Normalized Maximum Likelihood

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

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

NML code-length:

$$\ell(D) = \log_2 \frac{1}{p_{\text{nml}}(D)} = \log_2 \frac{1}{p_{\hat{\theta}}(D)} + \log_2 C$$
.

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood

Normalized Maximum Likelihood

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

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

NML code-length:

$$\ell(D) = \log_2 \frac{1}{p_{\text{nml}}(D)} = \log_2 \frac{1}{p_{\hat{\theta}}(D)} + \log_2 C$$
.

The more flexible (complex) the model class, the greater the normalizing constant. $(\Box \rightarrow \langle \overline{a} \rangle \land \overline{a})$

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood



Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood



Teemu Roos

Introduction to Information-Theoretic Modeling

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood



Teemu Roos

Introduction to Information-Theoretic Modeling

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Normalized Maximum Likelihood



Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal Models

We have seen three kinds of universal codes:

- two-part,
- Ø mixture,
- INML.

(日) (同) (三) (三)

3

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal Models

We have seen three kinds of universal codes:

- two-part,
- Ø mixture,
- Interpretation NML.

There are also universal codes that are not based on any (explicit) model class: Lempel-Ziv (WinZip, gzip)!

Two-Part Codes Mixture Codes Normalized Maximum Likelihood **Universal Prediction**

Uses of Universal Codes

So what do we do with them?

・ロト ・四ト ・ヨト ・ヨト

Э

SQC

Two-Part Codes Mixture Codes Normalized Maximum Likelihood **Universal Prediction**

Uses of Universal Codes

So what do we do with them?

We can use universal codes for (at least) three purposes:

イロト イポト イヨト イヨト

1

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Uses of Universal Codes

So what do we do with them?

We can use universal codes for (at least) three purposes:

compression,

イロト イポト イヨト イヨト

1

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Uses of Universal Codes

So what do we do with them?

We can use universal codes for (at least) three purposes:

- compression,
- 2 prediction,

◆ロ > ◆母 > ◆臣 > ◆臣 >

1

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Uses of Universal Codes

So what do we do with them?

We can use universal codes for (at least) three purposes:

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

◆ロ > ◆母 > ◆臣 > ◆臣 >

Two-Part Codes Mixture Codes Normalized Maximum Likelihood **Universal Prediction**

Universal Prediction

By the connection $p(D) = 2^{-\ell(D)}$, the following are equivalent:

• good compression: $\ell(D)$ is small,

◆ロ > ◆母 > ◆臣 > ◆臣 >

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal Prediction

By the connection $p(D) = 2^{-\ell(D)}$, the following are equivalent:

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

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal Prediction

By the connection $p(D) = 2^{-\ell(D)}$, the following are equivalent:

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

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

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Universal Prediction

By the connection $p(D) = 2^{-\ell(D)}$, the following are equivalent:

- good compression: $\ell(D)$ is small,
- good predictions: $p(D_i | D_1, ..., D_{i-1})$ is high 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**.

Two-Part Codes Mixture Codes Normalized Maximum Likelihood Universal Prediction

Model (Class) Selection

Since a model class that enables good compression of the data must be based on exploiting the **regular features in the data**, the code-length can be used as a **yard-stick** for comparing model classes.

< ロ > < 同 > < 三 > < 三 >

Two-Part Codes Mixture Codes Normalized Maximum Likelihood **Universal Prediction**

Time for a break?

Teemu Roos Introduction to Information-Theoretic Modeling

<ロ> <部> < 部> < き> < き> <</p>

3

990

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

1 Occam's Razor

2 MDL Principle

Oniversal Source Coding

- MDL Principle (contd.)
 - Modern MDL
 - Histogram Density Estimation
 - Clustering
 - Linear Regression
 - Wavelet Denoising



Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

MDL Principle

MDL Principle

"Old-style":

• Choose the model $p_{ heta} \in \mathcal{M}$ that yields the shortest *two-part* code-length

$$\min_{ heta,\mathcal{M}} \ \ell(\mathcal{M}) + \ell_1(heta) + \log_2 rac{1}{p_ heta(D)}.$$

Modern:

• Choose the model class \mathcal{M} that yields the shortest *universal* code-length

$$\min_{\mathcal{M}} \ell(\mathcal{M}) + \ell_{\mathcal{M}}(D).$$

naa

MDL Model Selection

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Intuitive Explanation of MDL

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

◆ロ > ◆母 > ◆臣 > ◆臣 >

MQ (P

MDL Model Selection

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Intuitive Explanation of MDL

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

We can only extract 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 > ◆母 > ◆臣 > ◆臣 >

MDL Model Selection

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Intuitive Explanation of MDL

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

We can only extract the structure that is "visible" to the used model class(es). For instance, the Bernoulli (coin flipping) model only sees the number of 1s.

When the model can express the structural properties pertaining to the data *but not more*, the total code-length is minimal.

MDL Model Selection

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Intuitive Explanation of MDL

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

We can only extract the structure that is "visible" to the used model class(es). For instance, the Bernoulli (coin flipping) model only sees the number of 1s.

When the model can express the structural properties pertaining to the data *but not more*, the total code-length is minimal.

Important: Too complex models lead to a long total code-length (Occam's Razor!).

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Multinomial Models

The multinomial model — the generalization of Bernoulli — is very simple:

$$\Pr(X = j) = \theta_j$$
, for $j \in \{1, \dots, m\}$.

イロト イポト イヨト イヨト

Э

SQA

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Multinomial Models

The multinomial model — the generalization of Bernoulli — is very simple:

$$\Pr(X = j) = \theta_j$$
, for $j \in \{1, \dots, m\}$.

Maximum likelihood:

$$\hat{\theta}_j = \frac{\#\{x_i = j\}}{n}.$$

Э

SQA

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Multinomial Models

The multinomial model — the generalization of Bernoulli — is very simple:

$$\Pr(X = j) = \theta_j$$
, for $j \in \{1, \dots, m\}$.

Maximum likelihood:

$$\hat{\theta}_j = \frac{\#\{x_i = j\}}{n}.$$

Two-part, mixture, and NML models readily defined.

-

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Multinomial Models

The multinomial model — the generalization of Bernoulli — is very simple:

$$\Pr(X = j) = \theta_j$$
, for $j \in \{1, \dots, m\}$.

Maximum likelihood:

$$\hat{\theta}_j = \frac{\#\{x_i = j\}}{n}.$$

Two-part, mixture, and NML models readily defined. \Rightarrow Exercises 7–9.

-

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Fast NML for Multinomials

The naïve way to compute the normalizing constant in the NML model

$$\frac{p_{\hat{\theta}}(x^n)}{C_n^m}, \qquad C_n^m = \sum_{y^n \in \mathcal{X}^n} p_{\hat{\theta}}(y^n),$$

takes exponential time $(\Omega(m^n))$.

イロト イポト イヨト イヨト

MQ (P
Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Fast NML for Multinomials

The naïve way to compute the normalizing constant in the NML model

takes exponential time $(\Omega(m^n))$.

However, there is a trick (Myllymäki & Kontkanen, 2007) to do it in linear time.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Fast NML for Multinomials

The naïve way to compute the normalizing constant in the NML model

takes exponential time $(\Omega(m^n))$.

However, there is a trick (Myllymäki & Kontkanen, 2007) to do it in linear time.

Kontkanen & Myllymäki, "A linear-time algorithm for computing the multinomial stochastic complexity", *Information Processing Letters* **103** (2007), 6, pp. 227–233

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Histogram Density Estimation

A histogram density is defined by

• The break-points between the bins,

イロト イポト イヨト イヨト

1

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Histogram Density Estimation

A histogram density is defined by

- The break-points between the bins,
- The heights of the bins.

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Histogram Density Estimation

A histogram density is defined by

- The break-points between the bins,
- The heights of the bins.

Choosing the number *and the positions* of break-points can be done by MDL.

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Histogram Density Estimation

A histogram density is defined by

- The break-points between the bins,
- The heights of the bins.

Choosing the number *and the positions* of break-points can be done by MDL.

The code-length is equivalent (up to additive constants) to the code-length in a multinomial model.

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Histogram Density Estimation

A histogram density is defined by

- The break-points between the bins,
- The heights of the bins.

Choosing the number *and the positions* of break-points can be done by MDL.

The code-length is equivalent (up to additive constants) to the code-length in a multinomial model.

 \Rightarrow Linear-time algorithm can be used.

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Histogram Density Estimation

MDL Histogram Density Estimation

Petri Kontkanen, Petri Myllymäki

Complex Systems Computation Group (CoSCo) Helsinki Institute for Information Technology (HIIT) University of Helsinki and Helsinki University of Technology P.O.Box 68 (Department of Computer Science) FIN-00014 University of Helsinki, Finland {Firstname}.{Lastname}Miit.fi

Abstract

We regard histogram density estimation as a model selection problem. Our approach is based on the information-theoretic minimum description length (MDL) principle, which can be applied for tasks such as data clustering, density estimation, image denoising and model calculation in gravel. MDL only on finding the optimal bin count. These regular histograms are, however, often problematic. It has been argued (Rissanen, Speed, & Yu, 1992) that regular histograms are only good for describing roughly uniform data. If the data distribution is strongly nonuniform, the bin count must necessarily be high if one wants to capture the details of the high density portion of the data. This in turn means that an unnecessary large amount of bins is wasted in the low density re-

1Q (2

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Histogram Density Estimation



Teemu Roos Introduction to Information-Theoretic Modeling

3

5900

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Clustering

Consider the problem of clustering vectors of (independent) multinomial variables.

◆ロ > ◆母 > ◆臣 > ◆臣 >

3

SQA

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Clustering

Consider the problem of clustering vectors of (independent) multinomial variables.

This can be seen as a way to encode (compress) the data:

イロト イポト イヨト イヨト

1

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Clustering

Consider the problem of clustering vectors of (independent) multinomial variables.

This can be seen as a way to encode (compress) the data:

If irst encode the cluster index of each observation vector,

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Clustering

Consider the problem of clustering vectors of (independent) multinomial variables.

This can be seen as a way to encode (compress) the data:

- I first encode the cluster index of each observation vector,
- then encode the observations using separate (multinomial) models.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

MQ (P

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Clustering

Consider the problem of clustering vectors of (independent) multinomial variables.

This can be seen as a way to encode (compress) the data:

- I first encode the cluster index of each observation vector,
- then encode the observations using separate (multinomial) models.

Again, the problem is reduced to the multinomial case, and the fast NML algorithm can be applied.

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Clustering

The clustering model can be interpreted as the **naïve Bayes** structure:



label = cluster index

 f_1, \ldots, f_n are *features*

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Clustering

The clustering model can be interpreted as the **naïve Bayes** structure:



label = cluster index f_1, \ldots, f_n are *features*

The structure is very restrictive. Generalization achieved by **Bayesian networks**.

イロト イポト イヨト イヨト

MQ (P

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Clustering

The clustering model can be interpreted as the **naïve Bayes** structure:



label = cluster index f_1, \ldots, f_n are *features*

The structure is very restrictive. Generalization achieved by **Bayesian networks**.

MDL criterion for learning Bayesian network structures again based on *fast NML for multinomials*.

MQ (P

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Bayesian Networks



Directed asyclic graph (DAG) describing conditional independencies (causality?).

イロト イポト イヨト イヨト

nar

3

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Bayesian Networks



Directed asyclic graph (DAG) describing conditional independencies (causality?).

Denote parents by $\operatorname{Pa}_i \subset \{X_1, \ldots, X_m\} \setminus X_i$.

◆ロ > ◆母 > ◆臣 > ◆臣 >

3

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Bayesian Networks



Directed asyclic graph (DAG) describing conditional independencies (causality?).

Denote parents by $\operatorname{Pa}_i \subset \{X_1, \ldots, X_m\} \setminus X_i$.

Joint probability

$$p(x_1,\ldots,x_m) = \prod_{j=1}^m p(x_i \mid \mathrm{pa}_i).$$

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Bayesian Networks

Problem with NML for Bayesian networks is the summation over all possible data-sets: exponential complexity.

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Bayesian Networks

Problem with NML for Bayesian networks is the summation over all possible data-sets: exponential complexity.

Approximation: normalization over smaller blocks at a time.



イロト イポト イヨト イヨト

MQ (P

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Bayesian Networks

Problem with NML for Bayesian networks is the summation over all possible data-sets: exponential complexity.

Approximation: normalization over smaller blocks at a time.



 \Rightarrow Factorized NML

イロト イポト イヨト イヨト

MQ (P

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising



Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising



Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising



Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising



Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

How to Encode Continuous Data?

In order to encode data using, say, the Gaussian density we face the question: How to encode continuous data?

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

How to Encode Continuous Data?

In order to encode data using, say, the Gaussian density we face the question: How to encode continuous data?

We already know how to encode using models with continuous *parameters*:

イロト イポト イヨト イヨト

MQ (P

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

How to Encode Continuous Data?

In order to encode data using, say, the Gaussian density we face the question: How to encode continuous data?

We already know how to encode using models with continuous *parameters*:

• two-part with optimal quantization $(\approx \frac{k}{2} \log_2 n)$,

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

How to Encode Continuous Data?

In order to encode data using, say, the Gaussian density we face the question: How to encode continuous data?

We already know how to encode using models with continuous *parameters*:

- two-part with optimal quantization $(\approx \frac{k}{2} \log_2 n)$,
- mixture code,

◆ロ > ◆母 > ◆臣 > ◆臣 >

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

How to Encode Continuous Data?

In order to encode data using, say, the Gaussian density we face the question: How to encode continuous data?

We already know how to encode using models with continuous *parameters*:

- two-part with optimal quantization $(\approx \frac{k}{2} \log_2 n)$,
- mixture code,
- NML.

◆ロ > ◆母 > ◆臣 > ◆臣 >

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

How to Encode Continuous Data?

In order to encode data using, say, the Gaussian density we face the question: How to encode continuous data?

We already know how to encode using models with continuous *parameters*:

- two-part with optimal quantization $(\approx \frac{k}{2} \log_2 n)$,
- mixture code,
- NML.

Obviously not possible to encode data with infinite precision. Have to **discretize**: encode x only up to precision δ .

< ロ > < 同 > < 回 > < 回 > < 回 > <

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

.

Э

SQC

Gaussian Density

Recall the Gaussian density function:

$$\phi_{\mu,\sigma^2}(x_1,\ldots,x_n) \stackrel{(i.i.d.)}{=} (2\pi\sigma^2)^{-n/2} e^{-\frac{\sum_{i=1}^n (x_i-\mu)^2}{2\sigma^2}}$$

<ロト <回ト < 回ト < 回ト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

.

3

DQ P

Gaussian Density

Recall the Gaussian density function:

$$\phi_{\mu,\sigma^2}(x_1,\ldots,x_n) \stackrel{(i.i.d.)}{=} (2\pi\sigma^2)^{-n/2} e^{-\frac{\sum_{i=1}^n (x_i-\mu)^2}{2\sigma^2}}$$

The code-length is then

$$\frac{n}{2}\log_2(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{i=1}^n (x_i - \mu)^2.$$

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

Ok, we have our Gaussian code-length formula:

$$\frac{n}{2}\log_2(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{i=1}^n (x_i - \mu)^2.$$

・ロト ・四ト ・ヨト ・ヨト

SQC

3
Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

Ok, we have our Gaussian code-length formula:

$$\frac{n}{2}\log_2(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{i=1}^n (x_i - \mu)^2.$$

Let's use the two-part code and plug in the maximum likelihood parameters:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu})^2.$$

イロト イポト イヨト イヨト

-

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

Ok, we have our Gaussian code-length formula:

$$rac{n}{2}\log_2(2\pi\hat{\sigma}^2) - rac{1}{2\hat{\sigma}^2}\sum_{i=1}^n (x_i - \hat{\mu})^2.$$

Let's use the two-part code and plug in the maximum likelihood parameters:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu})^2.$$

(日) (同) (三) (三)

-

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

Ok, we have our Gaussian code-length formula:

$$\frac{n}{2}\log_2(2\pi\hat{\sigma}^2) - \frac{1}{2\hat{\sigma}^2}\sum_{i=1}^n (x_i - \hat{\mu})^2.$$

Let's use the two-part code and plug in the maximum likelihood parameters:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu})^2.$$

イロト イポト イヨト イヨト

-

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

Ok, we have our Gaussian code-length formula:

$$\frac{n}{2}\log_2(2\pi\hat{\sigma}^2) - \frac{n}{2}$$

Let's use the two-part code and plug in the maximum likelihood parameters:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu})^2.$$

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

Ok, we have our Gaussian code-length formula:

$$\frac{n}{2}\log_2\hat{\sigma}^2 + constant.$$

Let's use the two-part code and plug in the maximum likelihood parameters:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu})^2.$$

< ロ > < 同 > < 回 > < 回 > < 回 > <

1

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

We get the total (two-part) code-length formula:

$$\frac{n}{2}\log_2\hat{\sigma}^2 + \frac{k}{2}\log_2 n + constant.$$

イロト イポト イヨト イヨト

3

SQA

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

We get the total (two-part) code-length formula:

$$\frac{n}{2}\log_2\hat{\sigma}^2 + \frac{k}{2}\log_2 n + constant.$$

Since we have two parameters, μ and σ^2 , we let k = 2.

◆ロ > ◆母 > ◆臣 > ◆臣 >

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Gaussian Density

We get the total (two-part) code-length formula:

$$\frac{n}{2}\log_2\hat{\sigma}^2 + \frac{2}{2}\log_2 n + constant.$$

Since we have two parameters, μ and σ^2 , we let k = 2.

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

A similar treatment can be given to linear regression models.

(日) (同) (三) (三)

3

SQA

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

A similar treatment can be given to *linear regression models*.

The model includes a set of regressor variables $x_1, \ldots, x_p \in \mathbb{R}$, and a set of coefficients β_1, \ldots, β_p .

◆ロ > ◆母 > ◆臣 > ◆臣 >

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

A similar treatment can be given to linear regression models.

The model includes a set of regressor variables $x_1, \ldots, x_p \in \mathbb{R}$, and a set of coefficients β_1, \ldots, β_p .

The dependent variable, Y, is assumed to be Gaussian:

• the mean μ is given as a linear combination of the regressors:

$$\mu = \beta_1 x_1 + \dots + \beta_p x_p = \beta' x,$$

• variance is some parameter σ^2 .

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

For a sample of size n, the matrix notation is convenient:

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad X = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \quad \epsilon = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

イロト イポト イヨト イヨト

3

SQR

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

For a sample of size *n*, the matrix notation is convenient:

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad X = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \quad \epsilon = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Then the model can be written as

$$Y = X\beta + \epsilon,$$

where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$.

イロト イポト イヨト イヨト

-

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

The maximum likelihood estimators are now

$$\hat{\beta} = (X'X)^{-1}X'Y, \quad \hat{\sigma}^2 = \frac{1}{n} \|Y - X\hat{\beta}\|_2^2 = \frac{\mathrm{RSS}}{n},$$

where RSS is the "residual sum of squares".

(日) (同) (三) (三)

3

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

The maximum likelihood estimators are now

$$\hat{\beta} = (X'X)^{-1}X'Y, \quad \hat{\sigma}^2 = \frac{1}{n} \|Y - X\hat{\beta}\|_2^2 = \frac{\mathrm{RSS}}{n},$$

where RSS is the "residual sum of squares".

Since the errors are assumed Gaussian, our code-length formula applies:

$$\frac{n}{2}\log_2\hat{\sigma}^2 + \frac{k}{2}\log_2 n + \text{constant.}$$

< ロ > < 同 > < 回 > < 回 > < 回 > <

nar

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

The maximum likelihood estimators are now

$$\hat{\beta} = (X'X)^{-1}X'Y, \quad \hat{\sigma}^2 = \frac{1}{n} \|Y - X\hat{\beta}\|_2^2 = \frac{\mathrm{RSS}}{n},$$

where RSS is the "residual sum of squares".

Since the errors are assumed Gaussian, our code-length formula applies:

$$\frac{n}{2}\log_2 \mathrm{RSS} + \frac{k}{2}\log_2 n + \mathrm{constant}.$$

イロト イポト イヨト イヨト

nar

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

The maximum likelihood estimators are now

$$\hat{\beta} = (X'X)^{-1}X'Y, \quad \hat{\sigma}^2 = \frac{1}{n} \|Y - X\hat{\beta}\|_2^2 = \frac{\mathrm{RSS}}{n},$$

where RSS is the "residual sum of squares".

Since the errors are assumed Gaussian, our code-length formula applies:

$$\frac{n}{2}\log_2 \mathrm{RSS} + \frac{k}{2}\log_2 n + \mathrm{constant}.$$

The number of parameters is now p + 1 (p of the β s and σ^2), so we get...

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Linear Regression

The maximum likelihood estimators are now

$$\hat{\beta} = (X'X)^{-1}X'Y, \quad \hat{\sigma}^2 = \frac{1}{n} \|Y - X\hat{\beta}\|_2^2 = \frac{\mathrm{RSS}}{n},$$

where RSS is the "residual sum of squares".

Since the errors are assumed Gaussian, our code-length formula applies:

$$\frac{n}{2}\log_2 \mathrm{RSS} + \frac{p+1}{2}\log_2 n + \mathrm{constant}.$$

The number of parameters is now p + 1 (p of the β s and σ^2), so we get...

イロト イポト イヨト イヨト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Subset Selection Problem

Often we have a large set of potential regressors, some of which may be irrelevant.

・ロト ・回ト ・ヨト ・ヨト

3

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Subset Selection Problem

Often we have a large set of potential regressors, some of which may be irrelevant.

The MDL principle can be used to select a subset of them by comparing the total code-lengths:

$$\min_{\mathcal{S}}\left[\frac{n}{2}\log_2 \mathrm{RSS}_{\mathcal{S}} + \frac{|\mathcal{S}|+1}{2}\log_2 n\right],\,$$

where RSS_S is the RSS obtained by using subset S of the regressors.

◆ロ > ◆母 > ◆臣 > ◆臣 >

MQ (P

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Subset Selection Problem

Often we have a large set of potential regressors, some of which may be irrelevant.

The MDL principle can be used to select a subset of them by comparing the total code-lengths:

$$\min_{\mathcal{S}}\left[\frac{n}{2}\log_2 \mathrm{RSS}_{\mathcal{S}} + \frac{|\mathcal{S}|+1}{2}\log_2 n\right],\,$$

where RSS_S is the RSS obtained by using subset S of the regressors.

 \Rightarrow Exercise 10.

イロト イポト イヨト イヨト

MQ (P

Occam's Razor Universal Source Coding MDL Principle (contd.)

Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Denoising

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

<ロト <回ト < 回ト < 回ト

Э

SQR

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Denoising

Complexity = Information + Noise

- = Regularity + Randomness
- = Algorithm + Compressed file

Denoising means the process of removing noise from a signal.

3

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

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.

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

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).

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Wavelet Denoising

One particularly useful way to obtain the regressor (design) matrix is to use **wavelets**.

・ロト ・四ト ・ヨト ・ヨト

3

SQR

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Wavelet Denoising

One particularly useful way to obtain the regressor (design) matrix is to use **wavelets**.



Image by Gabriel Peyré

5900

イロト イポト イヨト イヨト

Occam's Razor Universal Source Coding MDL Principle (contd.) Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Wavelet Denoising

IEEE TRANS, SIGNAL PROCESSING, VOL. ?, NO. ?, 2009

MDL Denoising Revisited

Teemu Roos Member, Petri Myllymäki, and Jorma Rissanen Fellow

Abstract-We refine and extend an earlier minimum description length (MDL) denoising criterion for wavelet-based denoising. We start by showing that the denoising problem can be reformulated as a clustering problem, where the goal is to obtain separate clusters for informative and non-informative wavelet coefficients, respectively. This suggests two refinements, adding a code-length for the model index, and extending the model in order to account for subband-dependent coefficient distributions. A third refinement is the derivation of soft thresholding inspired by predictive universal coding with weighted mixtures. We propose a practical method incorporating all three refinements, which is shown to achieve good performance and robustness in denoising both artificial and natural signals.

Index Terms-Minimum description length (MDL) principle, wavelets, denoising.

(both of which include the Gaussian and de densities as special cases).

A third approach to denoising is based description length (MDL) principle [16]-[20 ent MDL denoising methods have been su; [21]-[25]. We focus on what we consider MDL approach, namely that of Rissanen [24 is two-fold. First as an immediate result extending the earlier MDL denoising meth new practical method with greatly impro and robustness. Secondly, the denoising p to illustrate theoretical issues related to the involving the problem of unbounded paran and the necessity of encoding the model cl. $\neg \circ \circ$

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Example: Denoising



Teemu Roos Introduction to Information-Theoretic Modeling

<ロト < 団ト < 巨ト < 巨ト

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Example: Denoising



Teemu Roos Introduction to Information-Theoretic Modeling

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Example: Denoising



Teemu Roos Introduction to Information-Theoretic Modeling

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Example: Denoising



Teemu Roos Introduction to Information-Theoretic Modeling

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Example: Denoising



Teemu Roos Introduction to Information-Theoretic Modeling

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Example: Denoising



Teemu Roos Introduction to Information-Theoretic Modeling

Modern MDL Histogram Density Estimation Clustering Linear Regression Wavelet Denoising

Last Slide

Thanks for listening!

Teemu Roos Introduction to Information-Theoretic Modeling

◆ロ > ◆母 > ◆臣 > ◆臣 >

Э