

Knowledge: Testable Predictions, 可证伪性

Confirmation Bias (unwillingness to try negative example)

- 所得 data 可能 biased
- 可得 example 满足很多 rule



hypothesis 可能包含于 rule.

基于个人的 hypothesis, 尝试与 hypo 相违背的 example, 才有可能找到 rule.

Types of Data.

Records

tuples/vector o.g. (name, age) (Bob, 4)

Graphs

nodes, edges (directed/undirected)

Adjacency list/matrix

Image

由 pixels 组成 (red, green, blue)

Text

Strings

Unsupervised Learning

将 data 分类 (clustering)

Supervised Learning

Data Processing

↓ Can / can't be used? missing? inconsistent?

Explore Data → Extract Features → Create Model

Distance

compare data points 的方法

dissimilarity function

$d(x, y) = \begin{cases} \text{值大, 不相似} \\ \text{小, 相似} \end{cases}$

Minkowski Distance

$x, y \in \mathbb{R}^d, p \geq 1$

$$L_p(x, y) = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

当 $p=2$ 时为 Euclidean distance.

Cosine Similarity

算 2 个 vectors 的 $\cos(\theta)$

因 dissimilarity function 大时不相似, function 将为 $k - \text{sl}(x, y)$