

数据统计期末笔记清单

Made by: X-B.D.

一. 基础知识与重要分布

① 三大分布

• 卡方分布 $\chi^2(n)$

$$E(X) = \underline{n} \quad D(X) = \underline{2n}$$

ΔX 与 Y 相互独立, $X \sim \chi^2(m)$, $Y \sim \chi^2(n)$, 则 $X+Y \sim \underline{\chi^2(m+n)}$

② • T分布

$$T = \frac{X}{\sqrt{Y/n}} \sim t(n)$$

其中分子表示 $\underline{X \sim N(0,1)}$, 分母表示 $\underline{Y \sim \chi^2(n)}$

• F分布

$$F = \frac{X/m}{Y/n} \sim F(m,n)$$

其中分子表示 $\underline{X \sim \chi^2(m)}$, 分母表示 $\underline{Y \sim \chi^2(n)}$

② 常见分布与不常见分布

	EX	DX	表达式
(0,1)分布	p	$p(1-p)$	$p^x(1-p)^{1-x}$
几何分布	$\frac{1}{p}$	$\frac{1-p}{p^2}$	$(1-p)^{x-1}p$
二项分布 $B(n,p)$	np	$np(1-p)$	略
均匀分布 $U(a,b)$	$\frac{a+b}{2}$	$\frac{(a-b)^2}{12}$	x
指数分布 $Exp(\lambda)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	$\lambda e^{-\lambda x}$
正态分布	μ	σ^2	$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
泊松分布 $P(\lambda)$	λ	λ	$\frac{\lambda^k}{k!} e^{-\lambda}$
伽马分布 $\Gamma(\alpha, \beta)$	$\frac{\alpha}{\beta}$	α/β^2	$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$
逆伽马分布 $IG(\alpha, \beta)$	$\frac{\beta}{\alpha-1}$ (若 $\alpha > 1$)	$\frac{\beta^2}{(\alpha-1)^2(\alpha-2)}$	$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} e^{-\frac{\beta}{x}}$

$$\frac{\sum_{i=1}^n (x_i - \mu_0)^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = 1 + \frac{T^2}{n-1}$$

内部需包含一个自由度为 $n-1$ 的分布 T 的函数 $u(T)$

其中 $T = \frac{\sqrt{n(n-1)}(\bar{x} - \mu_0)}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \sim t(n-1)$

	EX	DX	表达式
贝塔分布 $B(a,b)$	$\frac{a}{a+b}$	$\frac{ab}{(a+b)^2(a+b+1)}$	$\frac{1}{B(a,b)} x^{a-1}(1-x)^{b-1}$

伽马函数: $\Gamma(x) = \int_0^{+\infty} t^{x-1} e^{-t} dt$

$\Gamma(x+1) = x \Gamma(x)$, $\Gamma(1) = 1$

贝塔函数: $B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$

$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$ (用 Γ 函数表示)

n 个泊松分布 $P(\lambda)$ 的和为 $\underline{P(n\lambda)}$

n 个指数分布 $Exp(\lambda)$ 的和为 $\underline{\Gamma(n, \lambda)}$

$X \sim \Gamma(n, \lambda)$, 则 $\alpha X \sim \Gamma(\underline{n}, \underline{\frac{\lambda}{\alpha}})$
 $\Gamma(n, \frac{1}{2}) = \underline{\chi^2(2n)}$

二、枢轴量与检验统计量

① 枢轴量

未知参数	条件	枢轴量及分布
μ	σ^2 已知	$(\bar{X} - \mu) / (\sigma / \sqrt{n}) \sim N(0, 1)$
	σ^2 未知	$(\bar{X} - \mu) / (S / \sqrt{n}) \sim t(n-1)$
$\mu_1 - \mu_2$	σ_1^2, σ_2^2 已知	$(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2) / \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \sim N(0, 1)$
	$\sigma_1^2 = \sigma_2^2$ 未知	$(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2) / S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \sim t(n_1 + n_2 - 2)$ ← 此处的 $S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1 + n_2 - 2}$
σ^2	μ 已知	$\sum_{i=1}^n (X_i - \mu)^2 / \sigma^2 \sim \chi^2(n)$
	μ 未知	$\sum_{i=1}^n (X_i - \bar{X})^2 / (n-1) S^2 \sim \chi^2(n-1)$
σ_1^2 / σ_2^2	μ_1, μ_2 未知	$\frac{S_1^2 / S_2^2}{\sigma_1^2 / \sigma_2^2} \sim F(n_1-1, n_2-1)$

② 检验统计量 (似然比)

	条件	$X \sim N(\mu, \sigma^2)$	μ, σ^2 未知	否定域及属于哪类分布
T ₀ 类	$H_0: \mu = \mu_0 \Leftrightarrow H_1: \mu \neq \mu_0$		$\left\{ T = \left \frac{\sqrt{n}(\bar{X} - \mu_0)}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}} \right > C \right\}$	$T \sim t(n-1)$
	$H_0: \mu \leq \mu_0 \Leftrightarrow H_1: \mu > \mu_0$		$\left\{ T = \frac{\sqrt{n}(\bar{X} - \mu_0)}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}} > C \right\}$	$T \sim t(n-1)$
T ₁ 类	$H_0: \sigma^2 = \sigma_0^2 \Leftrightarrow H_1: \sigma^2 \neq \sigma_0^2$		$\left\{ T_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma_0^2} > C_2 \text{ 或 } < C_1 \right\}$	$T_1 \sim \chi^2(n-1)$
	$H_0: \sigma^2 \leq \sigma_0^2 \Leftrightarrow H_1: \sigma^2 > \sigma_0^2$		$\left\{ T_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma_0^2} > C \right\}$	$T_1 \sim \chi^2(n-1)$
	条件	$X \sim N(\mu_1, \sigma_1^2)$ Y_1, Y_2, \dots, Y_{n_2}	$Y \sim N(\mu_2, \sigma_2^2)$ X_1, X_2, \dots, X_{n_1}	否定域及属于哪类分布
H ₀ 类	μ_1, μ_2 已知	$H_0: \sigma_1^2 = \sigma_2^2 \Leftrightarrow H_1: \sigma_1^2 \neq \sigma_2^2$	$\left\{ \frac{\sum (X_i - \mu_1)^2 / n_1}{\sum (Y_j - \mu_2)^2 / n_2} = F_0 > C_2 \text{ 或 } < C_1 \right\}$	$F_0 \sim F(n_1, n_2)$
		$H_0: \sigma_1^2 \leq \sigma_2^2 \Leftrightarrow H_1: \sigma_1^2 > \sigma_2^2$	$\left\{ \frac{\sum (X_i - \mu_1)^2 / n_1}{\sum (Y_j - \mu_2)^2 / n_2} = F_0 > C \right\}$	$F_0 \sim F(n_1, n_2)$
	μ_1, μ_2 未知	$H_0: \sigma_1^2 = \sigma_2^2 \Leftrightarrow H_1: \sigma_1^2 \neq \sigma_2^2$	$\left\{ \frac{\sum (X_i - \bar{X})^2 / (n_1-1)}{\sum (Y_j - \bar{Y})^2 / (n_2-1)} = F_1 > C_2 \text{ 或 } < C_1 \right\}$	$F_1 \sim F(n_1-1, n_2-1)$
		$H_0: \sigma_1^2 \leq \sigma_2^2 \Leftrightarrow H_1: \sigma_1^2 > \sigma_2^2$	$\left\{ \frac{\sum (X_i - \bar{X})^2 / (n_1-1)}{\sum (Y_j - \bar{Y})^2 / (n_2-1)} = F_1 > C \right\}$	$F_1 \sim F(n_1-1, n_2-1)$
$\sigma_1^2 = \sigma_2^2$ 未知	$H_0: \mu_1 = \mu_2 \Leftrightarrow H_1: \mu_1 \neq \mu_2$		$\left\{ T_2 = \left \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sum (X_i - \bar{X})^2 + \sum (Y_j - \bar{Y})^2}{n_1 + n_2}}} \right > C \right\}$	$T_2 \sim t(n_1 + n_2 - 2)$
	$H_0: \mu_1 \leq \mu_2 \Leftrightarrow H_1: \mu_1 > \mu_2$		$\left\{ T_2 = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sum (X_i - \bar{X})^2 + \sum (Y_j - \bar{Y})^2}{n_1 + n_2}}} > C \right\}$	$T_2 \sim t(n_1 + n_2 - 2)$
$\sigma_1^2 = \sigma_2^2$ 已知	$H_0: \mu_1 = \mu_2 \Leftrightarrow H_1: \mu_1 \neq \mu_2$		$\left\{ M_0 = \left \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \right > C \right\}$	$M_0 \sim N(0, 1)$
	$H_0: \mu_1 \leq \mu_2 \Leftrightarrow H_1: \mu_1 > \mu_2$		$\left\{ M_0 = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} > C \right\}$	$M_0 \sim N(0, 1)$

注: $\sigma_1^2 \neq \sigma_2^2$ 无法使用T义似然比得到UMP检验。

自强不息 厚德载物

三、回归分析与方差分析

① 一元线性回归与正比例回归

$$y_i = a + bx_i + \varepsilon \quad \varepsilon \sim N(0, \sigma^2)$$

$$l_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}, \quad l_{xx} = \frac{\sum (x_i - \bar{x})^2}{n}, \quad l_{yy} = \frac{\sum (y_i - \bar{y})^2}{n}, \quad u = \frac{\sum (\hat{y}_i - \bar{y})^2}{n}, \quad Q = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$

回归平方和 残差平方和

$$\hat{b} = \frac{l_{xy}}{l_{xx}}, \quad \hat{a} = \bar{y} - \hat{b}\bar{x}$$

平方和分解公式: $l_{yy} = Q + U$

用 $F = \frac{U}{Q/(n-2)}$ 衡量回归优差程度, 其服从 $F(1, n-2)$ 分布.

$$H_0: b=0 \Leftrightarrow H_1: b \neq 0$$

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}, \quad F(r) = \frac{r^2}{1-r^2} \cdot \frac{1}{n-2}$$

r^2 用 l_{xy}, l_{xx}, l_{yy} 表示: $r^2 = \frac{l_{xy}^2}{l_{xx} l_{yy}} = \frac{U}{l_{yy}}$

y_0 的置信区间枢轴量: $T = \frac{y_0 - \hat{y}_0}{\sqrt{dQ/(n-2)}} \sim (t(n-2))$, 其中 $d = 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}$

$$y = bx + e \quad e \sim N(0, \sigma^2)$$

$$\hat{b} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

检验 $H_0: b=0$

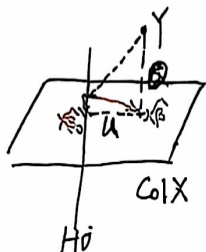
$$F = \frac{\hat{b}^2 \sum x_i^2}{Q/(n-1)} \sim F(1, n-1) \quad (Q = \sum (y_i - \hat{b}x_i)^2)$$

$F > \lambda$ 时否定 H_0

② 多元线性回归

$$Y = X\beta + \varepsilon \quad \hat{\beta} = (X^T X)^{-1} X^T Y, \quad \text{记 } \text{rank } X = r, \quad \text{rank } H = q$$

参数检验: $H_0: H\beta = 0 \Leftrightarrow H_1: H\beta \neq 0$



$$\lambda = \left(\frac{\|Y - X\hat{\beta}_0\|^2}{\|Y\|^2} \right)^{\frac{n}{2}} = \left(1 + \frac{r-q}{n-r} F \right)^{\frac{n}{2}}$$

其中 $F \sim F(r-q, n-r)$, $\lambda > \lambda_0$ 时拒绝 H_0 . 强不息 厚德载物

四. 贝叶斯估计

1. 已知样本 x_1, \dots, x_n , 估计参数 $\theta \in \Theta$, 似然函数为 L , θ 的先验分布为 $\pi(\theta)$

后验分布:

$$\pi(\theta | x_1, \dots, x_n) = \frac{\int_{\Theta} L \pi(\theta) d\theta}{\int_{\Theta} L \pi(\theta) d\theta}$$

① $X \sim B(1, p)$, $\pi(p) : \text{beta}(\alpha, \beta)$, x_1, x_2, \dots, x_n ,

则 $\pi(p | x_1, x_2, \dots, x_n) : \text{beta}(\alpha + \sum x_i, \beta + n - \sum x_i)$

② $X \sim \text{Poisson}(\lambda)$, $\pi(\lambda) : \text{gamma}(\alpha, \beta)$, x_1, x_2, \dots, x_n

则 $\pi(\lambda | x_1, x_2, \dots, x_n) : \text{gamma}(\alpha + \sum x_i, \beta + n)$

③ $X \sim \text{Exp}(\lambda)$, $\pi(\lambda) : \text{gamma}(\alpha, \beta)$, x_1, x_2, \dots, x_n

则 $\pi(\lambda | x_1, x_2, \dots, x_n) : \text{gamma}(\alpha + n, \beta + \sum x_i)$

④ $X \sim N(\mu, \sigma^2)$, σ^2 已知, μ 未知, $\pi(\mu) : N(\mu_0, \sigma_0^2)$, x_1, x_2, \dots, x_n

则 $\pi(\mu | x_1, \dots, x_n) : N\left(\frac{\mu_0 \sigma^2 + n \sigma_0^2 (\bar{x})^2}{\sigma^2 + n \sigma_0^2}, \frac{\sigma^2 \sigma_0^2}{\sigma^2 + n \sigma_0^2}\right)$

⑤ $X \sim N(\mu, \frac{1}{R})$, μ 已知, R 未知, $\pi(R) : \text{gamma}(\alpha, \beta)$, x_1, x_2, \dots, x_n

则 $\pi(R | x_1, \dots, x_n) : \text{gamma}\left(\alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2\right)$

2. 贝叶斯估计

定义损失函数 $L(\hat{\theta}, \theta)$, $\theta \in \Theta$. 似然函数为 L , θ 的先验分布为 π_1 , 后验分布为 π_2

贝叶斯估计量是积分

$$\int_{\Theta} L(\hat{\theta}, \theta) \pi_2 d\theta$$

取最小的 θ 值.

五、一些杂乱的定义与知识点.

1. Fisher 信息量与 C-R 不等式

$$I(\delta) = \underbrace{E\left[\frac{\partial \ln f(x, \delta)}{\partial \delta}\right]^2}_{\text{用一阶偏导}} = \underbrace{-E\left[\frac{\partial^2 \ln f(x, \delta)}{\partial \delta^2}\right]}_{\text{用二阶偏导}}, \quad \text{Var}_\theta(\psi(x_1, \dots, x_n)) \geq \boxed{\frac{g'(\theta)^2}{n I(\theta)}}$$

2. 相合: $P(\|\hat{\theta}_n - \theta\| < \varepsilon) = 1$; 强相合: $P(\lim_{n \rightarrow \infty} \hat{\theta}_n = \theta) = 1$

3. 充分统计量: $L(x_1, \dots, x_n | \theta) = g[\varphi(x), \theta] h(x)$, φ 为充分统计量

完全统计量: 对于每个可测函数 $u(\cdot)$, 若 $E_\theta u(\varphi(x)) = 0$ 可推出 $P(u(\varphi(x)) = 0) = 1$, 称 φ 为完全统计量

4. 指数型分布

① 通式: $f(x, \theta) = s(\theta) h(x) e^{\sum c(\theta) T(x)}$

充分统计量 φ 是完备的一个快速判据为 θ 有内点

② 单参数指数型分布: $f(x, \theta) = s(\theta) h(x) e^{c(\theta) T(x)}$

$H_0: \theta \leq \theta_1 \Leftrightarrow H_1: \theta > \theta_1$, 否定域 $\left\{ \sum_{i=1}^n T(x_i) > c \right\}$ 为检验水平为 α 的 UMP

其中临界值求法: $P(\sum T(x) > c | \theta = \theta_1) = \alpha$

5. 操作特性函数 $L_W(\theta) = P(\text{接受 } H_0 | \theta)$
功效函数 $P_W(\theta) = P(\text{拒绝 } H_0 | \theta)$ $\left. \begin{matrix} L_W(\theta) \neq P_W(\theta) = 1. \end{matrix} \right\}$

UMP 否定域: \forall 水平不超过 α 的 \tilde{w} 均有 $P_W(\theta) \geq P_{\tilde{w}}(\theta)$, 称 W 为 UMP 否定域 ($\forall \theta \in \Theta$)

无偏性: 正确时被拒绝的概率 小于 错误时被拒绝的概率

6. N-P 引理.

7. minimax 决策.

$$\sup_{\theta} R(\theta, \delta^*) \leq \sup_{\theta} R(\theta, \delta)$$