

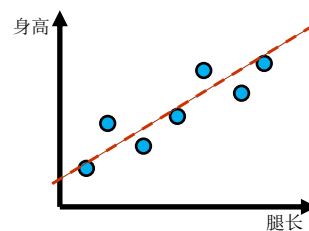
大数据分析

Statistics
Correlations
刘盛华

相关性分析 ——什么是相关

问题： 如果在一座古墓中发现一根腿骨, 那么通过这根腿骨的长度可以判断墓主的身高吗?

条件： 提供 348个成年男子的身高及腿长数据。



回答： 观察发现, 身高会随着腿长的变化而变化, 利用身高和腿长的关系可以解决此问题。

思考： 上述提到的关系就叫做相关关系

相关性分析 ——什么是相关



一个变量随着另一个变量变化。

一个变量随着另一个变量增加而增加。

一个变量随着另一个变量增加而减少。

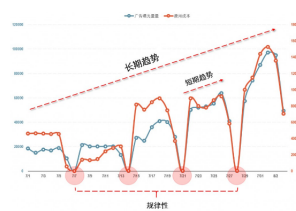
相关关系

正相关

负相关

相关性分析 ——相关的刻画方法

折线图



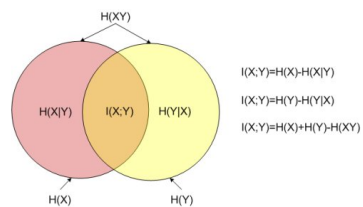
散点图



矩阵

10	-27	-2	-7	42	9	7	9	-23
40	100	-27	-18	28	11	27	32	26
-2	-27	100	-24	-12	20	29	10	3
-7	-18	-24	100	9	24	11	-27	-5
-42	28	-12	9	100	17	4	-3	-19
9	11	20	24	17	100	-51	8	9
7	27	29	11	4	-51	100	-37	27
9	32	10	-27	-3	8	-37	100	21
-23	26	3	-5	-19	9	27	21	100

互信息



相关性分析 ——相关性的量化

■ 传统统计相关性分析

□ 肯德尔相关系数

2个n维随机变量X和Y，它们之间的肯德尔系数 τ 定义为：

$$\tau = \frac{2P}{\frac{1}{2}n(n-1)} - 1 = \frac{4P}{n(n-1)} - 1$$

其中，P表示一致的序对个数

比如：

X	1	2	3	4	5	6	7	8
Y	3	4	1	2	5	7	8	6

➡ 序

$P = 5 + 4 + 5 + 4 + 3 + 1 + 0 + 0$
↑
8-max(1,3)

相关性分析 ——相关性的量化

■ 传统统计相关性分析

□ 皮尔森相关系数

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

□ 斯皮尔曼相关系数

$$\rho = \frac{\sum_{i=1}^n (r_i - \bar{r})(s_i - \bar{s})}{\sqrt{\sum_{i=1}^n (r_i - \bar{r})^2} \sqrt{\sum_{i=1}^n (s_i - \bar{s})^2}}$$

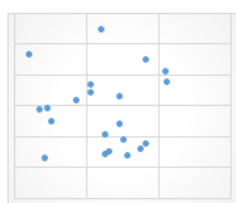
其中， r_i 和 s_i 分别是 X_i 和 Y_i 的秩 ← 序

形式一样
一个关注值
一个关注序

相关性分析

——相关性的量化

■ 对比



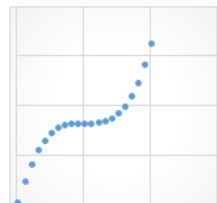
(a)

皮尔森相关系数: -0.0189
斯皮尔曼相关系数: -0.0208
肯德尔相关系数: -0.0095



(b)

皮尔森相关系数: 1.0000
斯皮尔曼相关系数: 1.0000
肯德尔相关系数: 1.0000



(c)

皮尔森相关系数: 0.9179
斯皮尔曼相关系数: 1.0000
肯德尔相关系数: 1.0000

相关性分析

——新的挑战

■ 传统统计相关性分析

- PEARSON相关系数
- Spearman相关系数
- KENDALL相关系数

非线性?
高维性?
海量性?

大数据

■ 大数据中的统计相关性分析

- 基于互信息的相关系数
- 基于协方差矩阵的相关系数
- 基于距离的相关系数

大数据中的统计相关性分析

■ 互信息

- 正式地，两个离散随机变量 X 和 Y 的互信息可以定义为

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right).$$

- 其中 $p(x, y)$ 是 X 和 Y 的联合概率分布函数，而 $p(x)$ 和 $p(y)$ 分别是 X 和 Y 的边缘概率分布函数。
- 在连续随机变量的情形下，求和被替换成了二重定积分

$$I(X; Y) = \int_Y \int_X p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right) dx dy,$$

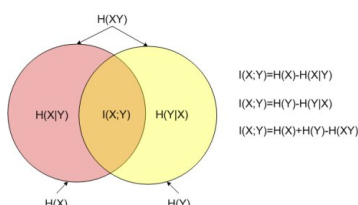
- 其中 $p(x, y)$ 当前是 X 和 Y 的联合概率密度函数，而 $p(x)$ 和 $p(y)$ 分别是 X 和 Y 的边缘概率密度函数。

大数据中的统计相关性分析

■ 互信息与其他量的等价关系

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X, Y) \\ &= H(X, Y) - H(X|Y) - H(Y|X) \end{aligned}$$

- 其中 $H(X)$ 和 $H(Y)$ 是边缘熵， $H(X|Y)$ 和 $H(Y|X)$ 是条件熵，而 $H(X, Y)$ 是 X 和 Y 的联合熵。注意到这组关系和并集、差集和交集的关系类似，用 Venn 图表示：



大数据中的统计相关性分析

■ 基于矩阵计算的相关系数

□ RV系数 (ROBERT P, 1976)

两个矩阵 A 和 B 的协方差定义为 $tr(AA' BB')$ ，方差分别为

$tr(AA')^2$, $tr(BB')^2$ (其中 $tr(\cdot)$ 是矩阵的迹，定义为矩阵主对角线元素的和)。鉴于上述定义，RV 系数以皮尔森相关系数的方式重新构造，即得

$$RV(A, B) = \frac{tr(AA' BB')}{\sqrt{tr(AA')^2 tr(BB')^2}}$$

RV是测量A和B协方差矩阵紧密程度的测度, 取值范围为[0,1]. 当RV约接近1, 说明用X (Y) 代替Y (X) 约合理

相关性分析 ——相关性的量化

■ 大数据中的统计相关性分析

□ 基于距离的相关系数

对于实数向量 $s = (s_1, s_2, \dots, s_p) \in R^p$ ，它的欧氏范数为 $\|s\| =$

$(s_1^2 + s_2^2 + \dots + s_p^2)^{1/2}$. 进一步定义 $i\langle s, X \rangle = s_1 X_1 + s_2 X_2 + \dots + s_p X_p$ 为 s 与 X 的内

积。同理，可以定义 $t = (t_1, t_2, \dots, t_q) \in R^q$ ， $\|t\|$ ， $i\langle t, Y \rangle$ 在此基础上，随机向量 (X, Y) 的联合特征函数定义为

$$f_{XY}(s, t) = E \exp[i\langle s, X \rangle + i\langle t, Y \rangle],$$

其中， i 为虚数单位， X, Y 各自的特征函数为

$$f_X(s) = f_{XY}(s, 0) = E \exp[i\langle s, X \rangle],$$

$$f_Y(t) = f_{XY}(0, t) = E \exp[i\langle t, Y \rangle].$$

相关性分析 ——相关性的量化

■ 大数据中的统计相关性分析

□ 基于距离的相关系数

定义随机向量 X 与 Y 的距离协方差 $V(X, Y)$ ，方差 $V^2(X)$ ， $V^2(Y)$ 。公式如下：

$$\begin{aligned} V^2(X, Y) &= \|f_{XY}(s, t) - f_X(s)f_Y(t)\|_{\omega}^2 \\ &= \int_{R^{p+q}} |f_{XY}(s, t) - f_X(s)f_Y(t)|^2 \omega(s, t) ds dt \\ V^2(X) &= V^2(X, X) = \|f_{XX}(s, t) - f_X(s)f_X(t)\|_{\omega}^2 \\ V^2(Y) &= V^2(Y, Y) = \|f_{YY}(s, t) - f_Y(s)f_Y(t)\|_{\omega}^2 \end{aligned}$$

其中， $\omega(s, t)$ 是权重函数，它的选择需要满足三个条件，即保证被积函数可积性； X 与 Y 独立时，相关系数为零； X 与 Y 同比例变化时，相关系数不变。在此基础上定义距离相关系数

$$R^2(X, Y) = \begin{cases} \frac{V^2(X, Y)}{\sqrt{V^2(X)V^2(Y)}} & V^2(X)V^2(Y) > 0 \\ 0 & V^2(X)V^2(Y) = 0 \end{cases}$$

此方法可以度量非线性相关性，且适用于度量任意两个不同维度的随机向量