

Beijing Forest Studio
北京理工大学信息系统及安全对抗实验中心



半监督聚类和患者相似性分析

数据挖掘组

硕士研究生 谢崇玮

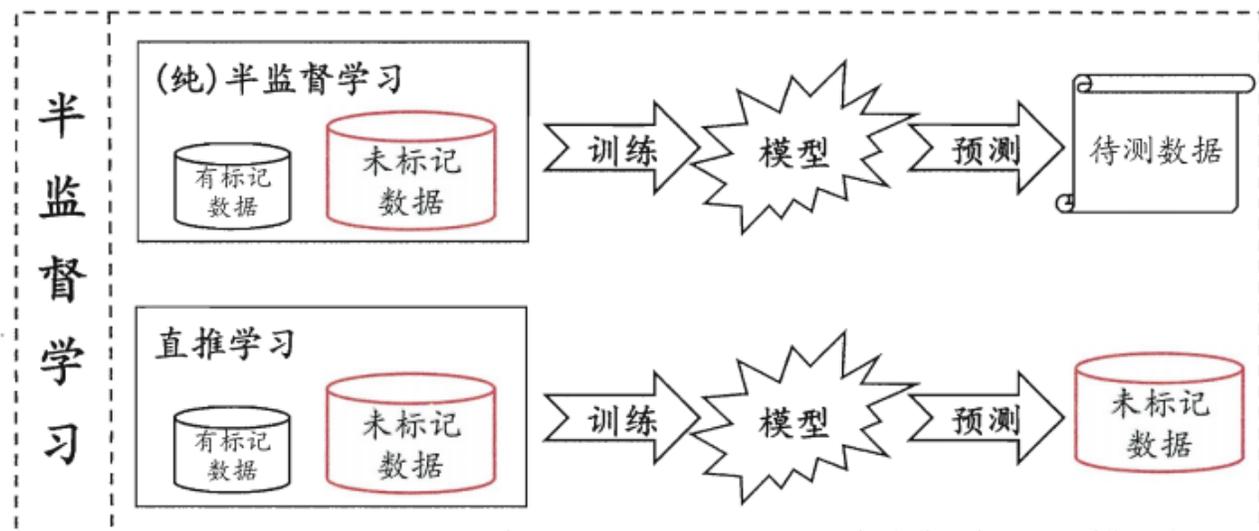
2022年9月4日



- 背景简介
- 基本概念
- 算法原理
- 总结分析
- 参考文献

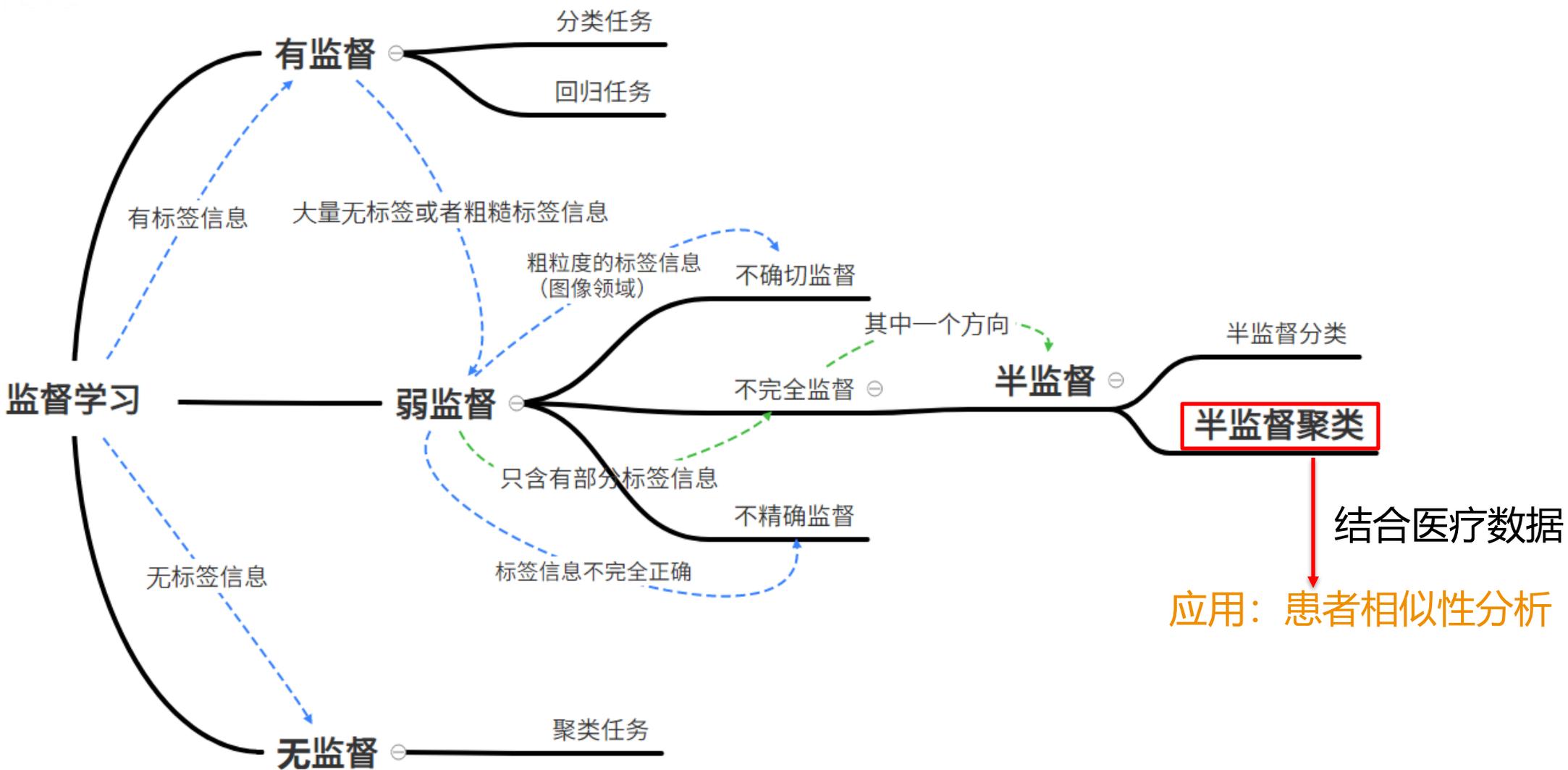
- 预期收获
 - 1.了解半监督聚类技术
 - 2.了解半监督聚类技术在患者相似性分析上的应用

- 半监督学习
 - 大量未标注数据，少量标注数据
 - 仅利用未标注数据聚类分析浪费了宝贵的标注资源
 - 标注数据的获取**耗时耗力**
- 患者相似性（主要利用电子医疗记录数据）
 - 对患者细粒度区分
 - 构建合适的**患者群落划分**
 - 为后续医生的**诊断和指导**提供帮助

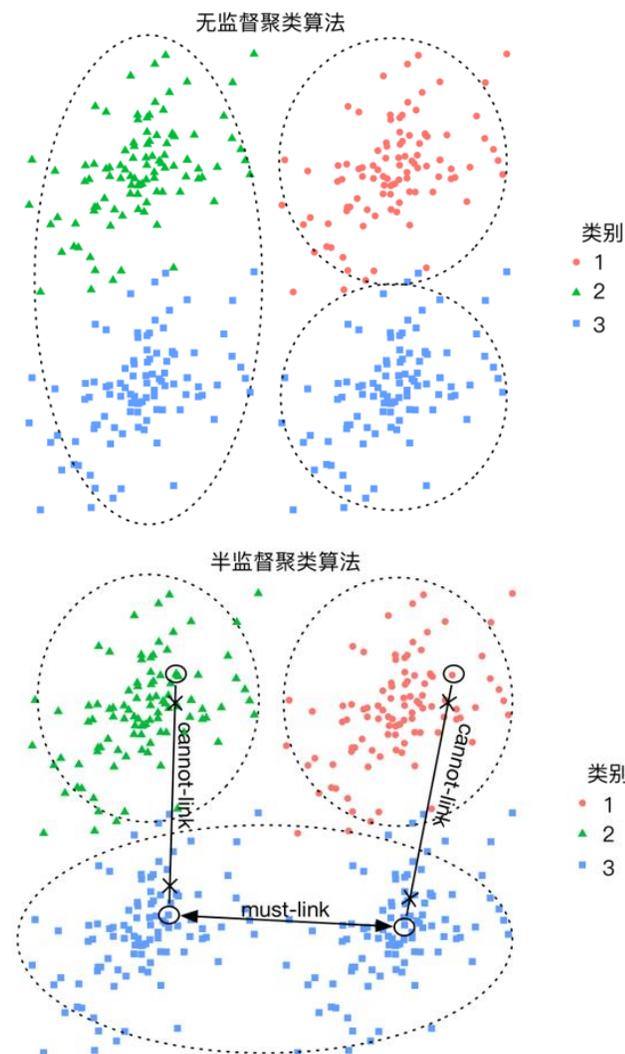




基本概念



- 半监督聚类
 - 结合半监督学习和聚类的方法（引入一些**监督信息**来指导**聚类**过程）
- 基于约束的方法
 - **成对约束**（Must-link/Cannot-link）
 - **正负样本约束**（A属于 Q_1 类/B不属于 Q_2 类）
 - **集群约束**（**簇大小约束**、**内部分布约束**等）
- 基于距离的方法
 - 首先训练**距离度量**以满足类别或限制信息，然后使用基于距离度量的聚类算法进行聚类
 - 凸优化的马氏距离、由最短路径算法改进的欧式距离、使用梯度下降算法的KL散度、谱聚类方法



- 聚类效果评判指标

- 外部指标 (监督)

$$RI = \frac{TP + TN}{TP + FP + TF + FN} = \frac{TP + TN}{C_N^2}$$

$$NMI(\Omega, C) = \frac{I(\Omega; C)}{(H(\Omega) + H(C))/2}$$

- 内部指标 (簇内相似度, 簇间分离度)

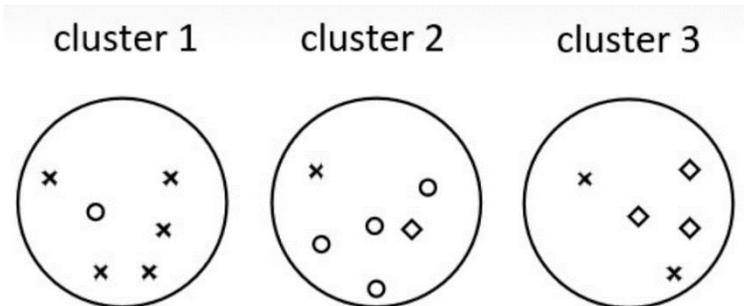
- 轮廓系数、SSE



$$SSE = \sum_{i=1}^r \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2, \bar{X}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij}, i = 1, 2, \dots, r$$

	同簇	非同簇
同类	TP	FN
非同类	FP	TN

	同簇	非同簇
同类	TP=20	FN=24
非同类	FP=20	TN=72



- 聚类中心K的选取标准

- 随着K值的增大, 数据集的划分会更加精确, 即SSE会逐渐变小。一般选取SSE随着K值的增大突然大幅下降的位置作为最优K值



- 患者相似性分析

- 选取用户电子病历记录中的特定属性
- 根据不同类型的特定属性选择不同的相似度计算规则，计算总相似度
- 根据相似度将患者聚类，找出患者自身与诊断和治疗方案之间的关系





算法原理

T	整合多源约束来提高聚类结果的有效性
I	聚类数据集
P	1.通过转换为成对关系矩阵, 统一地表示了不同类型的约束 2.整合多源约束指导聚类过程
O	聚类簇

P	算法大多基于单源约束, 很少考虑整合多源约束来提高聚类质量
C	多源包含不同类型的约束信息 (成对约束、正负标签约束)
D	统一不同类型约束
L	2021 SCI 1区

- 如何构造统一的表示来保存**不同类型**的约束?
- 如何**整合多源约束**来指导聚类过程?

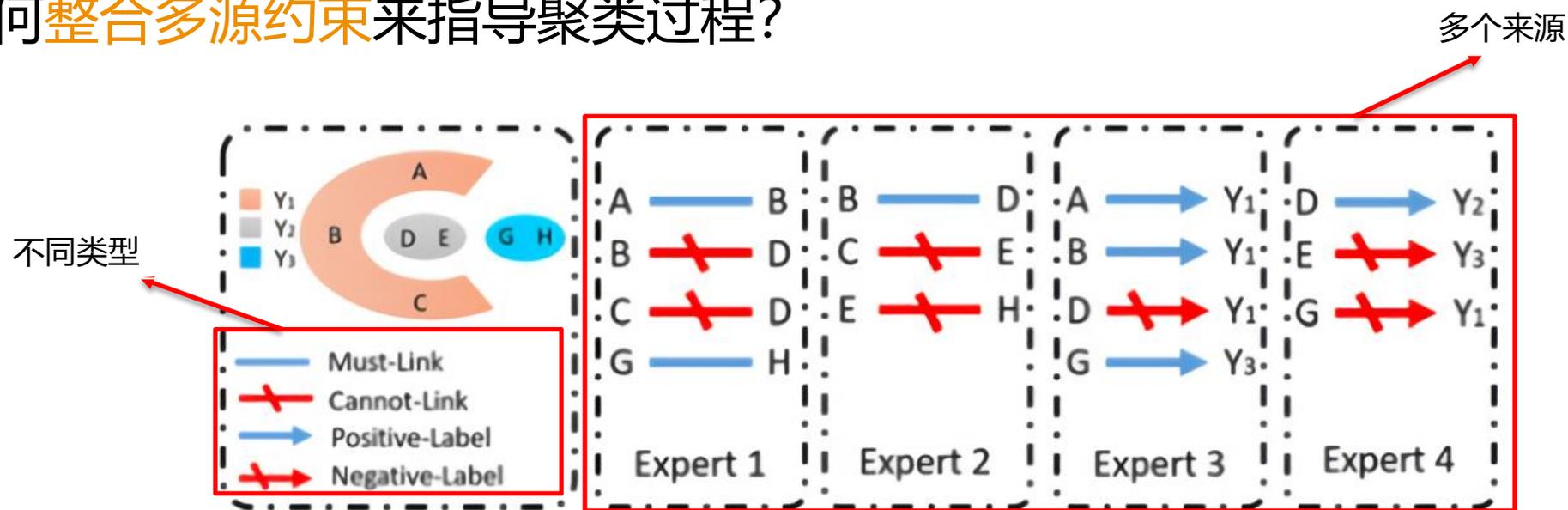


Fig. 1. Example of a data set with 3 clusters and constraints from 4 experts.

- 统一表示不同类型的约束

- 成对约束

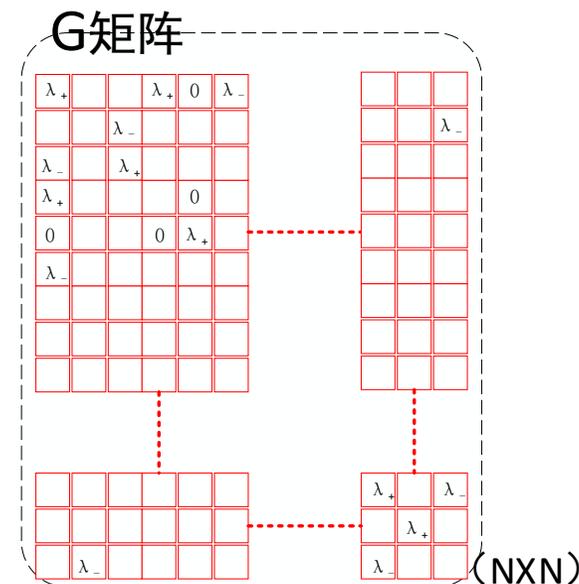
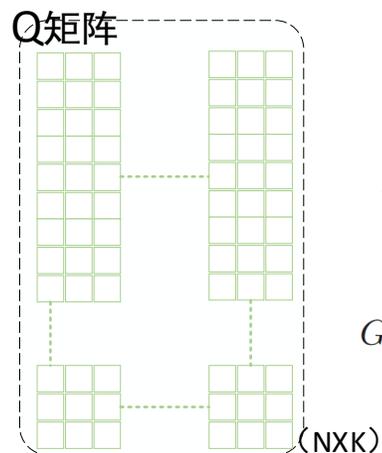
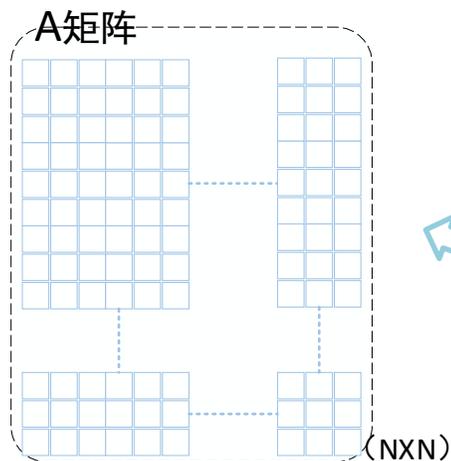
(ML $\rightarrow A_{ij} = \text{正}$, CL $\rightarrow A_{ij} = \text{负}$,

其它 $\rightarrow A_{ij} = 0$) $N \times N$ 矩阵

- 正负例样本

(正例 $\rightarrow Q_{il} = \text{正}$, 负例 $\rightarrow Q_{il} =$

负, 其它 $\rightarrow Q_{il} = 0$) $N \times K$ 矩阵



$$G_{(ij)} = \begin{cases} \lambda_+, & \exists l \langle \mathbf{x}_i, \mathbf{y}_l \rangle, \langle \mathbf{x}_j, \mathbf{y}_l \rangle \in P^+ \text{ or} \\ & |\{\mathbf{y}_l | \langle \mathbf{x}_i, \mathbf{y}_l \rangle, \langle \mathbf{x}_j, \mathbf{y}_l \rangle \in N^+\}| = k - 1, \\ \lambda_-, & \exists l \neq h \langle \mathbf{x}_i, \mathbf{y}_l \rangle, \langle \mathbf{x}_j, \mathbf{y}_h \rangle \in P^+ \text{ or} \\ & \exists l \langle \mathbf{x}_i, \mathbf{y}_l \rangle \in P^+, \langle \mathbf{x}_j, \mathbf{y}_l \rangle \in N^+, \\ 0, & \text{otherwise.} \end{cases}$$

- 整合多源指导聚类过程

- 获得m个专家（多源）提供的G矩阵，用 $U_{n \times k}$ 矩阵来示最终聚类结果

$$\min \Theta(U) = F(U) + \alpha E(U)$$

$\alpha \geq 0$; 此处设置=1

$$F(U) = \|\mathbb{X} - UZ\|_F^2$$

$$E(U) = \frac{1}{2} \sum_{t=1}^m \|G_t - \hat{U}\hat{U}^T\|_F^2 \quad (\text{Z是k*d维中心矩阵, } \hat{U} \text{表示U的归一化矩阵})$$

$$\max_{\hat{U}} \text{Tr} \left(\hat{U}^T \left(K + \alpha \sum_{t=1}^m G_t \right) \hat{U} \right), \text{ s.t. } \hat{U}^T \hat{U} = I_k.$$

$$\lambda_+ = \frac{1}{\gamma_{(i,i)}} (\max(K) - K_{(ij)}) \quad \text{优化目标} \quad \lambda_- = \frac{1}{\gamma_{(i,j)}} (\min(K) - K_{(ij)})$$

(γ_{ij} 是m源中 x_i 和 x_j 之间的约束数量, K_{ij} 表示其基于核的特征空间中的相似性)

- 实验设置

- 实验环境:

- 配备Intel i7-4710MQ CPU@2.5Hz和16GB的RAM的PC上进行

- 实验数据

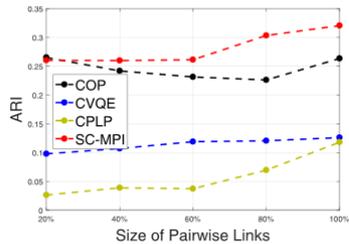
- 12个基准数据集

- 8种对比算法

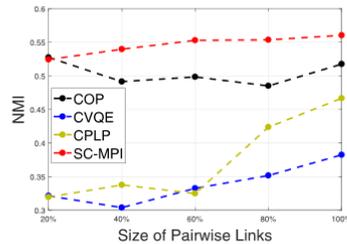
Data Sets	#Objects	#Features	#Classes
Yale [43]	165	1,024	15
ORL [43]	400	1,024	40
Banknote [42]	1,372	4	2
COIL20 [43]	1,440	1,024	20
Isolet [43]	1,560	617	26
OpticalDigits [42]	5,620	64	10
Statlog [42]	6,435	36	6
COIL100 [43]	7,200	1,024	100
MNIST [43]	10,000	784	10
PenDigits [43]	10,992	16	10
USPS [43]	11,000	256	10
Letters [42]	20,000	16	26

- 评估指标
 - ARI、NMI, 各自范围内越大聚类效果越好
- 聚类数K设置为给定数据集上的真实类数

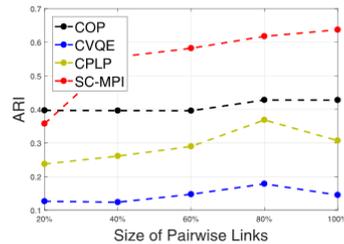
1、与单含成对约束算法对比



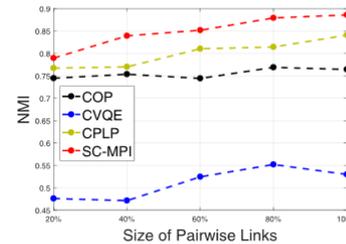
(a) ARI against size of pairwise links on data set Yale



(b) NMI against size of pairwise links on data set Yale

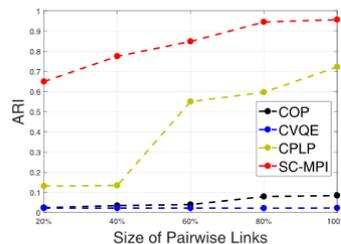


(c) ARI against size of pairwise links on data set ORL

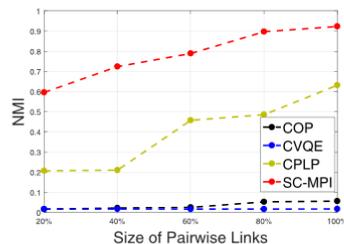


(d) NMI against size of pairwise links on data set ORL

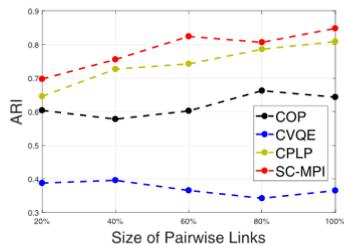
(ML数量=CL数量)



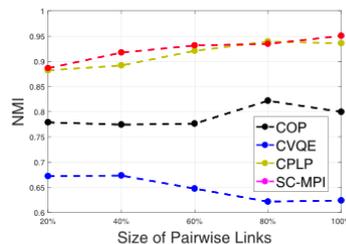
(e) ARI against size of pairwise links on data set Banknote



(f) NMI against size of pairwise links on data set Banknote



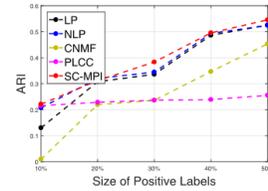
(g) ARI against size of pairwise links on data set COIL20



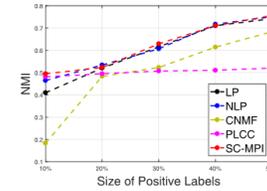
(h) NMI against size of pairwise links on data set COIL20

2、与给定正标签约束算法对比

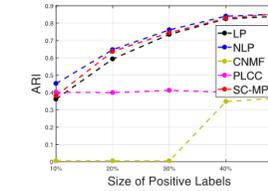
- 对于正标签约束，所提出的表示方法可以克服真实类标签和聚类标签之间的错位问题



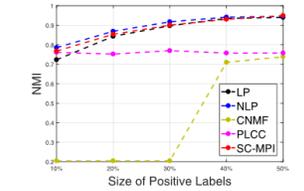
(a) ARI against size of positive labels on data set Yale



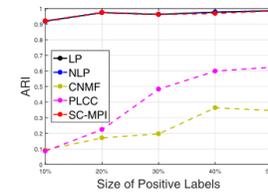
(b) NMI against size of positive labels on data set Yale



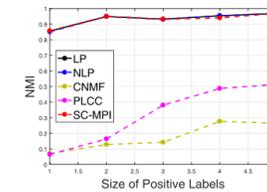
(c) ARI against size of positive labels on data set ORL



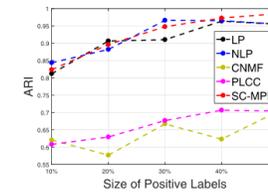
(d) NMI against size of positive labels on data set ORL



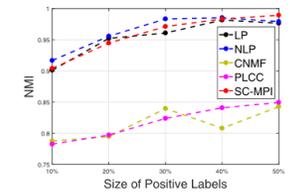
(e) ARI against size of positive labels on data set Banknote



(f) NMI against size of positive labels on data set Banknote



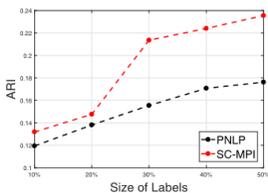
(g) ARI against size of positive labels on data set COIL20



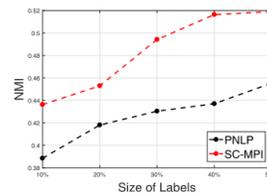
(h) NMI against size of positive labels on data set COIL20

3、与同时给定正、负标签约束算法对比

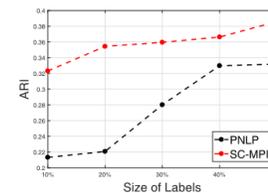
(正标签数量=负标签数量)



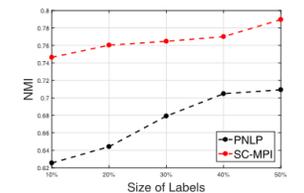
(a) ARI against size of labels on data set Yale



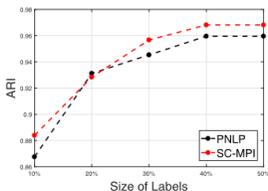
(b) NMI against size of labels on data set Yale



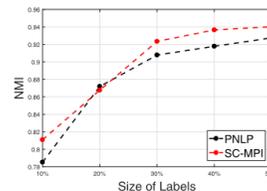
(c) ARI against size of labels on data set ORL



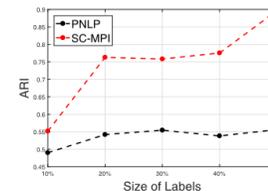
(d) NMI against size of labels on data set ORL



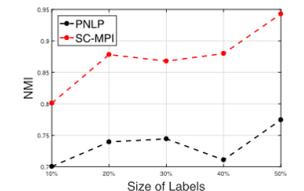
(e) ARI against size of labels on data set Banknote



(f) NMI against size of labels on data set Banknote

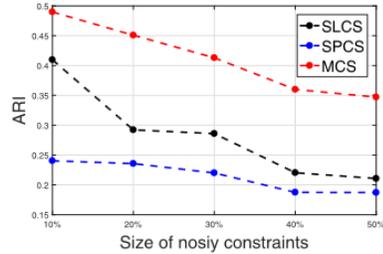


(g) ARI against size of labels on data set COIL20

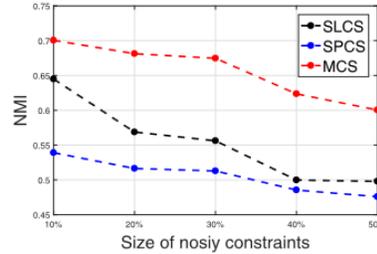


(h) NMI against size of labels on data set COIL20

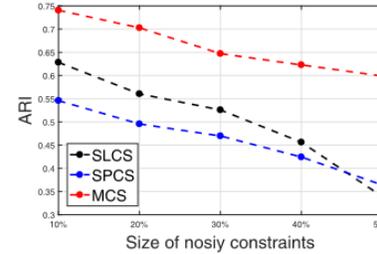
4、引入不正确约束信息的情况



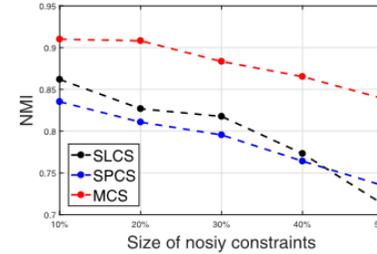
(a) ARI against size of noisy constraints on data set Yale



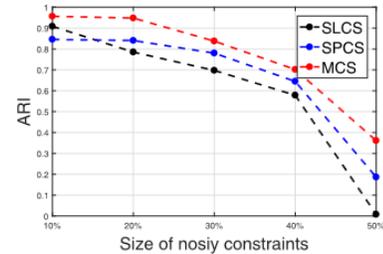
(b) NMI against size of noisy constraints on data set Yale



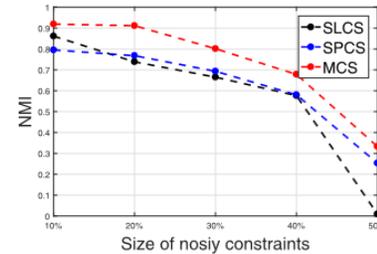
(c) ARI against size of noisy constraints on data set ORL



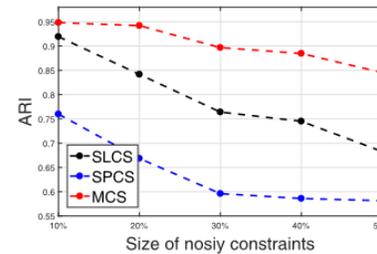
(d) NMI against size of noisy constraints on data set ORL



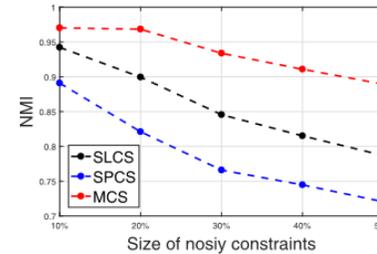
(e) ARI against size of noisy constraints on data set Banknote



(f) NMI against size of noisy constraints on data set Banknote



(g) ARI against size of noisy constraints on data set COIL20



(h) NMI against size of noisy constraints on data set COIL20

(整合多源约束, 降低不正确信息的影响)

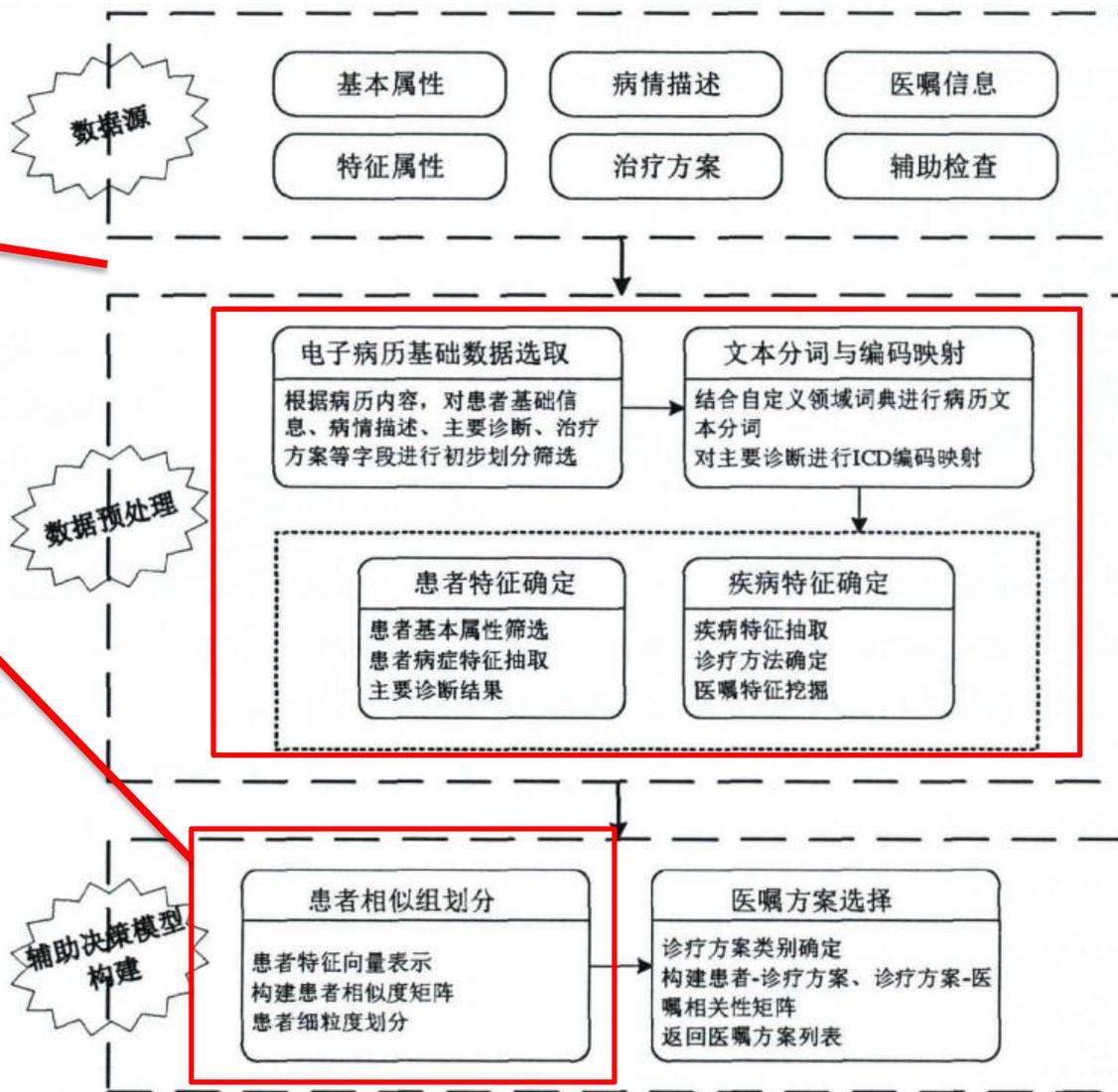


算法原理

T	结合 电子医疗记录 ，构建 患者配对约束 ，找到与目标患者 最相似的患者集
I	电子病历的入院记录
P	<ol style="list-style-type: none">1.从电子病历中获取患者自身的多种特征，并构建患者属性体系2.结合各类属性特征计算患者总相似度，获取成对监督信息3.利用成对信息指导聚类，实现患者相似组划分
O	各患者相似组

P	如何构建患者间约束信息指导患者相似集聚类
C	根据电子医疗记录可提取计算患者相似度
D	不同类别属性的相似度统一
L	2019 SCI二区

- 医嘱辅助决策模型框架
 - 数据获取与处理
 - 患者相似组划分
 - 诊疗医嘱推荐



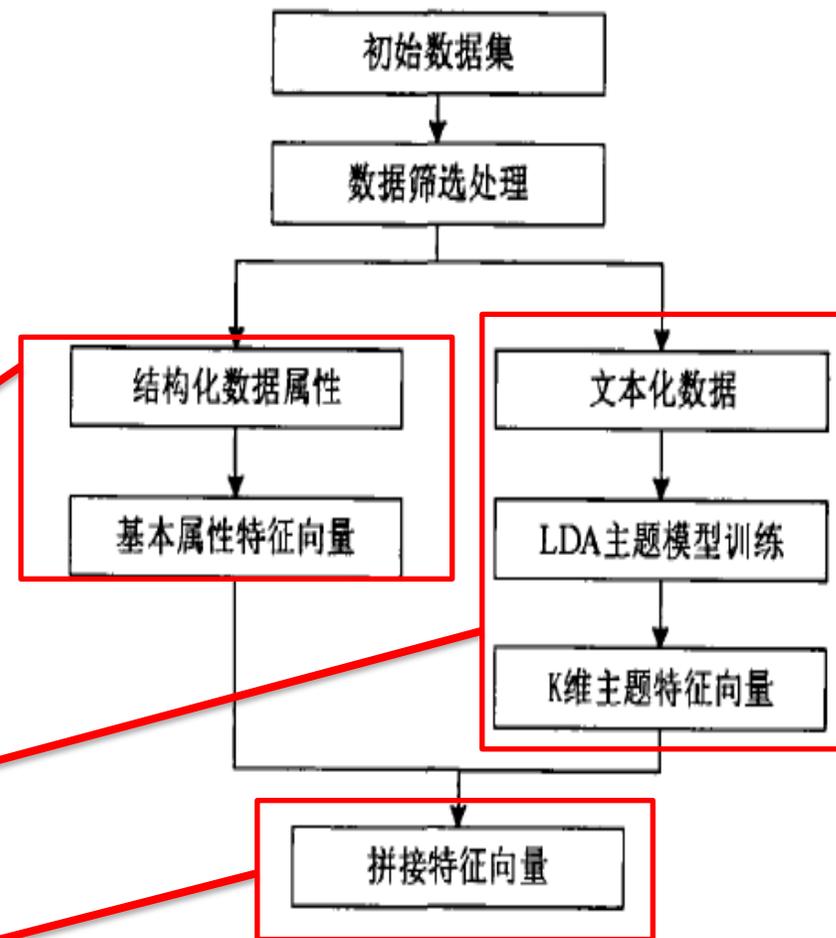
- 数据获取与处理

- 电子病历的入院记录

- 患者基本信息、患者自述信息、
 - 医生检查信息、辅助检查信息、初步诊断信息

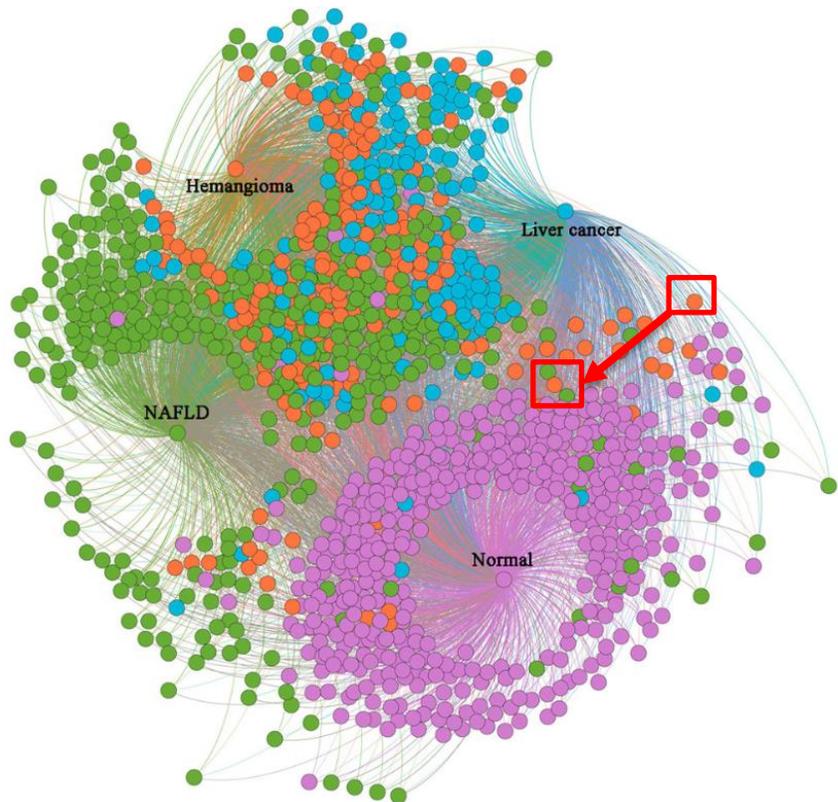
- 特征提取

- **结构化字段**：医生建议保留可能影响患者相似性的字段、同时为保证数据隐私，将病人ID、姓名等身份识别字段筛除。
 - **非结构化文本字段**：对N个字段中每个部分的关键词提取数或主题数定为K，共可产生N*K维属性值
 - 属性特征与文本主题特征进行拼接，形成最终**属性体系**



患者相似组划分

- 成对约束指导相似组构建
- 聚类评估指标: 轮廓系数



算法 4-2. 基于成对约束的患者相似组划分 K-means 聚类算法

输入: 文档数据集 $D = \{x_1, x_2, \dots, x_n\}$, must-link 集合 ML 和 cannot-link 集合 CL , 初始聚类中心数 K \longrightarrow 人工选取的影响

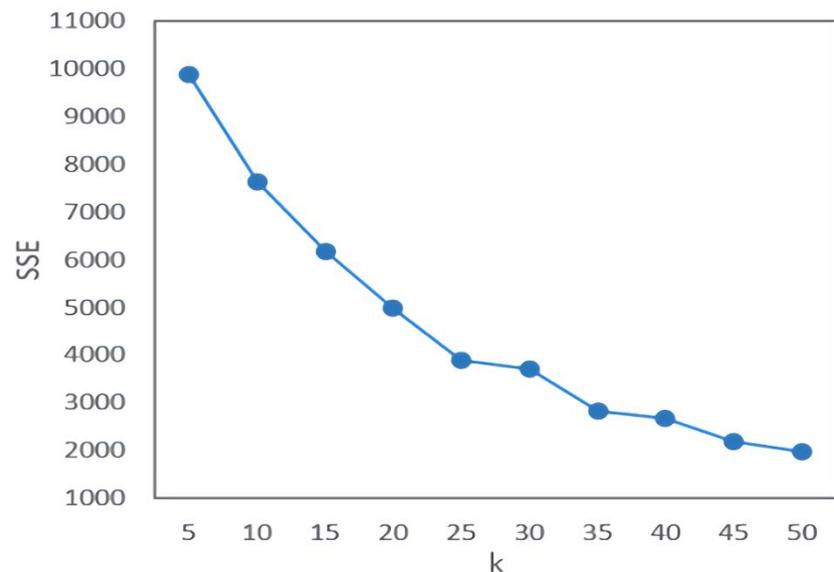
输出: K 个聚簇的集合

- 1) 随机选择样本数据集 D 中的 K 个数据样本 $C = \{C_1, C_2, \dots, C_k\}$ 作为初始聚类中心;
- 2) 重复执行下面的步骤。直到收敛为止
- 3) **for each** x_i **in** D **do**
- 4) 计算各样本 x_i 与均值向量间的距离, 找出与样本 x_i 距离最近的簇 C_r , 即使 x_i 满足 $Distance = \arg \min_k \|x_i - C_k\|^2$;
- 5) 检测将 x_i 划入当前 C_r 簇中是否违背 ML 与 CL 集合中的约束;
- 6) **if ! is_violated then:** $C_r = C_r \cup \{x_i\}$
- 7) **else:** 将当前目标簇舍去, 重新计算 $Distance$, 即找到两个聚类中心 C_r 和 C_l , 使 $\min(\|x_i - C_r\|^2 + \|x_l - C_l\|^2)$
- 8) 循环 5) 至 7) 步, 直到归入某类簇中
- 9) 重新计算 K 个类别中各自的聚类中心, 判断是否产生变化, 若变化, 则更新其均值向量, 使聚类中心 $C_i = (\sum_{j=1}^{n_i} x_j^i) / n_i$ 进行迭代计算, 直至聚类中心不产生变化, 或变化幅度小于所设阈值时停止迭代
- 3) **end**
- 4) **Return** K 个聚簇的集合

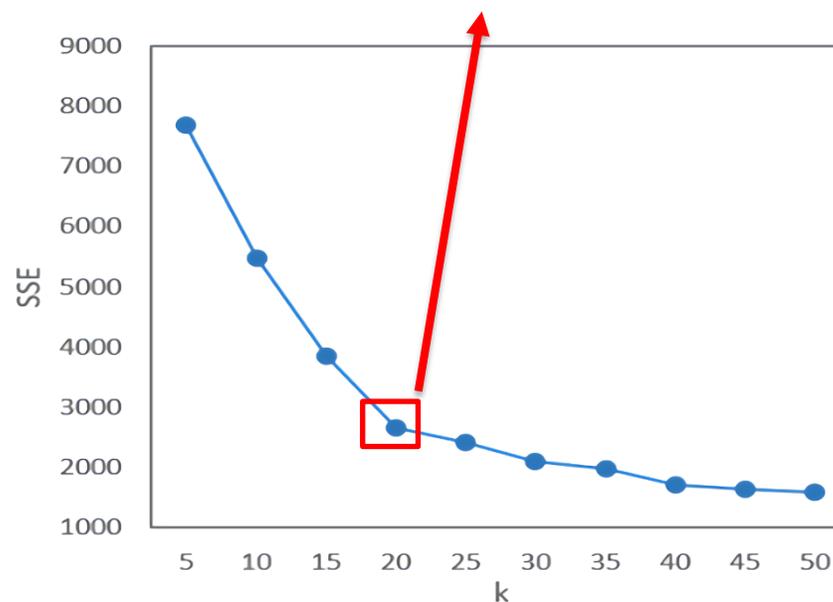
上一步所得约束集

- 聚类中心K的选取
 - 与K-Means对比
 - 成对约束信息更有针对性，因此聚类结果更清晰，样本数据分类更清晰。

K值的增大突然大幅下降的位置作为最优K值



无监督



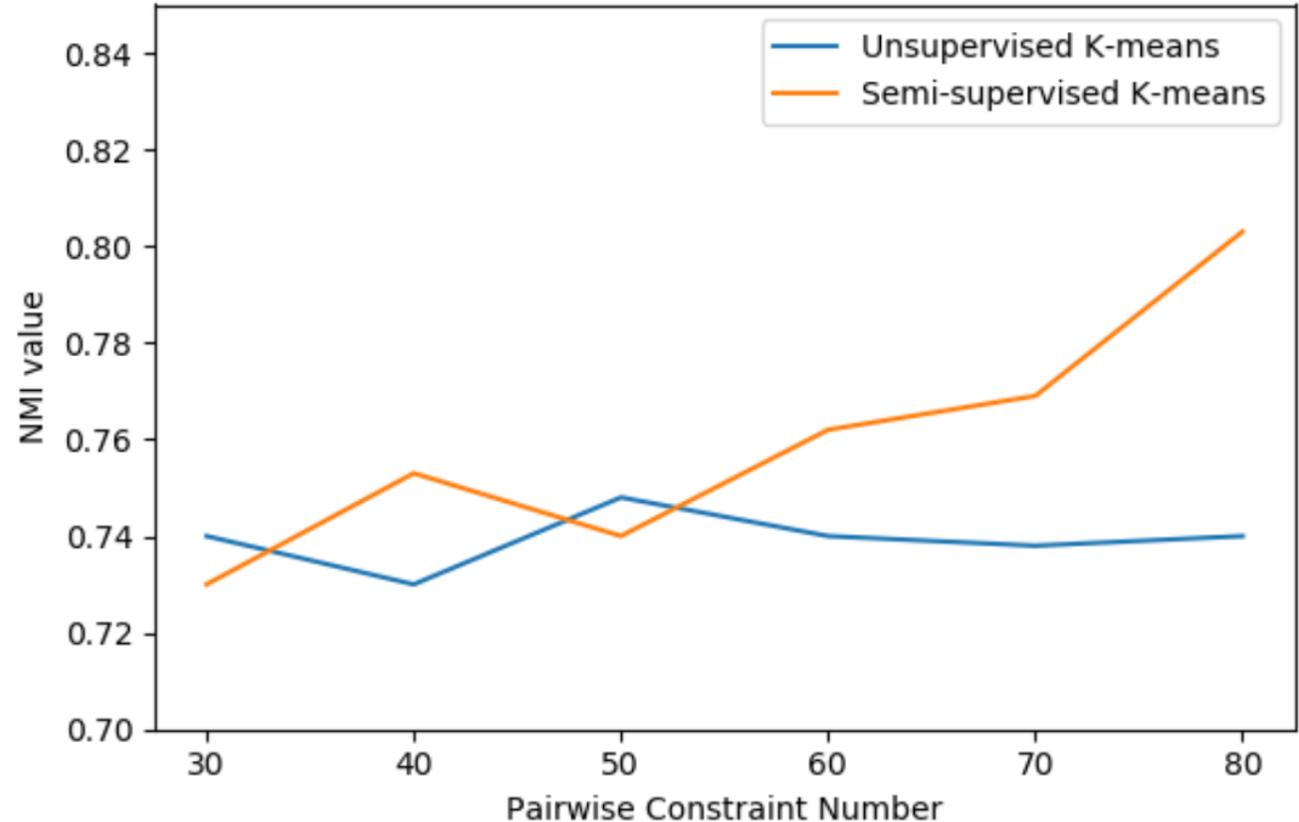
半监督

- 聚类效果

- NMI

- 半监督算法随着患者成对约束信息的增加，聚类效果不断优化。

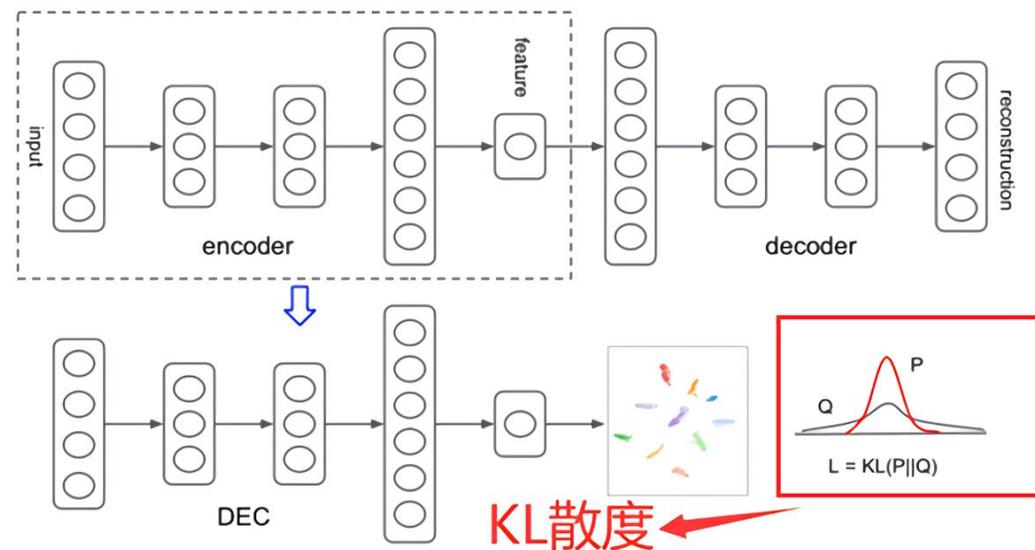
- 由于无监督K-means算法不受成对约束信息的影响，并且初始聚类中心是随机选择的，因此聚类效果略有不同，但总体波动不大。





总结分析

- 深度聚类网络结合
 - 这两种算法都是基于传统的划分聚类方法
 - 深度聚类网络主要贡献是面向高维数据，处理转换为低维特征后再进行聚类
- 老年人运动功能特征挖掘
 - 对老年人运动相关属性提取
 - 依据属性构建约束信息，对老年人群进行划分
 - 探究其簇中隐含的，能反映簇内相似性，簇间相离性的特性



- [1] Zhang J, Chang D. Semi-supervised patient similarity clustering algorithm based on electronic medical records[J]. *IEEE Access*, 2019, 7: 90705-90714.
- [2] Bai L, Liang J, Cao F. Semi-supervised clustering with constraints of different types from multiple information sources[J]. *IEEE Transactions on pattern analysis and machine intelligence*, 2021, 43(9): 3247-3258.
- [3] Huang D, Hu J, Li T, et al. Consistency regularization for deep semi-supervised clustering with pairwise constraints[J]. *International Journal of Machine Learning and Cybernetics*, 2022: 1-14.

大成若缺，其用不弊。
大盈若冲，其用不穷。
大直若屈。大巧若拙。
大辩若讷。静胜躁，寒
胜热。清静为天下正。

谢谢!

