

Ch 6.2 协方差和相关系数



回顾前一次课

若 (X, Y) 的联合概率密度为 $f(x, y)$, 则 $U = u(X, Y)$ 和 $V = v(X, Y)$ 的联合密度为 $f_{UV}(u, v) = f_{XY}(x(u, v), y(u, v))|J|$

n 维随机向量的联合分布函数、联合密度函数、边缘分布函数、边缘密度函数、性质、**随机向量 \mathbf{X} 和 \mathbf{Y} 相互独立**

多维正太分布: $X = (X_1, X_2, \dots, X_n) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$X \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ 分解 $\boldsymbol{\Sigma} = \mathbf{U}^\top \boldsymbol{\Lambda} \mathbf{U}$, 则 $Y = \boldsymbol{\Lambda}^{-1/2} \mathbf{U}(X - \boldsymbol{\mu}) \sim N(\mathbf{0}_n, I_n)$

设随机向量 $X \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 则有 $Y = \mathbf{A}X + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$

多维正太分布重要的性质

定理 5.13 设随机向量 $X = (X_1, X_2, \dots, X_n)^T$ 和 $Y = (Y_1, Y_2, \dots, Y_m)^T$, 以及

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \right),$$

- 随机向量 X 和 Y 的边缘分布分别为 $X \sim \mathcal{N}(\mu_x, \Sigma_{xx})$ 和 $Y \sim \mathcal{N}(\mu_y, \Sigma_{yy})$;
- 随机向量 X 与 Y 相互独立的充要条件是 $\Sigma_{xy} = (\mathbf{0})_{m \times n}$ (元素全为零的 $m \times n$ 矩阵);
- 在 $X = \mathbf{x}$ 的条件下随机向量 $Y \sim \mathcal{N}(\mu_y + \Sigma_{yx} \Sigma_{xx}^{-1}(\mathbf{x} - \mu_x), \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy})$;
- 在 $Y = \mathbf{y}$ 的条件下随机向量 $X \sim \mathcal{N}(\mu_x + \Sigma_{xy} \Sigma_{yy}^{-1}(\mathbf{y} - \mu_y), \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx})$.

期望

随机变量 $Z = g(X, Y)$ 的期望为 $E[Z] = \sum_{i,j} g(x_i, y_j) p_{ij}$

随机变量 $Z = g(X, Y)$ 的期望 $E[Z] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x, y) f(x, y) dx dy$

若随机变量 $X \geq Y$, 则有 $E[X] \geq E[Y]$

对任意随机变量 X, Y 有 $E[X + Y] = E[X] + E[Y]$

对独立随机变量 X 和 Y , 有 $E[XY] = E[X]E[Y]$

对任意随机变量 X 和 Y , 有Cauchy-Schwartz不等式

$$|E[XY]| \leq \sqrt{E[X^2]E[Y^2]}$$

对独立随机变量 X 和 Y , 有

$$\text{Var}(X \pm Y) = \text{Var}(X) + \text{Var}(Y)$$

协方差

随机变量的期望或方差仅涉及变量自身的统计信息, 没有刻画变量之间的统计信息

协方差: 描述随机变量 X 和 Y 相互关系的数字特征

定义: 设二维随机向量 (X, Y) 的期望 $E[(X - E(X))(Y - E(Y))]$ 存在, 则称其为 X 和 Y 的协方差, 记为

$$\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))] = E(XY) - E(X)E(Y)$$

协方差是两个随机变量与它们各自期望的偏差之积的期望, 由于偏差可正可负, 因此协方差可正可负.

协方差的性质

- $\text{Cov}(X, c) = 0$
- $\text{Cov}(X, X) = \text{Var}(X)$
- $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
- $\text{Var}(X \pm Y) = \text{Var}(X) + \text{Var}(Y) \pm 2\text{Cov}(X, Y)$

对任意常数 a 和 b , 随机变量 X 和 Y , 有

- $\text{Cov}(aX, bY) = ab\text{Cov}(X, Y)$
- $\text{Cov}(X + a, Y + b) = \text{Cov}(X, Y)$

协方差的性质

对任意常数 X_1, X_2 和 Y , 有

$$\text{Cov}(X_1 + X_2, Y) = \text{Cov}(X_1, Y) + \text{Cov}(X_2, Y)$$

对随机变量 X_1, X_2, \dots, X_n 和 Y_1, Y_2, \dots, Y_m , 有

$$\text{Cov}\left(\sum_{i=1}^n X_i, \sum_{j=1}^m Y_j\right) = \sum_{i=1}^n \sum_{j=1}^m \text{Cov}(X_i, Y_j)$$

以及进一步有

$$\text{Var}\left(\sum_i X_i\right) = \sum_i \text{Var}(X_i) + 2 \sum_{i < j} \text{Cov}(X_i, X_j)$$

协方差的性质

若随机变量 X 与 Y 独立, 则 $\text{Cov}(X, Y) = 0$, 反之不成立.

定理: 对任意随机变量 X 与 Y 有

$$(\text{Cov}(X, Y))^2 \leq \text{Var}(X)\text{Var}(Y)$$

等号成立的充要条件是 $Y = aX + b$ 几乎处处成立, 即 X 与 Y 之间几乎处处存在线性关系

正太分布

若随机向量 $(X, Y) \sim N(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2, \rho)$, 则 $\text{Cov}(X, Y) = \rho\sigma_x\sigma_y$

推论：若随机向量 $(X, Y) \sim N(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2, \rho)$, 则

X 和 Y 相互独立的**充要条件**是协方差 $\text{Cov}(X, Y) = 0$

例题

设随机变量 X_1, X_2, \dots, X_n 相互独立且服从正太分布, 方差为 σ^2 .
记 $\bar{X} = \sum_{i=1}^n X_i/n$, 讨论 \bar{X} 和 $\bar{X} - X_i$ 的独立性

例题

随机变量 (X, Y) 联合概率密度为

$$f(x, y) = \begin{cases} (x + y)/8 & 0 \leq x \leq 2, 0 \leq y \leq 2 \\ 0 & \text{其它} \end{cases}$$

求 $\text{Cov}(X, Y)$ 和 $\text{Var}(X + Y)$.

例题-匹配问题

有 n 对夫妻参加一次聚会, 将所有参会人员任意分成 n 组, 每组一男一女, 用 X 表示夫妻两人被分到一组的对数, 求 X 的期望和方差

相关系数

两个随机变量之间的关系：独立与非独立

非独立关系：线性关系和非线性关系。非线性关系较为复杂，无好办法来处理。线性相关程度可以通过线性相关系数来刻画

定义：随机变量 X 和 Y 的方差 $\text{Var}(X), \text{Var}(Y)$ 存在且不为0, 称

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

为 **X 与 Y 的相关系数**

- 若 $\rho_{XY} > 0$, 称 **X 与 Y 正相关**
- 若 $\rho_{XY} < 0$, 称 **X 与 Y 负相关**
- 若 $\rho_{XY} = 0$, 称 **X 与 Y 不相关**

相关系数性质

- 根据协方差性质可知 $|\rho_{XY}| \leq 1$
- $|\rho_{XY}| = 1$ 的充要条件为 X 与 Y 几乎处处有线性关系 $Y = aX + b$
- 若 X 与 Y 相互独立, 则 X 与 Y 不相关 ($\rho_{XY} = 0$), 但反之不成立

ρ_{XY} 刻画了 X 与 Y 的线性相关性, 又称线性相关系数。随机变量 X 与 Y 不相关, 仅表示 X 与 Y 之间不存在线性关系, 可能存在其他关系

例如, 设随机变量 $X \sim U(-1/2, 1/2)$ 和 $Y = \cos(X)$, 有

$$\text{Cov}(X, Y) = E[X \cos(X)] - E(X)E(\cos(X))$$

$$= E[X \cos(X)] = \int_{-1/2}^{1/2} x \cos(x) dx = 0$$

相关系数

对方差不为零的随机变量 X 和 Y , 下述条件相互等价:

- $\rho_{XY} = 0$
- $\text{Cov}(X, Y) = 0$
- $E(XY) = E(X)E(Y)$
- $\text{Var}(X \pm Y) = \text{Var}(X) + \text{Var}(Y)$

若随机向量 $(X, Y) \sim N(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2, \rho)$, 则 X 与 Y 的相关系数 $\rho_{XY} = \rho$, 因此 X 与 Y 独立的充要条件是 X 与 Y 不相关

例题

设随机变量 $X \sim N(\mu, \sigma^2)$ 和 $Y \sim N(\mu, \sigma^2)$ 相互独立，求 $Z_1 = \alpha X + \beta Y$ 和 $Z_2 = \alpha X - \beta Y$ 的相关系数 ($\alpha, \beta \neq 0$)

例题

设随机向量 (X_1, X_2, \dots, X_n) 服从多项分布 $M(m, p_1, p_2, \dots, p_n)$, 对任意 $i \neq j$, 求 X_i 和 X_j 的相关系数