

一、聚类任务

1. 无监督学习，训练样本的标记信息是未知的
2. 聚类仅能自动形成簇结构，但是所对应的概念语义需由使用者来把握和命名

二、性能度量

- 簇内相似度高、簇间相似度低
- 性能度量
 - 外部指标
 - 内部指标

三、距离计算

将属性划分为连续属性和离散属性

- 离散属性中
 - $\{1, 2, 3\}$ 这种属于有序属性
 - 不能从属性值上计算距离的为无序属性
- 各种计算距离的公式
 - 连续属性--》闵可夫斯基距离（也可用于有序属性）
 - 无序属性--》VDM
 - 混合属性--》闵可夫斯基距离 + VDM
 - 重要性不同的属性--》可以使用加权距离

四、原型聚类

4.1 k均值

k-均值算法的思路就是让同一个簇中的样本都是离本簇的中心较近，而离其他簇的中心较远。

输入: 样本集 $D = \{x_1, x_2, \dots, x_m\}$;

聚类簇数 k .

过程:

```
1: 从  $D$  中随机选择  $k$  个样本作为初始均值向量  $\{\mu_1, \mu_2, \dots, \mu_k\}$ 
2: repeat
3:   令  $C_i = \emptyset$  ( $1 \leq i \leq k$ )
4:   for  $j = 1, 2, \dots, m$  do
5:     计算样本  $x_j$  与各均值向量  $\mu_i$  ( $1 \leq i \leq k$ ) 的距离:  $d_{ji} = \|x_j - \mu_i\|_2$ ;
6:     根据距离最近的均值向量确定  $x_j$  的簇标记:  $\lambda_j = \arg \min_{i \in \{1, 2, \dots, k\}} d_{ji}$ ;
7:     将样本  $x_j$  划入相应的簇:  $C_{\lambda_j} = C_{\lambda_j} \cup \{x_j\}$ ;
8:   end for
9:   for  $i = 1, 2, \dots, k$  do
10:    计算新均值向量:  $\mu'_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$ ;
11:    if  $\mu'_i \neq \mu_i$  then
12:      将当前均值向量  $\mu_i$  更新为  $\mu'_i$ 
13:    else
14:      保持当前均值向量不变
15:    end if
16:  end for
17: until 当前均值向量均未更新
```

输出: 簇划分 $C = \{C_1, C_2, \dots, C_k\}$

为避免运行时间过长,
通常设置一个最大运行轮
数或最小调整幅度阈值,
若达到最大轮数或调整幅
度小于阈值, 则停止运行.

激活 Winc

https://blog.csdn.net/qq_34354651

4.2 学习向量量化(LVQ)

LVQ算法是聚类算法中的一个另类, LVQ(Learning Vector Quantization)算法需要样本的类别标记, 算是半个监督学习算法.

输入: 样本集 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$;

原型向量个数 q , 各原型向量预设的类别标记 $\{t_1, t_2, \dots, t_q\}$;

学习率 $\eta \in (0, 1)$.

过程:

```
1: 初始化一组原型向量  $\{p_1, p_2, \dots, p_q\}$ 
2: repeat
3:   从样本集  $D$  随机选取样本  $(x_j, y_j)$ ;
4:   计算样本  $x_j$  与  $p_i$  ( $1 \leq i \leq q$ ) 的距离:  $d_{ji} = \|x_j - p_i\|_2$ ;
5:   找出与  $x_j$  距离最近的原型向量  $p_{i^*}$ ,  $i^* = \arg \min_{i \in \{1, 2, \dots, q\}} d_{ji}$ ;
6:   if  $y_j = t_{i^*}$  then
7:      $p' = p_{i^*} + \eta \cdot (x_j - p_{i^*})$ 
8:   else
9:      $p' = p_{i^*} - \eta \cdot (x_j - p_{i^*})$ 
10:  end if
11: 将原型向量  $p_{i^*}$  更新为  $p'$ 
12: until 满足停止条件
```

输出: 原型向量 $\{p_1, p_2, \dots, p_q\}$

x_j 与 p_{i^*} 的类别相同.

x_j 与 p_{i^*} 的类别不同.

如达到最大迭代轮数.

https://blog.csdn.net/qq_34354651

4.3 高斯混合聚类

高斯分布就是正态分布

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)},$$

4.3.1 最大似然估计(MLE: Maxmium likelihood)

4.3.2 EM算法

EM算法可以解决机器学习中的“隐变量”问题，在真实数据场景中，有很多的样本会确实一些属性值，但是这些数据是有价值的，只是需要处理一下这些缺失的属性值，这些缺失的属性值就成为隐变量。

在训练学习器的过程中，我们是需要根据函数 $p(\mathbf{x}|\theta)$ 中的 \mathbf{x} 去获得 θ ，这样就可以获得完整的学习器：模型+参数。但是此时 \mathbf{x} 是不完整的，所以我们根据上面的描述： $p(\mathbf{x}|\theta)$ ，用 θ 去估计 \mathbf{x} ，但是 θ 也是未知的。就会陷入到一个死循环中，EM算法就是想采用一种逐步迭代的思路去打破这个循环：

- 初始化一个 θ ，用这个初始化参数去估计 \mathbf{x} 中的隐变量 x_c 。
- 用 x_c 放入样本 \mathbf{x} 中，去估计下一个 θ_{t+1}
- 不停的迭代循环，知道达到某个停止条件：最大迭代次数或者说是变动量小于阈值。

4.3.3 高斯混合聚类

然后假设样本是服从参数不同的高斯分布，总共有 k 个不同的高斯分布(看出来了吧，就是想往这个 k 上来靠)。那么就可以定义一个混合高斯分布

定义了高斯混合分布

$$p_M(\mathbf{x}) = \sum_{i=1}^k \alpha_i * p(\mathbf{x}|\mu_i, \Sigma_i), \text{ 其中 } \sum_{i=1}^k \alpha_i = 1$$

这个概率公式定义了每个样本的生成概率。那么我们已经有了样本，要求样本 x_j 的后验概率 $p(z_j = i | x_j)$ ，也就是这个样本 x_j 是从哪个分布中出来的概率。根据贝叶斯公式（

）可以得到：

$$\begin{aligned} p(z_j = i | x_j) &= \frac{P(z_j = i) * p_M(x_j | z_j = i)}{p_M(x_j)} \\ &= \frac{\alpha_i p(x_j | \mu_i, \Sigma_i)}{\sum_{l=1}^k \alpha_l p(x_j | \mu_l, \Sigma_l)} \end{aligned}$$

把这个复杂的一腿的公式记为： γ_{ij} 。

根据之前的聚类套路，就是计算每个样本 x_j 的这个后验概率 γ_{ij} ，哪个最大，就属于哪个类。但是计算这个东西的过程中，又碰到了手动增加进去的“隐变量”： α ，所以需要采用EM算法，将隐变量 α 和模型参数 μ, Σ 一起计算出来。模型参数 μ, Σ 的推导过程大家可以去看书。

输入：样本集 $D = \{x_1, x_2, \dots, x_m\}$;
高斯混合成分个数 k 。

过程：

- 1: 初始化高斯混合分布的模型参数 $\{(\alpha_i, \mu_i, \Sigma_i) \mid 1 \leq i \leq k\}$
- 2: **repeat**
- 3: **for** $j = 1, 2, \dots, m$ **do**
- 4: 根据式(9.30)计算 x_j 由各混合成分生成的后验概率，即
 $\gamma_{ji} = p_M(z_j = i | x_j) \ (1 \leq i \leq k)$
- 5: **end for**
- 6: **for** $i = 1, 2, \dots, k$ **do**
- 7: 计算新均值向量: $\mu'_i = \frac{\sum_{j=1}^m \gamma_{ji} x_j}{\sum_{j=1}^m \gamma_{ji}}$;
- 8: 计算新协方差矩阵: $\Sigma'_i = \frac{\sum_{j=1}^m \gamma_{ji} (x_j - \mu'_i)(x_j - \mu'_i)^T}{\sum_{j=1}^m \gamma_{ji}}$;
- 9: 计算新混合系数: $\alpha'_i = \frac{\sum_{j=1}^m \gamma_{ji}}{m}$;
- 10: **end for**
- 11: 将模型参数 $\{(\alpha_i, \mu_i, \Sigma_i) \mid 1 \leq i \leq k\}$ 更新为 $\{(\alpha'_i, \mu'_i, \Sigma'_i) \mid 1 \leq i \leq k\}$
- 12: **until** 满足停止条件
- 13: $C_i = \emptyset \ (1 \leq i \leq k)$
- 14: **for** $j = 1, 2, \dots, m$ **do**
- 15: 根据式(9.31)确定 x_j 的簇标记 λ_j ;
- 16: 将 x_j 划入相应的簇: $C_{\lambda_j} = C_{\lambda_j} \cup \{x_j\}$
- 17: **end for**

输出：簇划分 $C = \{C_1, C_2, \dots, C_k\}$

例如达到最大迭代轮数。

<https://blog.csdn.net/yogamer>

4.4 密度聚类

密度聚类实际上一个从某个点出发，找到自己足够近的邻居的过程，这些邻居在算法中用了更专业的词：邻域。一个经典的邻域算法是DBSCAN算法，算法中定义了两个邻域参数： $(\epsilon, \text{MinPts})$ 。由这两个参数确定一个邻域。

- ϵ : 定义了样本集中样本与某个样本 x_j 的最小距离, 即如果距离不大于 ϵ , 即样本位于样本 x_j 的邻域内。
- MinPts: 若样本 x_j 的邻域包含的样本数大于 MinPts, 则称 x_j 为核心对象。

另外, 还有几个其他的概念:

- 密度直达: 若样本 x_j 位于样本 x_i 的邻域内, 且 x_i 是核心对象, 则称 x_j 由 x_i 密度直达。
- 密度可达: 若 x_j 可以由一系列的中间节点 $p, p_1 \dots$ 由 x_i 直达, 则称 x_j 由 x_i 密度可达。
- 密度相连: 若 x_i 和 x_j 均通过 x_k 可达, 则称 x_i, x_j 相连。这里两个样本是距离是相互的, 但是不是每个样本点都是核心对象。

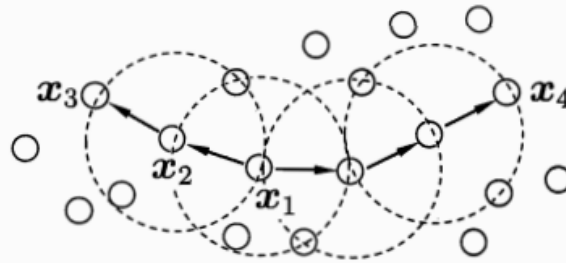


图 9.8 DBSCAN 定义的基本概念 ($MinPts = 3$): 虚线显示出 ϵ -邻域, x_1 是核心对象, x_2 由 x_1 密度直达, x_3 由 x_1 密度可达, x_3 与 x_4 密度相连。

<https://blog.csdn.net/pcgamer>

输入: 样本集 $D = \{x_1, x_2, \dots, x_m\}$;
邻域参数 $(\epsilon, MinPts)$.

过程:

```
1: 初始化核心对象集合:  $\Omega = \emptyset$ 
2: for  $j = 1, 2, \dots, m$  do
3:   确定样本  $x_j$  的  $\epsilon$ -邻域  $N_\epsilon(x_j)$ ;
4:   if  $|N_\epsilon(x_j)| \geq MinPts$  then
5:     将样本  $x_j$  加入核心对象集合:  $\Omega = \Omega \cup \{x_j\}$ 
6:   end if
7: end for
8: 初始化聚类簇数:  $k = 0$ 
9: 初始化未访问样本集合:  $\Gamma = D$ 
10: while  $\Omega \neq \emptyset$  do
11:   记录当前未访问样本集合:  $\Gamma_{old} = \Gamma$ ;
12:   随机选取一个核心对象  $o \in \Omega$ , 初始化队列  $Q = \langle o \rangle$ ;
13:    $\Gamma = \Gamma \setminus \{o\}$ ;
14:   while  $Q \neq \emptyset$  do
15:     取出队列  $Q$  中的首个样本  $q$ ;
16:     if  $|N_\epsilon(q)| \geq MinPts$  then
17:       令  $\Delta = N_\epsilon(q) \cap \Gamma$ ;
18:       将  $\Delta$  中的样本加入队列  $Q$ ;
19:        $\Gamma = \Gamma \setminus \Delta$ ;
20:     end if
21:   end while
22:    $k = k + 1$ , 生成聚类簇  $C_k = \Gamma_{old} \setminus \Gamma$ ;
23:    $\Omega = \Omega \setminus C_k$ 
24: end while
```

输出: 簇划分 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$

<https://blog.csdn.net/pcgamer>

1. 核心对象集合初始化为空
2. 第2行到第7行是确定每个样本 x_j 的邻域 $N_c(x_j)$ ，并从中挑出核心对象。
3. 从第10行开始到第24行，就是由核心对象来往外扩展生成簇的过程。
4. 从核心对象集合 Ω 中随机挑选一个核型对象，通过14-21行的代码将该核型对象的所有密度可达(注意密度可达的定义是需要对方也是核心对象)的样本全部加入到同一个簇内。在整个过程中，使用 $\Gamma = D$ 不停的去减去这些簇中的样本点，最后再第22行的时候，用 $\Gamma_{old} - \Gamma$ 就是这些样本点了，有点负负得正的意思。
 1. 刚开始 $\Gamma = D$ ， $\Gamma_{old} = \Gamma$
 2. 选出一个核心对象样本点 x_j 后， $\Gamma = \Gamma - x_j$
 3. 将核心对象样本点 x_j 的邻域(比如有 x_1, x_4, x_9)保存到 Δ 中，然后 $\Gamma = \Gamma - \Delta$
 4. 最后用 $\Gamma_{old} - \Gamma$ 就等于 (x_j, x_1, x_4, x_9) 了。
5. 形成完一个簇后，再重新从 Ω 中挑选核心对象进行扩展，直到 Ω 为空
6. 基于密度聚类算法会存在一些样本点没法成为核心对象，就会在算法结束之后成为一些孤立的样本点，成为噪声样本。

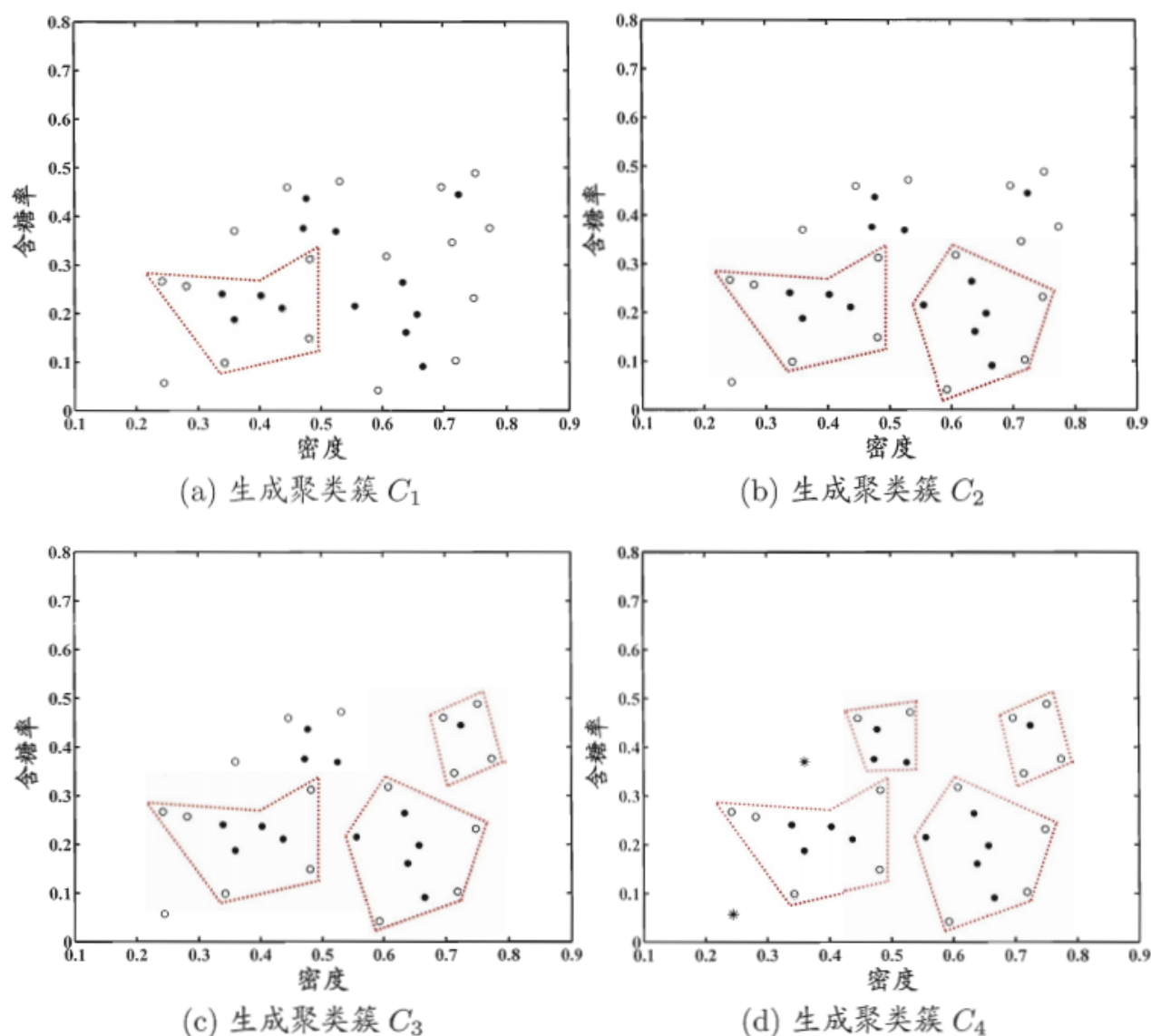
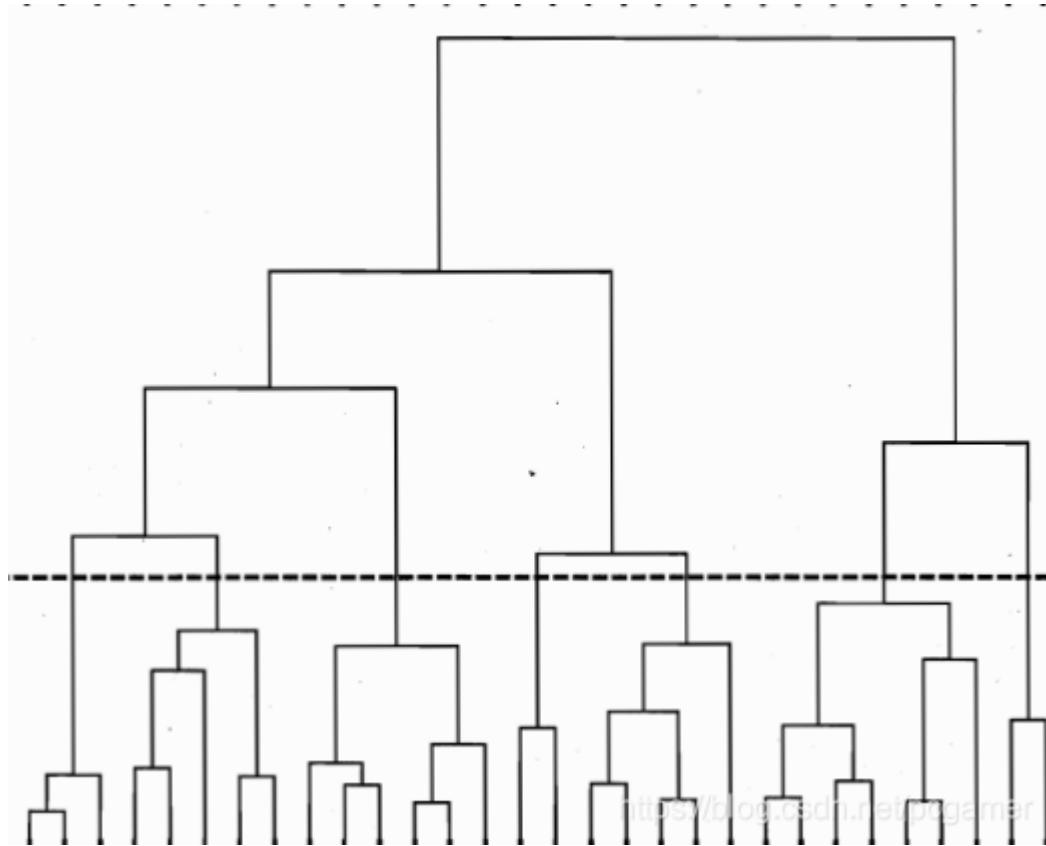


图 9.10 DBSCAN 算法($\epsilon = 0.11$, $MinPts = 5$)生成聚类簇的先后情况. 核心对象、非核心对象、噪声样本分别用“●” “○” “*”表示, 红色虚线显示出簇划分。g.csdn.net/pcgamer

4.5 层次聚类

层次聚类是将样本点分拆或者合并成类似树的结构，树的每一个分支就是一个簇。书中提到的AGNES是基于合并的自底向上的方法。用的思路也比较简单：

- 首先所有的样本都单独成为一个簇。
- 计算簇和簇之间的距离，距离足够近的形成一个新的簇，最终形成一个类似树的结构：



- 通过控制簇的个数k来控制树分支的合并情况。

这里就有一个问题，之前计算的距离是样本之间的距离，在层次聚类里需要计算簇之间的距离，所以需要定义簇的距离如何进行计算。簇实际上就是一个样本集合，所以定义一下关于集合的某种运算即可

- 最小距离: $dist(C_i, C_j) = \min_{x \in C_i, z \in C_j} dist(x, z)$, 将两个簇内样本之间的最小距离作为簇之间的距离。
- 最大距离: $dist(C_i, C_j) = \max_{x \in C_i, z \in C_j} dist(x, z)$, 将两个簇内样本之间的最大距离作为簇之间的距离。
- 平均距离: $dist(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{x \in C_i} \sum_{z \in C_j} dist(x, z)$, 将两个簇内样本之间的平均距离作为簇之间的距离。

算法过程

输入: 样本集 $D = \{x_1, x_2, \dots, x_m\}$;
聚类簇距离度量函数 d ;
聚类簇数 k .

过程:

```
1: for  $j = 1, 2, \dots, m$  do
2:    $C_j = \{x_j\}$ 
3: end for
4: for  $i = 1, 2, \dots, m$  do
5:   for  $j = 1, 2, \dots, m$  do
6:      $M(i, j) = d(C_i, C_j)$ ;
7:      $M(j, i) = M(i, j)$ 
8:   end for
9: end for
10: 设置当前聚类簇个数:  $q = m$ 
11: while  $q > k$  do
12:   找出距离最近的两个聚类簇  $C_{i^*}$  和  $C_{j^*}$ ;
13:   合并  $C_{i^*}$  和  $C_{j^*}$ :  $C_{i^*} = C_{i^*} \cup C_{j^*}$ ;
14:   for  $j = j^* + 1, j^* + 2, \dots, q$  do
15:     将聚类簇  $C_j$  重编号为  $C_{j-1}$ 
16:   end for
17:   删除距离矩阵  $M$  的第  $j^*$  行与第  $j^*$  列;
18:   for  $j = 1, 2, \dots, q - 1$  do
19:      $M(i^*, j) = d(C_{i^*}, C_j)$ ;
20:      $M(j, i^*) = M(i^*, j)$ 
21:   end for
22:    $q = q - 1$ 
23: end while
```

输出: 簇划分 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$
