

Artificial Intelligence — — Summary



Yanghui Rao

Assistant Prof., Ph.D

School of Data and Computer Science,

Sun Yat-sen University

raoyangh@mail.sysu.edu.cn

Probability

- **Product rule:**

$$P(A, B) = P(A)P(B | A) = P(B)P(A | B)$$

$$P(A, B_1, B_2, B_3) = P(A)P(B_1 | A)P(B_2 | A, B_1)P(B_3 | A, B_1, B_2)$$

- **Sum rule:** $P(A) = P(A, B) + P(A, B^c)$

$$P(A) = \sum_{i=1}^n P(A, B_i)$$

$$= \sum_{i=1}^n P(A | B_i)P(B_i)$$

Lec 2 Foundation of Mathematics

- There are two random variables X and Y . Which of the following is always true?
 - A. $\sum_X P(X|Y) = 1$
 - B. $\sum_Y P(X|Y) = 1$
 - C. All of the above
 - D. None of the above
- Is the statement True or False? Entropy of a discrete random variable is always non-negative.

Lec 2 Foundation of Mathematics

- There are two random variables X and Y. Which of the following is always true? (Answer: A)
 - A. $\sum_X P(X|Y) = 1$
 - B. $\sum_Y P(X|Y) = 1$
 - C. All of the above
 - D. None of the above
- Is the statement True or False? Entropy of a discrete random variable is always non-negative. (Answer: True)

Truth Tables

- Truth tables are used to define logical connectives and to determine when a complex sentence is true given the values of the symbols in it
- Note that \Rightarrow is a logical connective, so $P \Rightarrow Q$ is a logical sentence and has a truth value, i.e., is either true or false

Truth tables for the five logical connectives

P	Q	$\neg P$	$P \wedge Q$	$P \vee Q$	$P \Rightarrow Q$	$P \Leftrightarrow Q$
False	False	True	False	False	True	True
False	True	True	False	True	True	False
True	False	False	False	True	False	False
True	True	False	True	True	True	True

Quantifier Scope

- If a quantifier Q is followed by $($, then the scope of Q is to the matched $)$
 - $\forall x (F(x) \Leftrightarrow F(h))$
- If a quantifier Q is not followed by $($ or another quantifier, then the scope of Q is to the first connective
 - $\forall x F(x) \Leftrightarrow F(h)$
- If a quantifier $Q1$ is followed by another quantifier $Q2$, then the scope of $Q1$ is to the scope of $Q2$
 - $\forall x \exists y R(x, y)$
- F : ... can fly
- h : human being

False $\forall x (F(x) \Leftrightarrow F(h))$ \nleftrightarrow **True** $\forall x F(x) \Leftrightarrow F(h)$

Lec 3 Logic

- Fill in the following truth table:

P	Q	$(P \Rightarrow Q) \wedge (Q \Rightarrow P)$	$(\neg P \vee Q) \Leftrightarrow (P \Rightarrow Q)$
True	True		
True	False		
False	True		
False	False		

- If we represent “... is hot” by $H(\dots)$, and represent “fire” by f , what are the values of “ $H(f) \Rightarrow \forall x H(x)$ ” and “ $\exists x (H(f) \Leftrightarrow H(x))$ ”?

Lec 3 Logic

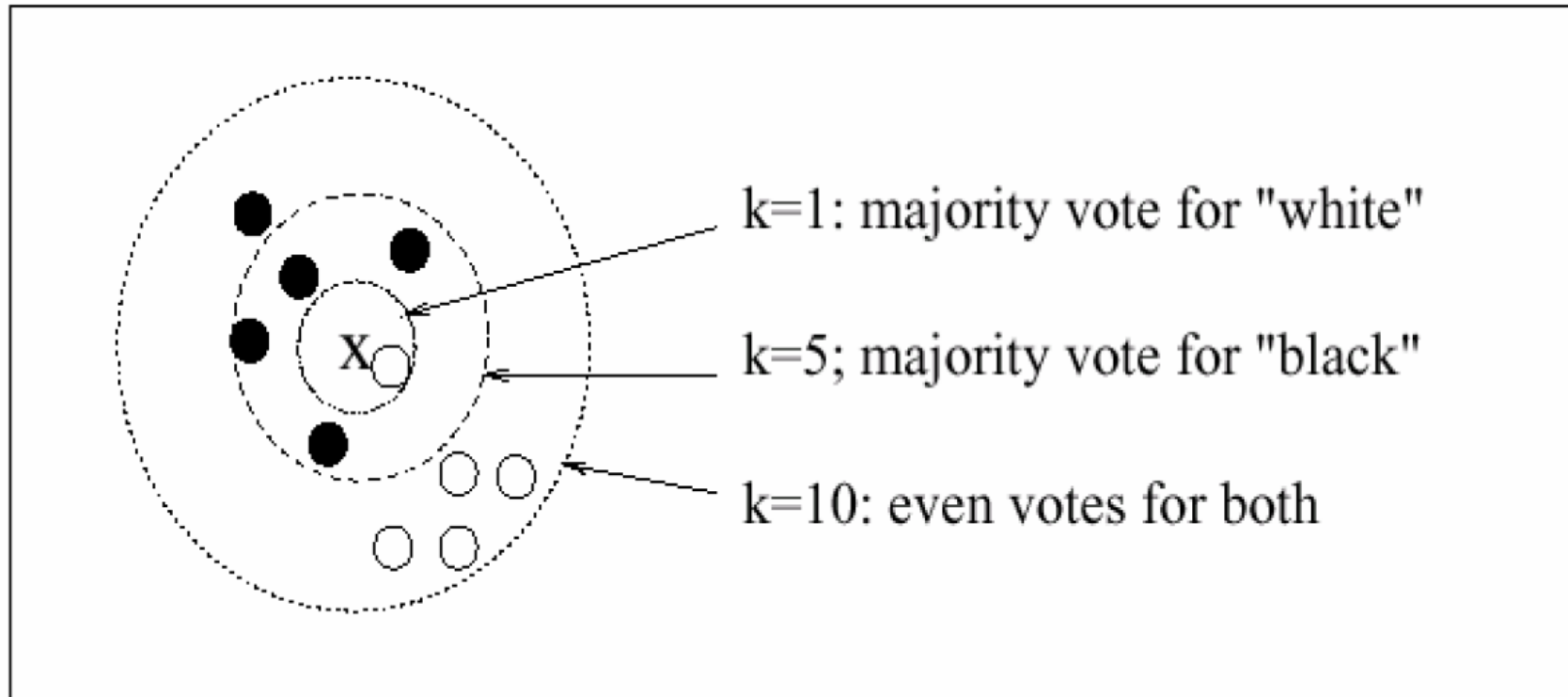
- Fill in the following truth table:

P	Q	$(P \Rightarrow Q) \wedge (Q \Rightarrow P)$	$(\neg P \vee Q) \Leftrightarrow (P \Rightarrow Q)$
True	True	True	True
True	False	False	True
False	True	False	True
False	False	True	True

- If we represent “... is hot” by $H(\dots)$, and represent “fire” by f , what are the values of “ $H(f) \Rightarrow \forall x H(x)$ ” and “ $\exists x (H(f) \Leftrightarrow H(x))$ ”? (Answer: False, True)

k -Nearest Neighbor

k -NN using a majority voting scheme



Naïve Bayesian Classifier

- This can be derived from Bayes' theorem

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i)P(C_i)}{P(\mathbf{X})}$$

- Since $P(\mathbf{X})$ is constant for all classes, only

$$P(C_i | \mathbf{X}) \propto P(\mathbf{X} | C_i)P(C_i)$$

needs to be maximized

- $P(C_i)$ can be obtained from training set s_i/s

Derivation

- **Assumption:** attributes are conditionally independent (i.e., no dependence relation between attributes):
$$P(\mathbf{X} | C_i) = \prod_{k=1}^n P(x_k | C_i)$$
- This greatly reduces the computation cost: Only counts the class distribution
- If A_k is categorical, $P(x_k | C_i) = s_{ik}/s_i$, count the distribution
- If A_k is continuous-valued, $P(x_k | C_i)$ can be computed based on Gaussian distribution

Lec 4 kNN and NB

- What is the meaning of “k” for the k-Nearest Neighbor (i.e., k-NN) and the k-Means clustering algorithm?
- If using k-NN for classification, what is the predicted class label when “ $x = 5$ ”? Is there any difference if based on City Block, Euclidean, or Supremum distance?

Note: Given a testing sample x , if there are multiple training samples' distances are the nearest, k-NN classifier will use the mode (众数) of the class labels of all nearest training samples as the predicted class label of x

X	Y
2	-
3	-
3	-
3	+
5	?

X	Y
2	+
3	-
3	-
3	+
5	?

Lec 4 kNN and NB

- What is the meaning of “k” for the k-Nearest Neighbor (i.e., k-NN) and the k-Means clustering algorithm?
 - Answer: a) The parameter “k” means the number of neighbors used to classify test examples for the k-NN. b) The parameter “k” specifies the number of clusters for the k-Means.
- Given a testing sample x , if there are multiple training samples' distances are the nearest, k-NN classifier will use the mode (众数) of the class labels of all nearest training samples as the predicted class label of x .
 - Answer: There is no difference if based on those distance measures. (1) Left table. The predicted class label is “-”; (2) Right table. The predicted class label is “-” for $k=1, 2, 3$ and “Unknown” for $k > 3$.

Information Gain (ID3)

- Class label: buy_computer="yes/no"
- 用字母 D 表示类标签，字母 A 表示每个属性
- $H(D)=0.940$ 14个训练样本中，9个买了电脑

$$H(D) = -\frac{9}{14} \log_2 \frac{9}{14} - \left(1 - \frac{9}{14}\right) \log_2 \left(1 - \frac{9}{14}\right)$$

- $H(D | A = "age") = 0.694$

$$H(D | A = "age") = \frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) \\ + \frac{4}{14} \times \left(-\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \right) + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right)$$

Information Gain (ID3)

- Class label: buy_computer="yes/no"
 - Compute the mutual information (互信息) between D and each attribute A
 - $H(D)=0.940$
 - $H(D|A="age")=0.694$
 - $g(D,A="age")=0.246$
 - $g(D,A="income")=0.029$
 - $g(D,A="student")=0.151$
 - $g(D,A="credit_rating")=0.048$
- “age”这个属性的条件熵最小（等价于信息增益最大），因而首先被选出作为根节点**
- | |
|--------------|
| $g(D, A)$ |
| $= H(D)$ |
| $- H(D A)$ |

Information Gain Ratio (C4.5)

- $\text{GainRatio}_A(D) = \text{Gain}_A(D) / \text{SplitInfo}_A(D)$

$$\text{SplitInfo}_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right)$$

- $\text{GainRatio}_{A=\text{"income"}}(D) = ?$

$$\text{SplitInfo}_{A=\text{"income"}}(D)$$

$$= -\frac{4}{14} \times \log_2 \left(\frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left(\frac{6}{14} \right) - \frac{4}{14} \times \log_2 \left(\frac{4}{14} \right)$$

$$= 0.926$$

- $\text{GainRatio}_{A=\text{"income"}}(D) = 0.029 / 0.926 = 0.031$

Gini Index (CART)

- D has 9 samples in `buys_computer` = “yes” and 5 in “no”

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

- The attribute *income* partitions D into 10 in D_1 : {medium, high} and 4 in D_2

$$gini_{income \in \{\text{medium, high}\}}(D) = \frac{10}{14} gini(D_1) + \frac{4}{14} gini(D_2)$$

$$= \frac{10}{14} \left(1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2 \right) + \frac{4}{14} \left(1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 \right)$$

$$= 0.450 = gini_{income \in \{\text{low}\}}(D)$$

Decision Tree

- But how can we compute the gini index, information gain of an attribute that is **continuous-valued**?
 - Given v values of A , then $v-1$ possible splits are evaluated. For example, the midpoint between the values a_i and a_{i+1} of A is $(a_i + a_{i+1}) / 2$

Incorporating model complexity

- In the case of a decision tree, let
 - L be the number of leaf nodes.
 - n_l be the l -th leaf node.
 - $m(n_l)$ be the number of training records classified by n_l .
 - $r(n_l)$ be the number of misclassified records by n_l .
 - $\zeta(n_l)$ be a penalty term associated with the node n_l .
- The resulting error e_c of the decision tree can be estimated as follows:

$$e_c = \frac{\sum_{l=1}^L (r(n_l) + \zeta(n_l))}{\sum_{l=1}^L m(n_l)}$$

Decision Tree

- We consider the training examples shown in the following table for a binary classification problem.
 - Calculate the respective changes in the Gini index value when a_1 and a_2 are used for partitioning the training set.
 - Calculate the respective changes in the classification (**training**) error when a_1 and a_2 are used for partitioning the training set.

a_1	a_2	a_3	Target Class
T	T	1	+
T	T	6	+
T	F	5	-
F	F	4	+
F	T	7	-
F	T	3	-
F	F	8	-
T	F	7	+
F	T	5	-

Decision Tree

- (1) The original Gini index is $1 - (\frac{4}{9})^2 - (\frac{5}{9})^2 = 0.494$

After splitting on a_1 , the Gini index becomes

$$\frac{4}{9}[1 - (\frac{3}{4})^2 - (\frac{1}{4})^2] + \frac{5}{9}[1 - (\frac{1}{5})^2 - (\frac{4}{5})^2] = 0.344$$

As a result, the change in Gini index is

$$\Delta G(a_1) = 0.494 - 0.344 = 0.15.$$

After splitting on a_2 , the Gini index becomes

$$\frac{5}{9}[1 - (\frac{2}{5})^2 - (\frac{3}{5})^2] + \frac{4}{9}[1 - (\frac{2}{4})^2 - (\frac{2}{4})^2] = 0.489$$

As a result,

$$\Delta G(a_2) = 0.494 - 0.489 = 0.005.$$

Decision Tree

- (2) The original classification error is $1 - \max(\frac{4}{9}, \frac{5}{9}) = \frac{4}{9}$

After splitting on a_1 , the classification error becomes

$$\frac{4}{9}[1 - \max(\frac{3}{4}, \frac{1}{4})] + \frac{5}{9}[1 - \max(\frac{1}{5}, \frac{4}{5})] = \frac{2}{9}$$

As a result, the change in classification error is

$$\triangle E(a_1) = 4/9 - 2/9 = 2/9.$$

After splitting on a_2 , the classification error becomes

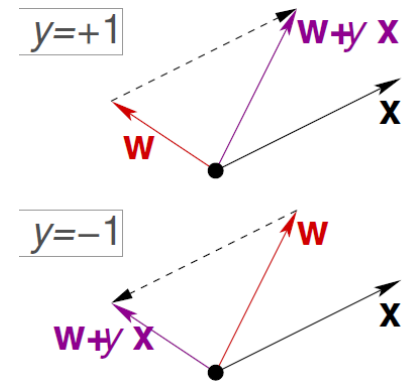
$$\frac{5}{9}[1 - \max(\frac{2}{5}, \frac{3}{5})] + \frac{4}{9}[1 - \max(\frac{2}{4}, \frac{2}{4})] = \frac{4}{9}$$

As a result,

$$\triangle E(a_2) = 4/9 - 4/9 = 0.$$

Perceptron Learning Algorithm

- Difficult: the set of $h(\mathbf{x})$ is of infinite size
- Idea: start from some initial weight vector $\mathbf{w}_{(0)}$, and “correct” its mistakes on D
- For $t = 0, 1, \dots$
 - find a mistake of $\mathbf{w}_{(t)}$ called $(\mathbf{x}_{n(t)}, y_{n(t)})$
 $\text{sign}(\mathbf{w}_{(t)}^T \mathbf{x}_{n(t)}) \neq y_{n(t)}$
 - (try to) correct the mistake by $\mathbf{w}_{(t+1)} \leftarrow \mathbf{w}_{(t)} + y_{n(t)} \mathbf{x}_{n(t)}$
 - until no more mistakes
- Return last \mathbf{W} (called \mathbf{W}_{PLA})



Perceptron Learning Algorithm

- Only if there exists an hyperplane that correctly classifies the data, the Perceptron procedure is guaranteed to converge; furthermore, the algorithm may give different results depending on the order in which the elements are processed, indeed several different solutions exist.

Perceptron Learning Algorithm

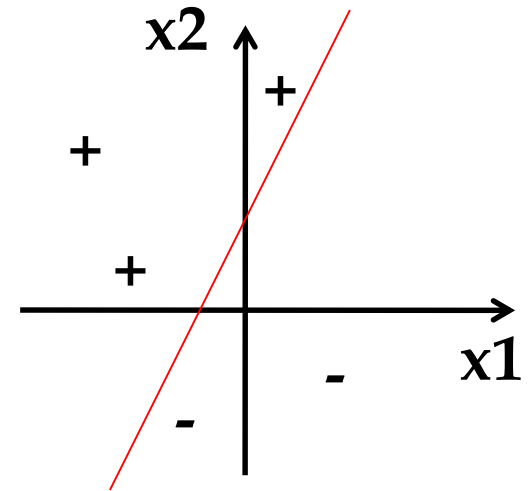
- What are the values of weights w_0 , w_1 , and w_2 for the perceptron whose decision surface is illustrated in the Figure? Assume the surface crosses the x_1 axis at -1, and the x_2 axis at 2.

- Answer:

$w_0 =$

$w_1 =$

$w_2 =$



Perceptron Learning Algorithm

- What are the values of weights w_0 , w_1 , and w_2 for the perceptron whose decision surface is illustrated in the Figure? Assume the surface crosses the x_1 axis at -1, and the x_2 axis at 2.

The surface crosses $(-1, 0)$ and $(0, 2)$

One surface: $-1 - x_1 + 0.5 \cdot x_2 = 0$

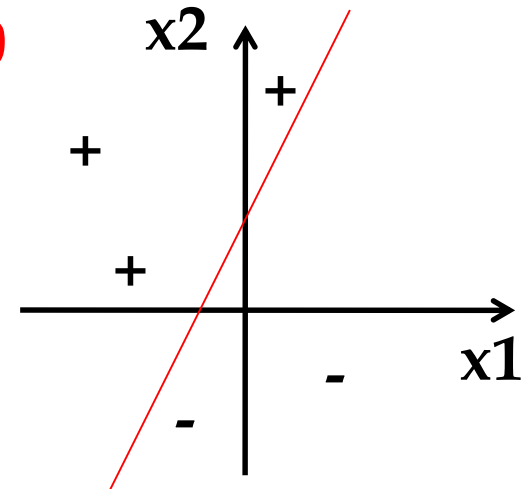
- Answer:

$$w_0 = -1 \cdot C$$

$$w_1 = -1 \cdot C$$

$$w_2 = 0.5 \cdot C$$

(where $C > 0$)



Logistic Regression Model

- Gradient Decent (梯度下降)

- Calculate the gradient vector
- Update the weighting in the opposite direction of the gradient vector at each surface point

- Repeat:
$$\begin{aligned}\tilde{\mathbf{W}}_{new}^{(j)} &= \tilde{\mathbf{W}}^{(j)} - \eta \frac{\partial C(\tilde{\mathbf{W}})}{\partial \tilde{\mathbf{W}}^{(j)}} \\ &= \tilde{\mathbf{W}}^{(j)} - \eta \sum_{i=1}^n \left[\left(\frac{e^{\tilde{\mathbf{W}}^T \tilde{\mathbf{X}}_i}}{1 + e^{\tilde{\mathbf{W}}^T \tilde{\mathbf{X}}_i}} - y_i \right) \tilde{\mathbf{X}}_i^{(j)} \right]\end{aligned}$$

- Until convergence

Apriori Algorithm

- **自连接**: 用 L_{k-1} 自连接得到 C_k
- **修剪**: 一个 k -项集, 如果他的一个 $k-1$ 项集 (他的子集) 不是频繁的, 那他本身也不可能是频繁的。
- pseudo code:

C_k : Candidate itemset of size k

L_k : frequent itemset of size k

$L_1 = \{\text{frequent items}\};$

for ($k = 1; L_k \neq \emptyset; k++$) **do begin**

C_{k+1} = candidates generated from L_k ;

for each transaction t in database **do**

 increment the count of all candidates in C_{k+1} that are contained in t

L_{k+1} = candidates in C_{k+1} with *minsup*

end

return $\cup_k L_k$;

Maximal Frequent Itemsets

- A maximal frequent itemset is defined as a frequent itemset for which **none of its immediate supersets are frequent**.
- We consider the itemset lattice shown in the following figure.
- The itemsets in the lattice are divided into two groups
 - Those that are frequent
 - Those that are infrequent

Closed Frequent Itemsets

- An itemset X is closed if none of its immediate supersets has exactly the same support count as X .
- In other words, X is not closed if at least one of its immediate supersets has the same support count as X .

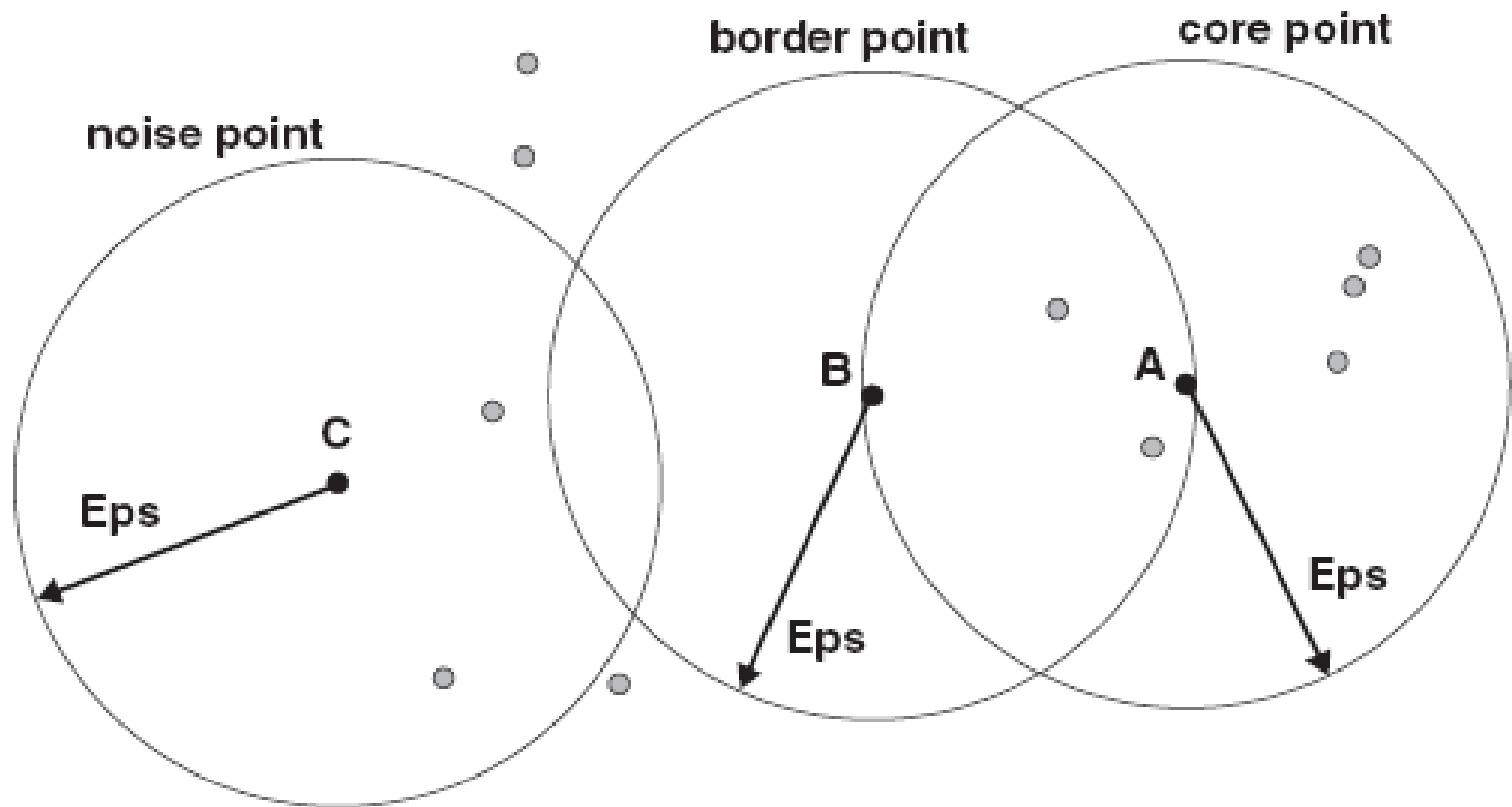
Partitional Clustering

- k -Means: Repeat...
 - Choose k arbitrary '**centroids**'
 - Assign each document to nearest centroid
 - Re-compute centroids
- Example of k -Means (划分法)
 - $x_1 = (0, 2)$, $x_2 = (0, 0)$, $x_3 = (1.5, 0)$, $x_4 = (5, 0)$,
 $x_5 = (5, 2)$
 - $k = 2$

DBSCAN

- We need to classify a point as being
 - In the interior of a dense region (a **core** point, 核心点).
 - At the edge of a dense region (a **border** point, 边界点)
 - In a sparsely occupied region (a **noise** or background point, 噪音点).
- The concepts of core, border and noise points are illustrated as follows.

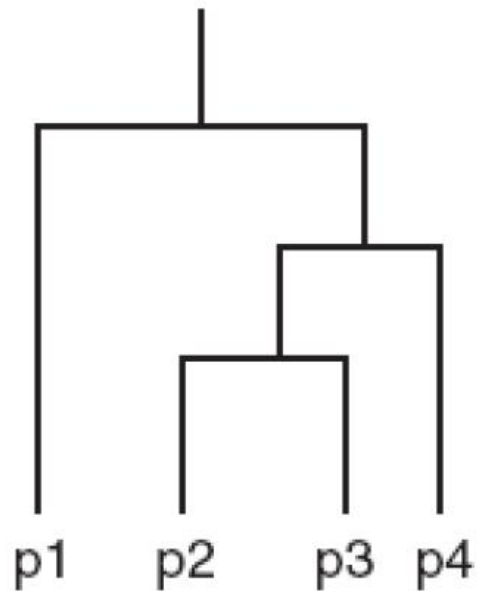
DBSCAN



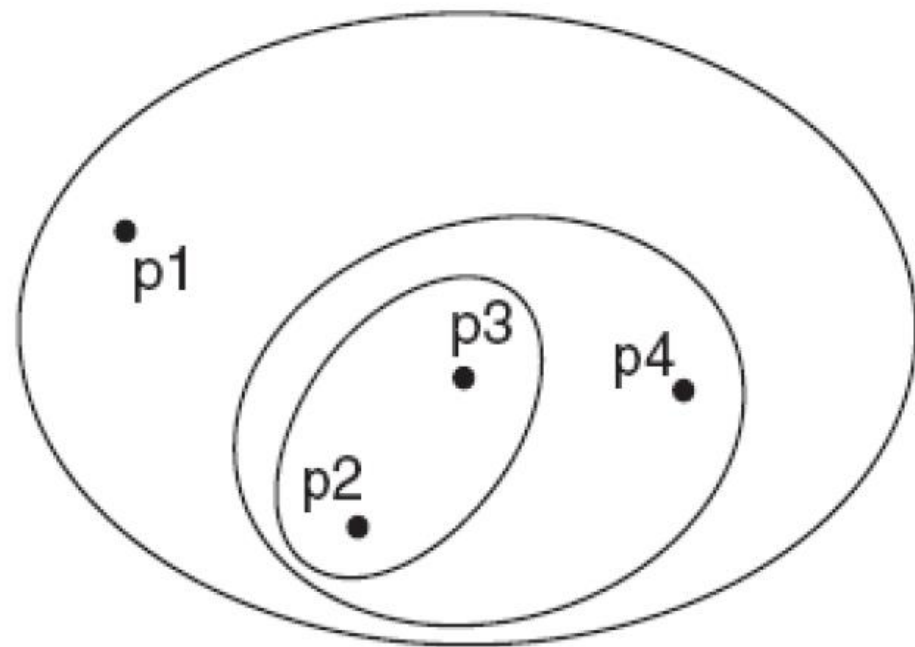
Hierarchical Clustering

- A hierarchical clustering is often displayed graphically using a tree-like diagram called the dendrogram (树状图).
- The dendrogram displays both
 - the cluster-subcluster relationships and
 - the order in which the clusters are merged (agglomerative) or split (divisive).
- For sets of 2-D points, a hierarchical clustering can also be graphically represented using a nested cluster diagram.

Hierarchical Clustering



(a) Dendrogram.



(b) Nested cluster diagram.

Hierarchical Clustering

- Different definitions of cluster distance leads to different versions of hierarchical clustering.
- These versions include
 - Single link (单连接) or MIN
 - Complete link (全连接) or MAX
 - Group average (组平均)

Single Link

- We now consider the single link or MIN version of hierarchical clustering.
- In this case, the distance of two clusters is defined as the minimum of the distance between any two points in the two different clusters.
- This technique is good at handling non-elliptical (非球状的) shapes.
- However, it is sensitive to noise and outliers.

Complete Link

- We now consider the complete link or MAX version of hierarchical clustering.
- In this case, the distance of two clusters is defined as the maximum of the distance between any two points in the two different clusters.
- Complete link is less susceptible (不敏感) to noise and outliers, but it tends to produce clusters with globular (球状) shapes.