

统计学

一、统计与数据基础

1.1 统计学 statistics 统计数据 分类数据 categorical data (有序/无序) qualitative; 数值数据 metric (离散/连续) quantitative

观测数据 observational ~ 实验数据 experimental ~ 截面数据 cross-sectional ~

时间序列数据 time series ~ 描述统计 descriptive ~ μ, σ, π 推断统计 inferential ~ \bar{x}, s, p

• 总体 population 样本 sample 样本量 sample size 参数 parameter 统计量 statistic
变量 variable (categorical ~ / metric ~)

1.2 数据来源: 间接 (二手)、直接 < 调查 实验

• 概率抽样 probability sampling: 简单随机抽样 ← 抽样框 sampling frame, ^{simple random}

分层抽样, 整群抽样, 系统抽样, 多阶段抽样
^{stratified} ^{cluster} ^{systematic} ^{multi-stage}

非概率抽样 non-~: 方便抽样, 判断抽样, 自愿样本, 滚雪球抽样

配额抽样 数据搜集方法: 自填, 面谈, 电话 (CATI) 及其特点

• 实验组 experiment group 对照组 control ~ 双盲法 统计方法 → 实验设计 (有效性)

• 数据误差: 抽样误差 (仅 probability 中) sampling error; 非抽样误差 non-~: 内部因果 外部推广

抽样框误差, 回答误差, 无回答误差, 调查员误差, 测量误差
(仅 prob 中)

1.3 数据预处理 data preprocessing: ~ 清洗 cleaning, ~ 集成 integration, ~ 归约 reduction, ~ 转换 transformation (审核)

• 数据整理与展示: 频数分布 frequency distribution 列联表 contingency table

条/柱形图 bar/column chart 帕累托图 Pareto chart 饼图 pie chart 环形图;

数据分组 组中值 class midpoint 直方图 histogram 箱图 box plot 须线 whiskers

线图 line plot 散点图 scatter diagram 雷达图 radar chart
inter/outer fence 25% 四分位距 quartile deviation
 $Q_{25\%} - 1.5 IQR$ (3) 内/外圈框 25% 回 分位数 中位数
whiskers
outlier 离群点
upper adjacent value 上相邻值
75% ~
quartiles median
IQR
inter/outer fence

1.4 数据的度量 (描述、概括): 集中趋势 central tendency, 离散程度, 分布形状

平均数/均值 mean $\bar{x} = \frac{\sum x_i}{n}$ simple mean $\bar{x} = \frac{\sum m_i f_i}{n}$ weighted mean 加权~

分位数 quantile 中位数 median 四分位数 quartile (SPSS 中 $\frac{n+1}{4}, \frac{3(n+1)}{4}$)

众数 mode M_o 三种数之选择.

全距 range 四分位距 inter-quartile range 平均差 (平均绝对离差) mean deviation

方差 variance 标准差 standard deviation $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$ $\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$

(变异) 离散系数 CV coefficient of variation = $\frac{s}{\bar{x}}$ 标准化 standardize $Z_i = \frac{x_i - \bar{x}}{s}$

偏度 skewness (偏度系数 SK) = $\sum_i \left(\frac{x_i - \bar{x}}{s}\right)^3$ 左偏/右偏 峰度 kurtosis ($\sim K$)
 $= \sum_i \left(\frac{x_i - \bar{x}}{s}\right)^4 - 3$

二、概率与抽样分布

2.1 试验事件 随机事件 (必然~不可能~) random event 基本事件 elementary event

概率 probability 古典定义 统计定义 主观概率定义

随机变量 random variable 概率函数 probability function 离散型/连续型

概率分布 probability distribution (密度) 分布函数 $F(x)$ $f(x)$ $f(x) \geq 0, \int_{-\infty}^{+\infty} f(x) dx = 1$

期望(值) expected value $E(X)$ 方差 $D(X)/V(X)$ $(E(X) = \int_{-\infty}^{+\infty} x f(x) dx, D(X) = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx)$

均匀分布 uniform ~ 二项分布 binomial ~ $P(X) = C_n^x p^x q^{n-x}$ ($p+q=1$), $E(X) = np, D(X) = npq$

泊松分布 Poisson ~ $P(X) = \frac{\lambda^x e^{-\lambda}}{x!}, E(X) = \lambda, D(X) = \lambda$ $f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$

正态分布 normal ~ $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, X \sim N(\mu, \sigma^2)$ 标准正态分布 (Z) $N(0, 1)$

"6σ 准则" ($6\sigma: 2 \times 10^{-9}, 3\sigma: 0.27\%$) "1.5σ 漂移" ($\rightarrow 3.4 \times 10^{-6}$)

2.2 统计量 $T(X_1, \dots, X_n)$ 样本 k 阶距 $m_k = \frac{1}{n} \sum_i x_i^k$ 样本 k 阶中心距 $v_k = \frac{\sum_i (x_i - \bar{x})^k}{n-1}$

抽样分布 sampling distribution χ^2 分布 Chi-square ~ X_1, \dots, X_n iid. $\sum_i X_i^2 \sim \chi^2(n)$

t 分布 $X \sim N(0, 1), Y \sim \chi^2(n), \frac{X}{\sqrt{Y/n}} \sim t(n)$ $E(t) = 0, D(t) = \frac{n}{n-2}$ $E(X^2) = n, D(X^2) = 2n$

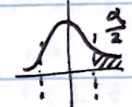
F 分布 $Y, Z \sim \chi^2(m), \chi^2(n), X = \frac{Y/m}{Z/n} \sim F(m, n)$ $E(X) = \frac{n}{n-2}$ $D(X) = \frac{2n^2(m+n-2)}{m(n-2)(n-4)}$

$X \sim t(n) \leftrightarrow X^2 \sim F(1, n)$ 分位数 $F_p(m, n) = \frac{1}{F_{1-p}(n, m)}$ (474)

• 马尔科夫不等式 $P(X \geq a) \leq \frac{E(X)}{a}$ 切比雪夫不等式 $P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$

中心极限定理 central limit theorem $\lim_{n \rightarrow \infty} \bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ (大样本: $30 <$)

三. 参数估计与假设检验

3.1 • 参数估计 parameter estimation 估计量 estimator $\hat{\theta}$ 估计值 estimated value 
 点估计 point ~ 区间估计 interval ~ 置信区间 confidence interval 置信上/下限
 置信度 (水平系数) confidence level $1 - \alpha$ 显著性水平 significance level α

评价估计量的标准: 无偏性 unbiasedness $E(\hat{\theta}) = \theta$, 有效性 efficiency $D(\hat{\theta}) < D(\hat{\theta}_0)$

• 均值估计: 1. N 总体, σ^2 已知, 大/小样本 一致性 consistency $\lim_{n \rightarrow \infty} \hat{\theta} = \theta$
 且非 N 总体, 大样本 $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$
 且大样本时 $s^2 \rightarrow \sigma^2$ $\mu \in \bar{X} \pm Z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} (\pm Z_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}})$

2. N 总体, σ^2 未知, 小样本 $t = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{(n-1)}$
 $\mu \in \bar{X} \pm t_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$
 比例估计: $\sigma_p^2 = \frac{\pi(1-\pi)}{n}$ $Z = \frac{p - \pi}{\sqrt{\pi(1-\pi)/n}} \sim N(0, 1)$ $\mu \in \bar{X} \pm t_{\frac{\alpha}{2}} \frac{s}{\sqrt{n}}$
 (大样本) $p \rightarrow \pi$ $\pi \in p \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{\pi(1-\pi)}{n}}$ (方差估计: $\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{(n-1)}$)

• 均值差估计: 独立样本 independent sample $\sigma^2 \in [\frac{(n-1)s^2}{\chi^2_{\frac{\alpha}{2}}}, \frac{(n-1)s^2}{\chi^2_{1-\frac{\alpha}{2}}}]$
 $Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$ $\mu_1 - \mu_2 \in \bar{X}_1 - \bar{X}_2 \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$ (或 $\pm Z_{\frac{\alpha}{2}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$)

1. 大样本 或 σ^2 已知 $\sigma_1^2 = \sigma_2^2$ 时: $\mu_1 - \mu_2 \in \bar{X}_1 - \bar{X}_2 \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2}}$
 2. 小样本, σ^2 未知 (N 总体) $\sigma_1^2 \neq \sigma_2^2$ 时: 自由度 $\nu = \frac{(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2})^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}}$
 $S_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}$
 $\mu_1 - \mu_2 \in \bar{X}_1 - \bar{X}_2 \pm t_{\frac{\alpha}{2}}(n_1 + n_2 - 2) \sqrt{S_p^2 (\frac{1}{n_1} + \frac{1}{n_2})}$
 $\mu_1 - \mu_2 \in \bar{X}_1 - \bar{X}_2 \pm t_{\frac{\alpha}{2}}(\nu) \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$

匹配样本 matched sample (信息要优于独立下)

$$M_d = \mu_1 - \mu_2 \in \bar{d} \pm Z_{\frac{\alpha}{2}} \frac{\sigma_d}{\sqrt{n}} (\pm t_{\frac{\alpha}{2}}(n-1) \frac{s_d}{\sqrt{n}})$$

比例差估计: $\pi_1 - \pi_2 \in p_1 - p_2 \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$ 方差比估计: $\frac{s_1^2}{s_2^2} \cdot \frac{\sigma_2^2}{\sigma_1^2} \sim F_{(n_1-1, n_2-1)}$

• 确定样本量 (改推): (均值) $n = \frac{(Z_{\frac{\alpha}{2}})^2 \sigma^2}{e^2}$ (或 s^2) $\frac{\sigma_1^2}{\sigma_2^2} \in [\frac{s_1^2/s_2^2}{F_{\alpha/2}}, \frac{s_1^2/s_2^2}{F_{1-\alpha/2}}]$
 (比例) $n = \frac{(Z_{\frac{\alpha}{2}})^2 \pi(1-\pi)}{e^2}$ (或 p) $\leq \frac{Z_{\frac{\alpha}{2}}^2}{4e^2}$

3.2 • 假设检验 hypothesis testing 原假设 null hypothesis 备择假设 alternative ~

α 错误 (弃真错误) β 错误 (取伪错误)

假设检验的流程: "H₀ 假定, 统计量是否落在拒绝域内?"

单侧(尾)/双侧(尾)检验 如 $|z| < |z_{\alpha/2}|$ 或 $|z| > |z_{\alpha/2}|$ 或利用 P 值 伴随(相伴)概率

左/右单侧检验 (下/上限检验)

检验统计量的确定: 样本量 n 大: z 小: t (未知 σ^2)

两个总体 N 均值差: z (大), t (小) 比例差: z 方差比: F

检验解释: H₀ ✓: "H₁ 为真出错的概率不超过 α "

H₀ ✗: "没有充足证据 (在 α 的显著水平下) 拒绝 H₀"

"H₀ 常为不轻易否定的 (原有、传统的) 命题, 要证明的常置于备择位 (H₁)"

α 可理解为苛刻的程度, 或相同 (H₀) 显著的程度

四、分类数据分析与方差分析

4.1 • χ^2 检验 < 拟合优度检验 goodness of fit test (适合性检验) $\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$

独立性检验 independence test (列联分析)

观察值频数 observed frequency f_o 期望值频数 expected frequency f_e

$f_e < 5$ 不宜使用

适合性检验: $df = n - 1$ 独立性检验: $df = (R - 1)(C - 1)$, $f_e = n \cdot \frac{RT}{n} \cdot \frac{CT}{n} = \frac{RT \cdot CT}{n}$

相关系数: 1. $\phi = \sqrt{\frac{\chi^2}{n}}$ (2x2 时 $\phi \in [0, 1]$, $n \uparrow \phi \uparrow$) 2. $c = \sqrt{\frac{\chi^2}{\chi^2 + n}}$ 列联系数 coefficient of contingency ($n \uparrow c_{max} \uparrow < 1$)

3. $V = \sqrt{\frac{\chi^2}{n \cdot \min[R-1, C-1]}}$ ($V \in [0, 1]$, R 或 $C = 2$ 时 $V = \phi$) 条件百分表之方向

4.2 • 方差分析 ANOVA analysis of variance 因素(因子) factor 水平(处理) treatment level

误差分解 → 分析检验 基本假定: 独立性、正态性、齐方差性

单因素方差分析 one-way ~ $H_0: \mu_1 = \dots = \mu_k$ 总平方和 SST sum of squares for total = $\sum_j \sum_i (X_{ij} - \bar{x})^2$

组间平方和 SSA ~ factor A = $\sum_i n_i (\bar{x}_i - \bar{x})^2$ 组内平方和/误差平方和 SSE ~ error = $\sum_j \sum_i (X_{ij} - \bar{x}_i)^2$

$SST = SSA + SSE$ ~ 均方/方差 MST/A/E

检验统计量 $F = \frac{MSA}{MSE} \sim F(k-1, n-k)$ (上限检验) 方差分析表 ANOVA table

"当因素仅 2 个水平时, one-way anova 等同于均值差 t 检验, 由于此时 $F = t^2$ "

关系强度(判定系数) $R^2 = \frac{SSA}{SST}$ "所解释的比例"

$n_1 (\frac{\mu_1 - \mu_2}{\sigma_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}})^2 = \frac{(\mu_1 - \mu_2)^2}{\sigma_p^2 (\frac{1}{n_1} + \frac{1}{n_2})}$

多重比较方法 multiple comparison procedures 最小显著差异方法 least significant difference

$|\bar{x}_i - \bar{x}_j| \leq LSD = t_{\frac{\alpha}{2}}(n-k) \sqrt{MSE (\frac{1}{n_i} + \frac{1}{n_j})}$ LSD difference

• 双因素方差分析 two-way ~ 无/有交互作用 (non-) interaction (无/可重复~)

1. $SST = SSR + SSC + SSE$, $F_R = \frac{MSR}{MSE}$, $F_C = \frac{MSC}{MSE}$, $R^2 = \frac{SSR + SSC}{SST}$
 df: $kr-1$ $r-1$ $k-1$ $(k-1)(r-1)$ $\sum_i \sum_j (X_{ij} - \bar{x}_i - \bar{x}_j + \bar{x})^2$ "注: 前两项难以区分是 error 还是 interaction."

2. $SST = SSR + SSC + SSRC + SSE$ $\sum_i \sum_j \sum_k (X_{ijk} - \bar{x}_{ij})^2$ "如果只对 R/C 单 anova, 则 C/R 被含入残差, 故双 anova 优."
 df: $krm-1$ $r-1$ $k-1$ $(k-1)(r-1)$ $\sum_i \sum_j (\bar{X}_{ij} - \bar{x}_i - \bar{x}_j + \bar{x})^2$ $F_{RC} = \frac{MSRC}{MSE}$ "一种联合对因变量的附加效应."

五. 一元与多元线性回归

5.1 • 相关分析 correlation 相关系数 (Pearson) correlation coefficient $r = \frac{\sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}} = \frac{\bar{xy} - \bar{x}\bar{y}}{s_x s_y}$ (总体 ρ)

r 的检验: $t = |r| \sqrt{\frac{n-2}{1-r^2}} \sim t(n-2)$ ($H_0: \rho = 0$) "大样本下 r 很容易通过检验, $(s_r = \frac{\sqrt{1-r^2}}{\sqrt{n}})$ "

• 回归分析 regression 因变量 dependent value 自变量 independent ~ 但实际相关性可能并不显著

回归模型 $y = \beta_0 + \beta_1 x + \varepsilon$, $E(y) = \beta_0 + \beta_1 x$, " ε 为期望为 0, 方差保持 σ^2 的正态误差项"

回归方程, 估计的回归方程 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ 最小二乘法 method of least squares

$Q = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$, $\frac{\partial Q}{\partial \hat{\beta}_0} \Big|_{\hat{\beta}_0 = \hat{\beta}_0} = 0$, $\frac{\partial Q}{\partial \hat{\beta}_1} \Big|_{\hat{\beta}_1 = \hat{\beta}_1} = 0 \rightarrow \hat{\beta}_1 = \frac{\bar{xy} - \bar{x}\bar{y}}{s_x^2}$, $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

1. 回归统计 multiple R / r, R^2 判定系数 coefficient of determination = $\frac{SSR}{SST} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$
 $SST = SSR + SSE$ (回归平方和 + 残差平方和) $R^2 \in [0, 1]$
 df: $n-1$ (1) k $(n-2)$ $n-k-1$ 一元时 $R=r$ " $R > R^2$, R^2 解释更优"

adjusted R^2 / R_a^2 调整的判定系数 $R_a^2 = 1 - (1 - R^2) \cdot \frac{n-1}{n-k-1} = 1 - \frac{MSE}{MST} < R^2$

se 估计标准误差 = $\sqrt{\frac{SSE}{n-2}} = \sqrt{MSE}$ "随着 $k \uparrow$ $SSE \downarrow$, R^2 无论怎样 \uparrow 趋近 1"
 (二乘法使 se 取 min, or $\sum \hat{y}_i (y_i - y_i) = 0$)

2. 方差分析 df, SS, MS, F, $\alpha_F(P)$ 3. 回归参数 $\hat{\beta}_0, \hat{\beta}_1, S_{\hat{\beta}_0}, S_{\hat{\beta}_1}, t_{\hat{\beta}_0}, t_{\hat{\beta}_1}, P, \pm t_{\frac{\alpha}{2}}$ ($ck, n-k-1$)

显著性检验 $\left\{ \begin{array}{l} \text{线性关系检验 } F = \frac{MSR}{MSE} \sim F(1, n-2) \text{ "回归方程显著"} \\ \text{回归系数检验 } S_{\hat{\beta}_1} = \frac{se}{\sqrt{n} s_x}, t = \frac{\hat{\beta}_1}{S_{\hat{\beta}_1}} \sim t(n-2), S_{\hat{\beta}_0} = \frac{se \sqrt{\sum x^2}}{n s_x} \end{array} \right.$

一元时 $\frac{\hat{\beta}_1}{S_{\hat{\beta}_1}} = \frac{\hat{\beta}_1 \sqrt{n} s_x}{se}$ $F = \frac{n \hat{\beta}_1^2}{se^2} = t^2$ 三者一致, 多元时 F 与 t 有别. ($F = \frac{\sum (\sum \beta_k (X_{ik} - \bar{X}_k))^2}{k se^2}$)

$t = |r| \sqrt{\frac{n-2}{1-r^2}} = \frac{r \sqrt{n} s_y}{se} = \frac{\hat{\beta}_1 \sqrt{n} s_x}{se}$
 $(r^2 = \hat{\beta}_1 x_1 \cdot \hat{\beta}_1 y_1)$ $se^2 = \frac{n s_y^2 (1-R^2)}{n-2} = \frac{n s_y^2 (1-r^2)}{n-2}$

4. 估计与预测 平均值(均值) $E(y_0) = \beta_0 + \beta_1 x_0$

点估计 $\left\{ \begin{array}{l} \text{平均值(均值)} E(y_0) = \beta_0 + \beta_1 x_0 \\ \text{个别值(单值)} \hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 \end{array} \right.$
 区间估计 $\left\{ \begin{array}{l} \text{均值} \rightarrow \text{置信区间 } \pm t_{\frac{\alpha}{2}} S_{\hat{y}_0} \\ \text{单值} \rightarrow \text{预测区间 } \text{prediction } \pm t_{\frac{\alpha}{2}} \sqrt{se^2 + S_{\hat{y}_0}^2} \end{array} \right.$
 $se^2 \rightarrow \sigma^2$ $E(y) = \hat{y}_0 = \frac{se \sqrt{\sigma_x^2 + (\alpha_0 - \bar{x})^2}}{\sqrt{n} s_x}$

5. 残差分析 residual analysis 残差 $\left\{ \begin{array}{l} \text{外生化 external} \sim / \text{学生化} \\ \text{内生化 (学生化)} \end{array} \right.$ $\left. \begin{array}{l} \text{标准化 standardize / Pearson} \\ \text{半学生化 semi-studentized} \end{array} \right.$

残差图验证 ϵ 的正态性. $\frac{\sum_{i=1}^n \epsilon_i^2}{n}$ $\frac{\epsilon_i}{S_{\epsilon_i}}$ $\frac{\epsilon_i}{S_{\epsilon_i}}$ (除 i 外拟合) $\sqrt{S_{\epsilon_i}^2 - \frac{\epsilon_i^2}{n}}$

• 多元回归模型 multiple \sim 估计的多元回归方程 $\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$ 最小二乘法

R^2 (多重 \sim), R_a^2 , S_e 同定义. $F_{(k, n-k-1)}$, $t_{(n-k-1)}$ 同定义. $(\bar{x}_j \bar{x}_k - \bar{x}_j \bar{x}_k) (\beta_j) =$

多重共线性 multicollinearity "回归线性关系显著且每个 X_i 与 y 都显著" $(\bar{x}_j \bar{y} - \bar{x}_j \bar{y})$

常见处理: 剔除相关的 X (变量选择) $x_1, x_2, x_3, \dots, x_k$
避免用 t 对单个 β 检验; 共线性
对 y 的估计限定在 X 的范围内.

常见表现: X_i 间显著相关;
 F 显著但许多 β_j 的 t_j 不显著 (或反之);
 β_j sign 与预期相反; 容忍度 tolerance

方差扩大因子 VIF variance. $1 - R_i^2$
inflation factor $\frac{1}{1 - R_i^2}$ 以 X_i 为因, 剩余 X 为自判定.

逐步回归 stepwise regression: X_i $F?$ (严重共线: > 10)
向前选择 forward selection + 向后剔除 backward elimination

六. 时间序列与指数

6.1 • 时间序列 平稳序列 stationary series 非平稳序列 non- \sim

加/乘法模型 additive/multiplicative model $Y_t = T_t \cdot S_t \cdot C_t \cdot I_t$

增长率 growth rate (环比/定基) 平均 \sim average $\sim \sqrt[n]{\frac{Y_n}{Y_0}} - 1$

预测方法: $T \times S \times X \rightarrow$ 平滑预测法 $\left\{ \begin{array}{l} \text{简单平均 } F_{t+1} = \frac{1}{k} \sum_{i=1}^k Y_i \\ \text{移动平均 } F_{t+1} = \bar{Y}_t = \frac{Y_t + \dots + Y_{t-k+1}}{k} \\ \text{指数平滑 } F_{t+1} = \alpha Y_t + (1-\alpha) F_t \\ \text{平滑系数 } = \alpha Y_t + \alpha(1-\alpha) Y_{t-1} + \dots + (1-\alpha)^k Y_{t-k} \end{array} \right.$

(不考虑 C) $T \times S \times X \rightarrow$ 趋势预测法 (LS) 季节指数 seasonal index

(评估: MSE, MAD, MPE, MAPE...) $S \checkmark \rightarrow$ 时序分解, 季节性回归...

6.2 • 指数 $\left\{ \begin{array}{l} \text{个体} \sim \\ \text{总} \sim \end{array} \right.$

数量指标 $\sim q$ (LS)
质量指标 $\sim p$ (PS)

简单 \sim 综合 拉氏/帕氏指数
加权 \sim 平均 Laspeyres/Pasche

指数体系 $\frac{\sum q_1 p_1}{\sum q_0 p_0} = \frac{\sum q_1 p_0}{\sum q_0 p_0} \cdot \frac{\sum q_0 p_1}{\sum q_0 p_0}$ 权数 (同度量因素, 媒介 (基期/报告期) 因素) $I = \frac{\sum iW}{\sum W}$

综合评价指数: 指标体系, 无量纲化处理, 确定权重, 计算分析.
 $(x_i, \frac{x_i}{x_0}, \frac{x_i - \min}{\max - \min}, \dots)$



概率论与数理统计 <陈希孺>

一、事件的概率

1.1 • 主观概率 试验 事件 基本事件 随机事件 偶/必然事件

古典概率 $P(E) = \frac{M}{N}$ 几何概率 ← “等可能性” (古典定义)

统计概率 $\lim_{n \rightarrow \infty} \frac{m}{n} = P(E)$ (统计定义) “并不是定义了概率, 而只是一种估计的方法”

公理化定义 柯氏公理体系 (Ω, \mathcal{F}) : \mathcal{F} 为 Ω 的一个 σ -代数 (子集族)

\forall 事件 $A \in \mathcal{F}$, 满足 ① $0 \leq P(A) \leq 1$ ② $P(\Omega) = 1, P(\emptyset) = 0$ ③ $A_i \cap A_j = \emptyset (i \neq j)$
则集函数 P 是 (Ω, \mathcal{F}) 上的概率 (测度), (Ω, \mathcal{F}, P) 概率空间 $P(\cup A_i) = \sum P(A_i)$

• 古典概率的计算 排列 P_n^n 组合 $C_n^r, \binom{n}{r}$ 分堆 $\frac{n!}{r_1! \dots r_k!} (r_1 + \dots + r_k = n)$

常用公式: $\sum_{i=0}^n \binom{n}{i} = 2^n, \sum_{i=0}^n (-1)^i \binom{n}{i} = 0$ (二次展开, 多项展开)

$$\binom{n}{m} + \binom{n}{m+1} = \binom{n+1}{m+1}, \sum_{i=0}^k \binom{m}{i} \binom{n}{k-i} = \binom{m+n}{k}, \sum_{i=0}^m \binom{n+i}{i} = \binom{n+m+1}{m}$$

1.2 • 事件的蕴含、包含、相等 $A \subset B, A=B$ 事件的互斥、对立 $P(A \cap B) = 0, B = \bar{A}$

事件的和 (并) $C = A+B$ 加法定理: A_i 两两互斥, $P(\sum A_i) = \sum P(A_i)$

事件的积 (交)、差 $C = AB, A-B = A\bar{B}$ (特别地, $P(\bar{A}) = 1 - P(A)$)
布尔运算律 (De Morgan 律)

条件概率 $P(B) \neq 0$ 时, $P(A|B) = \frac{P(AB)}{P(B)}$ 事件的独立、乘法定理:

\triangleright 一组独立事件任一部分, 或其组合, 对立 $\Rightarrow P(A) = P(A|B)$ 一般地, $P(A) = P(A|B) P(B) = P(A)P(B)$, A_i 相互独立.

• 全概率公式 $B_i B_j = \emptyset (i \neq j)$, $\sum B_i = \Omega$ (完备事件群) Δ 不是“两两”! $P(\cup A_i) = \sum P(A_i)$

$$P(A) = \sum_i P(B_i) P(A|B_i)$$

贝叶斯公式 $P(B_i|A) = \frac{P(B_i) P(A|B_i)}{\sum_j P(B_j) P(A|B_j)}$ ← “由结果推原因”

二、随机变量及概率分布

2.1 • 随机变量 离散型、连续型 概率分布 (离散) $p_i = P(X = a_i)$

分布函数 $F(x) = P(X \leq x)$ $X \sim F$ 概率密度函数 $f(x) = F'(x)$

\triangleright ① $f(x) \geq 0$ ② $\int_{-\infty}^{\infty} f(x) dx = 1$ ③ $\int_a^b f(x) dx = P(b \leq X \leq a) = F(b) - F(a)$

重要分布: 二项分布 $P_i = \binom{n}{i} p^i (1-p)^{n-i}$ ($X \sim B(n, p)$), 泊松分布

$P_i = \frac{e^{-\lambda} \lambda^i}{i!}$ ($X \sim P(\lambda)$) ← 二项 $n \rightarrow \infty, np \rightarrow \lambda$ 的推导, 超几何分布

负二项分布、几何分布: $\binom{i+r-1}{r-1} p^r (1-p)^i$, $p(1-p)^i$ (令 $r=1$) $\frac{\binom{M}{m} \binom{N-M}{n-m}}{\binom{N}{n}}$

正态分布 $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ ($X \sim N(\mu, \sigma^2)$), 标准正态分布 $N(0, 1)$ / Z

指数分布 $f(x) = \lambda e^{-\lambda x}$ ($x > 0$), 威布尔分布 $f(x) = \lambda \alpha x^{\alpha-1} e^{-\lambda x^\alpha}$ ($x > 0$)

均匀分布 $f(x) = \frac{1}{b-a}$ ($x \in [a, b]$) ($X \sim R(a, b)$) | $X \sim R(0, 1) \Rightarrow$ ($\alpha > 1$)

$F^{-1}(X) \sim F$

2.2 • n维随机变量 (向量) $X = (X_1, \dots, X_n)$, $f(x_1, \dots, x_n)$

多项分布 $P_{(k_1, \dots, k_n)} = \frac{N!}{k_1! \dots k_n!} p_1^{k_1} \dots p_n^{k_n}$ ($\sum k_i = N$), 二维正态分布 $f(x_1, x_2) =$

边缘分布 $P(X_1 = a_{1k}) = \sum_{j_2, \dots, j_n} P(k, j_2, \dots, j_n)$ $\frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} e^{-\frac{[(x_1-a_1)^2/\sigma_1^2 - 2\rho(x_1-a_1)(x_2-b_2)/\sigma_1\sigma_2 + (x_2-b_2)^2/\sigma_2^2]}{2(1-\rho^2)}}$

$f_1(x_1) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x_1, x_2, \dots, x_n) dx_2 \dots dx_n$ ← "派生但并不包含耦合 (相互作用) 的信息"

条件概率分布 $P(X_1 = a_i | X_2 = b_j) = \frac{P_{ij}}{\sum_k P_{kj}}$ $f(x_1 | a \leq x_2 \leq b) = \int_a^b f(x_1, t) dt$

(n维类似有 $f(x_1, \dots, x_n) = g(x_1, \dots, x_m) h(x_{m+1}, \dots, x_n | x_1, \dots, x_m)$) $f(x_1 | x_2) = \frac{f(x_1, x_2)}{f_2(x_2)} \int_a^b f_2(x_2) dx_2$

(连续全概率公式) $f_1(x_1) = \int_{-\infty}^{\infty} f(x_1 | x_2) f_2(x_2) dx_2$ (若 $f(x_1 | x_2)$ 不是 x_2 函数 (独立) 则能提出, $= f(x_1)$)

随机变量的独立性 $f(x_1, \dots, x_n) = f_1(x_1) \dots f_n(x_n)$

| X_1, \dots, X_n 独立, $\forall a_i < b_i, A_i = \{a_i \leq X_i \leq b_i\}$, A_1, \dots, A_n 也独立; 反之亦成立.

X_1, \dots, X_n 独立, $Y_1 = g_1(X_1, \dots, X_m), Y_2 = g_2(X_{m+1}, \dots, X_n)$ 则 Y_1, Y_2 独立.

若 $f(x_1, \dots, x_n)$ 能分表为 $g_1(x_1), \dots, g_n(x_n)$ n个函数之积, X_1, \dots, X_n 独立且 $f_1(x_1)$ 与 $g_2(x_2)$

事件的指示变量 $X_i = \begin{cases} 1 & A_i \text{ 发生} \\ 0 & A_i \text{ 不发生} \end{cases}$ (A 与 X 独立性一致) 仅差一因子 (倍数)

2.3 • 随机变量函数的概率分布: (离散) 对 Y 不同的取值情况求和.

$Y = g(X_1, \dots, X_n)$

(一元) $y = g(x)$ $l(y) = f(h(y)) |h'(y)|$; 不单调时可能需加和多个同值, 如

假设 g 严格单调, 反函数 $g^{-1} = h$. $Y = X^2$ 时 $l(y) = \frac{1}{2\sqrt{y}} \cdot [f(\sqrt{y}) + f(-\sqrt{y})]$

(= / 多元) $X_1 = h_1(Y_1, Y_2, \dots), X_2 = h_2(Y_1, Y_2, \dots) \dots$ (坐标变换)

$J_{X/Y} = \left| \frac{\partial h_i}{\partial y_j} \right| = \frac{1}{J_{Y/X}}$, $l(y_1, \dots, y_n) = |J(y_1, \dots, y_n)| f(h_1(y_1, \dots, y_n), \dots)$

| $2X \sim f(x)$, 则 $X \sim \frac{1}{2} f(\frac{x}{2})$

(注意系数的正倒) $X \sim 2f(2x)$ ✓

$$Y = X_1 + X_2: f(y) = \int_{-\infty}^{\infty} f(y-x, x) dx \stackrel{i}{=} \int_{-\infty}^{\infty} f_1(y-x) f_2(x) dx$$

• 和的密度函数: 理解 / $P(Y \leq y)$ 写出积分区域或求导 / 引入变量为新坐标系

如 $(X_1, X_2) \sim N(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$, $Y = X_1 + X_2 \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2 + 2\rho\sigma_1\sigma_2)$

卡方分布 $f_n(x) = \frac{1}{\Gamma(\frac{n}{2}) 2^{\frac{n}{2}}} e^{-\frac{x}{2}} x^{\frac{n}{2}-1}$ ($X \sim \chi_n^2$) $X_1, \dots, X_n \stackrel{iid}{\sim} N(0, 1)$, 则

$\triangleright X_1 \sim \chi_n^2, X_2 \sim \chi_m^2$, 则 $X_1 + X_2 \sim \chi_{n+m}^2$ $Y = X_1^2 + \dots + X_n^2 \sim \chi_n^2$

$X_1, \dots, X_n \stackrel{iid}{\sim} \lambda e^{-\lambda x}$, 则 $2\lambda(X_1 + \dots + X_n) \sim \chi_{2n}^2$.

△不要落1!!

高的密度函数: $\dots Y = \frac{X_2}{X_1}: f(y) = \int_0^{\infty} x_1 f(x_1, x_1 y) dx_1 \stackrel{i}{=} \int_0^{\infty} x_1 f_1(x_1) f_2(x_1 y) dx_1$

t分布 (学生t分布) $f_n(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} (1 + \frac{x^2}{n})^{-\frac{n+1}{2}}$ ($X \sim t_n$) $X_1 \sim \chi_n^2, X_2 \sim N(0, 1)$

F分布 $f_{m,n}(x) = m^{\frac{m}{2}} n^{\frac{n}{2}} \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2}) \Gamma(\frac{n}{2})} x^{\frac{m}{2}-1} (mx+n)^{-\frac{m+n}{2}}$ ($X \sim F_{m,n}$) $\stackrel{i}{=} \frac{X_2}{\frac{X_1}{n}} \sim t_n$

三大分布重要结论: ① $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, $(F_{m,n}(a) = F_{n,m}(1-a))$ $X_1 \sim \chi_n^2, X_2 \sim \chi_m^2$

$\bar{X} = \frac{\sum X_i}{n}, S^2 = \frac{\sum (X_i - \bar{X})^2}{n-1}$, 则 $\frac{\sum (X_i - \bar{X})^2}{\sigma^2} \stackrel{i}{=} \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$

② 同①假设, 则 $\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$. ③ 同①假设, 且 $Y_1, \dots, Y_m \stackrel{iid}{\sim} N(\mu_2, \sigma_2^2)$, 则 $\frac{S_Y^2}{\sigma_2^2} / \frac{S_X^2}{\sigma_1^2} \sim F_{m-1, n-1}$

④ 同③假设, 若 $\sigma_1^2 = \sigma_2^2$, 则

$$\frac{\sqrt{nm(n+m-2)} (\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\sum (X_i - \bar{X})^2 + \sum (Y_j - \bar{Y})^2}} \sim t_{n+m-2}$$

<2.4> • $\Gamma(x)$ 函数 $\Gamma(x) = \int_0^{\infty} e^{-t} t^{x-1} dt$ $\Gamma(x) = (x-1)\Gamma(x-1)$, $\Gamma(x) = (x-1)!$

B(x,y) 函数 $B(x,y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$ 余元 $\Gamma(p)\Gamma(1-p) = \frac{\pi}{\sin p\pi}$

$G_s(n)$ 函数 $G_s(n) = \int_0^{\infty} x^n e^{-\lambda x^2} dx = \lambda^{-\frac{n+1}{2}} \begin{cases} \frac{\sqrt{\pi}}{2} \cdot \frac{1}{2} \cdot \frac{3}{2} \dots n \text{ even} \\ \frac{1}{2} \cdot \frac{3}{2} \cdot \frac{5}{2} \dots n \text{ odd} \end{cases} \rightarrow (n-1)!!$
 $= \lambda^{-\frac{n+1}{2}} \cdot \frac{1}{2} \Gamma(\frac{n+1}{2})$

• 斯特林公式 $\lim_{n \rightarrow \infty} n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$

\triangleright 证明: Wallis 公式 $(I_{2n+1} < I_{2n} < I_{2n-1})$ 夹逼得 $\lim_{n \rightarrow \infty} \frac{(2n)!!}{(2n-1)!!} \frac{1}{2n+1} = \frac{\pi}{2}$

$\Sigma \rightarrow \int \ln$ 的面积, 鞍点法 $n! = \int_0^{\infty} e^{-t} t^n dt = \int_0^{\infty} e^{g(t)} dt = e^{g(n)} \int_0^{\infty} \frac{1}{2} g'(n)(t-n)^2 dt$

\triangleright 注意 $\frac{n!}{(n-i)!}$ 在 $n \rightarrow \infty$ 时 $\sim n^i$ 不带 e^i !
 $g(t) = -t + n \ln t$
 $g'(t) = -1 + \frac{n}{t}, g''(t) = -\frac{n}{t^2}, g'(n) = 0$ $= n^n e^{-n} \sqrt{2\pi n}$

• ① 的证明: 证 \bar{X} 与 $\sum (X_i - \bar{X})^2$ 独立. $Y_1 = \frac{\sum X_i}{\sqrt{n}}, \sum Y_i^2 = \sum X_i^2$
 ②③④ 顺推 $Y = aX, a = (\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}})$, 故 $\sum_{i=2}^n Y_i^2 = \sum (X_i - \bar{X})^2$

$(X_1, \dots, X_n) \sim \frac{1}{(\sqrt{2\pi}\sigma)^n} e^{-\frac{\sum (X_i - \mu)^2}{2\sigma^2}} \rightarrow \sum Y_i^2 - \mu\sqrt{n}Y_1 + n\mu^2$ 则 $\begin{cases} Y_1 \sim N(\sqrt{n}\mu, \sigma^2) \\ Y_2 \dots Y_n \sim N(0, \sigma^2) \text{ iid.} \end{cases}$

• 二项分布推出泊松分布与正态分布(可视为中心极限)

$$\binom{n}{i} p^i (1-p)^{n-i} \xrightarrow{i \rightarrow np + \epsilon, n \rightarrow \infty} \frac{e^{-\lambda} \lambda^i}{i!}$$
 (Stirling公式)

$$\frac{1}{\sqrt{2\pi np(1-p)}} \cdot \frac{1}{(1 + \frac{\epsilon}{np})^{np+\epsilon}} \cdot \frac{1}{(1 - \frac{\epsilon}{n(1-p)})^{n(1-p)-\epsilon}}$$
 与 $e^{\frac{x}{2}}$ 存在 $(-\frac{\epsilon}{2}x)$ 的一阶差.

$$= \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{\epsilon^2}{2np(1-p)} + \frac{\epsilon^2}{2np(1-p)}}$$
 这里必须考虑!

• 重要分布总结: Bernoulli (two-point), binomial, multinomial, Poisson.

(Pascal) negative binomial, geometric, hypergeometric, exponential, Weibull, normal (Gauss), ln-normal, logarithmic, Rayleigh, Chi-square, student, Fisher, Cauchy, Laplace, gamma, beta, multivariate beta (Dirichlet), uniform, maxmin...

对数正态分布: $Y \sim N(\mu, \sigma^2), X = e^Y \sim LN, f(x) = \frac{1}{\sqrt{2\pi\sigma}} \frac{1}{x} e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}}, x_m = e^{\mu - \sigma^2}$

$\sigma^2 \uparrow$ 偏态, 越明显.

$E(x) = e^{\mu + \frac{\sigma^2}{2}}, V(x) = e^{2\mu + \sigma^2} (e^{\sigma^2} - 1)$

对数级数分布: $-\frac{(1-p)^i}{i \ln p}, \quad \binom{n}{m} = (-1)^m \binom{m-n-1}{m}$. NB名称来自于 $(1-x)^{-r} = (1+x+x^2+\dots)^r$

NB (负二项) 分布 = $\sum_{r=1}^{\infty} \text{geometric dis.}$, $E(x) = \frac{r(1-p)}{p}, V(x) = \frac{r(1-p)}{p^2}$.

$\frac{r(1-p)}{p} = \lambda, \quad r \rightarrow \infty, p \rightarrow 1, \quad \text{Poisson. (e^{\lambda} \text{ 形式})}$

瑞利分布: 二维独立高斯分布的模长. $L(\mu, r), F(x) = \frac{1}{1 + e^{-\frac{x^\mu}{r}}}, f(x) = \frac{r}{x^2} (1 + e^{-\frac{x^\mu}{r}})^{-2}$

$x > 0, f(x) = \frac{x}{\sigma^2} e^{-\frac{x^2}{2\sigma^2}} (\sqrt{x^2} = x)$ ($F_{st}(x) = \frac{1}{1+e^{-x}}$) $V(x) = \frac{\pi^2}{3} (E(x)=0)$

$E(x) = \sqrt{\frac{\pi}{2}} \sigma, V(x) = (2 - \frac{\pi}{2}) \sigma^2$ $m(x) = \pi x \sec(\pi x)$

柯西分布: $f(x) = \frac{1}{\pi} \frac{y}{(x-x_0)^2 + y^2}$ ($\frac{1}{x} \cdot \frac{1}{1+x^2}$) 不存在期望, 各阶矩均不存在. 估计 x_0 适合用中位数.

(x_0, y) $X_1, X_2 \text{ iid. } N(0, 1), \text{ 则 } Y = \frac{X_1}{X_2} \sim C(0, 1)$ (自由度为1的t分布) (或 $U(-\frac{\pi}{2}, \frac{\pi}{2})$ 上 $y = \tan x$ 的分布)

拉普拉斯分布: $f(x) = \frac{1}{2\lambda} e^{-\frac{|x-\mu|}{\lambda}}$ $E(x) = \mu, V(x) = 2\lambda^2 (k.o.v=3)$ (或标准正态高的分布 $\frac{X_1}{X_2}$)

伽马 Γ 分布: $f(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}$. $Ga(\alpha, \lambda) + Ga(\beta, \lambda) = Ga(\alpha+\beta, \lambda)$

$\alpha=1$, 为指数分布; $\alpha=\frac{n}{2}, \lambda=\frac{1}{2}$, 为卡方分布 $X^2 \sim Ga(\frac{1}{2}, \frac{1}{2})$ $E(x) = \frac{\alpha}{\lambda}, V(x) = \frac{\alpha}{\lambda^2}$

$X_n^2 \sim Ga(\frac{n}{2}, \frac{1}{2})$

α : 形状参数 λ : 尺度参数 ($\lambda X \sim Ga(\alpha, 1)$)

贝塔 B 分布: $f(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$ $E(x) = \frac{\alpha}{\alpha+\beta}, V(x) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ Multinomial ~ Dirichlet, Binomial ~ Beta 为共轭先验.

狄利克雷分布: $f(x) = \frac{1}{B(\vec{\alpha})} \prod_i x_i^{\alpha_i-1}, B(\vec{\alpha}) = \frac{\prod \Gamma(\alpha_i)}{\Gamma(\sum \alpha_i)}$ 且 $\|X\|=1$. Beta 为其 marginal.

最大最小值: $F_{max}(x) = F(x)^n, F_{min}(x) = 1 - (1-F(x))^n$ 再求导.

三、随机变量的数字特征

3.1 数学期望 (均值) $E(\sum X) = \sum E(X)$, $X_i \text{ iid}, E(\pi X) = \pi E(X)$

要求绝对收敛 $\int_{-\infty}^{\infty} |x|f(x)dx < \infty$. $E(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx$ (故 $E(cX) = cE(X)$)

$\triangleright E(X_n^2) = n$. $E(\frac{1}{X_n^2}) = \frac{1}{n+2}$. $E(F_{m,n}) = \frac{n}{n+2}$. $E(Z^n) = (n-1)!!$ 偶/0 奇
 $E(X^n)$, N 分布, 只需对 $(X-\mu+\mu)^n$ 展开.

条件均值 $E(Y|X) = \int_{-\infty}^{\infty} yf(y|x)dy$. $E(Y) = E_y(E_y(Y|X))$ 中位数

方差 标准差 $V(X) = E(X^2) - E^2(X)$ (平行轴定理) $V(X+c) = V(X)$.

$X_i \text{ iid}, V(\sum X) = \sum V(X)$. $E((\hat{\theta} - \theta)^2) = V(\hat{\theta}) + (E(\hat{\theta}) - \theta)^2$ $V(cX) = c^2V(X)$.

标准化 $Z = \frac{X-\mu}{\sigma}$ $\triangleright V(X_n^2) = 2n$. $V(t_n) = \frac{n}{n-2}$. $V(F_{m,n}) = \frac{2n^2(m+n-2)}{m(n-2)^2(n-4)}$
 $(V(s^2) = 2(n-1)\sigma^4)$

矩 原点矩 中心矩 偏度 $\frac{\mu_3}{\mu_2^3}$ 峰度 $\frac{\mu_4}{\mu_2^2} (-3)$ 左偏, 右偏分布

协方差 $Cov(X,Y) = E((X-E(X))(Y-E(Y))) = E(XY) - E(X)E(Y)$

$X, Y \text{ iid}, Cov(X,Y) = 0$. $Cov^2(X,Y) \leq \sigma_x^2 \sigma_y^2$, $E(XY) \in E(X)E(Y) \pm D(X)D(Y)$

相关系数 $Corr(X,Y) = \rho_{XY} = \frac{Cov(X,Y)}{\sigma_x \sigma_y} \in [-1, 1]$ $\triangleright Cov \text{ or } Corr = 0$ 仅能称“不相关”,

最小二乘法 $E(Y-a-bX)^2 \rightarrow \min$ \uparrow 并不一定“独立”!
 线性关系, $(f(x,y) \neq f(x)f(y))$
 (二维正态, 是“独立”的)

$$= \sigma_y^2 + b^2 \sigma_x^2 - 2b \sigma_{XY} + (a - E(Y) + bE(X))^2$$

$$\min \sigma_y^2 (1 - \rho^2)$$

$$a_m = E(Y) - b_m E(X)$$

$$b_m = \frac{\sigma_{XY}}{\sigma_x^2} = \frac{\sigma_y}{\sigma_x} \rho_{XY}$$

\triangleright 柯西不等式的推论: $E(\frac{1}{X}) \geq \frac{1}{E(X)}$, $E(X) \leq \sqrt{E(X^2)}$

另一个推论: $V(X_1) = V(X_2)$, 则 $Cov(X_1+X_2, X_1-X_2) = 0$.

“均值趋于期望”

3.2 大数定理 $X_1, \dots, X_n \text{ iid}$. 且均值为 μ , 则 $\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| \geq \varepsilon) = 0$ 对 $\forall \varepsilon > 0$.

(存在方差 (σ^2))

“依概率收敛”

马尔科夫不等式 $P(Y \geq \varepsilon) \leq \frac{E(Y)}{\varepsilon}$, $Y \geq 0, \forall \varepsilon > 0$. (独立收敛)

契比雪夫不等式 若 $V(Y)$ 存在, 则 $P(|Y - E(Y)| \geq \varepsilon) \leq \frac{V(Y)}{\varepsilon^2}$

“依分布收敛”

中心极限定理 $X_i \text{ iid}$. $\forall x$ 有 $\lim_{n \rightarrow \infty} P(\frac{1}{\sqrt{ns}}(X_1 + \dots + X_n - n\mu) \leq x) = \phi(x)$

Lindeberg-Lévy CLT

特殊地, de Moivre-Laplace Th:

其中 ϕ 是标准正态分布函数 $\phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$.

X_i 二项分布, $\lim_{n \rightarrow \infty} P(\frac{\sum X - np}{\sqrt{npq}} \leq x) = \phi(x)$.

(有 $\phi(-x) = 1 - \phi(x)$)

\triangleright 正态分布 (同二项分布, 泊松分布一样) 是再生分布, 叠加, 分割仍为正态分布.

(线性变换稳定性: $X \sim N(\mu, \Sigma)$, $Y = CX \sim N(C\mu, C\Sigma C^T)$)

指数分布有无后效性, 且 $n \min X_i$ 与 X 同指数分布.

13.37 • 矩母函数 MGF、特征函数 类似于 Fourier, Laplace 变换, 可以把 Σ 的卷积化为乘积. 一个 MGF 决定了一个分布. (全部信息)

$$m(t) = E(e^{tx}) = \int_{-\infty}^{\infty} e^{tx} f(x) dx \quad (\Sigma e^{tx} f(x))$$

$$= E(1) + tE(x) + \frac{t^2}{2!} E(x^2) + \dots \quad \text{故 } n \text{ 阶矩 } E(x^n) = \frac{d^n}{dt^n} m(t)$$

指数分布 $m(t) = \frac{\lambda}{\lambda - t}$. 正态分布 $m(t) = e^{t\mu + \frac{t^2 \sigma^2}{2}}$ (多瓦 $e^{t\mu + \frac{t^2 \Sigma t}{2}}$)

卡方分布 $m(t) = \frac{1}{(1-2t)^{\frac{n}{2}}}$. $\varphi(t) = e^{-\frac{t^2}{2}}$ | $X \sim N(\mu, \Sigma)$. $E(x^T A x) = \text{tr}(A\Sigma) + \mu^T A \mu$.

二项分布 $m(t) = (1-p+pe^t)^n$. $m_{x^T A x}(t) = |K|^{-\frac{1}{2}} e^{-\frac{1}{2} \mu^T (I-K^{-1}) \Sigma^{-1} \mu}$

伽马分布 $m(t) = (\frac{\lambda}{\lambda-t})^\alpha$. $\varphi(t) = (\frac{\lambda}{\lambda-it})^\alpha$. 其中 $K = I - 2tA\Sigma$.

泊松分布 $m(t) = e^{\lambda(e^t-1)}$. $V(x^T A x) = 2\text{tr}(A\Sigma)^2 + 4\mu^T A \Sigma A \mu$.

$\varphi(t) = E(e^{itx})$ (以上 $m(t)$ 换为 $m(it)$ 即可). $\text{Cov}(Bx, x^T A x) = 2B\Sigma A \mu$.

$\Delta Y = ax + b$, $m_Y(t) = e^{bt} m_X(at)$. $X_i \text{ iid}$, $Y = \sum_{i=1}^n X_i$, $m_Y(t) = \prod_{i=1}^n \varphi_{X_i}(t)$. $(x-\mu)^T \Sigma^{-1} (x-\mu) \sim \chi^2(n)$.

• 正态分布相关来源: 1. central limit th. $Y = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{X_1 + \dots + X_n - n\mu}{\sigma/\sqrt{n}}$

$\varphi_Y(t) = e^{-i\mu t \frac{\sqrt{n}}{\sigma}} \varphi_X^n(\frac{t}{\sqrt{n}\sigma})$ 证明其形式为 $e^{-\frac{t^2}{2}}$. 令 $\eta = \frac{t}{\sqrt{n}\sigma}$, $\lim_{n \rightarrow \infty} \ln \varphi_Y(t) =$

$\lim_{\eta \rightarrow 0} \frac{t^2}{\sigma^2} \cdot \frac{\ln \varphi_X(\eta) - i\mu \eta}{\eta^2} \stackrel{L'H}{=} \lim_{\eta \rightarrow 0} \frac{t^2}{\sigma^2} \cdot \frac{\varphi_X'' \eta - \varphi_X'^2}{2\varphi_X^2} = -\frac{t^2}{2}$. 得证.

(已知 $\varphi_X(0) = 1$, $\varphi_X'(0) = i\mu$, $\varphi_X''(0) = -(\mu^2 + \sigma^2)$.)

2. Gauss (1809) error dis. $L(\theta) = f(x_1 - \theta) \dots f(x_n - \theta)$ ($f(x)$ 为误差函数)

极大似然估计 $\hat{\theta} = \bar{x} \rightarrow L(\theta)$: 令 $g = \frac{L'}{L}$, 则求导为 0 要求 $\sum g(x_i - \hat{\theta}) = 0$.

由 $\hat{\theta} = \bar{x}$ 发现 $g(x)$ 的性质 ($n=2$, $x_1 - \bar{x} = -(x_2 - \bar{x})$ 得 $g(x) = -g(-x)$, 且 $g(x) = cx$. $n=m+1$, $mg(-x) + gm(x) = 0$ 得 $g(cx) = cg(x)$)

故 $f(x) = C e^{cx^2}$. 归一化后得 $N(0, \sigma^2)$. 解释了最小二乘法.

3. Herschel & Maxwell symmetry of sp.

正交独立 + 旋转对称: $f(x)f(y) = p(x,y) = g(r,\theta) = g(r) = g(\sqrt{x^2+y^2})$

由于 $y=0$ 有 $g(x) = f(x) \cdot f(0)$. 令 $h(x) = \ln \frac{f(x)}{f(0)}$, 则 $h(x) + h(y) = h(\sqrt{x^2+y^2})$, $h(x) = cx^2$. (Maxwell 分布律类似) 故 $f(x) = \sqrt{\frac{\alpha}{\pi}} e^{-\alpha x^2}$.

4. Landon stable noise. 累加微小的随机噪声不改变分布模式 (仅层级/方差)

$X \sim p(x, \sigma^2)$, $\varepsilon \sim q(\varepsilon)$, $X' = X + \varepsilon$ 应有 $X' \sim p(x, \sigma^2 + V(\varepsilon))$.

X' 的 pdf $f(x) = \int p(x-\varepsilon, \sigma^2) q(\varepsilon) d\varepsilon \stackrel{\text{对 } p \text{ 展开}}{=} p - \frac{\partial p}{\partial x} \int \varepsilon q(\varepsilon) d\varepsilon + \frac{1}{2} \frac{\partial^2 p}{\partial x^2} \int \varepsilon^2 q(\varepsilon) d\varepsilon + o(\varepsilon^2)$.

把 p 对方差展开 $f(x) = p + \frac{\partial^2 p}{\partial \sigma^2} \cdot \bar{\varepsilon}^2$. 扩散方程 $\frac{1}{2} \frac{\partial^2 p}{\partial x^2} = \frac{\partial p}{\partial \sigma^2} \rightarrow p(x, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{x^2}{2\sigma^2}}$



5. maximum entropy. $H(p) = -\int p(x) \log p(x) dx$ 给定分布的 μ, σ^2 , 熵最大的为正态分布:

一个分布的熵总是小于相对熵 $H(p) \leq -\int p(x) \log q(x) dx$ ($\int p(x) \log \frac{q(x)}{p(x)} dx \leq 0$, 利用 $\log x \leq 1-x$)

代入 $q(x) = N(\mu, \sigma^2)$ 得 $H(p) \leq \frac{1}{2\sigma^2} \int p(x) (x-\mu)^2 dx + \log \sqrt{2\pi}\sigma = \frac{1}{2} + \log \sqrt{2\pi}\sigma$. (可在 $p(x) = N(\mu, \sigma^2)$ 时取等(max).)

6. Galton board. $\binom{n}{k} \frac{k-\frac{n}{2}+x}{2} \frac{2}{\sqrt{2\pi n}} e^{-\frac{x^2}{2n}}$ 与 binomial/poisson 逼近类似. 略.

四. 参数估计 Δ 参数/变量的范围!

4.1. 参数估计 假设检验 总体“无限总体” 样本 独立随机样本

统计量 样本方差 样本矩 样本均值 估计量

点估计 $X_1, \dots, X_n \stackrel{iid}{\sim} f(x; \theta_1, \dots, \theta_k)$ 估计 θ_i 或 $g(\theta_1, \dots, \theta_k)$, $\hat{\theta}_i(X_1, \dots, X_n) \rightarrow \theta_i$.

• 矩估计法 样本矩 \rightarrow 总体矩 $a_m(\theta_1, \dots, \theta_k) = a_m$ 解得 $\hat{\theta}_i = \hat{\theta}_i(X_1, \dots, X_n)$.
修正有 s 替代 $\sqrt{s^2}$ 等. 一般能用低阶矩不用高阶. $\hat{g} = \hat{g}(X_1, \dots, X_n) = g(\hat{\theta}_1, \dots, \hat{\theta}_k)$.

(MLE)

极大似然估计法 样本的分布 $L(x_1, \dots, x_n; \theta_1, \dots, \theta_k) = f(x_1; \theta_1, \dots, \theta_k) \dots f(x_n; \theta_1, \dots, \theta_k)$

$\frac{\partial \ln L}{\partial \theta_i} = 0$, $\ln L = \sum \ln f(x_i; \theta_1, \dots, \theta_k)$. “似然函数” 看作 θ_i 的函数. $L(x_1, \dots, x_n; \theta_1^*, \dots, \theta_k^*) = \max L(x_i; \theta_i)$

$N(\mu, \sigma^2)$ 矩和似然均给出 \bar{x}, m_2 . 一般更优良但个别情况会给出很不理 θ_i^* 想要的结果. 要求有分布具体的参数形式.

$R(\theta, \theta)$ 矩 \bar{x} , 似然 $\max(X_i)$.

贝叶斯法 先验分布 $h(\theta)$ 后验分布 $h(\theta | X_1, \dots, X_n) = \frac{h(\theta) f(X_1, \theta) \dots f(X_n, \theta)}{p(X_1, \dots, X_n)}$.
(广义先验密度 $\int h(\theta) d\theta \neq 1$) $p(X_1, \dots, X_n) = \int h(\theta) f(X_1, \theta) \dots f(X_n, \theta) d\theta$.

先验分布的选取: “同等无知原则” 之后常采用均值作为估计 $\hat{\theta} = \int \theta h(\theta | X_1, \dots, X_n) d\theta$.
(争议) $h(\mu) = 1, h(\sigma) = \frac{1}{\sigma}$ 等. (MLE 可以视为 $h(\theta | X) \max$ 且 $h(\theta) = 1$)

共轭先验分布 $h(\theta | X_1, \dots, X_n)$ 与 $h(\theta)$ 同一分布族, 则 $h(\theta)$ 共轭于似然函数 $f(X | \theta)$.
如 Bin: Beta, Poisson: Gamma, Uni: Pareto, 或称 θ 的共轭先验分布为 $h(\theta)$.
(A) Normal: Normal, (σ^2) Normal: IGamma...

4.2. 点估计的优良性准则 无偏性 $(X \sim U(a, \frac{1}{\lambda}) \sim I(Ua))$
(结合大数定理, 可以理解为 $\forall (\theta_1, \dots, \theta_k), E_{\theta_1, \dots, \theta_k}(\hat{g}(X_1, \dots, X_n)) = g(\theta_1, \dots, \theta_k)$ ($\forall \theta, E(\hat{g}) = g$)

$\frac{\sum_{i=1}^n \hat{g}_i}{N}$ 依概率收敛于 g . 如 $E(\bar{X}) = \mu, E(s^2) = \sigma^2, E(cs) < \sigma$ (可通过修正系数 $c_n s$, 对正态总体, $c_n = \frac{\sqrt{2} \Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})}$)

均方误差 $M_0(\theta) = E_0(\hat{\theta}(X_1, \dots, X_n) - \theta)^2 \stackrel{E(\hat{\theta}) = \theta \text{ 时}}{=} V_0(\hat{\theta})$.

最小方差无偏估计 (MVU) $\forall \theta, \theta_1, V_0(\hat{\theta}) \leq V_0(\theta_1)$ 且 $\hat{\theta}$ 为 MVE.

• Fisher 信息量 $0 = \frac{d}{d\theta} \int_{-\infty}^{\infty} f(x, \theta) dx = \int_{-\infty}^{\infty} \frac{\partial f(x, \theta)}{\partial \theta} dx = E\left(\frac{\partial \ln f(x, \theta)}{\partial \theta}\right)$

$$I(\theta) = V\left(\frac{\partial \ln f(x, \theta)}{\partial \theta}\right) = \int_{-\infty}^{\infty} \left(\frac{\partial \ln f(x, \theta)}{\partial \theta}\right)^2 f(x, \theta) dx = \int_{-\infty}^{\infty} \left(\frac{\partial f}{\partial \theta}\right)^2 \frac{1}{f} dx$$

(由 $0 = \frac{d}{d\theta} \int_{-\infty}^{\infty} f(x, \theta) dx$ 得 $\int_{-\infty}^{\infty} \frac{\partial^2 \ln f}{\partial \theta^2} f dx + \int_{-\infty}^{\infty} \left(\frac{\partial \ln f}{\partial \theta}\right)^2 f dx = 0$, 故 $I(\theta) = -E\left(\frac{\partial^2 \ln f(x, \theta)}{\partial \theta^2}\right)$)

样本的总信息量 $nI(\theta)$ ($V(\sum \frac{\partial \ln f}{\partial \theta}) \stackrel{iid}{=} \sum V(\frac{\partial \ln f}{\partial \theta})$) 亦即衡量似然函数的方差或 logL 的峰聚程度。

克拉美-劳不等式 $V(\hat{g}) \geq \frac{(g'(\theta))^2}{nI(\theta)}$ ($V\hat{g} \rightarrow g(\theta)$) (简写为 $\Delta^2 \hat{\theta} \geq \frac{1}{I(\theta)}$)

(CRLB: Cramer-Rao Lower Bound) | 证明: $Cov(\hat{g}, \sum \frac{\partial \ln f}{\partial \theta})^2 \leq V(\sum \frac{\partial \ln f}{\partial \theta}) V(\hat{g})$

• 相合性 (一致性) $\lim_{n \rightarrow \infty} P(|\hat{g}(X_1, \dots, X_n) - g(\theta_0)| \geq \epsilon) = 0$ $E(\hat{g} - \sum \frac{\partial \ln f}{\partial \theta}) = \int_{-\infty}^{\infty} \frac{\partial \ln \pi f}{\partial \theta} \hat{g} \pi f dx = \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} \hat{g} \pi f dx$

(样本大小无限扩大时, 估计依概率收敛于待估值) | 无偏有 $\frac{\partial g}{\partial \theta}$ 得法。
如样本的 $m_2, S^2 \rightarrow \sigma^2$ | 正态有 $I(\mu) = \frac{1}{\sigma^2}$, $I(\sigma^2) = \frac{1}{2\sigma^4}$

渐近正态性 $n \rightarrow \infty$ 时 \hat{g} 的分布 \rightarrow Normal 大/小样本性质/方法 (已知 μ, σ)
如又 (CLT) 即为 ($\rightarrow \infty$ 与否)

4.3 • 区间估计 $[\hat{\theta}_1(X_1, \dots, X_n), \hat{\theta}_2(X_1, \dots, X_n)]$ 置信系数 $1-\alpha$ 置信区间/界 (上, 下)

枢轴变量法: 1. 找一个与待估 $g(\theta)$ 有关 (置信水平 $\geq 1-\alpha, \forall \theta$)

的统计量 T (常为点估计). 2. 找一枢轴变量 $S = S(T, g(\theta))$, 其分布 F 应与

3. $S(T, g(\theta))$ 能进行改写 $a \leq S \leq b \rightarrow A(a, b, T) \leq g(\theta) \leq B(a, b, T)$. F 的上/下分位点, θ 无关.

贝伦斯-费歇尔 Behrens-Fisher 问题 $P(n_1 \frac{S}{\sigma_1} \leq S(T, g(\theta)) \leq n_2 \frac{S}{\sigma_2}) = 1-\alpha$ (见 5.1)

大样本法: 利用 n 很大时极限分布 (正态) 改写为 $g(\theta) \in [A, B]$ ($1-\alpha$)
(一般 $n > 30$, 其实大小区间估计本就不意义不大) (离散型枢轴法不易使用)

如 $\frac{X - np}{\sqrt{np(1-p)}}$ 近似 $\sim N(0, 1)$, 近似得 $\hat{p} \pm u_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ ($n \gg 1$)

贝叶斯法: "后验信度" $\int_{\theta_1}^{\theta_2} k(\theta | X_1, \dots, X_n) d\theta = 1-\alpha$ (上下界即改为 $-\infty / +\infty$)

常选择使 $\hat{\theta}_2 - \hat{\theta}_1$ 最小 (另一种为 $\int_{\theta_2}^{+\infty} = \frac{\alpha}{2}, \int_{-\infty}^{\theta_1} = \frac{\alpha}{2}$ 选择)

求解若易但先验分布的确定是问题, 而奈曼法确定分布很复杂.

奈曼理解: "抽取大量的区间, 有 $(1-\alpha)$ 的部分的区间里 (Neyman) 包含有真值. \leftarrow 真值有 $(1-\alpha)$ 的概率落在这一区间里" (Bayes) (大样本法有时无法做到很大难以控制误差)

| MVUE 的唯一性证明: 若 $\hat{\theta}_1, \hat{\theta}_2$ 均为 MVUE, 则 $\frac{\hat{\theta}_1 + \hat{\theta}_2}{2}$ (同无偏) 有

$$V\left(\frac{\hat{\theta}_1 + \hat{\theta}_2}{2}\right) \geq V_{\min} = V(\hat{\theta}_1) = V(\hat{\theta}_2) \quad V\left(\frac{\hat{\theta}_1 + \hat{\theta}_2}{2}\right) = \frac{1}{4} E(\hat{\theta}_1 - \theta + \hat{\theta}_2 - \theta)^2$$

$$= \frac{1}{4} E(\hat{\theta}_1 - \theta)^2 + \frac{1}{4} E(\hat{\theta}_2 - \theta)^2 + \frac{1}{2} E(\hat{\theta}_1 - \theta)(\hat{\theta}_2 - \theta)$$

故必须取等且 $\hat{\theta}_1$ 与 $\hat{\theta}_2$ 的任意线性组合也为 MVUE. $= \frac{1}{2} V_{\min} + \frac{1}{2} \sqrt{V_{\min} \cdot V_{\min}} = V_{\min}$

取等有 $\hat{\theta}_1 = a\hat{\theta}_2 + b$, 而 $E(\hat{\theta}_1) = E(\hat{\theta}_2)$. $b = (1-a)\theta$
故其同族均为 $a\hat{g} + (1-a)\theta$ 的带真值的平移, 不存在两个仅用样本的 MVUE.

五、假设检验

5.1

• 原假设 (零假设) H_0 对立假设 (备择假设) H_1

"统计上的显著性 (差异不能被随机误差解释) ≠ 现实中的重要性"

检验统计量 接受/否定域 临界值

简单/复合假设 赘余参数

"检验水平/显著性水平/苛刻度(α)/弃真错误"

• 功效函数 $\phi: \bar{x} \geq C \checkmark \text{ oth. } X, \beta_\phi(\lambda) = P_\lambda(\bar{X} < C)$

第一/二类错误 $\begin{cases} \beta_\phi(\theta), \theta \in H_0 \\ 1 - \beta_\phi(\theta), \theta \in H_1 \end{cases}$ 检验水平 α s.t. $\beta_\phi(\theta) \leq \alpha, \forall \theta \in H_0$

一致最优检验 对 $\forall \alpha$ 水平的检验 $\beta_\phi(\theta) \geq \beta_g(\theta), \forall \theta \in H_1$, 则 ϕ 为 α 水平 UBT. (α 尽量 min)

"一致即处处"的条件很高, 仅单参且单边或一些特例存在"

• 重要参数检验方法: 1° $H_0: \theta \geq \theta_0$, 2° $H_0: \theta \leq \theta_0$, 3° $H_0: \theta = \theta_0$

① 基于点估计

μ, σ^2 已知: 1° $\phi: \bar{x} \geq C = \theta_0 - \frac{\sigma}{\sqrt{n}} U_{\alpha}$ $\beta_\phi(\theta) = \Phi(\frac{\sqrt{n}(\theta_0 - \theta)}{\sigma} - U_{\alpha})$ (2, 3° 略)

$T(\theta) \leftrightarrow \theta_0, C$ (计算功效函数, "问题提法的倾向性" 3° 是双侧检验而非一致最优(不存在) (3°: $|\bar{x} - \theta_0| \leq C\sqrt{n}$) (1°, 2° 是)

由水平确定 C) μ, σ^2 未知: $\psi: \sqrt{n}(\bar{x} - \theta_0) \geq -t_{\alpha}(n-1) \checkmark \text{ oth. } X$

调整以适应 (t 检验) $\beta_\psi(\theta, \sigma) = P_{\theta, \sigma}(\frac{\sqrt{n}(\bar{x} - \theta_0)}{S} < -t_{\alpha}(n-1)) = F(\frac{\theta_0 - \theta}{\sigma} \frac{\sqrt{n}}{S})$

要求) $\beta_\psi(\theta) \geq 1 - \beta, \theta \leq \theta_1$ (Norm. 1. $\forall \theta < \theta_0$, 不存在 n s.t. $\theta \leq \theta_1$ 时 ψ 的 β 充分小的 β ($\sigma \rightarrow \infty$) 2. 除非 $\alpha \geq \frac{1}{2}$, oth. ψ 1°, 2° 也不是一致最优. (单边时可以指定一个 $\theta_1 < \theta_0$)

$\Delta \mu$; $U = \frac{\sqrt{nm}}{\sqrt{n+m}} \frac{\bar{x} - \bar{y} - \theta_0}{\sigma_0}$; $T = \frac{\sqrt{nm}}{\sqrt{n+m}} \frac{\bar{x} - \bar{y} - \theta_0}{S}$, $S^2 = \frac{\sum(X-\bar{X})^2 + \sum(Y-\bar{Y})^2}{n+m-2}$ (以约束二类错误.)

σ^2 相近: σ^2 : 1° $\phi: \frac{\sum(X-\bar{X})^2}{n-1} \geq \sigma_0^2 \chi^2_{\alpha}(n-1) \checkmark \text{ oth. } X$; 1° $\psi: \frac{S_1^2}{S_2^2} \geq a F_{\alpha}(n-1, m-1) (\frac{S_2^2}{S_1^2} \leq \frac{1}{a} F_{\alpha}(m-1, n-1)) \checkmark \text{ oth. } X$

exp. λ : 3° $\phi: \chi^2_{2n}(1-\alpha) \leq X \leq \chi^2_{2n}(\frac{\alpha}{2}) \checkmark \text{ oth. } X$ (1° 2° 一致最优而 3° 不是且没有见 5.3)

截尾寿命检验: 1. 定数截尾法 $r < n$, 至 r 个失效为止, $T = Y_1 + \dots + Y_r + (n-r)Y_r$

(由于部分 X 过长) $2\lambda T \sim \chi^2_{2r}$ $H_0: \lambda \leq \lambda_0$ 可由 $\phi: T \geq \frac{\chi^2_{2r}(1-\alpha)}{2\lambda_0}$

2. 定时截尾法 至 T_0 为止的总工作时间 (未失效以 T_0 计) $2\lambda T'$ 近似 $\sim \chi^2_{2r+1}$ (r 为失效个数)

3. 接替定时截尾 X 为总失效数, $X \sim \text{Poisson}(nT_0)$

$n=3$ T_0 ($X=7$) 对泊松分布参数进行检验 (见下)

Bin. p: 1° $H_0: p \leq p_0$, $\phi: X \leq C \checkmark \text{ oth. } X, \beta_\phi(p) = 1 - \sum_{i=0}^C \binom{n}{i} p^i (1-p)^{n-i}$, 要求 $\sum_{i=0}^C \binom{n}{i} p_0^i (1-p_0)^{n-i} = 1 - \alpha$ 得 C .

抽样验收 操作特征 (OC) ϕ $C_0 < C_0 + 1$ 可略放缩 α (n 很大时)

函数 $L_\phi(p) = \sum_{i=0}^C \binom{n}{i} p^i (1-p)^{n-i}$ 或随机化 (检验) 修改很小)

$L_\phi(p_0) = 1 - \alpha, L_\phi(p_1) = \beta$ (双约束) 以 $\frac{\beta_\phi(C_0) - \alpha}{\beta_\phi(C_0) - \beta_\phi(C_0 + 1)}$ 的比例通过 C_0 的产品.

符号检验 (如多人给产品打分, + - + ...) 检验 $H_0: p_0 = \frac{1}{2}$. $\varphi: C_1 \leq X \leq C_2$, $\sum_{i=0}^{C_1-1} \frac{\alpha}{2} = \frac{\alpha}{2}$, $\sum_{i=C_2+1}^n \frac{\alpha}{2} = \frac{\alpha}{2}$.
 非参数检验 (如果给分尺度无大差别 t 检验信息或许更多)

2. 总样本 $\hat{\beta} = \frac{\sum x_i}{n}$. $\approx \hat{\beta} \pm u_{\frac{\alpha}{2}} \sqrt{\frac{\beta(1-\beta)}{n}}$. 也是选一个估计, $n \rightarrow \infty$.
 Poisson λ : $H_0: \lambda \leq \lambda_0$. $\varphi: X \leq C$. $P_{\varphi}(\lambda) = 1 - \sum_{i=0}^C \frac{e^{-\lambda} \lambda^i}{i!} = K_{2C+2}(2\lambda)$ (K 为 χ^2 cdf.)
 C_0 满足 $2\lambda_0 = \chi_{2C_0+2}^2(1-\alpha)$. $\int_0^{\lambda} F(x) = \sum_{i=0}^{\lambda} \frac{e^{-x} x^i}{i!} = \frac{1}{\lambda!} \int_0^{\infty} t^{\lambda} e^{-t} dt = \frac{1}{2^{\lambda+1} \Gamma(\lambda+1)}$

② 大样本检验: Behrens-Fisher Problem:
 $X_1, \dots, X_n \sim N(\mu_1, \sigma_1^2)$, $Y_1, \dots, Y_m \sim N(\mu_2, \sigma_2^2)$. $H_0: \mu_1 = \mu_2$. χ_{2n+2}^2
 $\varphi: \frac{|\bar{X} - \bar{Y}|}{\sqrt{\frac{s_1^2}{n} + \frac{s_2^2}{m}}} \leq u_{\frac{\alpha}{2}}$. oth. x. $n, m \gg 1$ (使用极限分布) 但无法确知与 α 的偏差.

③ 贝叶斯法: $P(H_0 | X_1, \dots, X_n) \geq P(H_1 | X_1, \dots, X_n)$ (若 $\geq \frac{1}{2}$ 则思而不决) (不一定是 $\frac{1}{2}$: 统计决策)
 当 H_0 为单点时, 可以给该点 θ_0 一个先验概率 p_0 , 剩下 $1-p_0$ 以某种分布位于两边.

如 $X_1, \dots, X_n \sim N(\theta, \sigma^2)$. $H_0: a \leq \theta \leq b$. θ 广先验密度 $1, \sigma \sim \frac{1}{\sigma}$, 得 θ 的边缘后验密度
 Bayes: $P(H_0 | X) = T_{n-1}(\frac{\sqrt{n}(b-\bar{X})}{S}) - T_{n-1}(\frac{\sqrt{n}(a-\bar{X})}{S})$ | 贝叶斯的分母可以先不算而 $\circ \cdot (1 + \frac{t^2}{n-1})^{-\frac{n}{2}}$
 Neyman: $t \sim t_{n-1}$. 有 α . 构造否定域. $\geq \frac{1}{2}$? 确定后验分布的形式, 由归一 $t = \sqrt{n} \frac{\theta - \bar{X}}{S}$. 自然猜出系数 \circ .

Fisher: $\sqrt{n} \frac{(\bar{X} - \theta_0)}{S} \sim t_{n-1}$ 反过来决定了 θ 的信仰分布 信仰概率 (不存在先验分布, 无法可偏)

5.2 拟合优度检验 $Z = \sum \frac{(f_e - f_o)^2}{f_o}$ ($H_0: f_e$) F $F(b) - F(a)$
 $H_0: P(X=a_i) = p_i$ ($i=1, \dots, k$)

已知离散分布有限: $Z = \sum_{i=1}^k \frac{(np_i - n_i)^2}{np_i}$ ($\sum_{i=1}^k n_i = n$) 若 $H_0 \checkmark$, $n \rightarrow \infty$ 时 $Z \sim \chi_{k-1}^2$.

则 $\varphi: Z \leq \chi_{k-1}^2(\alpha) \checkmark$ oth. x. 拟合优度 $P(Z_0) = P(Z \geq Z_0 | H_0) \approx 1 - K_{k-1}(Z_0)$
 (注意样本量 n 对分辨率的影响) $P(Z_0)$ 小于 α 则差异愈稀奇, 拒绝 H_0 .

未知离散分布有限: $H_0: P(X=a_i) = p_i(\theta_1, \dots, \theta_r)$ ($i=1, \dots, k$) (θ 满足一定范围)
 (带参) 对某组 $(\theta_1^0, \dots, \theta_r^0)$ 成立.

多一步用极大似然估计 θ . $L = p_1^{n_1}(\theta_1) \dots p_k^{n_k}(\theta_r)$

② $p_i = p_i(\hat{\theta}_1, \dots, \hat{\theta}_r)$ 计算 Z . $\frac{\partial \ln L}{\partial \theta_j} = 0$ 即 $\sum_{i=1}^k \frac{n_i}{p_i(\theta)} \frac{\partial p_i(\theta)}{\partial \theta_j} = 0$ ($j=1, \dots, r$) 解出 $(\hat{\theta}_1, \dots, \hat{\theta}_r)$.

若 $H_0 \checkmark$, $n \rightarrow \infty$ 时 $Z \sim \chi_{k-r}^2$ ③ 计算 $p(Z) = 1 - K_{k-r}(Z)$ 比较 α 以决定.

特例: 列联表 齐一/独立性检验 | 高维列联表的压缩, 分层

($r=c=2$ 时称四格表) $Z = \frac{n(n_{11}n_{22} - n_{12}n_{21})^2}{n_1 n_2 n_{.1} n_{.2}}$
 似然估计得 $\hat{u}_i = \frac{n_{.i}}{n}$, $\hat{v}_j = \frac{n_{.j}}{n}$
 $\hat{p}_{ij} = \hat{u}_i \hat{v}_j = \frac{n_{.i} n_{.j}}{n^2}$
 $\chi_{(r-1)(c-1)}^2$ $Z = \sum_{i,j} \frac{(n_{ij} - n \hat{p}_{ij})^2}{n \hat{p}_{ij}}$ $rct - r - c - t + 2 \dots$ (\rightarrow 推出)
 高维独立性: (边缘表, 部分表) (辛普森悖论)
 $\left. \begin{matrix} A, BC \\ A, B, C \\ B, AC \dots \end{matrix} \right\} \left. \begin{matrix} BA, BC \\ CA, CB \\ \dots \end{matrix} \right\} C(r-1)(c-1)$
 (\hat{u}_i, \hat{v}_j 知, 若为 0 说明样本中不含 i, j 水平, 应视为没有)

一般(连续、无限)分布, 已知或半未知(特估): $H_0: F(x)$ 或 $F(x|\theta)$ 对一组 θ^0 .

多步划分区间. $-\infty = a_0 < a_1 < \dots < a_{k-1} < a_k = \infty$ 并统计各区间频数 v_i .
 $(a_0, a_1], \dots, (a_{k-1}, a_k)$ k 个区间, $p_i(\theta) = F(a_i; \theta) - F(a_{i-1}; \theta)$ ($i=1, \dots, k$)

($k \uparrow$ 更接近但每个区间太小(应 ≥ 5); $k \downarrow$ 使 $n(v_i)$ 与 χ^2 更接近;

一般 $10 \leq n \leq 100, k=6 \sim 8$; $100 \leq n \leq 200, k=9 \sim 12$; $200 < n, k \leq 20$)

有时 θ 很难解, 可用如 $\mu = \bar{x}$ 或 $\frac{1}{n} \sum m_i v_i$, $\sigma^2 = s^2$ 或 $(\frac{1}{n} \sum v_i (m_i - \mu)^2)^{\frac{1}{2}}$ 近似计算 Z .

(分组数据, m_i 为各区间中位数, 头尾可参考相邻区间长.)

5.3 \bullet UMPT 简单检验 $H_0: f_0(x), H_1: f_1(x)$ ($X=(X_1, \dots, X_n), f(x_1, \dots, f(x_n))$ 记为 $f(x)$)

奈-皮 $N-P$ 基本引理 α 的 UMPT φ 的否定域 Ω 应找 C , 使

$$\Omega = \{x \mid \frac{f_1(x)}{f_0(x)} > C\}, \text{ 且 } \int f_0(x) dx = \alpha.$$

复合检验 在 $H_0: H_1$ 中各取一值 $H_0': \theta_0; H_1': \theta_1$, 由 $N-P$ 引理求出其 φ .

若 1. φ 也是 $H_0: H_1$ 的水平 α 检验 2. φ 与 Ω 无关 则 φ 为 $H_0: H_1$ 的 UMPT (α).

(对于 $0 \leq a$ 形单侧原假设, θ_0 总取 a) 如 $N(\theta, \sigma^2)$ σ^2 已知, $H_0: \theta \leq a$.

同理可证 Norm. Exp. Bin. Poisson.

的单侧检验均为 UMPT.

证明一个 T 为 UMPT:

1. φ (即 Ω 或 A) 与 θ_1 无关 2. $\frac{f_1(x)}{f_0(x)} > \frac{f_1(y)}{f_0(y)} \quad \forall x \in \Omega, y \in A$.

\bullet 非中心 t 分布 $Z = \frac{X+\delta}{\sqrt{Y/n}} \sim t_{n,\delta} \quad E(t_{n,\delta}) = \delta \cdot \frac{\sqrt{n}}{2} \frac{\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n}{2})}$ (具体分布较复杂)

$\delta_2 > \delta_1$, 则 $F_{n,\delta_2}(x) \leq F_{n,\delta_1}(x)$.

$$V(t_{n,\delta}) = \frac{n}{n-2} (1+\delta^2) - E(t_{n,\delta})^2.$$

5.1 中 μ, σ^2 未知的 t 检验 $\beta_\varphi(\theta, \sigma) = F_{n-1, \frac{\sqrt{n}(\theta-\theta_0)}{\sigma}}(-t_{n-1}(\alpha)) = F(\frac{\theta-\theta_0}{\sigma})$ 可得其单调性及性质 1.

非中心 χ^2 分布 ($i=1, \dots, n$) $X_i \sim N(\mu_i, \sigma^2) \quad E(\chi_{n,\lambda}^2) = n + \lambda$

$$\lambda = \sum \mu_i^2, \chi_{n,\lambda}^2 = \sum X_i^2 \quad V(\chi_{n,\lambda}^2) = 2(n+2\lambda)$$

$X \sim \lambda e^{-\lambda x}$ 时

\bullet 截尾法证明: $\therefore n \cdot \min_n X \sim X$. 递推: $T_{k+1} = T_k + (n-k)(Y_{k+1} - Y_k)$

Δ 注意参数可取范围!

max, min! 如 $R(0, \theta)$, θ 下限 应为 $\max X_i$.

右后数性使 $Y_{k+1} - Y_k = \min_{n-k} X$. 由此 $\chi_{2k}^2 \rightarrow \chi_{2k+2}^2$.

(Norm. Chi-sq. Bin. Poisson. 均有可加性(再生性))

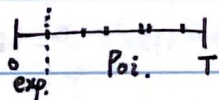
(Exp. \uparrow)

3: 设从单个开始替换

X 记 T 时替换决数.

$$X \sim \text{Poi}(\lambda T) \quad \text{递推: } \int_0^T \lambda e^{-\lambda(T-t)} \cdot \frac{e^{-\lambda t} (\lambda t)^n}{n!} dt$$

$$= \frac{e^{-\lambda T} (\lambda T)^{n+1}}{(n+1)!}$$



• 卡方检验的证明: 1. 一些验证 (E.V): $E(Z) = \sum_i \frac{E(np_i - v_i)^2}{np_i} = \sum_i (1 - p_i) = k - 1$
 $V(Z)$ 需 $E(Z^2)$, $E(np_i - v_i)^4$ 归为 Bin: $E(X^4)$, 二次、多项分布的构造 (由此计算 $E(X^4)$)
 而交叉项归为 Multi: $(p_i, p_j, 1 - p_i - p_j)$. Bin: $E\left(\frac{X(X-1)\dots(X-i+1)}{n(n-1)\dots(n-i+1)}\right) = p^i$ (Multi 同理)
 得 $V(Z) = 2(k-1)$. (方便化出求和)

2. 简单的理解: 每格均为 Poisson. $\frac{(X-EX)^2}{VX^2} \xrightarrow{D} N$, 加和为 χ^2 , 由于 $\sum_i v_i = n$ 而失去 1 个自由度.

3. 记 $Y_i = \frac{a_i - np_i}{\sqrt{np_i}}$. 向量 Y 为多元分布, $\Sigma = \text{Cov} Y = \begin{pmatrix} 1-p_1 & -\sqrt{p_1 p_2} \\ -\sqrt{p_1 p_2} & 1-p_2 \\ \vdots & \vdots \end{pmatrix}$
 (Y_i 的边缘为二项分布)

有 $\Sigma = I - PP^T$ 其中 $p = (\sqrt{p_1}, \dots, \sqrt{p_k})^T$ 且 $\|p\|=1$. ($E(a_i a_j) = n(n-1)p_i p_j$)

Σ 为投影阵 (幂等), 特征值 1 个 0, $k-1$ 个 1. ($\Sigma p = 0, \Sigma^2 = \Sigma$. 注意 Σ 不可逆)

$n \rightarrow \infty$ 时 $Y \xrightarrow{D} N(0, \Sigma)$. 现在通过一个旋转 $A \Sigma A^T = \begin{pmatrix} 0 & & \\ & 1 & \\ & & \ddots \end{pmatrix}$

则 $X = AY \xrightarrow{D} N(0, A \Sigma A^T)$ (把 $\sum_i p_i = 1$ 的缺失一维抓出来了)

即 $X = (0, X_1, \dots, X_{k-1})^T$, $X_i \stackrel{iid}{\sim} N(0, 1)$, Y_i 的形式无非是 a_i (多元 \rightarrow 正态的归一)

$Z = \sum_i Y_i^2 = \|Y\|^2 = \|X\|^2 \sim \chi_{k-1}^2$. 得证.

• 线代在统计中的简单应用: 协方差矩阵 $\Sigma = E((X-\bar{X})(X-\bar{X})^T)$

Σ 为正定 (半正定), 对称阵. 若 $|\Sigma| = 0$ 说明样本中有非独立变量. $\Sigma = \begin{pmatrix} V(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & V(X_2) \end{pmatrix}$

Σ 的变换 $Y = AX$, $\Sigma_Y = A \Sigma_X A^T$. Σ 的对角化: 找到一组独立的方向, $Y = Q^T X$.

(如二维正态下 $\sigma_2 X + \sigma_1 Y$)

但非唯一! $X \Lambda_1 X^T = \Lambda_2 \neq X$ 为对角阵.

($a, b, \sigma_1^2, \sigma_2^2, \rho$) $\sigma_2 X - \sigma_1 Y$ 独立, 但对角化并非此

(准确说为不相关)

Σ 的 PCA: 找到方差的最大解释 (特征向量)

正交方向. (见下)

六、回归、相关与方差分析

6.1 • 理论回归方程 (函数) $y = f(X_1, \dots, X_p) + e$, $E(e) = 0$. 线性回归分析
 自变量的选取、回归函数的形式 $\rightarrow \sigma^2$ 的控制 "模型"

应用: "描述", "估计回归函数", "预测", "控制" 自变量的随机性 回归设计

• 一元线性回归 中心化 $Y_i = \beta_0 + \beta_1 (X_i - \bar{X}) + e_i$ ($i=1, \dots, n$)

最小二乘法 $\sum_i \delta_i^2 = \sum_i (Y_i - \hat{Y}_i)^2 \rightarrow \min$ $e_i \stackrel{iid}{\sim}$, $E(e_i) = 0$, $V(e_i) = \sigma^2$ ($i=1, \dots, n$)

(亦即极大似然时)

$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 (X_i - \bar{X})$

$E(\hat{Y}_i) = \beta_0 + \beta_1 (X_i - \bar{X}) = \hat{Y}_i$

2. 偏导, 平方或线性代(以后)得 $\hat{\beta}_0 = \bar{Y}$, $\hat{\beta}_1 = \frac{\sum (X_i - \bar{X}) Y_i}{\sum (X_i - \bar{X})^2}$ $\sum (X_i - \bar{X}) Y_i = \sum (X_i - \bar{X}) (Y_i - \bar{Y})$

线性, 无偏 $E(\hat{\beta}_0) = \beta_0$, $E(\hat{\beta}_1) = \beta_1$, $V(\hat{\beta}_0) = \frac{\sigma^2}{n}$, $V(\hat{\beta}_1) = \frac{\sigma^2}{\sum (X_i - \bar{X})^2} = \frac{\sigma^2}{S_x^2}$

最小方差 ($I(\beta_0), I(\beta_1)$ 可得), 不相关 $Cov(\hat{\beta}_0, \hat{\beta}_1) = 0$ 当 $e_i \sim N$ 时两者独立.

残差方差 (估计方差): (利用 $E(e_i e_j) = \sigma^2 \delta_{ij}$) | 两个正态分布的联合分布只
(或与 $X - \bar{X}$ 正交) | 要不存在映射 (如 $Y = X$) 则为二维

$\sum \delta_i^2 = \sum (Y_i - \hat{Y}_i)^2 = \sum (Y_i - \bar{Y})^2 - \hat{\beta}_1^2 \sum (X_i - \bar{X})^2$ (只有计算) 正态分布.

区间估计: $\frac{\hat{\beta}_1 - \beta_1}{\hat{\sigma} / S_x} \sim t_{n-2}$ 或对函数 $\frac{\sum \delta_i^2}{\sigma^2} \sim \chi_{n-2}^2$ 回归诊断 ($\rho = 0$ id.)

$\hat{m}(x) = \hat{\beta}_0 + \hat{\beta}_1 (x - \bar{X}) \rightarrow m(x) \in \hat{m}(x) \pm \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x - \bar{X})^2}{S_x^2}} t_{n-2}(\frac{\alpha}{2})$ (不能“叠加”起来说直线都在此范围里)

假设检验: $H_0: \beta_1 = c$ (常 $c=0$) $\varphi: |\hat{\beta}_1 - c| \leq \hat{\sigma} S_x t_{n-2}(\frac{\alpha}{2}) \vee oth x$ ($n \rightarrow \infty$ 时 $m(x)$ 可以 $\rightarrow \hat{m}(x)$)

(一元回归中 β_1 的检验, 回归显著性检验 (F), 相关系数检验 三者一致.)

使用问题: 1. 解释回归方程 (变化区间, 人为控制) 2. 外推 (线性的把握)

3. 不可逆转使用 $x = \hat{c} + \hat{d}y$ (举例: $\rho \approx 0$ 时) 4. X 随机 即可使用 β , 但 $V(\hat{\beta})$ 不同, $\hat{\sigma}^2$ 应

视为固定 X 下 Y 方差

多元线性回归 $Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik} + e_i$ ($i=1, \dots, n$)

$\hat{\beta} = (X^T X)^{-1} X^T Y$ $\hat{\sigma}^2 = \frac{\sum \delta_i^2}{n-k-1}$, $\sum \delta_i^2 = \sum (Y_i - \hat{Y}_i)^2 = \sum Y_i^2 - \hat{\beta}_0 \sum Y_i - \dots - \hat{\beta}_k \sum X_{ki} Y_i \sim \sigma^2$

$V(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$ (元素 σ_{ij}) $\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \sqrt{\sigma_{jj}}} \sim t_{n-k-1}$, $\frac{\hat{m}(x) - m(x)}{\hat{\sigma} \sqrt{\frac{1}{n} + \sum_{ij} (X_i - \bar{X}_i)(X_j - \bar{X}_j) \sigma_{ij}}} \sim t_{n-k-1}$

1. 单个回归系数 $H_0: \beta_j = c$. 点 $y_0 \in \hat{m}(x_0) \pm \hat{\sigma} \sqrt{\lambda(x_0)} t_{n-k-1}(\frac{\alpha}{2})$ ($V(\Sigma X) = \sum_{ij} (\Sigma)_{ij}$)

2. 全体回归系数 (回归显著性) $H_0: \beta_1 = \dots = \beta_k = 0$ $R_1 = \sum \delta_i^2$, $R_2 = \sum (Y_i - \bar{Y})^2$ (只剩 $\hat{\beta}_0$)

(F 检验) $\frac{R_2 - R_1}{\sigma^2} \sim \chi_k^2$, $\frac{R_2 - R_1}{k \hat{\sigma}^2} \sim F_{k, n-k-1}$

3. 部分回归系数 $H_0: \beta_1 = \dots = \beta_r = 0$ $\frac{R_3 - R_1}{\sigma^2} \sim \chi_{k-r}^2$, $\frac{R_3 - R_1}{(k-r) \hat{\sigma}^2} \sim F_{k-r, n-k-1}$ (此重新拟合可通过 $X^T X = \begin{bmatrix} \dots & \dots \\ \dots & \dots \end{bmatrix}$ 分块消元得出)

使用问题: 1. X 交互作用 (X 随机时不应 2. 删去自变量, $R_3 - R_1 = \hat{\beta}^T (X^T X)^{-1} \beta$ 压低其他 X 单改一个) 复共线性 \rightarrow 稳定性 $|X^T X| \approx 0$ 时

可转为线性回归: 如多项式回归 $Y = b_0 + b_1 X + \dots + b_p X^p + e$ (高度相关)

$Y = b_0 + b_1 e^{tX}$, $Y = b_0 e^{b_1 X^2}$ (此处 mY 是否有 $e \sim N$ 待讨论.)

6.2

• 相关分析 样本相关系数 $r = \frac{\sum_i (X_{1i} - \bar{X}_1)(X_{2i} - \bar{X}_2)}{\sqrt{\sum_i (X_{1i} - \bar{X}_1)^2 \sum_i (X_{2i} - \bar{X}_2)^2}}^{\frac{1}{2}}$

$H_0: \rho = 0$ $\varphi: |r| \leq C \sqrt{0.4 \times 0.4}$. C 有 $\frac{\sqrt{n-2}r}{\sqrt{1-r^2}} \sim t_{n-2}$. $C = \frac{t_{n-2}(\frac{\alpha}{2})}{\sqrt{n-2 + t_{n-2}^2(\frac{\alpha}{2})}}$.

($n=100, C \approx 0.2$, 较微弱的 ($\rho=0$ 时) 相关需较大样本量) (正态, Y 对 X 回归可证: $\hat{\beta}_1 = r \frac{\sqrt{\sum(Y-\bar{Y})^2}}{\sqrt{\sum(X-\bar{X})^2}} \sim t_{n-2}$)
 $Y_i = b + \beta X_i + e_i$

偏相关 $X_2' = X_2 - L(X_3, \dots, X_p)$
 $X_1' = X_1 - L(X_3, \dots, X_p)$ 消除其他变量的影响 $\rho_{12 \cdot (3 \dots p)} = \text{Cov}(X_1', X_2')$

相关阵 $P = \begin{pmatrix} 1 & \rho_{12} & \dots \\ \rho_{12} & 1 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$, $\rho_{12 \cdot (3 \dots p)} = \frac{\rho_{12}}{\sqrt{1 - \rho_{13}^2 - \rho_{14}^2 - \dots}}$ (二元时 $\rho_{12 \cdot 3} = \frac{\rho_{12} - \rho_{13}\rho_{23}}{\sqrt{(1-\rho_{13}^2)(1-\rho_{23}^2)}}$)
 (ρ_{ij} 为 ij 元素) (Y 对 X 及 $X_3 \dots X_p$ 同时回归可证)

复相关 其他变量全体对一方的影响 $\rho_{12 \cdot (3 \dots p)} = \sqrt{1 - \frac{P_{11}}{P_{11}}}$ $\frac{\sqrt{n-p} \rho_{12 \cdot (3 \dots p)}}{\sqrt{1 - \rho_{12 \cdot (3 \dots p)}^2}} \sim t_{n-p}$
 $r_{12 \cdot (3 \dots p)}^2$ 为 Beta($\frac{p-1}{2}, \frac{n-p}{2}$) 分布, 或 $\frac{n-p}{p-1} \frac{r_{12 \cdot (3 \dots p)}^2}{1 - r_{12 \cdot (3 \dots p)}^2} \sim F_{p-1, n-p}$

6.3

• 方差分析 Fisher 试验设计三原则: 重复, 随机化 (完全/部分), 分区组

因素 水平 单因素完全随机化实验: $Y_{ij} = a_i + e_{ij}$ ($j=1, \dots, n_i; i=1, \dots, k$)

$H_0: a_1 = \dots = a_k$ (平方和分解) $SS = SS_A + SS_E$ $E(e_{ij}) = 0, V(e_{ij}) = \sigma^2, iid.$

$\varphi: \frac{MS_A}{MS_E} \leq F_{k-1, n-k}(\alpha)$ $SS_E = \sum_i \sum_j (Y_{ij} - \bar{Y}_i)^2$ $SS = \sum_i \sum_j (Y_{ij} - \bar{Y})^2$ (Normal)
 $SS_A = \sum_i n_i (\bar{Y}_i - \bar{Y})^2$

该因素的显著性: \bar{Y}_i 与 \bar{Y}_j 的差异 $MS_E = \frac{SS_E}{n-k} = \sigma^2$ $MS_A = \frac{SS_A}{k-1}$

显著不代表 α 中没有相同的, \leftrightarrow 随机误差 MS_E

就其中一对而言 $a_u - a_v \in \bar{Y}_u - \bar{Y}_v \pm \sqrt{\frac{n_u + n_v}{n_u n_v}} \sigma t_{n-k}(\frac{\alpha}{2})$ 若包含0则达不到显著.

区间同时成立的概率更小,

可采用多重比较法)

双因素完全随机实验: $Y_{ij} = \mu + a_i + b_j + e_{ij}$ ($i=1, \dots, k; j=1, \dots, l$)

$Y_{..} = \frac{\sum_i \sum_j Y_{ij}}{kl}$, $Y_{i.} = \frac{\sum_j Y_{ij}}{l}$, $Y_{.j} = \frac{\sum_i Y_{ij}}{k}$, $\hat{a}_i = Y_{i.} - Y_{..}$, $\sum_i a_i = 0, \sum_j b_j = 0$.

$SS = SS_A + SS_B + SS_E$, $SS = \sum_i \sum_j (Y_{ij} - Y_{..})^2$ $b_j = Y_{.j} - Y_{..}$ $SS_E = \sum_i \sum_j (Y_{ij} - Y_{i.} - Y_{.j} + Y_{..})^2$


$H_{0A}: a_1 = \dots = a_k = 0$, $H_{0B}: b_1 = \dots = b_l = 0$ $SS_A = \sum_i l (Y_{i.} - Y_{..})^2$ (注意不是 $Y_{ij} - \frac{Y_{i.} Y_{.j}}{Y_{..}}$!)

$MS_E = \frac{SS_E}{(k-1)(l-1)}$ $\varphi: \frac{MS_A}{MS_E}, \frac{MS_B}{MS_E} \sim F$ 检验同上. 交互作用 SS_{AB} (SS_E 可能偏大)

单因素随机区组实验: b_j 区组效应 SS_B 区组平方和 区组大小等于水平数.

划分不当 (区组不显著时), 甚至 $MS_B < MS_E$, 加之自由度损失, 发现A效应更难 (SS_E 更大)

正交表 部分实施水平组合保证¹ 仍可对 a_i, b_j, C_{n_i}, d_v 等估计². 总平方和 SS 仍可分解.

$SS_A = \frac{n(a_1^2 + \dots + a_k^2)}{k}$ $L_9(2^4)$  因素/区组 (为 n 因素)

(两区组/两水平 F 检验与双样本 t 检验一致)

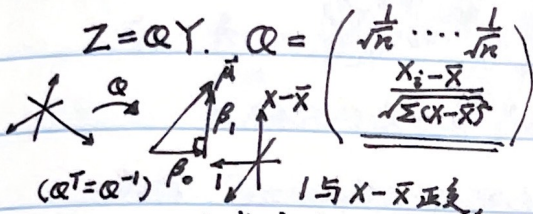
(多总体 (联合) 均值 t 检验)

6.4

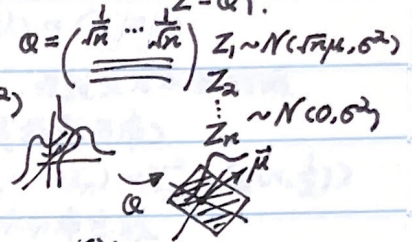
• 矩阵变换视角

$Y \sim N(\beta_0 + \beta_1(X-\bar{x}), \sigma^2 I)$

另一个例子: $Y \sim N(\mu, \sigma^2 I)$
 $Z = QY$



$Z_1 \sim N(\sqrt{n}\beta_0, \sigma^2)$
 $Z_2 \sim N(\beta_1 \cdot \sqrt{\sum(X-\bar{x})^2}, \sigma^2)$
 $Z_3 \sim N(0, \sigma^2)$
 \vdots
 $Z_n \sim N(0, \sigma^2)$



为消除未知(待估)参数,
将坐标架转至该方向, 剩余均值为0,
为真正(残差的)自由度.

$\frac{\sum Y}{\sqrt{n}} = Z_1 \rightarrow \sqrt{n}\beta_0, \hat{\beta}_0 = \frac{\sum Y}{\sqrt{n}}$
 $\frac{\sum(X-\bar{x})Y}{\sqrt{\sum(X-\bar{x})^2}} = Z_2 \rightarrow \beta_1 \sqrt{\sum(X-\bar{x})^2}, \hat{\beta}_1 = \frac{\sum(X-\bar{x})Y}{\sum(X-\bar{x})^2}$

多元时 $Y \sim N(X\beta, \sigma^2 I)$

$\beta_0 + \beta_1(X-\bar{x}) + \beta_2(G-\bar{G}) + \dots$

注意 $X-\bar{x}$ 与 $G-\bar{G}$ 一般不垂直(相关, $\sum(X-\bar{x})(G-\bar{G}) \neq 0$). 需正交化(如Schmidt法).

$EY = \beta_0 \cdot 1 + \beta_1(X-\bar{x}) + \beta_2(G-\bar{G}) + \dots$ (无的"中心化"即正交化?)

$(G-\bar{G}) = \frac{\sum(X-\bar{x})(G-\bar{G})}{\sqrt{\sum(X-\bar{x})^2 \sum(G-\bar{G})^2}}$ 同理可由 $\sum G^T Y = \dots$ 解得 $\hat{\beta}_1, \hat{\beta}_2$.

残差方和 $\sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2 = \|Y - \hat{Y}\|^2 = \|QY - Q\hat{Y}\|^2 = \|Z - Q\hat{Y}\|^2$ (与2偏导一致)

注意 $Q\hat{Y} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \end{pmatrix}$, $\hat{Y} = X\hat{\beta}$, 就是用 Z_1, Z_2 估计出的 $\hat{\beta}$
 $\hat{\beta} \sim \chi_{n-2}^2$ (χ^2 已达到 $\sum \delta_i^2 \min$)

• 帽子矩阵 $X\hat{\beta} = \hat{Y} \sim N(EY, \sigma^2 I)$, $EY = X\beta$ 其中设计矩阵 $X = (1, X_1, X_2, \dots)$.



$X^T X \hat{\beta} = X^T Y$, $\hat{\beta} = (X^T X)^{-1} X^T Y \sim N(\beta, \sigma^2 (X^T X)^{-1})$

$X^T X = \begin{pmatrix} n & 0 & \dots & 0 \\ 0 & \sum X_1^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sum X_{k-1}^2 \end{pmatrix}$
 $(X^T X)^{-1} = \begin{pmatrix} \frac{1}{n} & 0 & \dots & 0 \\ 0 & \dots & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sum X_{k-1}^2} \end{pmatrix}$

(=无时例: $\begin{pmatrix} \sum X^2 & \sum XG \\ \sum XG & \sum G^2 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} \sum XY \\ \sum GY \end{pmatrix}$ $\hat{Y} = HY$)

则 $E\hat{\beta} = \beta$, $V\hat{\beta} = \sigma^2 (X^T X)^{-1}$ (H-元时例: $\frac{\hat{\beta}_i - \beta_i}{\sigma \sqrt{h_{ii}}} \sim t_{n-k-1}$, $h_{ij} = \frac{1}{n} + \frac{(X_i - \bar{X})(X_j - \bar{X})}{\sum(X-\bar{x})^2}$)

$\frac{1}{n} \leq h_{ii} \leq 1$, $H = X(X^T X)^{-1} X^T$, 对称, 幂等, 半正定, $HX = X$, $H1 = 1$.

$\sum_i h_{ii} = \text{tr} H = rH = k+1$
 $\sum_{i,j} h_{ij} = n$ ($X^T X$ 满秩)

$e = Y - \hat{Y} \sim N((I-H)EY, \sigma^2(I-H))$ (投影阵) $(I-H)X = 0$.

注意残差不独立, 但与估计值独立且正交.

$V(e_i) = E(e_i^2) = \sigma^2(1-h_{ii})$

$\text{Cov}(\hat{Y}, e) = H^T \text{Cov} Y (I-H) = \sigma^2 H^T (I-H) = 0$

$\sigma^2(n-k-1) = \sum e_i^2 = e^T e = (X\hat{\beta} + \varepsilon)^T (I-H)(X\hat{\beta} + \varepsilon)$
 $(= \hat{\beta}^T X^T X \hat{\beta} = n\bar{Y}^2 + \sum \hat{\beta}_i^2 \sum X_i^2 \text{ 或 } \|Q\hat{Y}\|^2 = \varepsilon^T (I-H)\varepsilon \rightarrow \chi_{n-k-1}^2$ (该二次型经过旋转可以化为 $\lambda = (0, 1)$))

$\|Y\|^2 = \|\hat{Y}\|^2 + \|Y - \hat{Y}\|^2 = Z_1^2 + \dots + Z_{k+1}^2$ (ε 为误差项, $V(\varepsilon) = \sigma^2$)

$\sigma^2 H$, 由于 $H1 = 1$ 对角化后有一个1属

$n\bar{Y}^2 + \sum_{k-r}^2 + \sum_{r+1}^2 + \sum_{n-k-1}^2$

(解释: 1. Q正交化 2. 用H直接写出 $\hat{Y}, Y - \hat{Y}$ 分布) 予 β_0 ($EY = \beta_0 1$, Q第一行 $\frac{1}{\sqrt{n}}$ 应删去)

$n\bar{Y}^2 + \sum_{k-r}^2 + \sum_{r+1}^2 + \sum_{n-k-1}^2$

残差图诊断: 线性, 方差齐性, 不相关性, 正态性;
学生化残差 $\hat{\varepsilon}_i = \frac{e_i}{\sigma \sqrt{h_{ii}}}$ (| t_i | > 2)
异常点, 强影响点/高杠杆点 (h_{ij}).
标准化残差 $\frac{e_i}{\sigma}$.

伽马、贝塔、指数、泊松分布相关:

Gamma: $f(x) = \frac{\lambda^\alpha x^{\alpha-1}}{\Gamma(\alpha)} e^{-\lambda x} = \frac{(\lambda x)^{\alpha-1} e^{-\lambda x}}{(\alpha-1)! \cdot \lambda}$ $\sum_{n=0}^{\infty} \text{id. 指数分布 } (\lambda) = \text{Gamma}(n, \lambda)$
 $\int_0^{\infty} t \text{Gamma}(n, \lambda) dx = \sum_{i=0}^{n-1} \text{Poi}(\lambda t)$ (泊松过程, 每一段都是指数分布)
 与泊松分布“互逆”, 为发生 n 次的时间.

Beta: $\text{Beta}(\alpha, \beta) \times \text{Bin}(n, k) \rightarrow \text{Beta}(\alpha+k, \beta+n-k)$ (由此可解释 $2\lambda(X_1 + \dots + X_n) \sim \chi_{2n}^2 = \text{Gamma}(n, \frac{1}{2})$)
 $f(x) = \frac{1}{B(\alpha, \beta)} x^\alpha (1-x)^\beta$ 且 $X \sim \text{Beta}(\alpha, \beta)$ 时, $\frac{\beta X}{\alpha(1-X)} \sim F_{2\alpha, 2\beta}$
 $T_1 \sim \text{Gamma}(k_1, \lambda)$ $T_2 \sim \text{Gamma}(k_2, \lambda)$ 且 $X \sim \text{Beta}(\alpha, \beta)$ 时, $\frac{\beta X}{\alpha(1-X)} \sim F_{2\alpha, 2\beta}$

比例的分布: $T = T_1 + T_2 \sim \text{Gamma}(k_1 + k_2, \lambda)$ (甚至只需 iid. 即可, 可视作 $\text{Poi}(\lambda)$)
 则 $\frac{T_1}{T} \sim \text{Beta}(k_1, k_2)$ $(\frac{\chi_{2\alpha}^2}{2\alpha} / \frac{\chi_{2\beta}^2}{2\beta} = \text{Gamma}(\alpha, \frac{1}{2}) / \text{Gamma}(\beta, \frac{1}{2}))$

$\hat{\beta}$ 为 BLUE 证明: 1. $\text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$ 若另一估计量 $\beta' = A^T Y$ (A 与 X 同形状)

应满足 $E(\beta') = \beta \rightarrow A^T X \beta = \beta, \forall \beta! A^T X = I$, 则 $V(\beta') = \sigma^2 A^T A \geq \sigma^2 (X^T X)^{-1}$

故 $V(\beta') \geq V(\hat{\beta})$ 始终 $\geq V(\hat{\beta})$, 因 $A^T A - (X^T X)^{-1} = A^T (I - H) A \geq 0$ (半正定)
 $V(\beta_i)$ 组合均有 $\hat{\beta}$ 最小方差.

2. (假设 $e \sim \text{Normal}$) $Y \sim N(X\beta, \sigma^2 I)$, $I(\beta) = \int \frac{(\frac{\partial f}{\partial \beta})^2}{f} dy$ 其中
 故 CRL 下界 $V_{\min}(\hat{\beta}) = \frac{1}{I(\beta)} = \sigma^2 (X^T X)^{-1}$ $\frac{\partial f}{\partial \beta} = \frac{f}{\sigma^2} \cdot \frac{1}{\sigma^2} \cdot X^T (Y - X\beta)$ (平方应视作 $X^T X$)
 $(\frac{\partial f}{\partial \beta})^T = \frac{1}{\sigma^2} e^{-\frac{(Y-X\beta)^T (Y-X\beta)}{2\sigma^2}}$ $I(\beta) = \int \frac{1}{\sigma^2} X^T (Y-X\beta) (Y-X\beta)^T X \frac{1}{\sigma^2} dy = X^T \frac{1}{\sigma^2} X$ (向量求导)

组间方差分布证明: $\bar{Y}_i \sim N(a_i, \frac{\sigma^2}{n_i})$, $\sqrt{n_i} \bar{Y}_i \sim N(\sqrt{n_i} a_i, \sigma^2 I)$ (将 $X^T O X$ 提出归一) (参考线代)
 $\frac{n_i \bar{Y}_i}{n} = \bar{Y}$, 则 α 第一行取 $(\frac{\sqrt{n_i}}{n})^T \rightarrow \sqrt{n} \bar{Y}$ 得证.

全书小结:

概率论 { 概率(几何)与分布 古典概率, 贝叶斯, 随机变量, 分布(五大), E.V.
 极限理论 大数, 中心极限
 随机过程与 Markov 过程
 鞅论等

数理统计 { 参数估计与假设检验(统计推断) 奈曼, 贝叶斯, 大样本
 非参数统计 符号检验, χ^2 拟合优度 极大似然
 相关与回归分析 H.r...
 统计决策
 试验设计与分析, 抽样理论等 ANOVA
 时间序列分析

(补) 概率导论

1 样本空间与概率

条件独立 计算概率的方法: '计数法(排列/组合/分割)(古典概型)

2. 序贯树形图 3. 全概率公式

- 邦费罗尼不等式 $P(\prod_{i=1}^n A_i) \geq \sum_{i=1}^n P(A_i) - (n-1)$ 条件概率的全概率公式
- 博雷尔-坎泰利引理 $\{P_i\} \ i=1, \dots, \infty$. N 为序列事件均未发生, I 为有无穷多次发生.
若 $\sum_{i=1}^{\infty} P_i < \infty$ 则 $P(I) = 0$.
- 生日问题

2 离散随机变量

条件期望 全期望定理 $E(X) = \sum_y P_Y(y) E(X|Y=y)$ (或称全)

· 巴拿赫火柴问题 · 圣彼得堡悖论 容斥恒等式(示性函数证明)

· 熵与信息量 $H(X) = -\sum P_i \log P_i$. $H(X) \leq -\sum P_i \log q_i$ 及叉熵

$$\begin{aligned} \text{互信息 } I(X, Y) &= H(X) - H(X|Y) \quad (\sum P_i = 1) \\ &= H(X) + H(Y) - H(X, Y) = \sum_{x,y} P(x,y) \log \frac{P(x,y)}{P(x)P(y)} \geq 0 \end{aligned}$$

3 一般随机变量

· 布丰投针 $\frac{2l}{\pi d}$ 联合累积分布函数 · 混合随机变量 $F_X(x) = pF_Y(x) + (1-p)F_Z(x)$

随机个随机独立变量和 $T = X_1 + \dots + X_N$ $E(T), V(T) = V(X) \cdot E(N) + E^2(X) \cdot V(N)$
($X_i, N \text{ i.i.d.}$) $= E(X) \cdot E(N)$

4 随机变量高级主题

导出分布 · 拉普拉斯分布 和的方差 $V(\sum_i X_i) = \sum_i V(X_i) + \sum_{i \neq j} \text{Cov}(X_i, X_j)$

重期望法则 $E(E(X|Y)) = E(X)$ 用条件期望作为估计量: $\hat{X} = E(X|Y)$, $\tilde{X} = \hat{X} - X$.
(外面的 E 为 E_Y) (如 Bayes 估计) 则 $E(\tilde{X}) = 0 = E(\tilde{X}|Y)$, (含 X 和 Y)

条件方差 全方差定理 $V(X) = E(V(X|Y)) + V(E(X|Y))$. $\text{Cov}(\hat{X}, \tilde{X}) = 0$. $V(X) = V(\hat{X}) + V(\tilde{X})$.
(代入 $E(\tilde{X}^2) = E(V(X|Y))$)

矩母函数 $M_Y(s) = M_N(\ln M_X(s))$. 可得 $V(X|Y) = E^2(X) V(Y) + E^2(Y) V(X) + V(X) V(Y)$ (得全方差式)
 $M_X(s) = \sum_x e^{sx} p(x)$ ($Y = X_1 + \dots + X_N$) ($X, Y \text{ i.i.d.}$)
 $= E(e^{sX})$ (常见分布的矩母函数, 由形式判断分布)

5 极限理论

弱大数定律 $\lim_{n \rightarrow \infty} P(|Y_n - a| \geq \varepsilon) = 0, \forall \varepsilon > 0$ 依概率收敛
 强大数定律 $P(\lim_{n \rightarrow \infty} Y_n = a) = 1$ 以概率1收敛 (几乎处处收敛)

切诺夫界 利用马尔可夫不等式. 中心极限定理的证明: $Z_n = \frac{X_1 + \dots + X_n}{\sigma\sqrt{n}}$

收敛矩母: 针对 $X_i \in [0, 1]$ (如伯努利), e^{sX} 为凸函数. 利用了 $M_{Z_n} \rightarrow M_Z$ 则 $F_{Z_n} \rightarrow F_Z$.
 $P(X \geq a) \leq \frac{M_X(s)}{e^{sa}}$ 令 $\sum X_i = X$ 就先 $M_{Z_n}(s) = (M_X(\frac{s}{\sigma\sqrt{n}}))^n$
 $Ece^{sX} \leq (e^s - 1)E(X) + 1 \leq e^{E(X)(e^s - 1)}$ (可取等号) 对0展开 $M_X(s) = 1 + 0 \cdot s + \frac{\sigma^2}{2} s^2 + o(s^2)$
 $M_X(s) = E(e^{sX}) \leq \prod_i e^{E(X_i)(e^s - 1)} = e^{\mu(e^s - 1)}$ 泊松分布. 得到 $\lim_{n \rightarrow \infty} M_{Z_n}(s) = e^{\frac{s^2}{2}}$ 标准正态.
 o/p. $p \rightarrow 0$

则 $\forall s > 0$ 有 $P(X \geq a) \leq e^{-sa} M_X(s) = e^{-sa} e^{\mu(e^s - 1)}$ (适当取)
 取 $a = \mu + \frac{\mu}{s}$, $s = \ln(1 + \frac{\mu}{a - \mu})$ 得 $= e^{-\mu(\frac{\mu}{a - \mu} - (1 + \frac{\mu}{a - \mu}) \ln(1 + \frac{\mu}{a - \mu}))} \leq e^{-\frac{\mu^2}{a}}$
 同理 $s < 0$ 由 $P(X \leq a) \leq e^{-sa} M_X(s)$ 得 $P(X \leq (1 - \frac{\mu}{a})\mu) \leq e^{-\mu(\frac{\mu}{a} + (1 - \frac{\mu}{a}) \ln(1 - \frac{\mu}{a}))} \leq e^{-\frac{\mu^2}{a}}$

琴生不等式 f 凸 ($\frac{\partial^2 f}{\partial x^2} > 0$): $f(E(X)) \leq E(f(X))$. 依均方收敛 $\lim_{n \rightarrow \infty} E(X_n - a)^2 = 0$
 (期望的) (代入 $f(x) = (x - a)^2$) $(a \leq f(x))$ C 依概率收敛 (切比雪夫不等式) (强于)

6 伯努利过程和泊松过程

伯努利过程 无记忆性与独立性 '相邻到达的间隔时间: $(1-p)^{t-1}$ 几何分布

1. 分裂与合并 $\begin{matrix} \uparrow \\ \text{分裂} \\ \downarrow \end{matrix}$ $\begin{matrix} \downarrow \\ \text{合并} \\ \uparrow \end{matrix}$ (p, q) 2. k 次到达时间: $T_k \sim \text{Geo}(p)$, $T = \sum_{i=1}^k T_i \sim k$ 阶帕斯卡分布

3. t 时间内的到达次数: $\binom{t}{k} p^k (1-p)^{t-k}$ 二项分布. (k 次到达前的失败次数 \sim 负二项分布, 可代入 $t = t - k$)
 (下面均可如此类推) (p) \downarrow \downarrow \downarrow $(p+q-pq)$
 $p = \lambda\delta, \delta \rightarrow 0$ (q) \uparrow \uparrow \uparrow $(= 1 - (1-p)(1-q))$ (Δ 注意2/3的互逆关系)

泊松过程 (点过程) 到达率 λ 时间同质性与独立性 小区间概率

1. 相邻到达的间隔时间: $\lambda e^{-\lambda t}$ 指数分布 (无记忆性) $P(0, t) = 1 - \lambda t + o(t)$
 2. k 次到达时间 $T_k \sim \text{Exp}(\lambda)$, $T = \sum_{i=1}^k T_i \sim k$ 阶 爱尔朗分布. Δ 求概率分布: $P(1, t) = \lambda t + o(t)$
 3. t 时间内的到达次数: $\frac{\lambda^k t^k e^{-\lambda t}}{(k-1)!}$ (泊松分布) $P(k, t) = o(t)$ ($k \geq 2$)
 4. 分裂与合并 $(\sum_{k=1}^{\infty} G_k(t) = \lambda)$ $\textcircled{1}$ 分析事理) 分析该点/事件发生概率 $(f(x)dx)$
 (λ) \downarrow (λ_1) \downarrow $(\lambda_1 + \lambda_2)$ $\textcircled{2}$ 累积分布函数求导. 矩母函数等. $(P(X \leq x))$

随机个随机独立变量 (λ) \downarrow (λ_1) \downarrow $(\lambda_1 + \lambda_2)$ $X_1 \sim \text{Exp}(\lambda_1), X_2 \sim \text{Exp}(\lambda_2)$. id.
 则 $\min\{X_1, X_2\} \sim (\lambda_1 + \lambda_2)$ (串联系统) $E(\min) = \frac{1}{\lambda_1 + \lambda_2}$
 和 一随机过程 但 $\max\{X_1, X_2\}$ 不是 Exp ! $E(\max) = \frac{1}{\lambda_1 + \lambda_2} + \frac{1}{\lambda_2} \cdot \frac{\lambda_1}{\lambda_1 + \lambda_2} + \frac{1}{\lambda_1} \cdot \frac{\lambda_2}{\lambda_1 + \lambda_2}$
 对称性: 倒流的伯努利/泊松过程没有区别.

随机插入 t_1 t_2 t_1 t_2 无记忆性导致应为两段独立. (注意调查的固定/参考点)
 (k 阶帕斯卡/爱尔朗) 无限服务队列 $\pi(\lambda, E0)$

7 马尔可夫过程


马尔可夫链模型 转移概率矩阵 n 步转移概率. K-C方程

状态: 可达的 $i \rightarrow j \in A_{ci}$ 常返/非常返的 $\forall j \in A_{ci}, i \in A_{cj}$. (a存在性: $\forall i, \exists j \in A_{ci}$ 且常返)

常返类 R 周期的 $i \in S_k$ 且 $P_{ij} > 0$ 则 $j \in S_{k \bmod T+1}$. / 非周期的 $\exists n, s.t. \forall i, j \in R, P_{ij}(n) > 0$.
 $(U_S=R, \cap S=\emptyset)$

除多类链(看初态)可能有进入过程, 周期类 $(P_{ij}(n))$ 摆动外, \exists 稳态概率 $\pi_j = \lim_{n \rightarrow \infty} P(X_n=j)$
 (非无穷链) 平衡分布. \sim 方程组 $\begin{cases} \pi_j = \sum_i \pi_i P_{ij} \\ 1 = \sum_i \pi_i \end{cases}$ (与初始 X_0 无关)

· 门口没有伞

· 生灭过程  局部平衡方程 $\pi_i b_i = \pi_{i+1} d_{i+1}$. 长期行为的频率解释 (大数 T_n)

吸收概率 $a_i = P(X_n=s | X_0=i)$, s 为吸收的. \sim 期望时间方程组 $\begin{cases} \mu_i = 0 & \text{常返} \\ \mu_i = 1 + \sum_j P_{ij} \mu_j & \text{非常返} \end{cases}$
 (期望) 首次时间与回访时间 $(i \rightarrow s) \begin{cases} P_{ss} = 1 \\ t_s = 0 \\ t_i = 1 + \sum_j P_{ij} t_j \end{cases}$. $t_0^* = 1 + \sum_j P_{0j} t_j = \frac{1}{\pi_0}$.

连续时间马尔可夫链 $q_{ij} = \nu_i p_{ij}$

X_{ct} $\xrightarrow[\text{Exp}(\nu_i)]{i}$ $(q_{ij} \delta + 0 \delta)$

转移速率矩阵 同理有稳态收敛定理
 $\lim_{t \rightarrow \infty} P(X_{ct}=j | X_0=i) = \pi_j$
 π 为 $\begin{cases} \pi_j \sum_k q_{jk} = \sum_k \pi_k q_{kj} \\ 1 = \sum_k \pi_k \end{cases}$ 的解.

· (排队论) 缓冲器存信息 (丢包概率)

局部 \sim (时逆性) $\pi_j q_{ji} = \pi_i q_{ij}$

其中所有非常返态 j 有 $\pi_j = 0$.

抽样马尔可夫链 $Y_n = X_{ln}$. $r_{ij}(l), \pi_i$ 一致.

8 贝叶斯统计推断

最大后验概率 (MAP) 最小均方 (LMS) 线性最小均方 (LLMS)

· 匹配的滤波器 $X_i = a_i + W_i$
 $\|a\|^2 = \sum_i a_i^2$. (正态噪声模型 $\rightarrow \hat{\theta} = E(\theta) + \rho \frac{\sigma_\theta}{\sigma_x} (X - E(\theta)) = E(E(\theta|X))$)
 下, LMS 即 LLMS. 此时方差为 $(1-\rho^2)\sigma_\theta^2$. $\rho = \frac{Cov(\theta, X)}{\sigma_\theta \sigma_x}$

9 经典统计推断

贝叶斯线性回归 (最大后验概率, $\frac{\sigma^2}{2\sigma^2} + \frac{\sigma_1^2}{2\sigma^2}$)
 y LLMS. $\sum_i (y_i - \theta_0 - \theta_1 x_i)^2 \rightarrow \min$. 在 $\sigma^2, \sigma_1^2 \rightarrow \infty$ 不提供先验信息时
 (为最小二乘法提供了合理性) 与经典回归一致)

回归方法的考虑: 1. 异方差性 (加权 \sim) 2. 非线性 3. 多重共线性 4. 过度拟合 5. 因果关系

(f 为对参数取 \max , 此处为两点式)

似然比检验 $L(x) = \frac{f(x; H_1)}{f(x; H_0)}$ 临界值 ξ s.t. $P(L(x) > \xi; H_0) = \alpha$,

奈-皮引理: $P(L(x) > \xi; H_0) = \alpha$ 当观察 x 有 $L(x) > \xi$ 时拒绝 H_0 .



$P(L(x) \leq \xi; H_1) = \beta$ 若 \exists 其他检验拒绝域为 R ,

$P(X \in R; H_0) \leq \alpha$, 则 $P(X \in R; H_1) \geq \beta$.

(证明: $\xi = \frac{P_0(\theta_0)}{P_0(\theta_1)}$ 此时 MAP 其

犯错概率 $\leq \frac{\xi}{1+\xi} (\leq \alpha) + \frac{1}{1+\xi} (\geq \beta)$

(H_0, H_1 试作两点)

(左面计算 $L(x)$ 及其分布)

广义似然比检验 最大似然 $P_x(x; \theta)$ 得 $\hat{\theta}$, 计算 $\frac{P_x(x; \hat{\theta})}{P_x(x; \theta_0)}$ 根据 α 判断 H_0 .

χ^2 检验推导: $P_x(x; \theta) = c \theta_1^{n_1} \dots \theta_m^{n_m} \rightarrow \hat{\theta}_i = \frac{n_i}{n}$ ($n = \sum_i n_i$).

$\varphi: \prod_i \frac{(\frac{n_i}{n})^{n_i}}{(\theta_{0i})^{n_i}} > \xi$ v. 即 $\sum_i n_i \ln \frac{n_i}{n \theta_{0i}} > \ln \xi$ 当 $n \rightarrow \infty$ 时 (大数) $\frac{n_i}{n} \rightarrow \theta_{0i}$ (H_0 下),

\approx 所展开得 $\approx \frac{1}{2} \sum_i \frac{(n_i - n \theta_{0i})^2}{n \theta_{0i}}$

(分布证明见前)