

# 清华大学

## 1.2 条件概率与独立性

乘法公式:  $P(AB) = P(B)P(A|B)$   $P(A_1 \cdots A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1A_2) \cdots P(A_n|A_1 \cdots A_{n-1})$

拿球问题: 有放回/无放回.

全概率公式:  $P(A) = \sum_{i=1}^n P(B_i)P(A|B_i)$

彩票问题: 结论: 公平机会均等. 证明: 首肯分析法 (设  $A_k$  为  $n$  个人中的  $k$  个人中奖, 求概率公式)

Bayes 公式:  $P(B_i|A) = \frac{P(B_i)P(A|B_i)}{\sum_{j=1}^n P(B_j)P(A|B_j)}$

独立性  $P(AB) = P(A)P(B)$

相关系数 (一种特例)  $\tau(A, B) = \frac{P(AB) - P(A)P(B)}{\sqrt{P(A)P(A)P(B)(1-P(B))}}$

## 1.4 随机变量

离散 vs 连续:   
 取值有限 or 可列多  $\Rightarrow$  离散   
 存在密度函数  $\Rightarrow$  连续   
 既连续也不离散: 广泛存在.

常见随机变量.

负二项和泊松都满足可加性.

分布	分布列 or 密度函数	期望	方差	实际意义 & 和其他分布关系 & 性质.
二项 $b(n, p)$	$P_k = C_n^k p^k (1-p)^{n-k} \quad k=0, 1, \dots, n$	$np$	$np(1-p)$	
负二项 $nb(r, p)$	$P_k = C_{k-1}^{r-1} (1-p)^{k-r} p^r \quad k=r, r+1, \dots$	$\frac{r}{p}$	$\frac{r(1-p)}{p^2}$	第 $r$ 次成功时总次数, 可以和几何分布串联 (都是成功线)
泊松分布 $P(\lambda)$	$P_k = \frac{\lambda^k}{k!} e^{-\lambda} \quad k=0, 1, 2, \dots$	$\lambda$	$\lambda$	$\lambda$ 实际也理解为 "均分成很多小段后, 平均每段出现事件次数" 重复伯努利试验中 $np_n \rightarrow \lambda$ 则 $b(n, p) \rightarrow P(\lambda)$
几何分布 $Ge(p)$	$P_k = (1-p)^{k-1} p \quad k=1, 2, \dots$	$\frac{1}{p}$	$\frac{1-p}{p^2}$	无记忆性 (需要)
正态分布 $N(\mu, \sigma^2)$	$P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad x \in (-\infty, +\infty)$	$\mu$	$\sigma^2$	标准化 $\frac{x-\mu}{\sigma} \sim N(0, 1)$ 可加性
均匀分布 $U(a, b)$	$P(x) = \frac{1}{b-a} \quad a < x < b$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	

指数 $E(\lambda)$	$p(x) = \lambda e^{-\lambda x} \quad (x \geq 0)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	无记忆性(重要) 几何分布

期望与方差.

$$E(X) = \begin{cases} \sum x_i p(x_i) \\ \int_{-\infty}^{+\infty} x p(x) dx \end{cases}$$

$$E(ax+b) = aE(X) + b$$

$$E(x_1 + \dots + x_n) = E(x_1) + \dots + E(x_n)$$

$$Var(X) = E((X-E(X))^2)$$

$$Var(aX+b) = a^2 Var(X)$$

$$Var(X) = E(X^2) - E^2(X)$$

期望和方差可能不存在 - Cauchy 分布.

性质: ①  $E(X)$  是  $E(X-C)$  取 min 时的  $C$  (中位数是  $E(X-d)$  min 时的  $C$ )

$$② \text{切比雪夫不等式 } \forall \varepsilon > 0 \quad P(|X-E(X)| \geq \varepsilon) \leq \frac{Var(X)}{\varepsilon^2}$$

切比雪夫估计离散量时一般取 0.5 估计  
而且这都还是大致的估计.

$$③ X, Y \text{ 独立, 则 } E(XY) = E(X)E(Y) \quad Var(XY) = Var(X) + Var(Y)$$

多元随机变量

$$F(x,y) = \int_{-\infty}^x \int_{-\infty}^y p(u,v) du dv \quad p(x,y) = \frac{\partial^2}{\partial x \partial y} F(x,y)$$

$$\text{边际分布 } F_X(x) = \int_{-\infty}^{+\infty} F(x,y) dy$$

$$\text{边际密度 } p_X(x) = \int_{-\infty}^{+\infty} p(x,y) dy$$

$$\text{独立性: } P(X_1=x_1, \dots, X_n=x_n) = P(X_1=x_1) P(X_2=x_2) \dots P(X_n=x_n)$$

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P_{X_i}(x_i)$$

$$\Leftrightarrow F(X_1, \dots, X_n) = \prod_{i=1}^n F_{X_i}(x_i)$$

验证是否独立: ① 把边缘密度求出来看乘积是否为联合密度

② 联合密度能拆成  $f_{X_1}(x_1) f_{X_2}(x_2)$  的独立.

多元随机变量求方差: 多元积分. 利用对称性

$$E(X-Y)^2 = Var(X-Y) \text{ 当 } E(X-Y)=0 \text{ 时.}$$



# 清华大学

## L7 协方差与相关系数

$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y) = E[(X - EX)(Y - EY)] \quad \text{双线性: } \text{Cov}(C_1X + C_2Y + C_3Z, C_4X + C_5Y + C_6Z) = C_1C_4\text{Cov}(X, Y) + \dots$$

$$\text{和方差公式: } \text{Var}(X+Y) = \text{Var}X + \text{Var}Y + 2\text{Cov}(X, Y)$$

$$[\text{Cov}(X, Y)]^2 \leq \text{Var}X \text{Var}Y \quad \text{证明用 Cauchy 方法, 设 } E[(X-EX)^2 + 2tE[(X-EX)(Y-EY)] + t^2E[(Y-EY)^2]] = g(t) \text{ 求导}$$

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}X}\sqrt{\text{Var}Y}} \in [-1, 1] \quad \text{不相关} \Leftrightarrow \text{独立}$$

$$\text{二元正态分布} \quad p(x, y) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left[-\frac{1}{2} (x-\bar{\mu})^T \Sigma^{-1} (x-\bar{\mu})\right] \quad \text{元来说 } n=2 \quad \Sigma = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$$(X, Y) \sim N(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$$

$$\text{二元正态} \Leftrightarrow \text{边缘是一元正态} \quad (\text{反之不成立, 如 } W \sim \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} X \sim N(0, 1) \quad Y = WX \sim N(0, 1) \quad X, Y \text{ 不独立且相关})$$

二元正态下 独立  $\Rightarrow$  不相关 (正态分布间仅存线性关系)

$$\text{条件分布: } X \sim P(\lambda_1), Y \sim P(\lambda_2) \text{ 则 } X|X+Y \sim b(n, \frac{\lambda_1}{\lambda_1+\lambda_2}) \quad X+Y \sim P(\lambda_1+\lambda_2)$$

$$P_{X|Y}(x|y) = \frac{P(x, y)}{P_Y(y)}$$

## L8

条件期望与随机变量函数

$$E(X|Y) = g(Y)$$

$$\text{离散型} \quad E(X|Y) \sim \frac{E(X|Y=y_i)}{P(Y=y_i)} \quad \dots \quad \frac{E(X|Y=y_n)}{P(Y=y_n)} \leftarrow Y=y_n \text{ 下的期望}$$

$$\text{连续型} \quad Z = E(X|Y) \sim Z = E(X|Y=y) \text{ 时 } p_Y(y) \rightarrow p_Z(z)$$

$$\text{重期望} \quad E(X) = E(E(X|Y)) \Rightarrow \text{应用递归方法求期望方差}$$

$$\frac{\lambda_1^{x_1} \lambda_2^{x_2}}{x_1! x_2!} e^{-\lambda_1 - \lambda_2} \quad \frac{1}{n!} \sum_{i=0}^n \frac{n!}{i! (n-i)!} \lambda_1^i \lambda_2^{n-i} e^{-\lambda_1 - \lambda_2} = \frac{\lambda_1^n}{n!} e^{-\lambda_1}$$

$$\frac{(\lambda_1 + \lambda_2)^n}{n!} e^{-(\lambda_1 + \lambda_2)}$$



# 随机变量函数

方法一: 先求F再求P.  $\Rightarrow$  注意定义域, 分类讨论.

方法二: 对连续型用求导直接替换

$$P_Y(y_1, \dots, y_n) = \begin{cases} \sum P_X(x_1^{(1)}(y_1, \dots, y_n) \dots x_n^{(1)}(y_1, \dots, y_n)) \cdot |J^{(1)}| & \text{当 } (y_1, \dots, y_n) \text{ s.t. } y_i = g_i(x) \text{ 有解时} \\ 0 & \text{其他} \end{cases}$$

$$J^{(1)} = \frac{\partial x_1^{(1)} \dots x_n^{(1)}}{\partial y_1 \dots y_n}$$

卷积公式:  $(X, Y) \sim P_{XY}(x, y)$   $Z = X + Y$  满足  $P_Z(z) = \int_{-\infty}^{\infty} P_{XY}(z-w, w) dw$

次序统计量  $P_k(x) = \frac{n!}{(k-1)!(n-k)!} (F(x))^{k-1} p(x) (1-F(x))^{n-k}$  如:  $Z = \min\{X_1, \dots, X_n\} \sim \text{Exp}(\lambda_1 + \dots + \lambda_n)$   
 $X_i \sim \text{Exp}(\lambda_i)$

## L9. 大数定律

LLN: 伯努利.  $S_n$  为  $n$  次伯努利试验中  $A$  发生次数,  $p$  为每次  $A$  发生的概率, 则  $\forall \varepsilon > 0$   $\lim_{n \rightarrow \infty} P(|\frac{S_n}{n} - p| < \varepsilon) = 1$

切比雪夫:  $X_1, \dots, X_n, \dots$  两两不相关, 方差有界,  $\forall \varepsilon > 0$  有  $\lim_{n \rightarrow \infty} P(|\frac{1}{n} \sum_{k=1}^n X_k - \frac{1}{n} \sum_{k=1}^n E(X_k)| < \varepsilon) = 1$

依根号收敛  $\{X_n\}$  为随机变量序列,  $\forall \varepsilon > 0$   $\lim_{n \rightarrow \infty} P(|X_n - \mu| < \varepsilon) = 1$  则  $X_n \xrightarrow{p} \mu$

马尔可夫条件  $\frac{1}{n} \text{Var}(\sum_{k=1}^n X_k) \rightarrow 0$  满足条件即满足LLN一般形式

CLT:  $\{X_n\}$  独立同分布  $E(X_1) = \mu$   $\text{Var}(X_1) = \sigma^2$ , 则  $\forall y$  有  $\lim_{n \rightarrow \infty} P(\frac{X_1 + \dots + X_n - n\mu}{\sqrt{n}\sigma} \leq y) = \Phi(y)$

i.e.  $\frac{X_1 + \dots + X_n - n\mu}{\sqrt{n}\sigma} \rightarrow N(0, 1)$  更一般  $X_1 + \dots + X_n \sim N(\sum_{k=1}^n E(X_k), \sum_{k=1}^n \text{Var}(X_k))$

0.5校正: 对随机变量(离散型)的期望. 一般易如用于

Stat

## L1 统计量

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad \bar{X}, S^2 \text{ 是对 } E(X) \text{ Var}(X) \text{ 的无偏估计}$$

$$Y = X_1^2 + \dots + X_n^2 \sim \chi^2(n) \quad X_1, \dots, X_n \sim N(0, 1) \text{ i.i.d.}$$

$$T = \frac{X_1}{\sqrt{\frac{X_2^2 + \dots + X_n^2}{n-1}}} \sim t(n) \quad X_1, Z_1, \dots, Z_n \sim N(0, 1) \text{ i.i.d.}$$

$$F = \frac{X_1/m}{X_2/n} \sim F(m, n) \text{ 没有开根. } X_1 \sim \chi^2(m) \quad X_2 \sim \chi^2(n) \quad X_1, X_2 \text{ 独立.}$$

统计抽样定理  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$   
 $\bar{X}, S^2$  独立 i.i.d.

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$



# 清华大学

点估计

矩估计: 用样本矩代替理论矩  $\bar{X} \rightarrow E(X)$   $S^2 \rightarrow \text{Var}(X)$

一般先列出  $E(X)$ ,  $\text{Var}(X)$  关于  $\theta$  的表达式, 然后用  $\bar{X}, S^2$  代替后反解  $\theta$

极大似然 MLE:  $L(\theta) = L(\theta; X_1, \dots, X_n) = \prod_{k=1}^n p(X_k; \theta)$  称为似然函数  $l(\theta) = \ln L(\theta)$  为对数似然.

$\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$  s.t.  $L(\hat{\theta}) = \max_{\theta \in \Theta} L(\theta)$  则  $\hat{\theta}$  为极大似然估计

求  $\hat{\theta}$  ①求导 ②取值范围 ☆

点估计评价指标: ①相合性 ②无偏性 ③有效性  $\text{Var}(\hat{\theta}) \leq \text{Var}(\hat{\theta}_0)$  且对  $\theta$  s.t. 不等式成立  $\Leftarrow$  增大样本量

正态分布矩估计与 MLE:  $\begin{cases} \text{矩估计} & \hat{\mu} = \bar{X}, \hat{\sigma}^2 = S^2 \\ \text{MLE} & \hat{\mu}_{\text{MLE}} = \bar{X}, \hat{\sigma}_{\text{MLE}}^2 = \frac{n-1}{n} S^2 \end{cases}$

区间估计

置信水平: 随机区间至少以  $1-\alpha$  覆盖 置信数

枢轴量法: 构造  $T = g(X, \theta)$   $T$  分布与  $\theta$  无关 选取  $C, d$  满足  $P(C \leq T \leq d) \geq 1-\alpha$  再求解.

$X \sim N(\mu, \sigma^2)$   $\begin{cases} \sigma^2 \text{已知} & \frac{\sqrt{n}(\bar{X}-\mu)}{\sigma} \sim N(0,1) \\ \sigma^2 \text{未知} & \frac{(\bar{X}-\mu)/\frac{1}{\sqrt{n}}}{\sqrt{\frac{S^2}{n}}} = \frac{\sqrt{n}(\bar{X}-\mu)}{S} \sim t(n-1) \end{cases}$

双总体下.

$X \sim N(\mu_1, \sigma_1^2)$   $Y \sim N(\mu_2, \sigma_2^2)$ ,  $\sigma_1^2 = \sigma_2^2 = \sigma^2$  未知时  $t = \frac{\sqrt{mn(m+n-2)} (\bar{X}-\bar{Y}) - (\mu_1-\mu_2)}{\sqrt{(m-1)S_1^2 + (n-1)S_2^2}} \sim t(m+n-2)$  可用于估计  $\mu_1 - \mu_2$ .

方差估计:  $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$  作为枢轴量.

双总体下  $F = \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \sim F(m-1, n-1)$

估计区间长度: 不给  $X$  下可以用均值不等式优化.

## L4 假设检验

$H_0: \theta \in \Theta_0 \Leftrightarrow H_1: \theta \in \Theta_1$ ,  $W$  拒绝域 显著性水平  $\alpha = \max P(\text{拒 } H_0 | H_0 \text{ 为真})$

第I类错误: 拒真 第II类: 受伪  $\alpha$  起到限制第I类错误, 保护原假设作用.

U检验:  $\sigma^2$  已知  $H_0: \mu = \mu_0$   $H_1: \mu \neq \mu_0$   $W = \{ \bar{x} | |\bar{x} - \mu_0| > \frac{\sigma}{\sqrt{n}} U_{1-\frac{\alpha}{2}} \}$   
 $H_0: \mu \geq \mu_0$   $H_1: \mu < \mu_0$   $W = \{ \bar{x} | \bar{x} - \mu_0 < -\frac{\sigma}{\sqrt{n}} U_{1-\frac{\alpha}{2}} \}$

T检验:  $\sigma^2$  未知:  $H_0: \mu = \mu_0$   $H_1: \mu \neq \mu_0$   $W = \{ \bar{x} | |\bar{x} - \mu_0| > \frac{\sigma}{\sqrt{n}} t_{1-\frac{\alpha}{2}}(n-1) \}$   
 $H_0: \mu \geq \mu_0$   $H_1: \mu < \mu_0$   $W = \{ \bar{x} | \bar{x} - \mu_0 < -\frac{\sigma}{\sqrt{n}} t_{1-\frac{\alpha}{2}}(n-1) \}$

P值: 原假设成立条件下, 检验统计量比观测值更异常. 越小越好

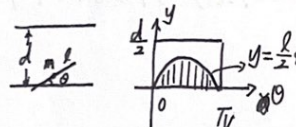
雨课堂:

快排:  $X_{ij}$  表示  $y_i, y_j$  是否被比较,  $X_{ij} = 1 \Leftrightarrow y_i, y_j$  是第一次比较的 (对于  $i, j$  之间的元素)

$$P(X_{ij}=1) = \frac{2}{j-i+1} \quad \text{比较次数} = E\left(\sum_{i=1}^n \sum_{j=i+1}^n X_{ij}\right) = \sum_{i=1}^n \sum_{j=i+1}^n E(X_{ij}) \sim n \ln n + O(n)$$

二元正态分布:  $X, Y \sim N(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$   $(ax+by, cx+dy)$  也服从二元正态.  
 $\Rightarrow$  构造线性组合, 用  $\text{Cov}(X, Y) = 0$  推  $X, Y$  独立.

1. 为随机数生成:  $\forall x \in \mathbb{R}$   $F(x)$  单增连续, 且  $F(-\infty) = 0 \leq F(x) \leq 1 = F(+\infty)$  则  $X = F^{-1}(U)$  ( $U \sim U(0,1)$ )  
 的累积分布函数为  $F(x)$

布丰投针问题.   $p = \frac{\int_0^{\frac{\pi}{2}} \frac{l}{2} \sin \theta d\theta}{\frac{\pi}{2}} = \frac{2l}{\pi d}$

贝特朗悖论: 弦长大于1的根号. 分析1-弦中心唯一确定  $\frac{1}{4}$  分析2-固定弦端点, 在圆弧上选取另弦端点  $\frac{1}{3}$   
 分析3-选直径  $\frac{1}{2}$  等可能假设不同.

Monty Hall问题: 换门:  $\frac{2}{3}$  不换:  $\frac{1}{3}$

康托尔分布: 既非连续也非离散



# 清华大学

4.5 拟合优度检验.

$$\chi^2 = \sum_{i=1}^n \frac{(N_i - N p_i)^2}{N p_i} \sim \chi^2(k-s-1) \quad s \text{ 为用样本估计的参数数, } k \text{ 为取值数.}$$

独立性检验:

$$\chi^2 = \sum_{i=1}^s \sum_{j=1}^t \frac{(n_{ij} - n \cdot \frac{c_i}{n} \cdot \frac{d_j}{n})^2}{n \cdot \frac{c_i}{n} \cdot \frac{d_j}{n}} \sim \chi^2((s-1)(t-1))$$

A \ B	1, 2, ..., t	行合计
1	$n_{11} \dots n_{1t}$	$C_1$
2	$\vdots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$
s	$n_{s1} \dots n_{st}$	$C_s$
列合计	$d_1 \dots d_t$	$N$