# Webmining 2

Bayesian Methods.

Other Classification Approaches.

Accessing Model Performance

Introduction of Clustering

K- Means Clustering

## BAYESIAN Methods

$x =$ variables

$y =$ target class

$P(y|x) \Rightarrow P(y=C_1|x)$

$\qquad P(y=C_2|x)$

$\qquad P(y=C_3|x)$

Baye's theorem (贝叶斯定理)

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)}$$

$\arg\max\limits_{y \in \{C_1, C_2, C_3\}} P(y|x)$

$P(x|y) P(y) = P(y) P(x_1|y) P(x_2|y, x_1) \cdots P(x_n|y, x_1, x_2, \cdots x_{n-1})$

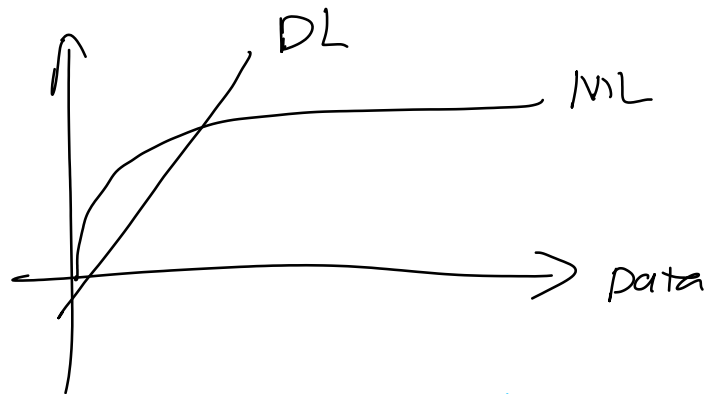## Naïve BAYES CLASSIFIERS

先验概率:

$P(y=Yes) = \frac{9}{14}$    $P(y=N) = \frac{5}{14}$

(conditional independence assumption)

(Naïve Bayes)

## K- nearest neighbour.

ML: Training & Testing

decision tree

neural network



## MODEL EVALUATION

## VALIDATION SET Approach.

验证集. 测试集 (50% + 50%)

## K-FOLD CROSS VALIDATION

↳ K个 测试集 → 验证

supervised learning

$Accuracy = \frac{1}{n} \sum\limits_{i=1}^{n} I(y_i = \hat{y_i})$

$I(y_i = \hat{y_i}) = \begin{cases} 1 & if (y_i = \hat{y_i}) \\ 0 & otherwise \end{cases}$

## RESAMPLING METHODS

$S_{train} \longrightarrow S_{test}$

## LEAVE-ONE-OUT CROSS- VALIDATION

使用一个测试集, 其余有 n-1 个测试集

## K- MEANS CLUSTERING

minimize $\left\{ \sum\limits_{k=1}^{k} W(C_k) \right\}$

$C_1 \cdots C_k$

$W(C_k) = \frac{1}{|C_k|} \sum\limits_{i,i' \in C_k} \sum\limits_{j=1}^{P} (x_{ij} - x_{i'j})^2$