

日期: /

数据

1. 数据

1. 数据集: 一组含有特征的数据对象的集合

2. 样本: 由一组特征所描述的一个对象

3. 特征: 描述样本或对象在某方面的表现或性质的事项 (对模型有用的部)

		特征 (Features)			
		Refund	Marital Status	Taxable Income	Cheat
样本 (Samples)	Tid	1 Yes	Single	125K	No
		2 No	Married	100K	No
		3 No	Single	70K	No
		4 Yes	Married	120K	No
		5 No	Divorced	95K	Yes
		6 No	Married	60K	No
		7 Yes	Divorced	220K	No
		8 No	Single	85K	Yes
		9 No	Married	75K	No
		10 No	Single	90K	Yes

例: 个人身高, 北京某时刻气温
属性因对象而异
或随时间变化

4. 数据的数学表示

$D = \{x_1, x_2, \dots, x_m\}$ 表示 m 个样本的数据集
每个样本由 d 个属性描述; 样本 $x_i = (x_{i1}, x_{i2}, \dots, x_{in})$ 为 d 维样本空间 X 中的一个向量, $x_{ij} \in X$
 x_{ij} 为 x_i 在第 j 个属性上的取值; (x_i, y_i) 表示第 i 个样例, y_i 是样本 x_i 的标记

日期: /

如下表, $m=17$, $x_1 = (\text{青绿}; \text{蜷缩}; \text{浊响}; \text{清晰}; \text{凹陷}; \text{硬滑})$

$x_{21} = \text{“乌黑”}$ 是 x_2 在第1个属性上的取值

描述西瓜的6个属性

编号	色泽	根蒂	敲声	纹理	脐部	触感	好瓜	标记 (label)
1	青绿	蜷缩	浊响	清晰	凹陷	硬滑	是	
2	乌黑	蜷缩	沉闷	清晰	凹陷	硬滑	是	
3	乌黑	蜷缩	浊响	清晰	凹陷	硬滑	是	
6	青绿	稍蜷	浊响	清晰	稍凹	软粘	是	
7	乌黑	稍蜷	浊响	稍糊	稍凹	软粘	是	
10	青绿	硬挺	清脆	清晰	平坦	软粘	否	
14	浅白	稍蜷	沉闷	稍糊	凹陷	硬滑	否	
15	乌黑	稍蜷	浊响	清晰	稍凹	软粘	否	
16	浅白	蜷缩	浊响	模糊	平坦	硬滑	否	
17	青绿	蜷缩	沉闷	稍糊	稍凹	硬滑	否	

二、特征类型

1. 特征取值

① 定义: 给特征指派一个特征的数字或符号

② 相同的特征可以有不同取值

不同的特征可以有相同取值

• 样本-->姓名: 张三、年龄: 18、性别: 男

• 特征: {姓名、年龄、性别}, 取值: {张三、18、男}

2. 特征类型: 定性, 定量

① 定性属性: 不具有数的大部分性质

② 定量属性: 用数表示, 具有数的大部分性质
可以是连续或整数

日期: /

属性类型	描述	例子	操作
分类的 (定性的)	标称 标称属性的值只是不同的名字, 即标称值只提供足够的信息以区分对象(=, ≠)	邮政编码、雇员 ID 号、眼球颜色、性别	众数、熵、列联相关、 χ^2 检验
	序数 序数属性的值提供足够的信息确定对象的序(<, >)	矿石硬度(好, 较好, 最好)、成绩、街道号码	中值、百分位、秩相关、游程检验、符号检验
数值的 (定量的)	区间 对于区间属性, 值之间的差是有意义的, 即存在测量单位(+, -)	日历日期、摄氏或华氏温度	均值、标准差、皮尔逊相关、 t 和 F 检验
	比率 对于比率变量, 差和比率都是有意义的(*, /)	绝对温度、货币量、计数、年龄、质量、长度、电流	几何平均、调和平均、百分比变化

③ 标称类型 Nominal

- 只能区分样本之间的不同, 例如: 学号、籍贯、邮政编码

序数类型 Ordinal

- 能够对样本之间的顺序进行区分, 例如: 排名、年级、衣服的号码(S, M, L, XL, XXL)

区间类型 Interval

- 能够对样本在坐标系上的相对距离进行度量, 例如: 日历上的日期、摄氏或华氏温度等

比率类型 Ratio

- 能够对样本在坐标系上的绝对位置进行标定, 例如: 开尔文温度、长度、时间、质量、货币单位等

④ 本质区别: 对应的操作不同

{ 相异性: =, ≠ 加法: +, -
序: <, > 乘法: *, /

标称: 相异性 (=, ≠)

序数: 相异性 序 (=, ≠, <, >)

日期:

值之间若有意义。

区间 相异性, 序, 加法 ($=, \neq, <, >, +, -$)

比率 相异性, 序, 加法, 乘法

值之间差与比率均有意义。

• 摄氏温度?

- 冰水混合物的温度定为0摄氏度, 沸水的温度定为100摄氏度

• 华氏温度?

- 氯化铵和水的混合物的冰点温度 (即氨水结冰的温度) 为温度计的零度, 人体温度为温度计的100度。(华氏温度将水的冰点定为32度, 沸点定为212度)

• 开尔文温度 (绝对温度)?

- 以绝对零度作为计算起点的温度, 即将水三相点的温度准确定义为273.15K。
- 20K是10K的2倍

摄氏温度或华氏温度的零度是硬性规定的, 其比率是无物理意义的

⑤ 区间类型和比率类型的特征

• **区间类型 (interval-scaled) 属性特点:**

- 例: 温度属性, 一般表示: $10^{\circ}\text{C} \sim 15^{\circ}\text{C}$ 。
- 1. 用相等的单位尺度度量, 区间属性的值有序, 可以为正、0、负。(值的秩评定)
- 2. 允许比较与定量评估值之间的差。
- 3. 区间标度属性是数值的, 中心趋势度量中位数和众数, 还可以计算均值。

• **比率类型 (ratio-scaled) 属性特点:**

- 1. 具有固有零点的数值属性。(也就是该种属性中会有固有的为0的值)
- 2. 一个值是另一个的倍数 (或比率)。
- 3. 值是有序的。(可以计算差、均值、中位数、众数)
- 例: 度量重量、高度、速度和货币量 (例如 100 元是 1 元的 100 倍) 的属性。

⑥ 坐标变换:

日期:

/

属性类型		变 换	注 释
分类属性 (定性属性)	标称	任何一对一变换, 例如值的一个排列	邮政如果所有雇员的ID号都重新赋值, 不会导致任何不同
	序数	值的保序变换, 即 $new_value = f(old_value)$ 其中 f 是单调函数	包括好, 较好, 最好的属性可以完全等价地用值{1, 2, 3}或{0.5, 1, 10}表示
数值属性 (定量属性)	区间	$new_value = a * old_value + b$ 其中 a, b 是常数	华氏和摄氏温度标度零度的位置和1度的大小(单位)不同
	比率	$new_value = a * old_value$	长度可以用米或英尺度量

3. 特征类型: 离散与连续 (值的个数).

① 离散 Discrete Feature.

{ 具有有限或无限种可能值. (邮政编码...)
常常为整数变量.
可以为分类的, 也可以是数值的.

② 连续 Continuous Feature.

{ 取实数值的特征. (温度, 高度...)
有限精度测量与表示.
浮点变量表示.

③ 二值特征 Binary Feature.

{ 仅有2个值的特征. (0或1) (性别, 对错)
是离散的特殊情形

④ 非对称特征 Asymmetric Binary Feature.

日期:

/

{ 出现非0值才重要的特征, 状态结果不重要
非0才1, = 0 特征: 只有非0值重要的二元特征.
[体检: 阳性(1) 阴性(0)]

三. 样本的相似性和相异性

1. 相似度和相异度基本概念

① 相似度: similarity

{ 两个样本相似程度的数值化度量
2个样本越相似, 相似性越高
非负, 取值通常在 $[0, 1]$

② 相异度: dissimilarity

{ 两个样本相异程度的数值化度量
2个样本越相似 相异性越低.
非负, 取值通常在 $[0, 1]$ 和 $[0, \infty)$

通常将距离作为相异度的同义词

2. 简单属性的相似度和相异度

属性类型	相异度	相似度
标称的	$d = \begin{cases} 0 & \text{如果 } x=y \\ 1 & \text{如果 } x \neq y \end{cases}$	$s = \begin{cases} 1 & \text{如果 } x=y \\ 0 & \text{如果 } x \neq y \end{cases}$
序数的	$d = \frac{ x-y }{(n-1)}$ <p>值映射到整数 0 到 $n-1$, 其中 n 是值的个数</p>	$s = 1 - d$
区间或比率的	$d = x - y $	$s = -d, s = \frac{1}{1+d}, s = e^{-d},$ $s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

日期: /

3. 距离计算

① 欧几里德距离

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

n 为对象数据维度 (属性的个数), x_k 和 y_k 分别是数据对象 x 和 y 的第 k 个属性.

② 闵可夫斯基距离

欧氏距离是闵可夫斯基的特例

$$d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{\frac{1}{r}}$$

$\left\{ \begin{array}{l} r=1, \text{ 曼哈顿距离, } L_1 \text{ 范数 } L_1\text{-norm} \\ r=2, \text{ 欧氏距离, } L_2 \text{ 范数, } L_2\text{-norm} \\ r=\infty, \text{ 上确界距离, } L_\infty \text{ 范数 } L_\infty\text{-norm} \end{array} \right.$

$$d(x, y) = \sum_{k=1}^n |x_k - y_k|$$

$$d(x, y) = \sqrt{\sum_{k=1}^n |x_k - y_k|^2}$$

$$d(x, y) = \|x - y\|_\infty = \lim_{r \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{\frac{1}{r}}$$

4. 距离特性

① 非负性: 对于任意 p 和 q , 存在 $d(p, q) \geq 0$;
当且仅当 $p = q$ 时 $d(p, q) = 0$

日期: /

② 对称性: 对于任意 p 和 q , $d(p, q) = d(q, p)$

③ 三角不等式: 对于任意 p, q 和 r ,
 $d(p, r) \leq d(p, q) + d(q, r)$

满足以上3条列为一种度量(metric)

5. 相似性

① 若 $s(x, y)$ 为数据点 x 和 y 之间的相似度

$\begin{cases} \text{仅当 } x=y \text{ 时, } s(x, y)=1 & (0 \leq s \leq 1) \text{ [非负]} \end{cases}$

$\begin{cases} \text{对于所有 } x \text{ 和 } y, s(x, y) = s(y, x) & \text{[对称性]} \end{cases}$

② 二元数据的相似度度量

对2个二元向量 x 和 $y \in \mathbb{R}^n$, 值取0或1

f_{ij} 表示 $x=i, y=j$ 时的属性数目. ($i, j=0, 1$)

11, SMC: 简单匹配系数

$$SMC = \frac{\text{值匹配的属性个数}}{\text{属性个数}}$$

$$= \frac{f_{11} + f_{00}}{f_{01} + f_{00} + f_{10} + f_{11}}$$

$\begin{cases} \text{通常 } SMC \in [0, 1] & 0 \text{ 表示不相似, } 1 \text{ 表示相似} \end{cases}$

适用: 对称二值离散型属性

0和1同样重要 (例如性别)

日期: /

2, Jaccard 系数.

$$J = \frac{\text{匹配个数}}{\text{匹配中涉及的属性个数}}$$

- 如果每个非对称的二元属性对应于商店的一种商品，则1表示该商品被购买，0表示该商品未被购买。由于未被顾客购买的商品数远大于被购买的商品数，因此，简单匹配系数(SMC)会判定所有的事务都是类似的。Jaccard值越大说明相似度越高。
- Jaccard系数来处理仅包含非对称的二值离散型属性。不对称的二值离散型属性是指属性取值为1或者0不是同等重要，例如：是否是癌症的结果，因此通常用1来表示阳性结果，而用0来表示阴性结果。

32

$$x = (1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0)$$

$$y = (0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1)$$

$$\begin{aligned} \text{SMC} &= (F_{11} + F_{00}) / (F_{01} + F_{10} + F_{11} + F_{00}) \\ &= (0+7) / (2+1+0+7) = 0.7 \end{aligned}$$

$$J = (F_{11}) / (F_{01} + F_{10} + F_{11}) = 0 / (2 + 1 + 0) = 0$$

③ 余弦相似度 Cosine Similarity

x, y 为 文档向量

$$\cos(x, y) = \frac{\langle x, y \rangle}{\|x\| \|y\|} = \frac{\langle x', y' \rangle}{\|x'\| \|y'\|}$$

'表示能量, $\langle x, y \rangle$ 表示向量内积,

日期: /

$$\langle x, y \rangle = \sum_{k=1}^n x_k y_k = x^T y$$

且 $\|x\|$ 是向量 x 的长度, $\|x\| = \sqrt{\sum_{k=1}^n x_k^2} = \sqrt{\langle x, x \rangle}$

适用: 非对称、非归一化, 可以处理非二元向量。

欧氏距离衡量空间点的直线距离, 余弦距离衡量点在空间的方向差异

• Example:

$$d_1 = 3205000200$$

$$d_2 = 1000000102$$

$$d_1 \cdot d_2 = 3 \cdot 1 + 2 \cdot 0 + 0 \cdot 0 + 5 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 + 2 \cdot 1 + 0 \cdot 0 + 0 \cdot 2 = 5$$

$$\|d_1\| = (3^2 + 2^2 + 0^2 + 5^2 + 0^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2)^{0.5} = (42)^{0.5} = 6.481$$

$$\|d_2\| = (1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150$$

6. 相关性

① 应用: 用于测量两个变量之间或两个对象之间的关系。

② 若2个数据对象中的值来自不同属性, 使用相关性来度量属性之间相似度

Features 相关

相似/相异

Samples

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

• 皮尔森相关系数 Pearson's Correlation

- 度量两个变量之间的线性相关性

$$\text{corr}(\mathbf{x}, \mathbf{y}) = \frac{\text{covariance}(\mathbf{x}, \mathbf{y})}{\text{standard_deviation}(\mathbf{x}) \times \text{standard_deviation}(\mathbf{y})} = \frac{s_{xy}}{s_x s_y}$$

协方差与标准差之比

$$\text{covariance}(\mathbf{x}, \mathbf{y}) = s_{xy} = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$$

$$\text{standard_deviation}(\mathbf{x}) = s_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2}$$

$$\text{standard_deviation}(\mathbf{y}) = s_y = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2}$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k \text{ 是 } \mathbf{x} \text{ 的均值} \quad \bar{y} = \frac{1}{n} \sum_{k=1}^n y_k \text{ 是 } \mathbf{y} \text{ 的均值}$$

- 对于非线性的相关性难以建模

$$\bullet X = (-3, -2, -1, 0, 1, 2, 3)$$

$$\bullet Y = (9, 4, 1, 0, 1, 4, 9)$$

$$Y = X^2$$

负相关

正相关

- Mean(X) = 0, Mean(Y) = 4

- 皮尔森相关系数 Correlation

$$= (-3)(5) + (-2)(0) + (-1)(-3) + (0)(-4) + (1)(-3) + (2)(0) + 3(5)$$

$$= -15 + 0 + 3 + 0 - 3 + 0 + 15$$

$$= 0 \text{ (即相关度为0)}$$