

数据准备

数据挖掘方法: (1)分类: 有监督学习, 面向预定义的类别 (2)聚类: 没有预先定义的类别, 借助相似度度量自动生成 (3)关联规则 (4)异常检测

事务型数据 Transaction data: 一条记录对应一个项目 (Item)集合(无序)

关联规则: 分析事务型数据, 从而根据一部分项的存在记录来判断另一部分项目是否存在于事务中
基本形式: A->B, A、B 均为集合
支持度为 (A+B) 在全体事务中的比重 $s(A>B) = (|A \cup B|) / N$, 置信度 (A+B) 占 A 出现的事务中的比重 $c(A>B) = (|A \cup B|) / |A|$ 。
例: 考虑(Diaper,Milk->Beer)关联规则

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

$s = \frac{(Milk, Diaper, Beer)}{|T|} = \frac{2}{5} = 0.4$
 $c = \frac{(Milk, Diaper, Beer)}{(Milk, Diaper)} = \frac{2}{3} = 0.67$

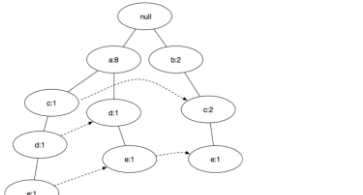
频繁项集: 支持度高于阈值的项集集合 A
频繁项集生成方法 (1)最基本的: 穷举所有可能集合, 计算支持度, 复杂度过高!
(2)Apriori (剪枝思想, 但仍是“生成-测试”范式): 逐步减去所有的非频繁项集, 然后基于频繁项集生成其超集。(先验原理: 如果一个项目集是频繁的, 那么它的所有子集也是频繁的。则非频繁项集的所有超集也是非频繁的。)



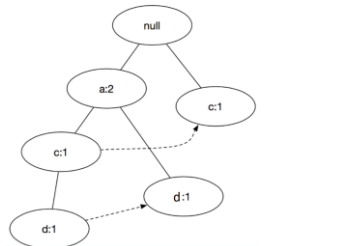
(3)FP-Growth: 本质是输入数据的压缩表示, 通过逐个读入事务, 并将事务映射到 FP 树中的某条路径来构造。
建树: 对各个项按支持度排序, 将排序后的项集逐步读入并建立树状结构, 对相同项节点采用指针连接, 方便快速访问。



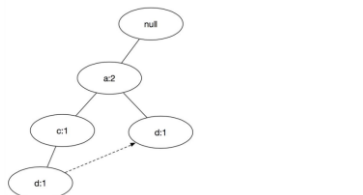
基于 FP-树, 生成频繁项集: FP-Growth 本质上是自底向上的探索。首先查找以 e 为结尾的频繁项集, 其次是 d/c/b, 最后是 a, 包含 e 的路径如下:



判定 e 本身是否为频繁项集 (此处设阈值为 2, 高于阈值), 将 e 的前缀路径转化为条件 FP 树 (需要更新路径上的支持度计数, 只有包含 e 的事务会被需要), 删除那些非频繁的项(例如 b)



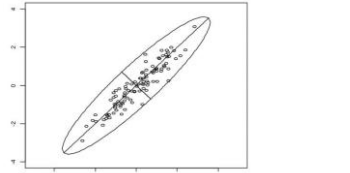
考虑增长结尾的频繁项集的问题, 以{de}为结尾的频繁项集判定为例, 在前一张图上统计与 d 相关的支持度求和为 2, 为频繁项集; 以 de 为结尾, 得到前缀路径如下, 通过其条件 FP 树发现{ade}支持度为 2, 也频繁。



异常检测: 异常数据 ≠ 错误数据, 而是包含不同寻常规律的数据。

(1)基于分布: 前提是识别数据集的具体分布, 错误识别会导致错误检测; 基于一元正态分布的离群点判定: 已知参数的前提下根据正态分布判定离群的概率, 奥卡姆剃刀。
(2)基于度量: 基于距离或密度的方式, 检测空间中远离大多数数据点的离群点。●求数据点到 K 最近邻的平均距离: 高于阈值判定为异常点 ●基于密度: 采用 K 近邻距离的倒数作为密度或给定半径内点的个数
(3)基于聚类 ●抛弃远离其他族的小簇, 但簇的个数将影响结果。●先聚类所有对象, 再评估对象属于簇的程度 (点到簇中心的距离或相对距离)
数据预处理 数据质量问题: (1)数据测量、采集等过程中出现的错误(2)噪声、离群点 (3)缺失数据(-)删除和填补(并重) (4)重复数据(多源数据归并-实体歧义, 多马甲账号)

数据聚合 多源数据归并时的问题: 不同属性名称、单位尺度、属性统计方式、不同数据源的统计不一致性
解决: (1)换算和汇总 (2)实体对齐
数据采样 采样代表性; 启发式采样规模确定方法 (分组采样, 组内数据高度相似, 不同组对象差异大, 每组至少取一个)
数据规约 维度规约: 删除不具有区分度的特征, 可能降低噪声, 避免维度灾难的同时使模型更容易理解, 也可可视化
降维方法 1 主成分分析 PCA: 通过正交变换将一组可能存在相关性的变量转换为一组线性不相关的变量, 转换后的这组变量叫主成分。多维降维: 找出主轴与几个最长的轴作为新维度。选择投影方差最大的轴。



最大特征值对应的特征向量可以最大化投影方差。
完成第一个方向基选择后, 第二个方向基应选择使其不存在相关性, 故优化目标为协方差为 0, 即 $\sum_{i=1}^m a_i b_i = 0$ 。等价于协方差矩阵 $C = \frac{1}{n} X X^T$ 对角化, 即除对角线 (方差) 外的其它元素化为 0。故优化目标转变为寻找矩阵 P, 满足 $P C P^T$ 是一个对角矩阵。
最大 K 个特征值的特征向量对应的线性组合是 K 个新指标 (K 个特征值的比重反映了主成分的信息量, 一般 > 0.85)

输入: n 维样本集 $X = (x_1, x_2, \dots, x_m)^T$, 要降维到 n'
输出: 降维后的样本集 Y

- (1) 对所有样本进行中心化 $x_i = x_i - \frac{1}{n} \sum_{j=1}^n x_j$
- (2) 计算样本的协方差矩阵 $C = \frac{1}{n} X X^T$
- (3) 求出协方差矩阵的特征值及其对应的特征向量
- (4) 将特征向量按特征值大小从上到下按行排列成矩阵, 取前 k 行组成矩阵 P
- (5) $Y = P X$ 即为降维到 k 维后的数据

PCA 特点 (1)依赖原始变量也只能反映原始变量 (2)PCA 内在假设之一是原始变量直接存在一定关联。若原始变量相互独立则降维失败, 数据越相关, 降维效果越好 (3)PCA 的结果未必清晰可解释

降维方法 2 特征子集选择 (而不是归纳新特征) 去除冗余特征和不相关特征, 或为特征赋予不同权重
数据离散 (将连续属性转换为分类属性)
(1)二元化将连续或离散属性转化为多个二元属性 0/1
(2)非监督离散化 不用类别信息 (等宽/等频率/等深)
(3)有监督离散化 (基于熵)
先进行二分, 选择熵最小的点进行分割。对其中具有较大方差 (即纯度不高, 信息混乱) 的部分再下一轮分割。

基于规则的分类 规则分类器的基本形式: 规则 condition (属性组合, 前提) -> 标签
规则分类器的基本原理: 互斥原理 (每条记录至多被一条规则所覆盖) & 穷举原理 (每条记录至少被一条规则所覆盖)
规则分类器的有序性 基于规则的排序方案: 按照规则的质量 (如准确性) 进行排序; 基于类的排序方案: 同类规则排在一起, 相对顺序被忽略。
如何制定规则分类器:

(1)直接从数据中自动学习规则: RIPPER, CN2
●顺序覆盖: 贪心
算法开始时, 决策表为空, 即不包含任何规则
每一步针对某个类 y, 提取覆盖当前训练集的最佳规则
●什么是最好的规则? 覆盖的样本尽可能多, 同时样本类别尽可能一致
如果规则覆盖大多数正例 (即 y), 而没有或覆盖极少负例 (非 y), 则保留
●将该规则加入决策表的末端, 同时删除该规则覆盖的所有训练样本
重复上述过程, 直至满足终止条件 (例如: 某个增益的阈值, 如准确/准确率)
规则增长策略:
“从一般到特殊”: 初始规则条件为空, 给定目标标签 y, 逐步加入合取项 (AND 相连) 来提高规则质量
“从特殊到一般”: 随机选择一个正样本作为初始种子, 逐步删除规则中的合取项, 来覆盖更多的同类别正例
(2)间接借助其他分类模型学习规则: 决策树
基于监督学习的分类
决策树: 从根节点到叶节点的一条路径都对应一类分类规则。
特征选择 (有较强区分能力的特征) -> 生成决策树 -> 决策树剪枝, 避免过拟合

信息熵 $Ent(D) = -\sum_{k=1}^{|I|} p_k \log_2 p_k$
信息增益 $Gain(D, a) = Ent(D) - Ent(D|a)$
 $= Ent(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} Ent(D^v)$
 $Gain_ratio(D, a) = \frac{Gain(D, a)}{IV(a)}$
 $IV(a) = -\sum_{v=1}^V \frac{|D^v|}{|D|} \log_2 \frac{|D^v|}{|D|}$

●特征选择准则 (1)信息增益 Gain (偏好取值较多的属性) (2)信息增益率 Gain_ratio (引入惩罚项, 但倾向于选择取值较少的属性) 先从候选特征中找到信息增益高于平均水平的集合, 再从中找到信息增益率最大特征 (3)基尼指数 (一个随机样本被分错的概率, Gini 指数低则信息纯度低)
K 个类, 样本点属于第 k 类的概率为 p_k :
 $Gini(p) = \sum_{k=1}^K p_k(1-p_k) = 1 - \sum_{k=1}^K p_k^2$
 $Gini(D, A) = \sum_{i=1}^{|D|} Gini(D_i)$
例如: $Gini = \frac{1}{2} \times (1 - \frac{1}{2}) + \frac{1}{2} \times (\frac{1}{2} - \frac{1}{2}) = 0.375$

●决策树生成 计算当前节点各个属性的信息增益率, 基于最大信息增益率属性, 迭代式对节点进行分类, 直到每个节点上的样本类别统一
●剪枝 预剪枝: 每个节点划分前衡量当前节点的划分能否提高决策树的泛化能力 (验证集上精度)
后剪枝: 自底向上考虑每个非叶子节点被替换成叶子节点后能否提高泛化性能
最近邻分类 计算与未知样本与其他样本的距离找到 K-最近邻, 基于 K-最近邻的类别确定分类结果。动机: 表征空间上相似的文档是相似的; 基于实例的学习, 不需要对数据进行抽象 (如提取特征); 消极学习, 不需要模型, 但开销大; 基于局部信息判断, 受噪声影响大。
距离度量: 欧式距离, 汉明距离 (0/1 向量 统计多少维数字不同), 余弦相似度, 马氏距离, 无穷范数 (分量最大值) K 小则易受噪声干扰, K 大则错误涵盖其他类别样本。
支持向量机 SVM: 二分类转化为寻找最大间隔超平面, 实现对高维空间中节点进行有效分割, 使得超平面和支持向量 (离该超平面最近的点) 的间隔最大化

求解对偶问题即求极值时: (1)内部函数均为 0, 解方程组, 检查函数是否非负 (2)边界上也可取极值 (比如单一分量为 0) 核函数: 将高维空间下的内积运算转化为低维空间下的核函数计算。
软间隔: 允许少数样本不满足超平面约束 (引入惩罚项 C)

原问题 $\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i(w^T \phi(x_i) + b))$
对偶问题 $\max_{\alpha} \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) - \sum_{i=1}^m \alpha_i$
s.t. $\sum_{i=1}^m \alpha_i y_i = 0, 0 \leq \alpha_i \leq \frac{1}{C}, i = 1, 2, \dots, m$

不平衡分类问题解决方案: (1)代价敏感学习: 代价矩阵衡量将一个类别分到另一个类的代价, 优化目标由准确/召回变为加权后的代价 (2)抽样: 少数类过采样, 多数类欠采样, K-最近邻用少数类样本生成新样本

聚类
基本问题: 聚类依据? 相似性度量? 簇的数量?
方法: 层次聚类/划分聚类
K 均值聚类 (K-means)
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 $\{\mu_1, \mu_2, \dots, \mu_k\}$
2: repeat
3: $C_i = \emptyset (1 \leq i \leq k)$
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: $d_{ij} = \|x_j - \mu_i\|_2$
6: 根据距离最近的均值向量确定 x_j 的簇标: $\lambda_j = \arg \min_{i \in \{1, 2, \dots, k\}} d_{ij}$
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 μ_i 更新为 μ'_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: 簇标 $C = \{C_1, C_2, \dots, C_k\}$
 $SSE = \sum_{i=1}^k \sum_{x \in C_i} \text{dis}^2(x, \mu_i)$ (mi 为簇 Ci 的中心)

初始中心的选择 (1)K 均值的后处理: 清除较小的可能代表离群点的簇, 对 SSE 高的簇拆分, 对 SSE 低的合并 (2)二分 K 均值聚类: 先分 2 个簇, 再不断选择一个分裂空簇的处理: 最大 SSE 的簇拆分, 或选择一个最远样本点新生成一个簇。
缺点: 易受离群点干扰, 当簇存在不同规模、密度及不规则形状的情况下, K 均值聚类效果较差。

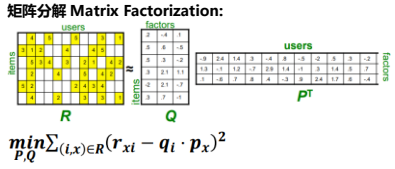
层次聚类
优点: 不需预设簇的数量, 结果有意义对应到分类学目录上; 缺点: 局部最优, 每步合并决策都是最终的
凝聚式聚类: (1)合并邻近度最高的两个簇 (2)基于更新的簇重新计算邻近度, 更新邻近度矩阵
邻近度 (1)单链 MIN: 不同簇最近的点之间的邻近度 (擅长处理非椭圆形状的区域, 但对噪声比较敏感) (2)全链 MAX: 不同簇最远的点之间的邻近度 (对噪声不太敏感, 但可能使得较大的簇变得支离破碎) (3)组平均: 所有来自不同簇的两点之间的平均邻近度 (4)中心距离: 两个簇中心之间的邻近度 (5)沃德法: 合并后簇中各点到新中心的距离平方和。
分裂式聚类: (1)二分 K 均值聚类 (2)最小生成树聚类: 由差异矩阵生成一颗最小生成树 (节点之间权重最小), 每步断开差异最大的一条边, 从而创建一个新的簇。
MST 聚类结果与单链凝聚聚类的结果相同。

基于密度聚类
密度: 样本一定半径的样本数量
DBSCAN 算法 (1)将所有点区分为核心点 (半径内样本数超过阈值的点), 边界点 (非核心点但处于稠密区域边界内/上) 和噪声点 (处于稀疏区域的点) (2)删除噪声点 (3)将所有边界点在预定半径内的核心点之间连一条边 (4)连通的点形成核心簇 (5)将所有的边界点指派到一个与之关联的核心点所在的簇中
优点: 对噪声鲁棒; 缺点: 密度变化大的簇受影响
模糊聚类 (计算归属度为计算从属度)
 $SSE = \sum_{j=1}^k \sum_{i=1}^N w_{ij} (x_i - c_j)^2, \sum_{j=1}^k w_{ij} = 1$
从属度由 $w_{ij} \in [0, 1]$ 改为 $w_{ij} \in [0, 1]$ 。

聚类问题的评估 非监督评估: (1)基于邻近度矩阵 (理想结果为“簇内点邻近度全为 1, 簇间为 0”, 呈对角模式); (2)凝聚度 (簇内两点的邻近度之和 / 簇内各点到簇中心的邻近度之和) 与分离度 (簇间两点的邻近度之和 / 簇中心到其他簇中心的邻近度之和) 有监督评估: (1)分类度量: 熵、纯度 (簇在多大程度上包含单个对象, 以最多数量的比例计算)、准确率、召回率、F 值 (2)相似性度量: 分类对应的矩阵, 同一类的样本对应的元素为 1, 不同类为 0, 比较两矩阵相关性
推荐系统
推荐评估: (1)评分均方根误差 RMSE $\sqrt{\sum_{i=1}^n (r_{xi} - \hat{r}_{xi})^2}$
(2)分类 (推荐正确 / 错误: Precision/Recall/F; TopN: Pre@N, Rec@N) (3)排序评估
基于内容的推荐: 物品画像 & 用户画像 (基于用户评分进行加权) 采用余弦相似性度量进行评分
优点: 用户推荐过程独立, 个性化, 推荐结果可解释
缺点: 难以提取物品特征, 难建立新用户画像
多样化评估: 最大边界相关性 MMR
基于路径推荐: 知识图谱指向量化画像, 基于图谱上的游走实现推荐, 路径可作为推荐的依据。
偏见 Bias 问题: 位置 (第一更受关注), 模态 (与众不同的模态吸引关注), 关键效应 (标题党的威力)
双向选择问题: 用户 <-> 物品, 稳定匹配
协同过滤 本质: 矩阵补全
基于内存 (Memory-based): User-based: (预测时忽略空值)

$sim(a, b) = \frac{\sum_{i \in \text{product}(P)} (r_{ai} - \bar{r}_a)(r_{bi} - \bar{r}_b)}{\sqrt{\sum_{i \in \text{product}(P)} (r_{ai} - \bar{r}_a)^2} \sqrt{\sum_{i \in \text{product}(P)} (r_{bi} - \bar{r}_b)^2}}$
 $pred(a, p) = \bar{r}_a + \frac{\sum_{i \in \text{Neighbors}(a)} sim(a, b_i) \cdot (r_{bi} - \bar{r}_b)}{\sum_{i \in \text{Neighbors}(a)} sim(a, b_i)}$
Item-based $r_{ix} = \frac{\sum_{j \in N(i, x)} s_{ij} \cdot r_{jx}}{\sum_{j \in N(i, x)} s_{ij}}$ (无需平均修正)

相似度计算时, 如果未评分, 直接设为 0, 不用减平均数。
基于物品的推荐效果更好? 物品属性单一, 用户偏好多样。
基于用户优点: 不受多模态非结构化信息表征与特征选取困扰; 缺点: 冷启动、稀疏性、热度偏差。
冷启动问题: (1)先非个性化推荐 (2)借助个人信息或其它网站浏览信息 (3)诱导式推荐迭代收集用户反馈 (4)基于内容混合 (5)Side information: 众包文本、知识图谱
基于模型 (Model-based): 潜在因子
矩阵分解 Matrix Factorization:



$\min_{P, Q} \sum_{(i, x) \in R} (r_{xi} - q_i \cdot p_x)^2$
解决过高的 K 带来的过拟合? 引入正则项, 避免过大参数值。
 $\min_{P, Q} \sum_{(i, x) \in \text{training}} (r_{xi} - q_i \cdot p_x)^2 + \left[\lambda_1 \sum_i \|p_i\|^2 + \lambda_2 \sum_i \|q_i\|^2 \right]$
“error” “length”

(1)非负矩阵分解 (2)概率矩阵分解 (3)社交约束
社会网络
基本元素: 节点 (网络中实体), 有向/无向边 (关系), 邻居/出入度, 连通性/连通组件。
节点角色: 意见领袖、结构洞 (作用: 为组织引入外部信息, 衡量方式: 聚集系数 (任意两好友也互为好友的概率) 低)
链接预测 三元组包
(1)两节点间存在边的概率, 与它们共同好友的个数成正比 (如何削弱好友个数影响? 对指标进行正则化)
(2)共同好友作为中介的引导力: 还要考虑好友的好友数
(3)考虑“共同好友”的好友: 基于多条关系的链接预测
 $s_{xy}^{katz} = \sum_{i=1}^{\infty} \beta^i \cdot \text{paths}_{xy}^{<i>} = \beta A_{xy} + \beta^2 (A^2)_{xy} + \beta^3 (A^3)_{xy} + \dots$

社团挖掘:
边介数: 网络中任意两点的最短路径, 有多少条会通过该边
(1)基于层次聚类
Girvan-Newman 算法: 计算网络中所有边的边介数; 去除边介数最高的边; 重新计算去除边后网络中所有边的边介数; 跳至步骤二, 直至网络中没有边存在
(2)基于划分聚类: K-means, 谱聚类 (基于邻域, 面向图)
关系抽取方法 (1)基于规则: 纯手工定值规则, 文本匹配 “模板” 和 “槽”; 难制定规则, 领域专家构筑大规模知识库, 领域移植困难。
(2)基于模式: 从种子关系中获得模式, 再由模式寻找更多种子, 迭代优化。 (基于字面匹配, 移植性差, 适合特定关系)

双重语义模式关系抽取 DIPRE:
模式的 Order 和 Middle, 即为 Occurrence 集合的 Order 和 Middle
模式的 URL Prefix、Prefix、Suffix, 分别为 Occurrence 集合中最长的公共 URL 前缀与后缀。
●其他部分采用通配符填充
Snowball: 仅信任支持度 (满足每个模式的元组数量) 和置信度高 (符合该模式的元组确实符合相应关系的模式)
(3)基于机器学习: 转化为分类问题, 训练模型求解
开放关系抽取 (1)基于知识图谱 (2)基于句法
远程链接 如果某个实体对之间具有某种关系, 那么所有包含这个实体对的句子都是用于描述这种关系; 目的: 获取足够数量的、高质量的标注
语义迁移: 不是所有包含该实体的句子都表达该关系且错误不断放大 优化: 动态转移矩阵描述各类相互标错的概率; 规则学习设计相应否定模式列表去除错误的标签; 注意力机制
事件抽取的模板: 选定模板模板后, 通过事件元素 (事件参与者) 与事件元素角色 (事件元素在事件中充当的角色) 的识别, 将相应的元素填入模板合适的槽 (描述命名实体基本信息, 内容可包括名称/类别/种类) 内