

# Multi-Interval Discretization of Continuous-Valued Attributes for Classification Learning

# MDL : Fundamental idea

## Learning as Data Compression

We assume that each sequence is 10000 bits long :

- 00010001000100010001 ... 0001000100010001000100010001
- 01110100110100100110 ... 1010111010111011000101100010
- 00011000001010100000 ... 0010001000010000001000110000

# MDL : Fundamental idea

## Learning as Data Compression

We assume that each sequence is 10000 bits long :

- for i = 1 to 2500; print "0001"; next; halt
- 01110100110100100110 ... 1010111010111011000101100010
- 00011000001010100000 ... 0010001000010000001000110000

# MDL : Fundamental idea

## Learning as Data Compression

We assume that each sequence is 10000 bits long :

- `for i = 1 to 2500; print "0001"; next; halt`
- `print "011101001101000010101010...1010111010111011000101100010"; halt`
- `00011000001010100000 . . . 0010001000010000001000110000`

# MDL : Fundamental idea

## Learning as Data Compression

We assume that each sequence is 10000 bits long :

- for  $i = 1$  to 2500; print "0001"; next; halt
- print "011101001101000010101010...1010111010111011000101100010"; halt
- can be compressed to some length  $\alpha n$ , with  $0 < \alpha < 1$

# MDL and model selection

- The goal of statistical inference may be cast as trying to find regularity in the data.
- For a given set of hypotheses  $H$  and data set  $D$ , we should try to find the hypothesis or combination of hypotheses in  $H$  that compresses  $D$  most.

# MDL and model selection

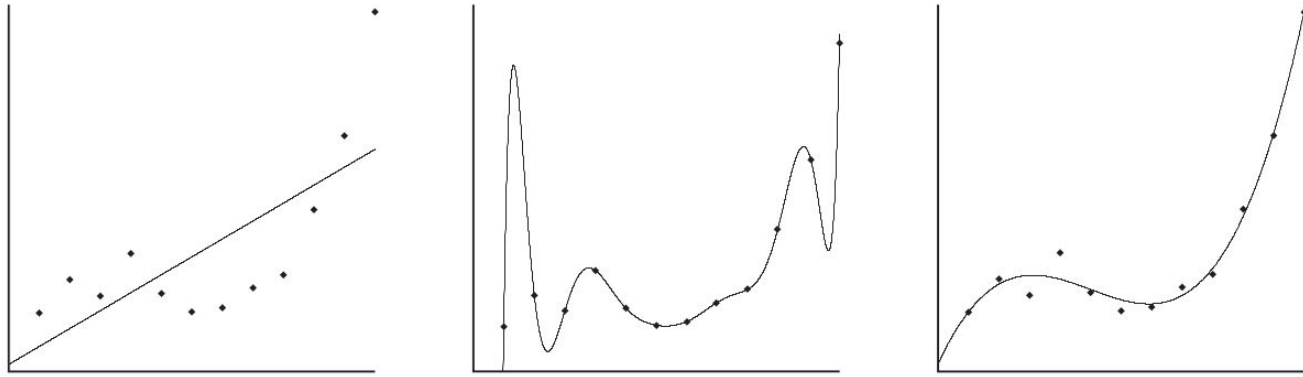
## Crude , Two-Part Version of MDL principle (Informally Stated)

Let  $\mathcal{H}^{(1)}, \mathcal{H}^{(2)}, \dots$  be a list of candidate models (e.g.,  $\mathcal{H}^{(k)}$  is the set of  $k$ th-degree polynomials), each containing a set of point hypotheses (e.g., individual polynomials). The best point hypothesis  $H \in \mathcal{H}^{(1)} \cup \mathcal{H}^{(2)} \cup \dots$  to explain the data  $D$  is the one which minimizes the sum  $L(H) + L(D|H)$ , where

- $L(H)$  is the length, in bits, of the description of the hypothesis; and
- $L(D|H)$  is the length, in bits, of the description of the data when encoded with the help of the hypothesis.

The best *model* to explain  $D$  is the smallest model containing the selected  $H$ .

# MDL and model selection

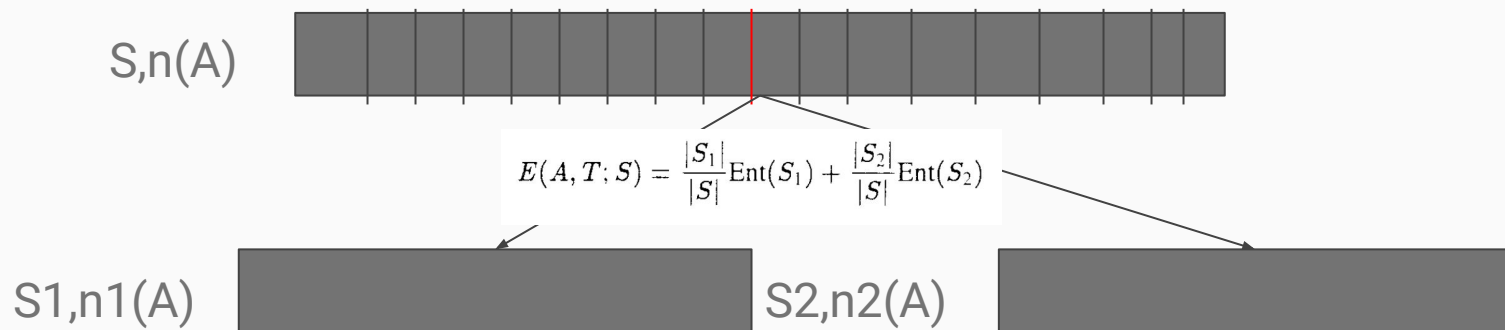


A simple, complex and tradeoff (third-degree) polynomial.



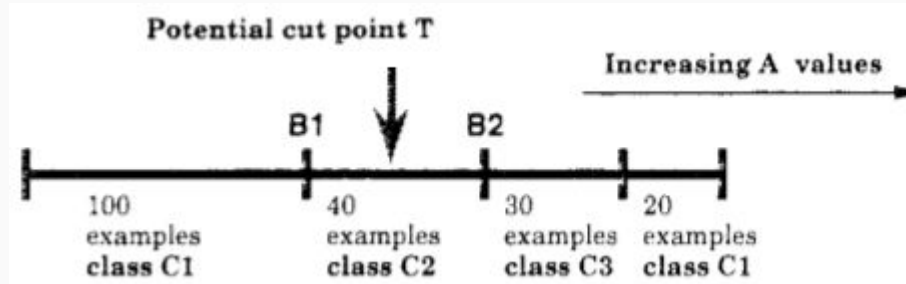
# Fayyad & Irani

- Most real-world applications of classification learning involve continuous-valued attributes.
- The discretization process - crucial - often uses heuristics
- ID3, C4, CART use entropy minimization for binary discretization



# Why multiple ranges rather than binary ranges ?

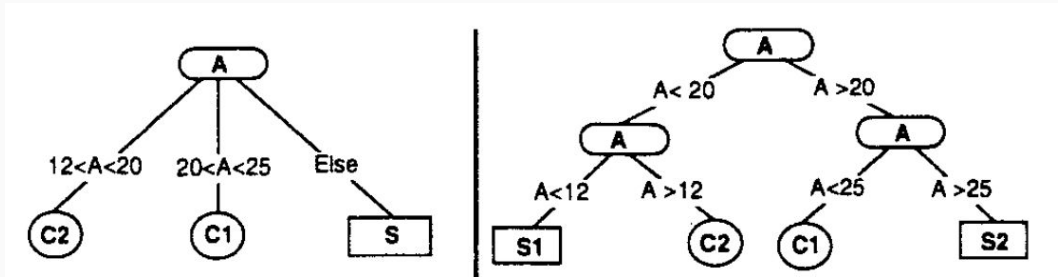
- Although polynomial in complexity, the selection criterion must be calculated  $N-1$  times.
- Entropy minimizing cutpoints are always class boundaries



=>Motive for generalizing the algorithm to generate multiple intervals.

# Why multiple ranges rather than binary ranges ?

- Unnecessary tree growth and irrelevant values problem.



- To cut or not to cut, that is the question !



# Applying MDL, a coding problem.

- Coding the Null Theory:
  - Transmit the classes of the examples in  $S$
  - Huffman( $N$  messages)  $\Rightarrow$  total cost of transmitting  $S$  + overhead
- Coding the Partion Theory:
  - Specify the cut point,  $\log_2(N-1)$
  - Sequence  $S_1$  then Sequence  $S_2$
- Cut if  $\text{Cost}(\text{HT}) < \text{Cost}(\text{NT})$ , the solution is readily available from Information theory.

```
>>> from my_module import MDLP
>>> from sklearn.datasets import
load_iris
>>> iris = load_iris()
>>> X = iris.data
>>> y = iris.target
>>> mdlp = MDLP()
>>> conv_X = mdlp.fit_transform(X, y)
```

What's my micro-research ?  
Implementation =)