

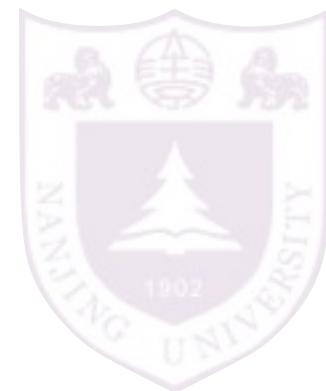
Lecture 6:

Bayesian Methods and Lazy Methods

http://cs.nju.edu.cn/yuy/course_dm12.ashx



Bayes rule



classification using posterior probability

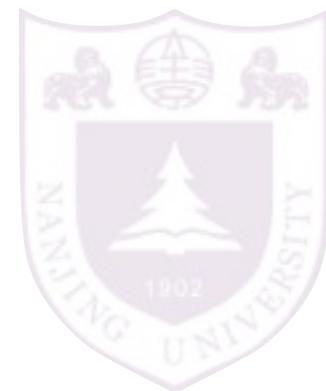
for binary classification

$$f(x) = \begin{cases} +1, & P(y = +1 | \mathbf{x}) > P(y = -1 | \mathbf{x}) \\ -1, & P(y = +1 | \mathbf{x}) < P(y = -1 | \mathbf{x}) \\ \text{random,} & \textit{otherwise} \end{cases}$$

in general

$$f(x) = \arg \max_y P(y | \mathbf{x})$$

Bayes rule



classification using posterior probability

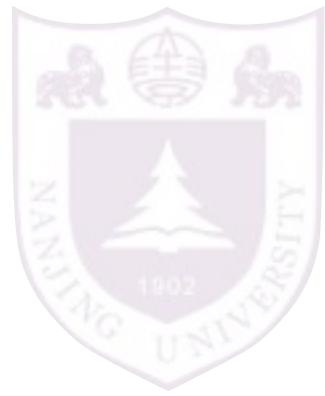
for binary classification

$$f(\mathbf{x}) = \begin{cases} +1, & P(y = +1 | \mathbf{x}) > P(y = -1 | \mathbf{x}) \\ -1, & P(y = +1 | \mathbf{x}) < P(y = -1 | \mathbf{x}) \\ \text{random,} & \textit{otherwise} \end{cases}$$

in general

$$\begin{aligned} f(\mathbf{x}) &= \arg \max_y P(y | \mathbf{x}) \\ &= \arg \max_y P(\mathbf{x} | y)P(y)/P(\mathbf{x}) \\ &= \arg \max_y P(\mathbf{x} | y)P(y) \end{aligned}$$

how the probabilities be estimated



Naive Bayes

$$f(x) = \arg \max_y P(\mathbf{x} | y)P(y)$$

estimation the a priori by frequency:

$$P(y) \leftarrow \tilde{P}(y) = \frac{1}{m} \sum_i I(y_i = y)$$

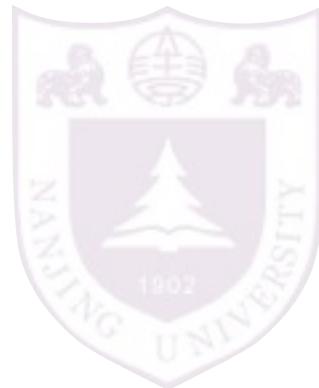
assume features are conditional independence given the class (**naive assumption**):

$$\begin{aligned} P(\mathbf{x} | y) &= P(x_1, x_2, \dots, x_n | y) \\ &= P(x_1 | y) \cdot P(x_2 | y) \cdot \dots \cdot P(x_n | y) \end{aligned}$$

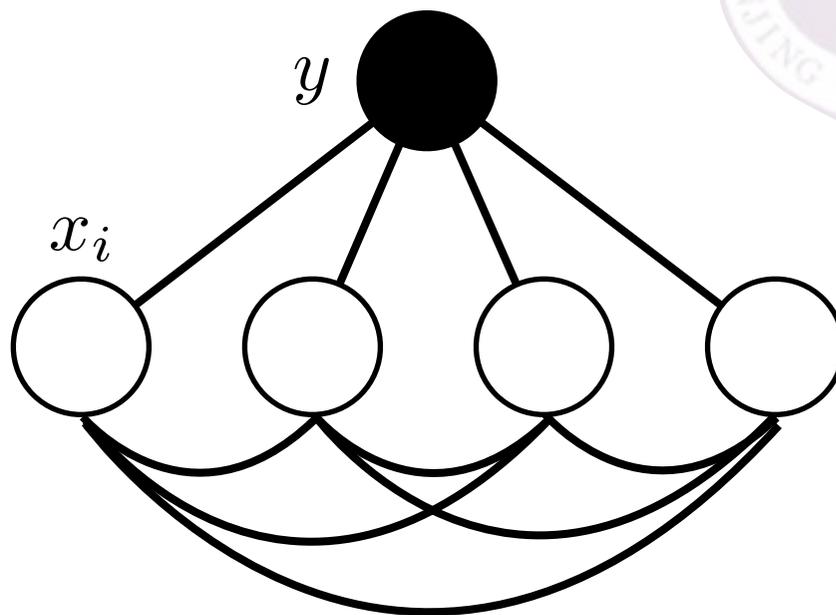
decision function:

$$f(x) = \arg \max_y \tilde{P}(y) \prod_i \tilde{P}(x_i | y)$$

Naive Bayes

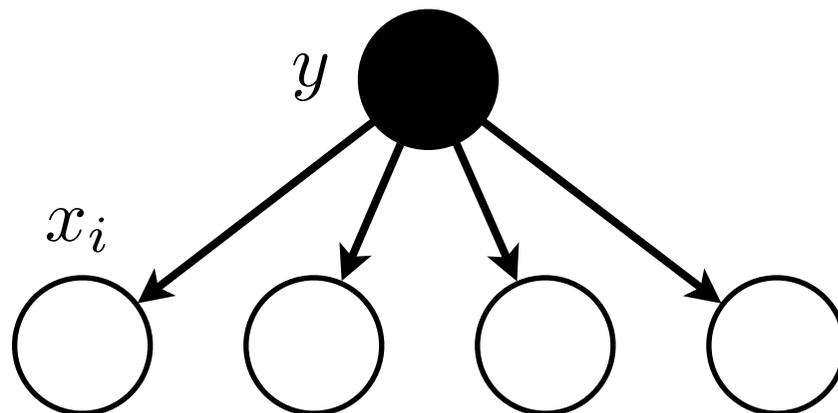


graphic representation
no assumption:

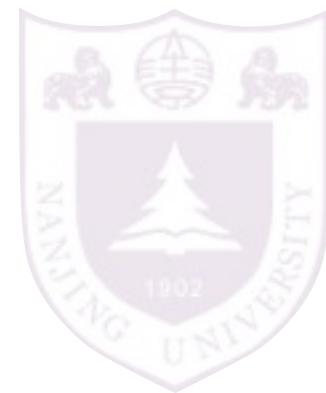


naive Bayes assumption:

$$P(\mathbf{x} | y) = \prod_i P(x_i | y)$$



Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

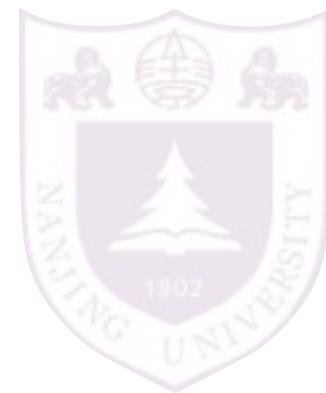
$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

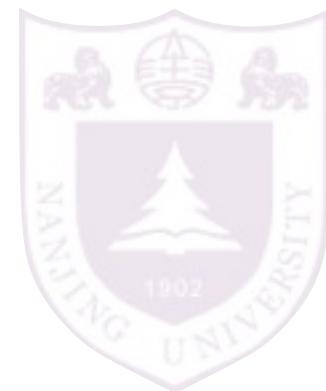
$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$

Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

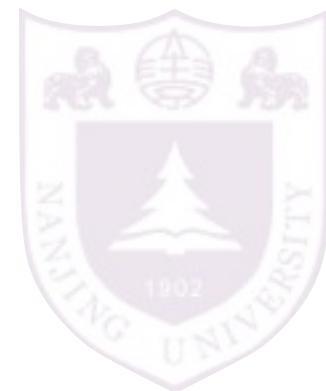
...

$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$

$$P(\text{color} = 3 \mid y = \text{yes})P(\text{weight} = 3 \mid y = \text{yes})P(y = \text{yes}) = 0.5 \times 0.5 \times 0.4 = 0.1$$

$$P(\text{color} = 3 \mid y = \text{no})P(\text{weight} = 3 \mid y = \text{no})P(y = \text{no}) = 0.33 \times 0.33 \times 0.6 = 0.06$$

Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

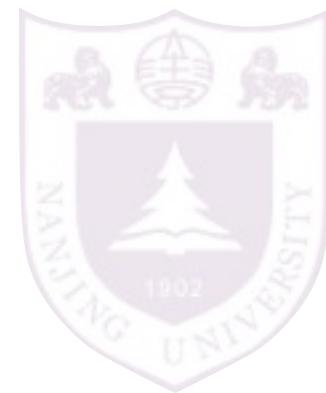
$$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$$

$$P(\text{color} = 3 \mid y = \text{yes})P(\text{weight} = 3 \mid y = \text{yes})P(y = \text{yes}) = 0.5 \times 0.5 \times 0.4 = 0.1$$

$$P(\text{color} = 3 \mid y = \text{no})P(\text{weight} = 3 \mid y = \text{no})P(y = \text{no}) = 0.33 \times 0.33 \times 0.6 = 0.06$$

$$f(y \mid \text{color} = 0, \text{weight} = 1) \rightarrow$$

Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

$$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$$

$$P(\text{color} = 3 \mid y = \text{yes})P(\text{weight} = 3 \mid y = \text{yes})P(y = \text{yes}) = 0.5 \times 0.5 \times 0.4 = 0.1$$

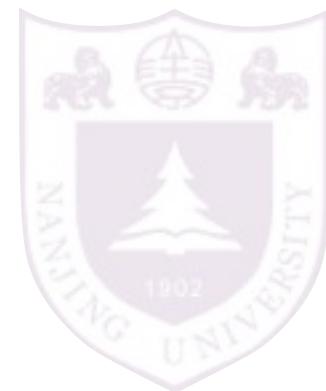
$$P(\text{color} = 3 \mid y = \text{no})P(\text{weight} = 3 \mid y = \text{no})P(y = \text{no}) = 0.33 \times 0.33 \times 0.6 = 0.06$$

$$f(y \mid \text{color} = 0, \text{weight} = 1) \rightarrow$$

$$P(\text{color} = 0 \mid y = \text{yes})P(\text{weight} = 1 \mid y = \text{yes})P(y = \text{yes}) = 0$$

$$P(\text{color} = 0 \mid y = \text{no})P(\text{weight} = 1 \mid y = \text{no})P(y = \text{no}) = 0$$

Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

+

color	sweet?
0	yes
1	yes
2	yes
3	yes

smoothed (Laplacian correction) probabilities:

$$P(\text{color} = 0 \mid y = \text{yes}) = (0 + 1) / (2 + 4)$$

$$P(y = \text{yes}) = (2 + 1) / (5 + 2)$$

for counting frequency,
assume every event
has happened once.

$$f(y \mid \text{color} = 0, \text{weight} = 1) \rightarrow$$

$$P(\text{color} = 0 \mid y = \text{yes})P(\text{weight} = 1 \mid y = \text{yes})P(y = \text{yes}) = \frac{1}{6} \times \frac{1}{7} \times \frac{3}{7} = 0.01$$

$$P(\text{color} = 0 \mid y = \text{no})P(\text{weight} = 1 \mid y = \text{no})P(y = \text{no}) = \frac{2}{7} \times \frac{1}{8} \times \frac{4}{7} = 0.02$$

Naive Bayes



advantages:

very fast:

scan the data once, just count: $O(mn)$

store class-conditional probabilities: $O(n)$

test an instance: $O(cn)$ (c the number of classes)

good accuracy in many cases

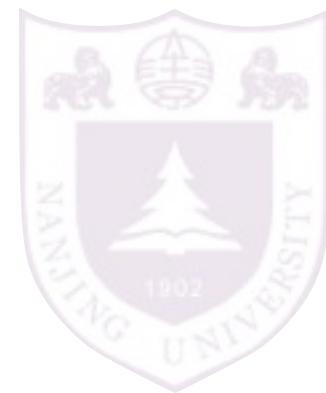
parameter free

output a probability

naturally handle multi-class

disadvantages:

Naive Bayes



advantages:

very fast:

scan the data once, just count: $O(mn)$

store class-conditional probabilities: $O(n)$

test an instance: $O(cn)$ (c the number of classes)

good accuracy in many cases

parameter free

output a probability

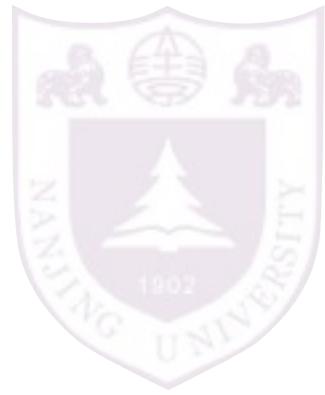
naturally handle multi-class

disadvantages:

the strong assumption may harm the accuracy

does not handle numerical features naturally

Relaxation of naive Bayes assumption



assume features are conditional
independence given the class

if the assumption holds, naive Bayes
classifier will have excellence performance

if the assumption does not hold ...

Relaxation of naive Bayes assumption



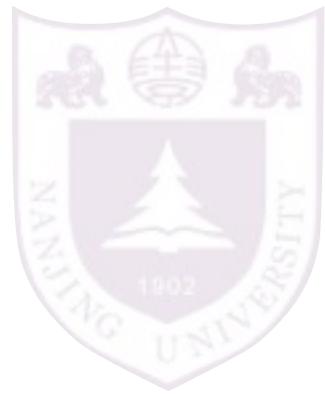
assume features are conditional independence given the class

if the assumption holds, naive Bayes classifier will have excellence performance

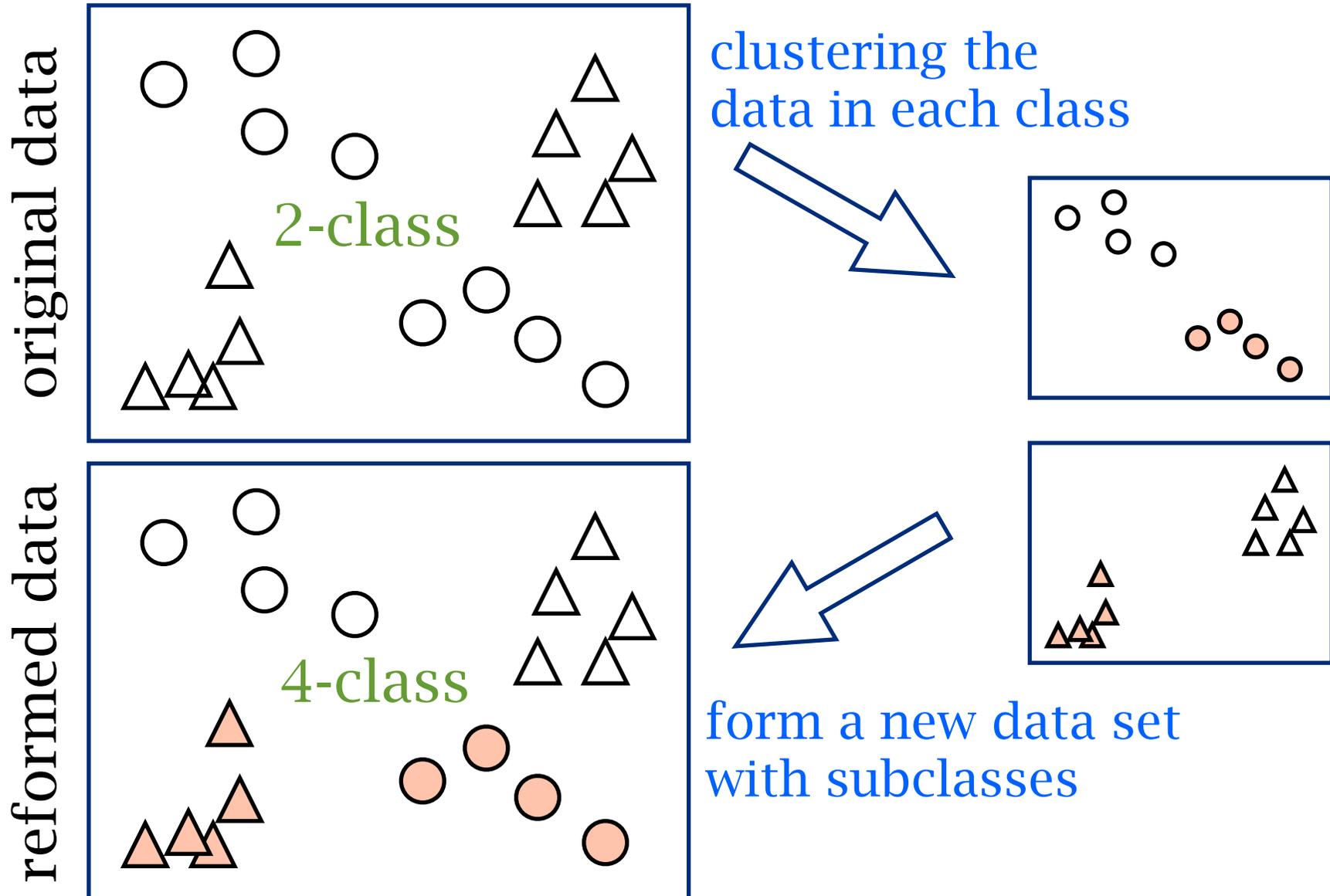
if the assumption does not hold ...

- ▶ Naive Bayes classifier may also have good performance
- ▶ Reform the data to satisfy the assumption
- ▶ Invent algorithms to relax the assumption

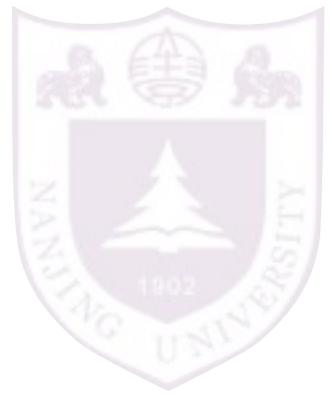
Reform the data



clustering to generate data with subclasses

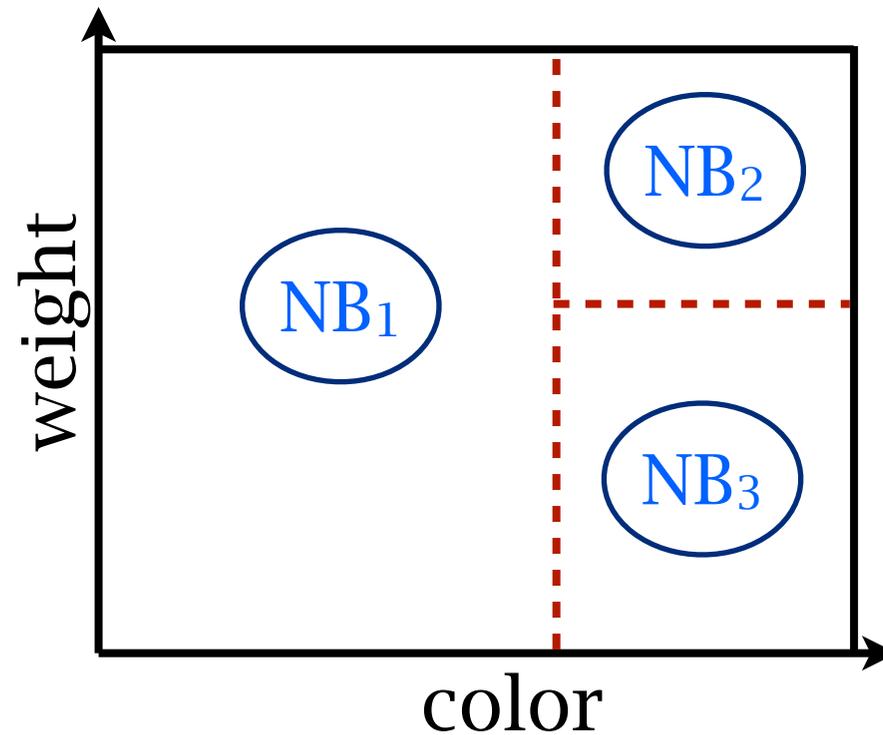


Semi-naive Bayes classifiers

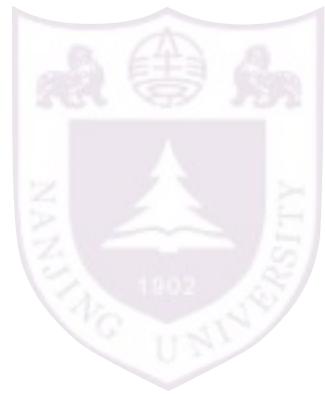


TreeNB

train an NB classifier in each leaf node of a rough decision tree

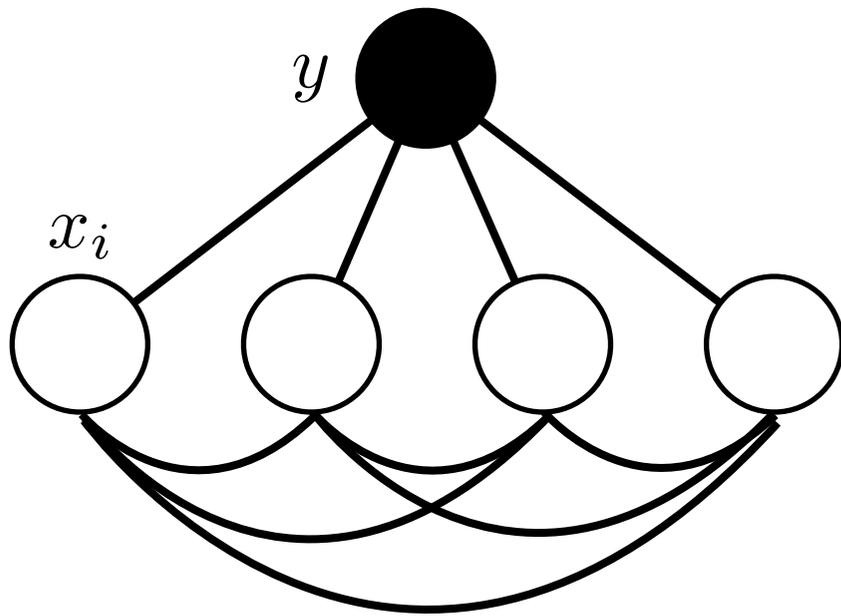


Semi-naive Bayes classifiers

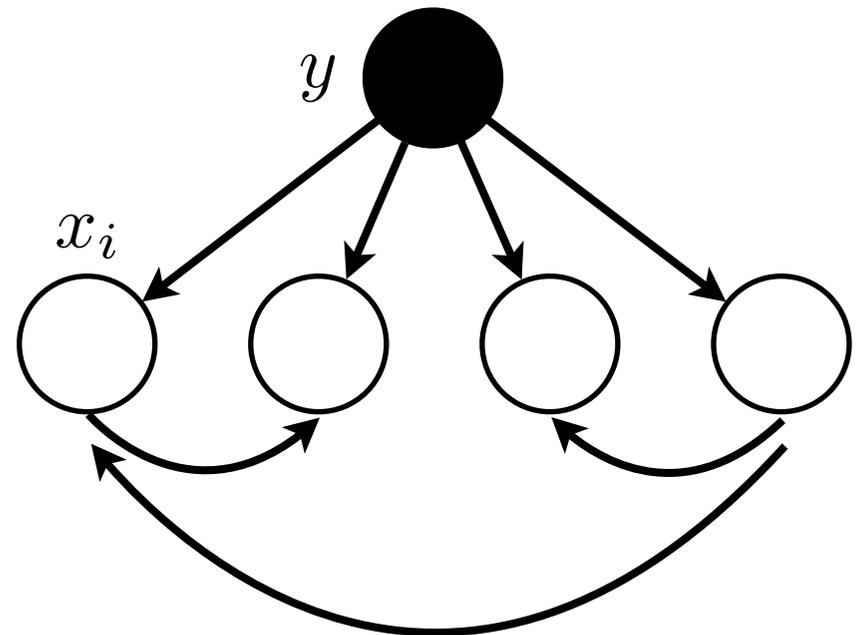


TAN (Tree Augmented NB)

extends NB by allowing every feature to have one more parent feature other than the class, which forms a tree structure



fully connected

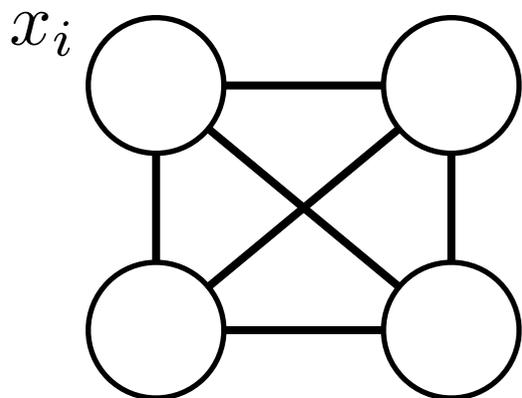


TAN

Semi-naive Bayes classifiers



TAN (Tree Augmented NB)

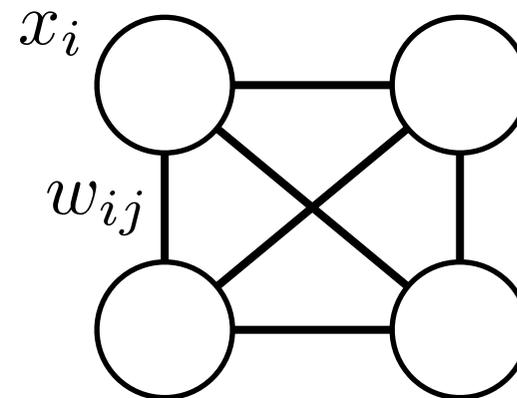


fully connected graph
among features

mutual information
for every node pair



$$\begin{aligned} I(X_i, X_j | Y) &= \mathbb{E}_Y [I(X_i; X_j) | Y] \\ &= \mathbb{E}_Y [H(X_i) - H(X_i | X_j) | Y] \\ &= \sum_{x_i, x_j, y} P(x_i, x_j, y) \log \frac{P(x_i, x_j | y)}{P(x_i | y)P(x_j | y)} \end{aligned}$$

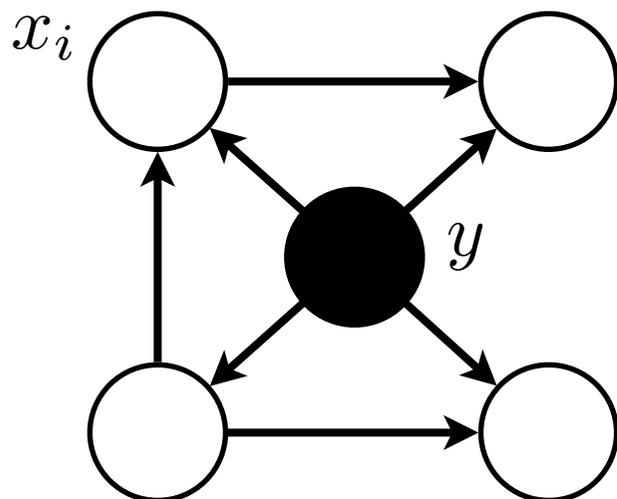


weights assigned

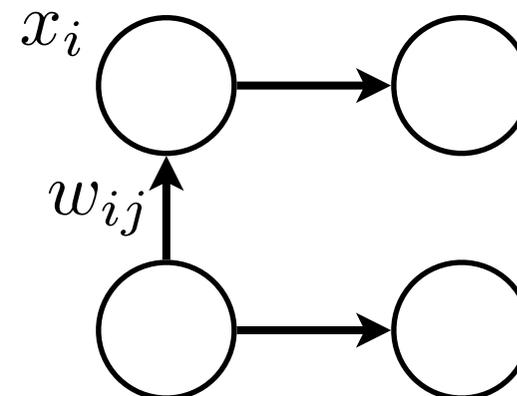
maximum
weighted
spanning tree



and
choose
a root



connect to the
class node





Semi-naive Bayes classifiers

AODE (average one-dependent estimators)

expand a posterior probability
with one-dependent estimators

$$\begin{aligned} P(\mathbf{x} | y) &= P(x_2, \dots, x_n | x_1, y) P(x_1 | y) \\ &= P(x_1 | y) \prod_i P(x_i | x_1, y) \end{aligned}$$

compare with NB:

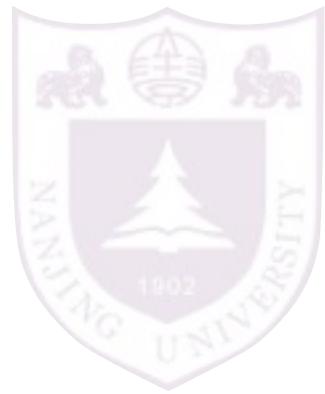
$$P(\mathbf{x} | y) = \prod_i P(x_i | y)$$

- ▶ the conditional independency is less important
- ▶ harder to estimate (fewer data)

AODE: average ODEs

$$f(x) = \arg \max_y \sum_i I(\text{count}(x_i \geq m)) \cdot \tilde{P}(y) \cdot \tilde{P}(x_i | y) \cdot \prod_j \tilde{P}(x_j | x_i, y)$$

Handling numerical features



Discretization

recall what we have talked about in Lecture 2

Estimate probability density ($P(X) \rightarrow p(x)$)

Gaussian model:

$$p(x) = \frac{1}{\sqrt{2\pi\delta^2}} e^{-\frac{(x-\mu)^2}{2\delta^2}}$$

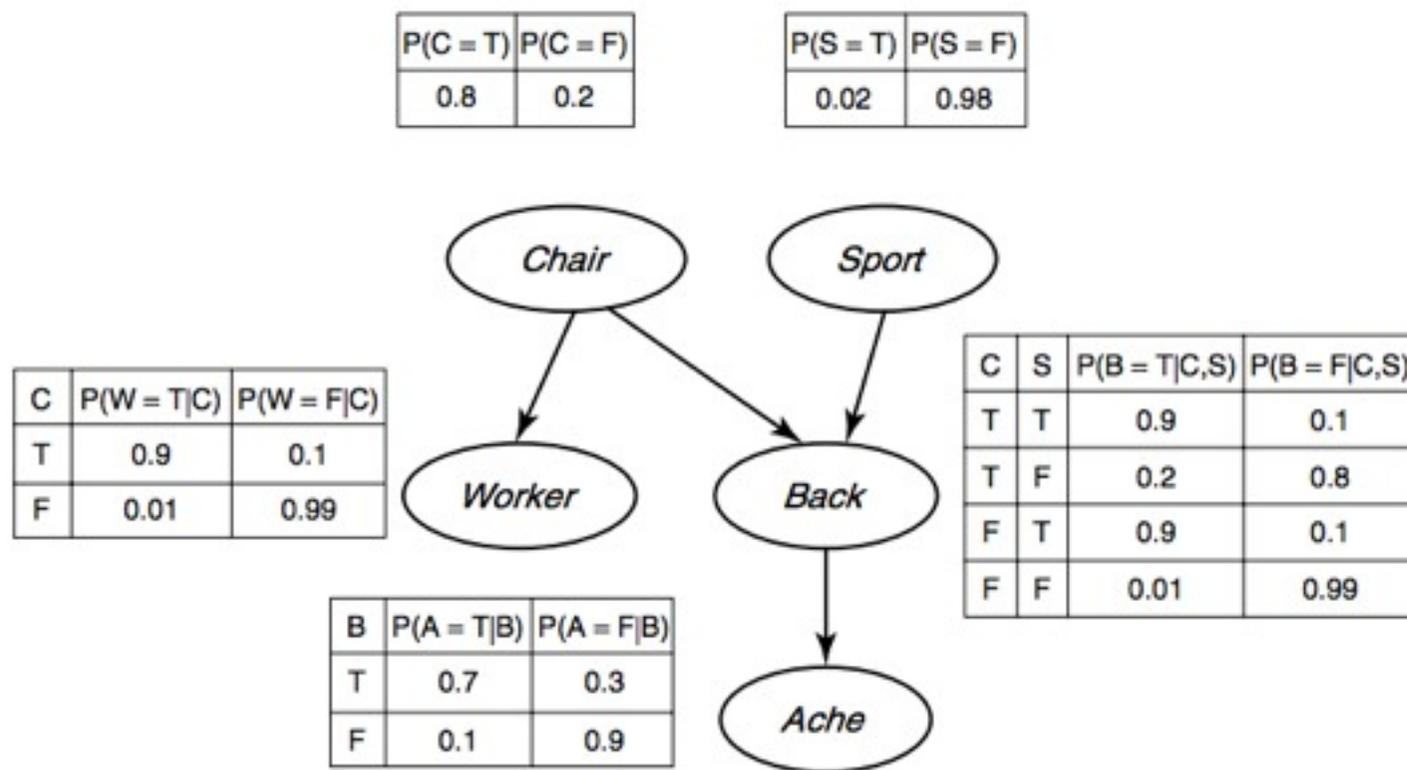
$$p(x_1, \dots, x_n) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

training: calculate mean and covariance
test: calculate density



Bayesian networks

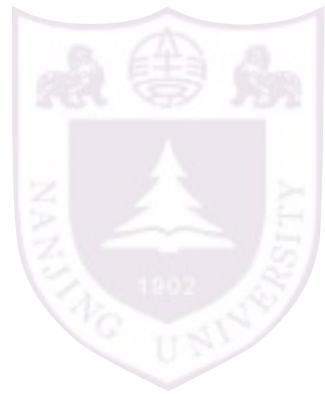
inference in a graphic model representation
a model simplified by conditional independence
a clear description of how things are going



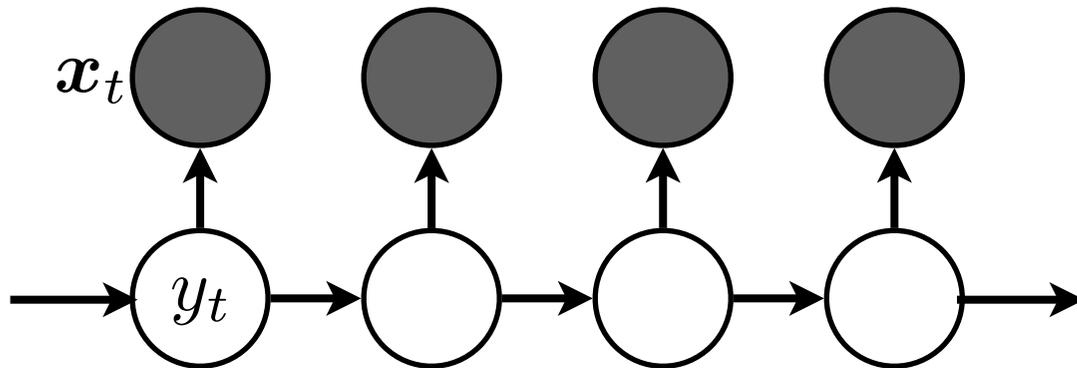
Judea Pearl
Turing Award 2011

“for fundamental contributions to artificial intelligence through the development of a calculus for probabilistic and causal reasoning”

Bayesian networks/Graphic models



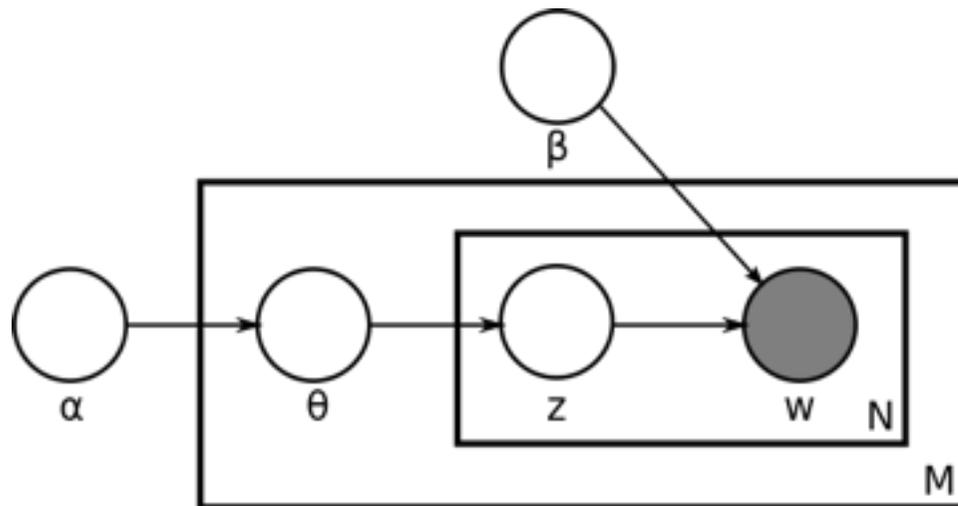
Hidden Markov Model (HMM)



voice

words

Topic Model: Latent Dirichlet Allocation



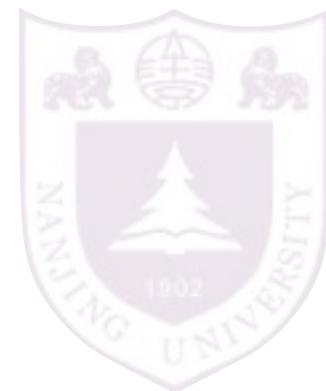
α, β parameters

θ document

z topic

w words

Lazy methods



similarity function $S(\mathbf{x}_1, \mathbf{x}_2)$

training data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$

no model is built until meet a test instance \mathbf{x}

to predict the label of \mathbf{x}

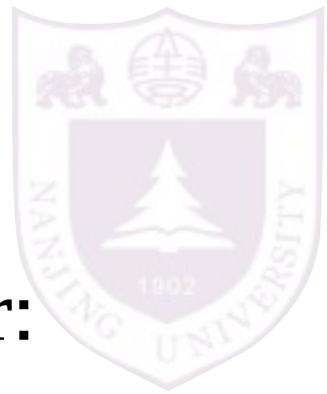
objects that look similar are indeed similar

find similar training instances S

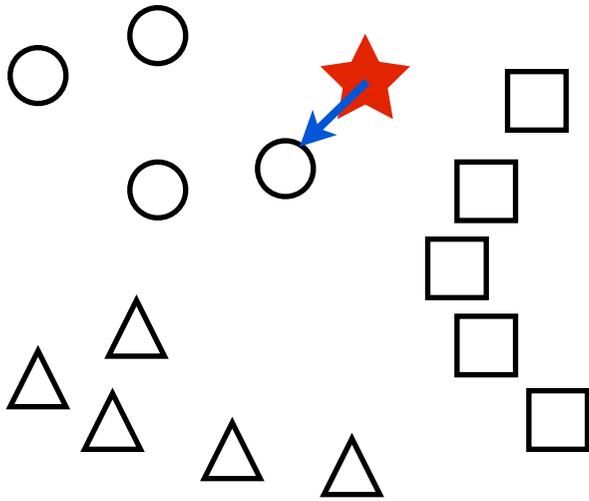
build a model on S

use the model to predict the label of \mathbf{x}

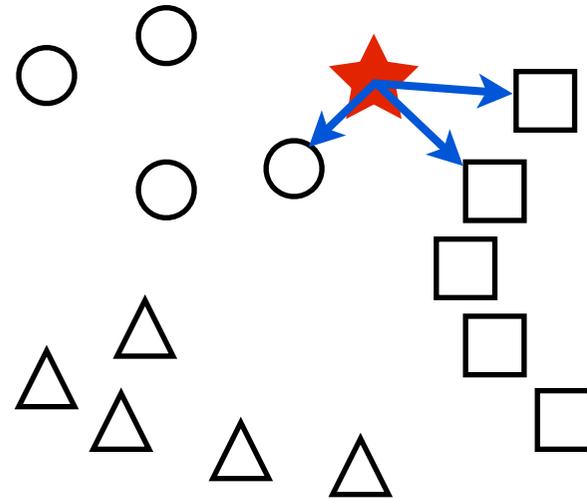
Nearest neighbor classifier



1-nearest neighbor:



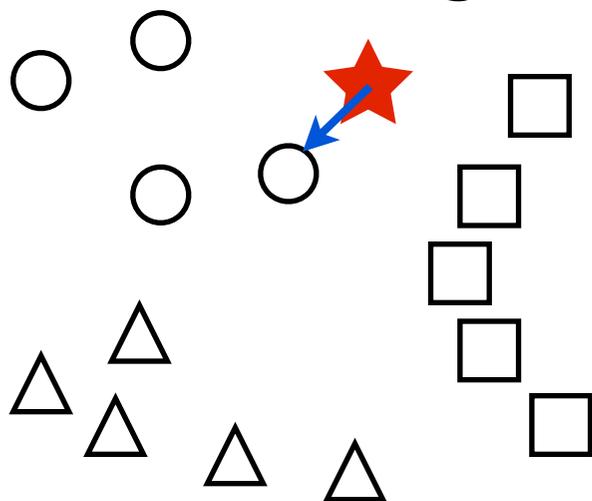
k -nearest neighbor:



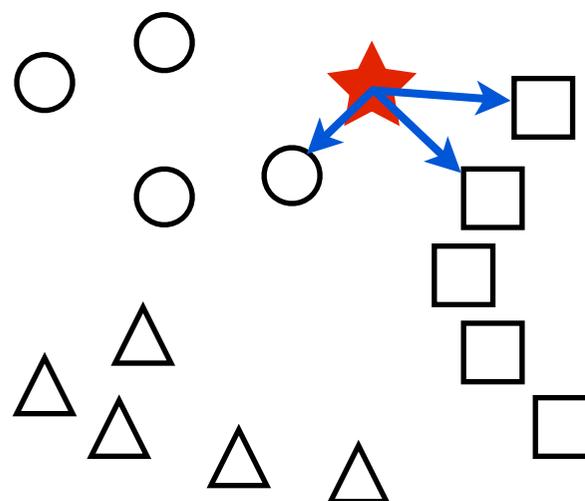
Nearest neighbor classifier



1-nearest neighbor:



k -nearest neighbor:

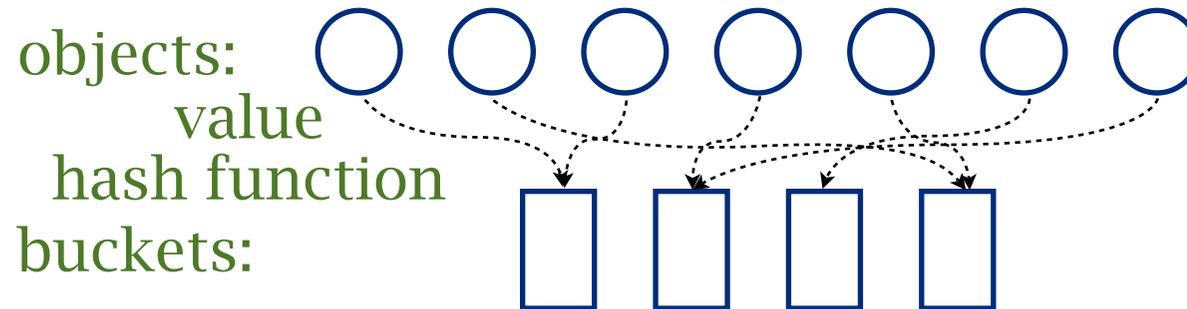


- ▶ asymptotically less than 2 times of the optimal Bayes error
- ▶ naturally handle multi-class
- ▶ no training time
- ▶ nonlinear decision boundary
- ▶ slow testing speed for a large training data set
- ▶ have to store the training data
- ▶ sensitive to similarity function

Locality sensitive hashing



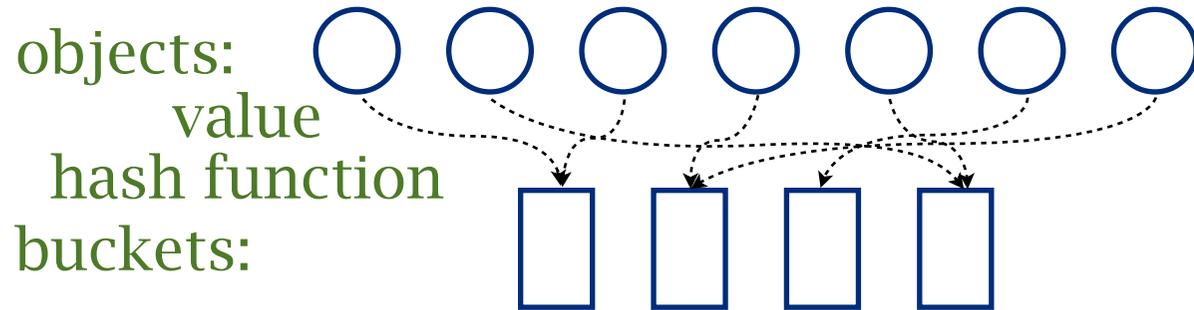
hashing



Locality sensitive hashing



hashing



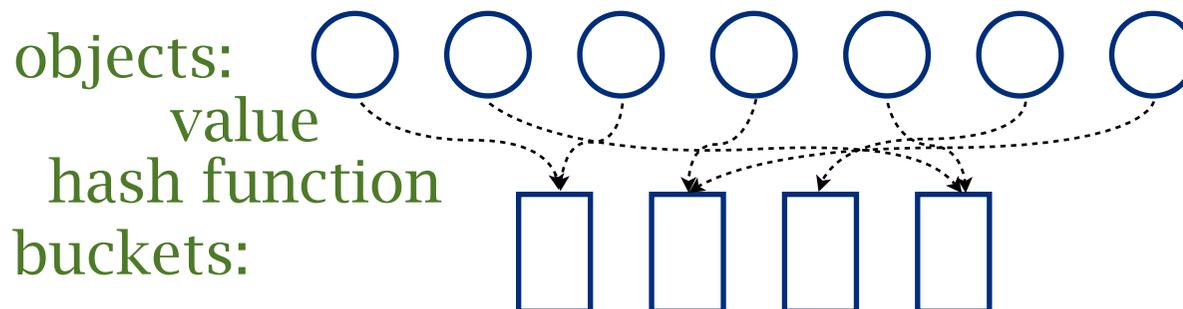
locality sensitive hashing:

similar objects in the same bucket

Locality sensitive hashing



hashing



locality sensitive hashing:

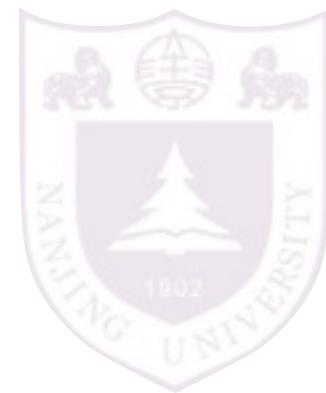
similar objects in the same bucket

A LSH function family $\mathcal{H}(c, r, P_1, P_2)$ has the following properties for any $\mathbf{x}_1, \mathbf{x}_2 \in S$

if $\|\mathbf{x}_1 - \mathbf{x}_2\| \leq r$, then $P_{h \in \mathcal{H}}(h(\mathbf{x}_1) = h(\mathbf{x}_2)) \geq P_1$
similar objects should be hashed in the same bucket with high probability

if $\|\mathbf{x}_1 - \mathbf{x}_2\| \geq cr$, then $P_{h \in \mathcal{H}}(h(\mathbf{x}_1) = h(\mathbf{x}_2)) \leq P_2$
dissimilar objects should be hashed in the same bucket with low probability

Locality sensitive hashing



Binary vectors in Hamming space

objects: (1100101101)

Hamming distance: count the number of positions with different elements

$$\|110101001, 110001100\|_H = 3$$

Locality sensitive hashing



Binary vectors in Hamming space

objects: (1100101101)

Hamming distance: count the number of positions with different elements

$$\|110101001, 110001100\|_H = 3$$

LSH functions: $\mathcal{H} = \{h_1, \dots, h_n\}$ where $h_i(\mathbf{x}) = x_i$

	h_2	h_5	h_9
110101001	1	0	1
110010100	1	1	0
000110110	0	1	0
111001001	1	0	1
000011101	0	1	1



Locality sensitive hashing

Binary vectors in Hamming space

objects: (1100101101)

Hamming distance: count the number of positions with different elements

$$\|110101001, 110001100\|_H = 3$$

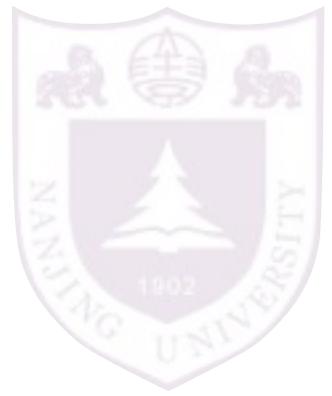
LSH functions: $\mathcal{H} = \{h_1, \dots, h_n\}$ where $h_i(\mathbf{x}) = x_i$

	h_2	h_5	h_9
110101001	1	0	1
110010100	1	1	0
000110110	0	1	0
111001001	1	0	1
000011101	0	1	1

$$P(h_i(\mathbf{x}_1) = h_i(\mathbf{x}_2)) = 1 - \frac{\|\mathbf{x}_1 - \mathbf{x}_2\|}{d}$$



frequency in the same bucket for a sample of hashing functions

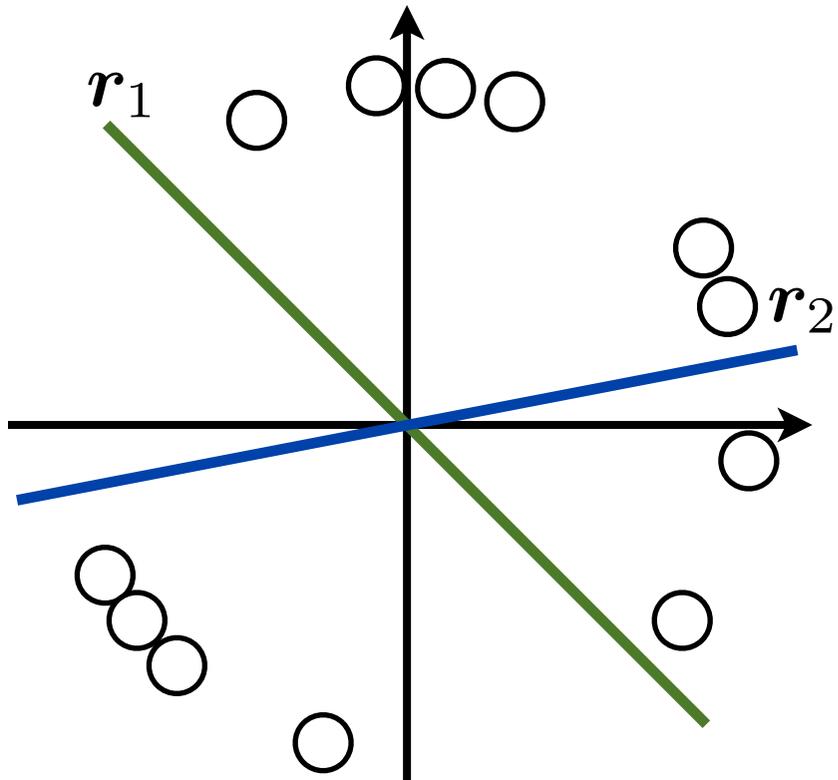


Locality sensitive hashing

Real vectors with angle similarity

$$\theta(\mathbf{x}_1, \mathbf{x}_2) = \arccos \frac{\mathbf{x}_1^\top \mathbf{x}_2}{\|\mathbf{x}_1\| \|\mathbf{x}_2\|}$$

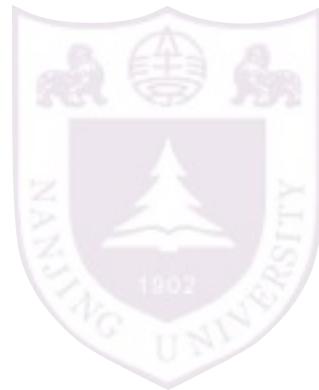
LSH functions: $\mathcal{H} = \{h_r\} (r \in \mathbb{B}^n)$ where $h_r(\mathbf{x}) = \text{sign}(r^\top \mathbf{x})$



$$P(h_r(\mathbf{x}_1) = h_r(\mathbf{x}_2)) = 1 - \frac{\theta(\mathbf{x}_1, \mathbf{x}_2)}{\pi}$$

↑
frequency in the same bucket for
a sample of hashing functions

习题



朴素贝叶斯假设是指数据的属性之间相互独立？

朴素贝叶斯假设不满足时，朴素贝叶斯的性能一定不好？

k近邻分类算法是否需要训练预测模型？