

一、相似度：

- 两个样本相似程度的数值化度量
- 两个样本越相似，他们之间的相似性就越高
- 相似度是非负的，通常取值范围在[0, 1]

如果 $s(x, y)$ 是数据点 x 和 y 之间的相似度

- 1) 仅当 $x=y$ 时 $s(x, y)=1$ 。($0 \leq s \leq 1$) (非负性)
- 2) 对于所有 x 和 y , $s(x, y)=s(y, x)$ 。(对称性)

简单匹配系数 SMC: 如果样本的属性都是 对称的二值离散型属性，则样本间的距离可用简单匹配系数计算。 对称的二值离散型属性是指属性取值为 1 或者 0 同等重要

• **简单匹配系数 (Simple Matching Coefficient, SMC)** 定义如下：

$$\begin{aligned} \text{SMC} &= \text{值匹配的属性个数} / \text{属性个数} \\ &= (f_{11} + f_{00}) / (f_{01} + f_{10} + f_{11} + f_{00}) \end{aligned}$$

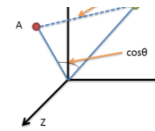
Jaccard 系数: 应用于非对称二元属性; 不对称的二值离散型属性是指属性取值为 1 或者 0 不是同等重要

$$\begin{aligned} J &= \frac{\text{匹配个数}}{\text{匹配中不涉及的属性个数}} \\ &= \frac{f_{11}}{f_{01} + f_{10} + f_{11}} \end{aligned}$$

余弦相似度: $\cos\langle x, y \rangle$: 衡量点在空间的方向差异, 适应于非对称属性, 还可以处理非二元向量

如果 x 和 y 是两个文档向量, 则

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



其中'表示向量或者矩阵的转置, $\langle x, y \rangle$ 表示两个向量的内积:

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

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

二、相异度: (距离)

- 两个样本之间差异程度的数值化度量
- 两个样本越相似，他们之间的相异性越低
- 相异度是非负的，取值在[0,1]和[0, ∞)均有

• 距离（如欧几里得距离）满足以下三个特性：

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 $d(p, q) = 0$ 。
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 ， $d(p, r) \leq d(p, q) + d(q, r)$ 。

• 其中 $d(p, q)$ 是 p 和 q 之间的距离。

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	$d = \frac{ x-y }{(n-1)}$ 值映射到整数 0 到 $n-1$ ，其中 n 是值的个数	$s = 1 - d$
区间或比率的	$d = x - y $	$s = -d, s = \frac{1}{1+d}, s = e^{-d},$ $s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

• 欧几里德距离

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

• 闵可夫斯基距离

• 闵氏距离的欧式距离的一种**泛化**，欧式距离是闵氏距离的一种特例

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



• $r=1$, 曼哈顿距离, **L1范数, L1-norm**

$$d(x, y) = \|x - y\|_1 = \sum_{k=1}^n |x_k - y_k|$$

• $r=2$, 欧氏距离, **L2范数, L2-norm**

$$d(x, y) = \|x - y\|_2 = \sqrt{\sum_{k=1}^n |x_k - y_k|^2}$$

• $r=\infty$, 上确界距离, **∞ 范数, L-norm**

• 对象各个属性之间的最大距离，即上确界。

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

三、相关性：测量两个变量(高度和重量)之间或两个对象之间的关系，若两个数据对象中的值来自不同的属性，使用相关性来度量属性之间的相似度

• 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

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

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

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

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

29

- 对于非线性相关性难以建模

$$\begin{aligned} \bullet X &= (-3, -2, -1, 0, 1, 2, 3) \\ \bullet Y &= (9, 4, 1, 0, 1, 4, 9) \end{aligned} \quad Y = X^2$$

负相关

正相关

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

- 皮尔森相关系数 Correlation

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

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

$$= 0 \text{ (即相关度为0)}$$