



第六章 样本及抽样分布

目录

1. 随机样本
2. 直方图和箱线图
3. 抽样分布



- 第1~5章内容讲述概率论的基本内容

在概率论中，所研究随机变量，它的分布都假设是已知的。在此前提下，去研究它的性质、特点及规律性，如求数字特征、变量的分布等；

- 第6~8章内容讲述数理统计部分

数理统计是具有广泛应用的数学分支，它以概率论为理论基础。在数理统计中，研究对象（随机变量）分布往往是未知的（或部分可知），我们需要根据实验或观测数据，对研究对象的客观规律性作出合理估计判断。

数理统计包括：如何收集、整理数据资料；如何对所得的数据资料进行分析、研究，从而对所研究的对象性质、特点作出推断。



学习统计无须把过多时间花在计算上，可以更有效地把时间用在基本概念、方法原理的正确理解上. 国内外著名的统计软件包： SAS， SPSS， MATLAB， STAT等，都可以让你快速、简便地进行数据处理和分析.

数理统计学是一门应用性很强的学科. 它关于数据资料收集、整理、分析、和推断的一门学科。对所考察的问题作出推断和预测, 直至为采取一定的决策和行动提供依据和建议.

数理统计学 { 合理收集数据-**试验设计、抽样调查**等
整理分析数据-**统计推断**



几个实际问题:

1. 估计产品寿命问题: 根据用户调查获得某品牌洗衣机50台的使用寿命为, 5, 5.5, 3.5, 6.2,。根据这些数据希望得到如下推断:

- A. 可否认为产品的平均寿命不低于4年?
- B. 保质期设为多少年, 才能保证有95%以上的产品过关?



2. 商品日投放量问题：如草莓的日投放量多少合理？如何安排银行各营业网点的现金投放量？快餐食品以什么样的速度生产最为合理等等。

与概率论一样, 数理统计也是研究大量随机现象的统计规律的一门数学学科, 它以概率论为理论基础, 根据试验或观察得到的数据, 对研究对象的客观规律性作出种种合理的估计和科学的推断.



数理统计的基本概念

数理统计的分类

描述统计学

对随机现象进行观测、试验，以取得有代表性的观测值

推断统计学

对已取得的观测值进行整理、分析，作出推断、决策，从而找出所研究的对象的规律性



推断统计学

推断 统计学

参数估计

假设检验

方差分析

回归分析





1. 随机样本



随机样本

一个统计问题总有它明确的研究对象.

总体: 研究对象的全体



研究某批灯泡的质量

该批灯泡寿命的全体
就是总体



考察国产轿车的油耗

所有国产轿车每公里耗油
量的全体就是总体

随机样本

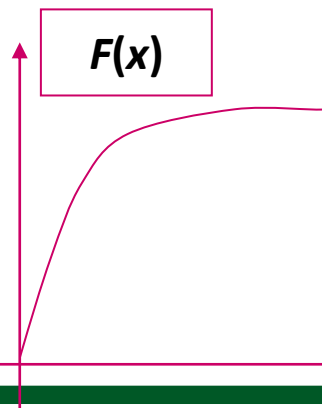
对某数量指标进行试验或观测，将试验的全部可能的观测值称为**总体**，对应一个随机变量；

总体可以用一个随机变量 X 或其分布来描述

如：研究某批灯泡的寿命时，我们关心的就是**寿命**，那么，寿命这个总体就可以用随机变量 X 表示，或用其分布函数 $F(x)$ 表示。



寿命可用一概率
(指数)分布来刻画



随机样本

个体

总体中每个对象称为个体.

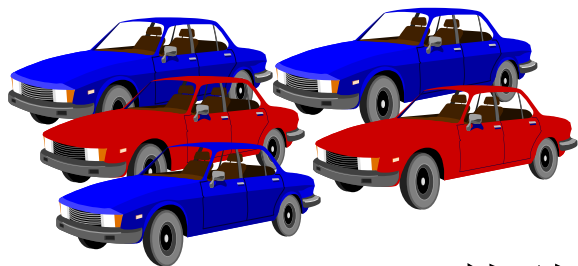
总体中所包含个体的个数称为总体的**容量**，容量为有限的称为**有限总体**，容量为无限的称为**无限总体**。

样本

为推断总体分布及各种特征, 按一定规则从总体中抽取若干个体进行观察试验以获得有关总体的信息. 所抽取的部分个体称为样本. 样本中所包含的个体数目称为**样本容量**.



随机样本



从国产轿车中
抽5辆进行耗
油量试验

样本容量为 5

抽到哪 5 辆是随机的！

样本是随机变量

容量为 n 的样本可以看作 n 维随机变量 (X_1, X_2, \dots, X_n) .

一旦取定一组样本，得到的是 n 个具体的数 x_1, x_2, \dots, x_n ,

称为样本 (X_1, X_2, \dots, X_n) 的一组观测值，简称样本值.



随机样本

实际中，总体的分布往往是未知的，或部分已知（类型或参数），数理统计中，我们需要通过从总体中抽取一部分个体（即样本），根据个体对总体分布作出推断。

在相同条件下，对总体 X 进行 n 次重复、独立的观察，将 n 次观察结果按试验次序记为 X_1, X_2, \dots, X_n ，可以认为它们相互独立，且服从相同分布的随机变量。



随机样本

简单随机样本

抽取的样本 X_1, X_2, \dots, X_n 满足下面两点:

1. 独立性: X_1, X_2, \dots, X_n 是相互独立的随机变量; 独立同分布;
2. 代表性: $X_i (i = 1, 2, \dots, n)$ 与所考察的总体 X 具有相同的分布.

简单随机样本是应用中最常见的情形, 今后, 说到

“ X_1, \dots, X_n 是来自某总体的样本” 时, 若不特别说明, 就指简单随机样本.



随机样本

定义： 设 X 是具有分布函数 F 的随机变量，若 X_1, X_2, \dots, X_n 是具有同一分布函数 F 的、相互独立的随机变量，则称 X_1, X_2, \dots, X_n 为从分布函数 F （或总体 F 、或总体 X ）得到的**容量为 n 的简单随机样本**，简称**样本**。它们的观察值 x_1, x_2, \dots, x_n 称为**样本值**，又称为 X 的 **n 个独立的观察值**。

将样本写成一个随机向量 (X_1, X_2, \dots, X_n) ， X_1, X_2, \dots, X_n 独立同分布。

所以随机向量 (X_1, X_2, \dots, X_n)

分布函数为：
$$F^*(x_1, x_2, \dots, x_n) = \prod_{i=1}^n F(x_i)$$

概率密度为：
$$f^*(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i)$$



随机样本

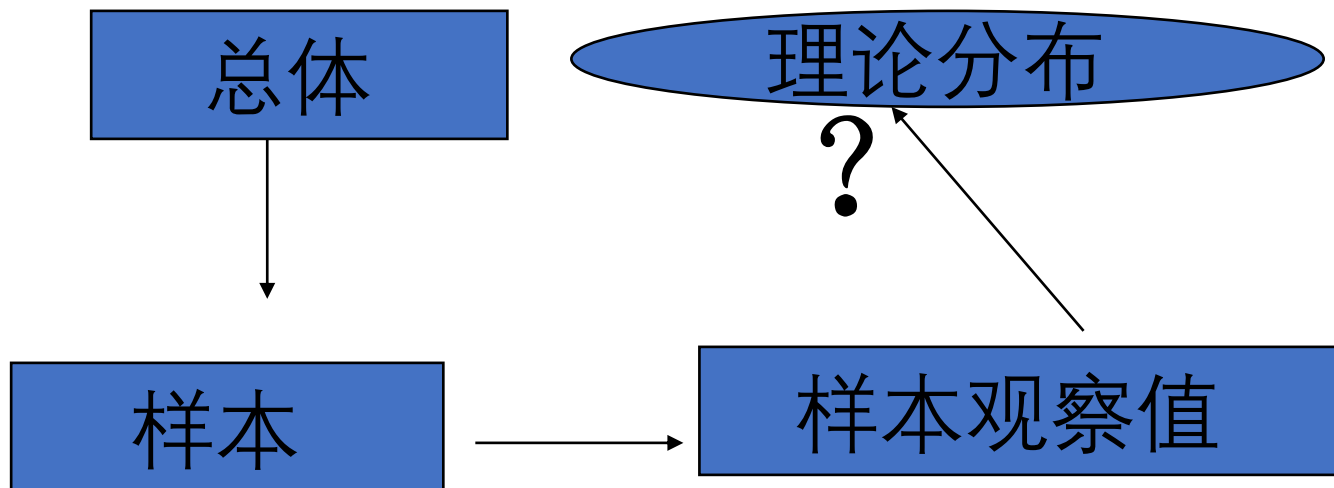
总体、样本、样本值的关系

事实上我们抽样后得到的资料都是具体的、确定的值. 如我们从某班大学生中抽取10人测量身高, 得到10个数, 它们是样本取到的值而不是样本. 我们只能观察到随机变量取的值而见不到随机变量.



随机样本

总体、样本、样本观察值的关系



样本空间 —— 样本所有可能取值的集合.





2. 直方图和箱线图



直方图和箱线图

为了研究总体分布的性质，人们通过试验得到许多观察值，一般来说这些数据是杂乱无章的。为了利用它们进行统计分析，将这些数据加以整理，还常借助于表格或图形对它们加以描述。

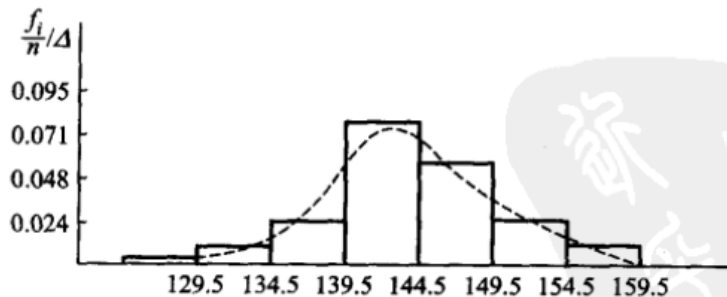
例：给出84个Etruscan人男子头颅的最大宽度(mm)，需要画出对应的频率直方图

141	148	132	138	154	142	150	146	155	158
150	140	147	148	144	150	149	145	149	158
143	141	144	144	126	140	144	142	141	140
145	135	147	146	141	136	140	146	142	137
148	154	137	139	143	140	131	143	141	149
148	135	148	152	143	144	141	143	147	146
150	132	142	142	143	153	149	146	149	138
142	149	142	137	134	144	146	147	140	142
140	137	152	145						

直方图和箱线图

解：进行整理。所有数据落在 $[126, 158]$ ，取区间 $[124.5, 159.5]$ 。将区间等分为7个小区间，小区间长度 $\Delta = \frac{159.5 - 124.5}{7} = 5$ ， Δ 称为组距。小区间端点称为组限。统计出每个小区间的频数、频率。

组 限	频 数 f_i	频率 f_i/n	累积频率
124.5~129.5	1	0.011 9	0.011 9
129.5~134.5	4	0.047 6	0.059 5
134.5~139.5	10	0.119 1	0.178 6
139.5~144.5	33	0.392 9	0.571 5
144.5~149.5	24	0.285 7	0.857 2
149.5~154.5	9	0.107 1	0.952 4
154.5~159.5	3	0.035 7	1



直方图的外廓曲线
接近于总体 X
的概率密度曲线

直方图和箱线图

定义： 设有容量为 n 的样本观察值 x_1, x_2, \dots, x_n ，样本 p 分位数($0 < p < 1$)记为 x_p ，它具有以下性质：

1. 至少有 np 个观测值小于或等于 x_p ；
2. 至少有 $n(1 - p)$ 个观测值大于或等于 x_p ；

样本 p 分位数可按以下法则求得，将 x_1, x_2, \dots, x_n 从小到大排序成 $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$.

$$x_p = \begin{cases} x_{([np]+1)}, np \notin Z \\ \frac{1}{2} [x_{(np)} + x_{(np+1)}], np \in Z \end{cases}$$



直方图和箱线图

特别，当 $p=0.5$ 时，0.5分位数 $x_{0.5}$ 也记为 Q_2 或 M ，称为样本中位数，即有

$$x_{0.5} = \begin{cases} x_{([n/2]+1)}, n \text{ 是奇数} \\ \frac{1}{2}[x_{(n/2)} + x_{(n/2+1)}], n \text{ 是偶数} \end{cases}$$

0.25分位数 $x_{0.25}$ 称为**第一四分位数**，又记为 Q_1 ；

0.75分位数 $x_{0.75}$ 称为**第三四分位数**，又记为 Q_3 ；



直方图和箱线图

例： 设有一组容量为18的样本值如下(已排序)

122	126	133	140	145	145	149	150	157
162	166	175	177	177	183	188	199	212

求样本分位数： $x_{0.2}$ ， $x_{0.25}$ ， $x_{0.5}$ 。

解：

(1). 因为 $np = 18 \times 0.2 = 3.6$ ， $x_{0.2}$ 位于 $[3.6] + 1 = 4$ 处，即有 $x_{0.2} = x_{(4)} = 140$.

(2). 因为 $np = 18 \times 0.25 = 4.5$ ， $x_{0.25}$ 位于 $[4.5] + 1 = 5$ 处，即有 $x_{0.25} = x_{(5)} = 145$.

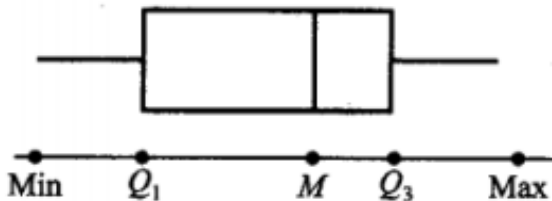
(3). 因为 $np = 18 \times 0.5 = 9$ ， $x_{0.5}$ 是这组数中间两个数的平均值，即有 $x_{0.5} = \frac{1}{2}(157 + 162) = 159.5$.



直方图和箱线图

箱线图:

确定5个点: Min , Q_1 , M , Q_3 , Max

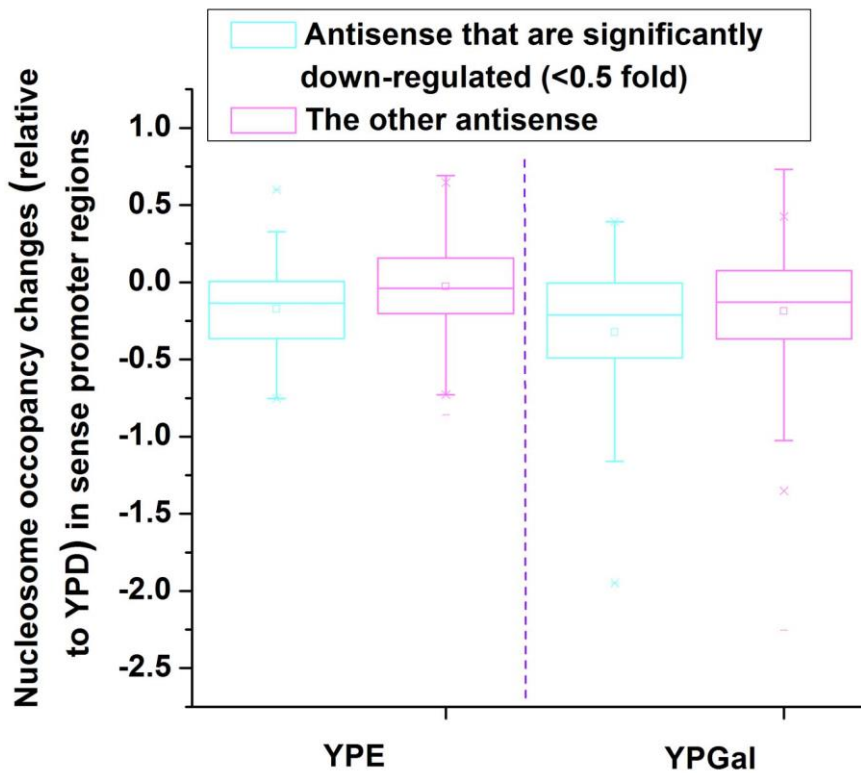


从箱线图可看出数据的分布特性:

1. 中心位置
2. 散布程度
3. 对称性

直方图和箱线图

箱线图:





3. 抽样分布



抽样分布

经验分布函数：我们还可以作出与总体分布函数 $F(x)$ 相应的统计量——经验分布函数。也称作**样本分布函数**用 $S(x)$, $-\infty < x < +\infty$ 表示 X_1, X_2, \dots, X_n 不大于 x 的随机变量的个数。定义经验分布函数 $F_n(x)$ 为

$$F_n(x) = \frac{1}{n} S(x), -\infty < x < +\infty$$

一般，设 x_1, x_2, \dots, x_n 是总体 F 的一个容量为 n 的样本值，将 x_1, x_2, \dots, x_n 从小到大次序排序，并重新编号，设为

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

经验分布函数 $F_n(x)$ 的观察值为

$$F_n(x) = \begin{cases} 0 & , x < x_{(1)} \\ \frac{k}{n}, & x_{(k)} \leq x < x_{(k+1)} \\ 1, & x \geq x_{(n)} \end{cases}$$



抽样分布

注1: $F_n(x)$ 为分布函数, 因为满足分布函数的特征性质

注2: $F_n(x)$ 实际上是累积频率直方图曲线。

注3:

由伯努利大数定律:

只要 n 相当大, $F_n(x)$ 依概率收敛于 $F(x)$ 。



抽样分布

对于经验分布函数，格里汶科 (Glivenko) 在1933年证明了：
对于任一实数 x ，当 $n \rightarrow \infty$ 时 $F_n(x)$ 以概率1一致收敛于分布函数 $F(x)$ ，即

$$P\{\lim_{n \rightarrow \infty} \sup_{-\infty < x < \infty} |F_n(x) - F(x)| = 0\} = 1$$

因此对于任一实数 x 当 n 充分大时，经验分布函数的任一个观察值 $F_n(x)$ 与总体分布函数 $F(x)$ ，只有微小的差别，从而在实际上当可当作 $F(x)$ 来使用。



抽样分布

统计量:

由样本推断总体特征, 需要对样本值进行“加工”, “提炼”. 这就需要构造一些样本的函数, 它把样本中所含的信息集中起来.

定义: 设 X_1, X_2, \dots, X_n 是来自总体 X 的一个样本, $g(X_1, X_2, \dots, X_n)$ 是 X_1, X_2, \dots, X_n 的函数, 若 g 中不含未知参数, 则称 $g(X_1, X_2, \dots, X_n)$ 是一**统计量**。



抽样分布

实例：

设 X_1, X_2, X_3 是来自总体 $N(\mu, \sigma^2)$ 的一个样本，其中 μ 为已知， σ^2 为未知，判断下列各式哪些是统计量，哪些不是？

$$T_1 = X_1,$$

$$T_2 = X_1 + X_2 e^{X_3},$$

$$T_3 = \frac{1}{3}(X_1 + X_2 + X_3),$$

是

$$T_4 = \max(X_1, X_2, X_3),$$

$$T_5 = X_1 + X_2 - 2\mu,$$

$$T_6 = \frac{1}{\sigma^2}(X_1^2 + X_2^2 + X_3^2).$$

不是



抽样分布

常见的统计量:

样本平均值: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

样本方差: $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right)$

样本标准差: $S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$

样本k阶(原点)矩: $A_k = \frac{1}{n} \sum_{i=1}^n X_i^k, k = 1, 2, \dots$

样本k阶中心矩: $B_k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k, k = 1, 2, \dots$



抽样分布

X_1, X_2, \dots, X_n 为来自 X 的样本
样本数字特征 (随机变量):

1. 样本均值 $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

2. 样本方差

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

3. 样本 k 阶原点矩

$$A_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

4. 样本 k 阶中心矩

$$B_k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k$$

x_1, x_2, \dots, x_n 为样本值
样本数字特征观察值:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

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

$$a_k = \frac{1}{n} \sum_{i=1}^n x_i^k$$

$$b_k = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^k$$



抽样分布

若总体 X 的 k 阶矩 $E(X^k) \stackrel{\text{记成}}{=} \mu_k$ 存在,

则当 $n \rightarrow \infty$ 时, $A_k \xrightarrow{P} \mu_k, k = 1, 2, \dots$ 。

证: X_1, X_2, \dots, X_n 独立且与 X 同分布, 所以 $X_1^k, X_2^k, \dots, X_n^k$ 独立且与 X^k 同分布, 故有

$$E(X_1^k) = E(X_2^k) = \dots = E(X_n^k) = \mu_k$$

由辛钦大数定理

$$A_k = \frac{1}{n} \sum_{i=1}^n X_i^k \xrightarrow{P} \mu_k, k = 1, 2, \dots$$

进而由依概率收敛的序列的性质知

$$g(A_1, A_2, \dots, A_k) \xrightarrow{P} g(\mu_1, \mu_2, \dots, \mu_k)$$

其中 g 为连续函数, 这是下一章矩估计法理论基础。



抽样分布

来自正态总体的几个常用统计量的分布

统计学上的三大分布： χ^2 分布、 t 分布和 F 分布。

正态分布

定理：若 X_1, X_2, \dots, X_n 相互独立， $X_i \sim N(\mu_i, \sigma_i^2)$ ，则

$$\sum_{i=1}^n a_i X_i \sim N\left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2\right)$$

特别地，若 $X_1, X_2, \dots, X_n \stackrel{i.i.d}{\sim} N(\mu, \sigma^2)$

则 $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

$$U = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \sim N(0, 1)$$

在已知总体 μ ， σ^2 时，可用本定理计算样本均值 \bar{X} 。



抽样分布

标准正态分布的 α 分位数

定义

若 $P(X > u_\alpha) = \alpha$, 则称 u_α 为标准正态分布的上 α 分位数.

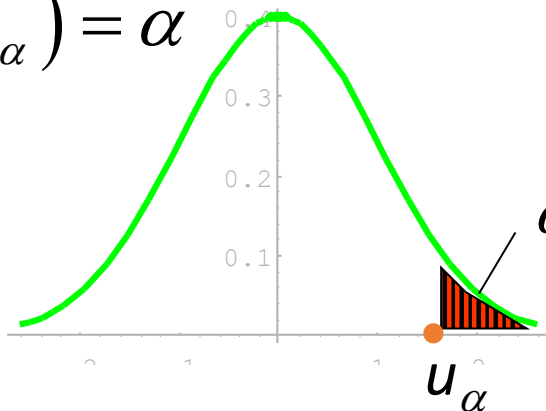
若 $P(|X| > u_{\frac{\alpha}{2}}) = \alpha$, 则称 $u_{\frac{\alpha}{2}}$ 为标准正态分布的双侧 α 分位数.



抽样分布

标准正态分布的 α 分位数图形

$$P(X > u_\alpha) = \alpha$$



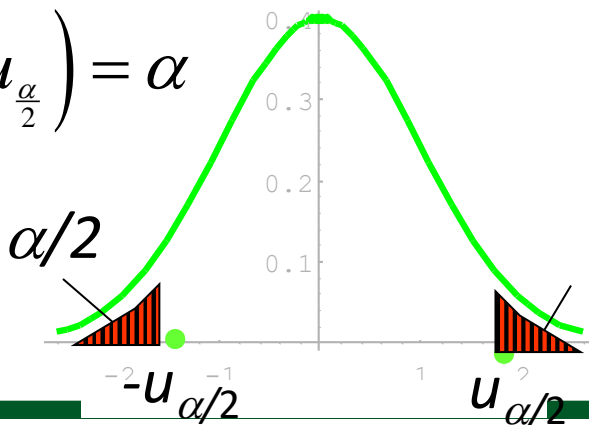
$$u_{0.05} = 1.645$$

$$u_{0.025} = 1.96$$

$$u_{0.005} = 2.575$$

常用
数字

$$P(|X| > u_{\frac{\alpha}{2}}) = \alpha$$



$$-u_{\alpha/2} = u_{1-\alpha/2}$$

抽样分布

(一) χ^2 分布

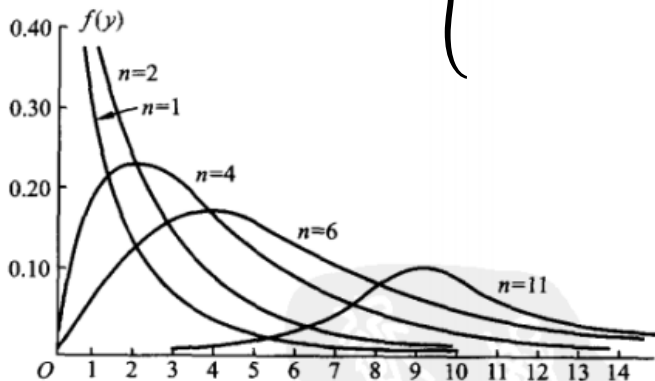
设 X_1, X_2, \dots, X_n 是来自总体 $N(0,1)$ 的样本, 则称统计量

$$\chi^2 = X_1^2 + X_2^2 \dots + X_n^2$$

服从自由度为 n 的 χ^2 分布, 记为 $\chi^2 \sim \chi^2(n)$.

概率密度为

$$f(y) = \begin{cases} \frac{1}{2^{n/2} \Gamma(n/2)} y^{n/2-1} e^{-y/2}, & y > 0 \\ 0, & y \leq 0 \end{cases}$$



抽样分布

Gamma Function

对于 $\alpha > 0$, Gamma Function $\Gamma(\alpha)$ 定义为

$$\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} e^{-x} dx$$

重要性质

1. 对 $\alpha \geq 1$, $\Gamma(\alpha + 1) = \alpha \Gamma(\alpha)$;
2. 对正整数 n , $\Gamma(n) = (n - 1)!$
3. $\Gamma\left(\frac{1}{2}\right) = \pi^{\frac{1}{2}}$

Γ 分布见第三章第5节例3 (P78)

例 3 设随机变量 X, Y 相互独立, 且分别服从参数为 $\alpha, \theta; \beta, \theta$ 的 Γ 分布 (分别记成 $X \sim \Gamma(\alpha, \theta), Y \sim \Gamma(\beta, \theta)$). X, Y 的概率密度分别为

$$f_X(x) = \begin{cases} \frac{1}{\theta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\theta}, & x > 0, \\ 0, & \text{其他,} \end{cases} \quad \alpha > 0, \theta > 0.$$
$$f_Y(y) = \begin{cases} \frac{1}{\theta^\beta \Gamma(\beta)} y^{\beta-1} e^{-y/\theta}, & y > 0, \\ 0, & \text{其他,} \end{cases} \quad \beta > 0, \theta > 0.$$

试证明 $Z = X + Y$ 服从参数为 $\alpha + \beta, \theta$ 的 Γ 分布, 即 $X + Y \sim \Gamma(\alpha + \beta, \theta)$.

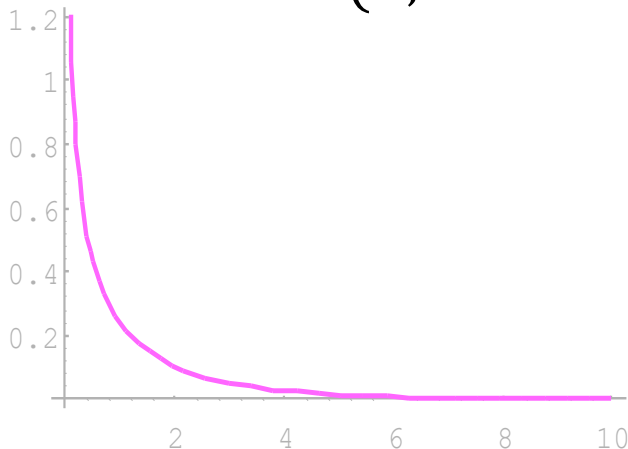


抽样分布

对于 χ^2 分布

注1 $n = 1$ 时, 其密度函数为

$$f(x) = \begin{cases} \frac{1}{\sqrt{2\pi}} x^{-\frac{1}{2}} e^{-\frac{x}{2}}, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

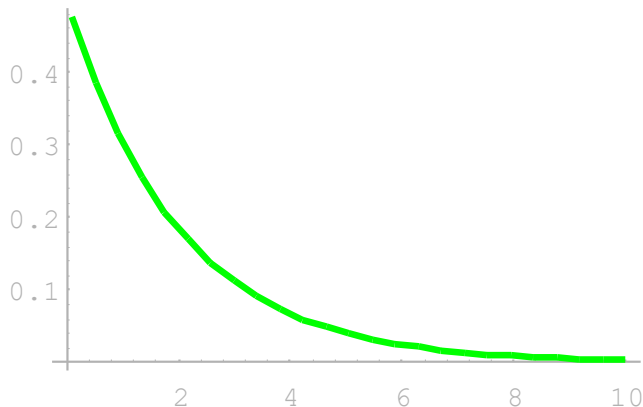


抽样分布

注2 $n = 2$ 时, 其密度函数为

$$f(x) = \begin{cases} \frac{1}{2}e^{-\frac{x}{2}}, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

为参数为1/2的指数分布.



抽样分布

由第二章第五节例3知 $\chi^2(1)$ 分布即为 $\Gamma\left(\frac{1}{2}, 2\right)$ 分布, 由

X_1, X_2, \dots, X_n 的独立性知 $X_1^2, X_2^2, \dots, X_n^2$ 相互独立, 由 Γ 分布的可加性知

$$\chi^2 = \sum_{i=1}^n X_i^2 \sim \Gamma\left(\frac{n}{2}, 2\right)$$

即可得到 χ^2 的概率密度。



抽样分布

χ^2 分布的性质

1° χ^2 分布的可加性：设 $\chi_1^2 \sim \chi^2(n_1)$, $\chi_2^2 \sim \chi^2(n_2)$, 且 χ_1^2, χ_2^2 相互独立

$$\chi_1^2 + \chi_2^2 \sim \chi^2(n_1 + n_2)$$

2° χ^2 分布的数学期望和方差：若 $\chi^2 \sim \chi^2(n)$ 则

$$E(\chi^2) = n, D(\chi^2) = 2n$$

3° $n \rightarrow \infty$ 时, $\chi^2(n) \rightarrow$ 正态分布

4° $\chi^2(n)$ 分布的上 α 分位数有表可查



抽样分布

证2: 设 $\chi^2(n) = \sum_{i=1}^n X_i^2$ $X_i \sim N(0,1)$ $i = 1, 2, \dots, n$,

X_1, X_2, \dots, X_n 相互独立

则 $E(X_i) = 0$, $D(X_i) = 1$, $E(X_i^2) = 1$

$$E(\chi^2(n)) = E\left(\sum_{i=1}^n X_i^2\right) = n$$

$$E(X_i^4) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^4 e^{-\frac{x^2}{2}} dx = 3$$

$$D(X_i^2) = E(X_i^4) - E^2(X_i^2) = 2$$

$$D(\chi^2(n)) = D\left(\sum_{i=1}^n X_i^2\right) = 2n$$



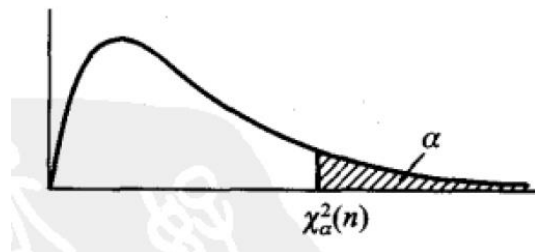
抽样分布

χ^2 分布的分位点

对于给定的 α , $0 < \alpha < 1$, 满足条件

$$P\{\chi^2 > \chi_{\alpha}^2(n)\} = \int_{\chi_{\alpha}^2(n)}^{\infty} f(y)dy = \alpha$$

的点 $\chi_{\alpha}^2(n)$ 为 $\chi^2(n)$ 分布的上 α 分位点。



对于不同的 α , n , 上 α 分位点的值已制表, 可以查用。

例如, 对于 $\alpha = 0.1, n = 25$, 查得 $\chi_{0.1}^2(25) = 34.382$ 。

抽样分布

但该表只详列到 $n = 40$ 为止

费希尔 (R. A. Fisher) 证明, 当 n 充分大时

$$\chi^2_{\alpha}(n) \approx \frac{1}{2} (z_{\alpha} + \sqrt{2n-1})^2$$

z_{α} 是标准正态分布的上 α 分位点, 利用该式子可求得 $n > 40$ 时 $\chi^2(n)$ 的上 α 分位点近似值。



抽样分布

(二) t 分布

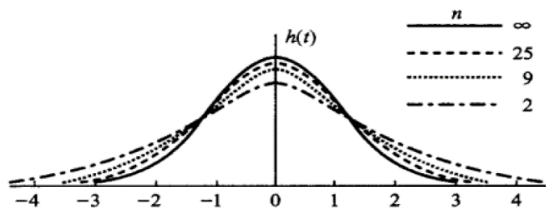
设 $X \sim N(0,1)$, $Y \sim \chi^2(n)$, 且 X , Y 相互独立, 称随机变量

$$t = \frac{X}{\sqrt{Y/n}}$$

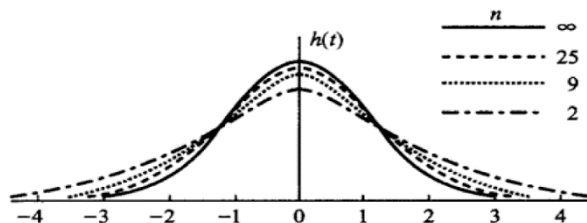
服从自由度为 n 的 t 分布, 记为 $t \sim t(n)$

t 分布又称学生(Student)分布, $t(n)$ 分布概率密度函数为

$$h(t) = \frac{\Gamma[(n+1)/2]}{\sqrt{\pi n} \Gamma(n/2)} \left(1 + \frac{t^2}{n}\right)^{-(n+1)/2}, \quad -\infty < t < +\infty$$



抽样分布



$h(t)$ 的图形关于 $t = 0$ 对称，当 n 充分大时其图形类似于标准正态概率密度图形。

具有自由度为 n 的 t 分布 $t \sim t(n)$,其数学期望与方差为:

$$E(t) = 0, D(t) = n/(n-2) \quad (n > 2)$$

利用 Γ 函数性质可得 $\lim_{n \rightarrow \infty} h(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}$

n 足够大时， t 分布近似于 $N(0,1)$

抽样分布

t 分布的分位点

对于给定的 α , $0 < \alpha < 1$, 满足条件

$$P\{t > t_{\alpha}(n)\} = \int_{t_{\alpha}(n)}^{+\infty} h(t)dt = \alpha$$

的点 $t_{\alpha}(n)$ 称为 $t(n)$ 分布的上 α 分位点。

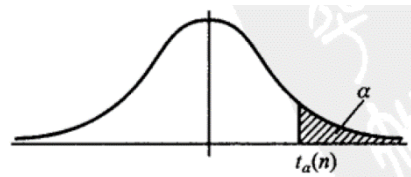
根据 t 分布的对称性及上 α 分位点定义可知

$$t_{1-\alpha}(n) = -t_{\alpha}(n)$$

t 分布的上 α 分位点可自附表4查得。

当 $n > 45$, 对于常用的 α 的值, 就用正态分布近似

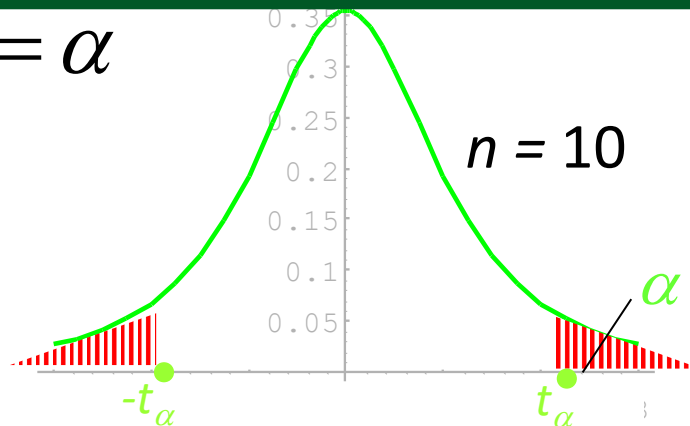
$$t_{\alpha}(n) \approx z_{\alpha}$$



抽样分布

$$P(T > t_{\alpha}) = \alpha$$

$$-t_{\alpha} = t_{1-\alpha}$$



$$P(T > 1.8125) = 0.05 \Rightarrow t_{0.05}(10) = 1.8125$$

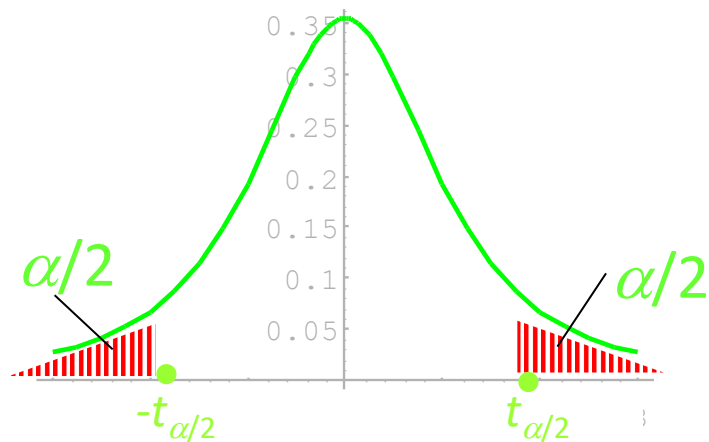
$$P(T < -1.8125) = 0.05, \quad P(T > -1.8125) = 0.95$$

$$\Rightarrow t_{0.95}(10) = -1.8125$$

抽样分布

$$P(T > t_{\alpha/2}) = \frac{\alpha}{2}$$

$$P(|T| > t_{\alpha/2}) = \alpha$$



$$P(T > 2.2281) = 0.025$$

$$P(|T| > 2.2281) = 0.05$$

$$\Rightarrow t_{0.025}(10) = 2.2281$$



抽样分布

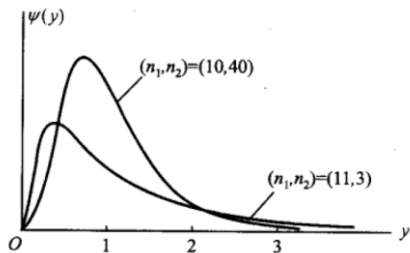
(三) F 分布

设 $U \sim \chi^2(n_1)$, $V \sim \chi^2(n_2)$, 且 U, V 相互独立, 则称随机变量

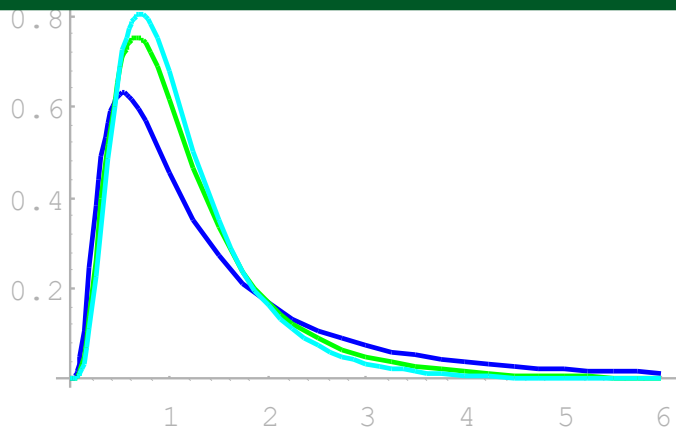
$$F = \frac{U/n_1}{V/n_2}$$

服从自由度为 (n_1, n_2) 的 F 分布, 记为 $F \sim F(n_1, n_2)$, 概率密度函数为

$$\varphi(y) = \begin{cases} \frac{\Gamma[(n_1 + n_2)/2](n_1/n_2)^{n_1/2} y^{(n_1/2)-1}}{\Gamma(n_1/2)\Gamma(n_2/2)[1 + (n_1 y/n_2)]^{(n_1+n_2)/2}}, & y > 0 \\ 0, & \text{others} \end{cases}$$



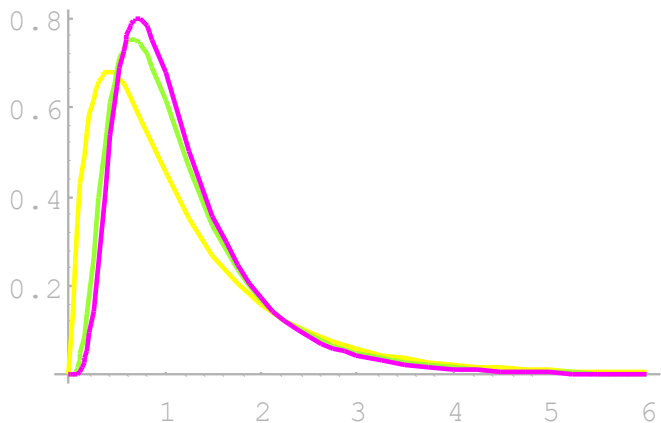
抽样分布



$$m = 10, n = 4$$

$$m = 10, n = 10$$

$$m = 10, n = 15$$



$$m = 4, n = 10$$

$$m = 10, n = 10$$

$$m = 15, n = 10$$



抽样分布

从定义可知，若 $F \sim F(n_1, n_2)$ ，则

$$\frac{1}{F(n_1, n_2)} \sim F(n_2, n_1)$$

F 分布的数学期望为：

$$E(F) = \frac{n_2}{n_2 - 2}, \quad n_2 > 2$$

即它的数学期望并不依赖于第一自由度 n_1 .



抽样分布

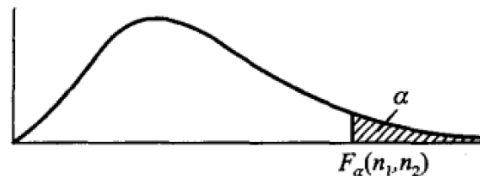
F分布的分位点

对于给定的 α , $0 < \alpha < 1$, 满足条件

$$P\{F > F_{\alpha}(n_1, n_2)\} = \int_{F_{\alpha}(n_1, n_2)}^{+\infty} \varphi(y) dy = \alpha$$

的点 $F_{\alpha}(n)$ 称为 $F(n_1, n_2)$ 分布的上 α 分位点。

查询见附表6



F分布的上 α 分位点重要性质

$$F_{1-\alpha}(n_1, n_2) \sim \frac{1}{F_{\alpha}(n_2, n_1)}$$

常用来求F分布表中未列出的常用的上 α 分位点。

抽样分布

(四) 正态总体的样本均值与样本方差的分布

设总体 X (可为任何分布, 只要均值和方差存在) 的均值为 μ , 方差为 σ^2 , X_1, X_2, \dots, X_n 是来自 X 的一个样本, 则其样本均值 \bar{X} 和方差 S^2 有如下性质

$$\begin{aligned} E(\bar{X}) &= \mu, D(\bar{X}) = \frac{\sigma^2}{n} \\ E(S^2) &= E\left(\frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right)\right) = \frac{1}{n-1} \left[\sum_{i=1}^n E(X_i^2) - nE(\bar{X}^2) \right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n (\sigma^2 + \mu^2) - n(\sigma^2/n + \mu^2) \right] = \sigma^2 \end{aligned}$$



抽样分布

进而，设 $X \sim N(\mu, \sigma^2)$ ，知正态分布的线性组合 $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

也服从正态分布，得到以下定理

定理一： 设 X_1, X_2, \dots, X_n 是来自正态总体 $N(\mu, \sigma^2)$ 的样本，则其样本均值

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

定理二： 设 X_1, X_2, \dots, X_n 是来自正态总体 $N(\mu, \sigma^2)$ 的样本，则

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

2: \bar{X} 与 S^2 相互独立

证明见本章末附录。



抽样分布

定理三： 设 X_1, X_2, \dots, X_n 是来自正态总体 $N(\mu, \sigma^2)$ 的样本，则

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$$

证：由定理一、定理二可知

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1), \quad \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

且两者独立，则由 t 分布定义可知

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} / \sqrt{\frac{(n-1)S^2}{\sigma^2(n-1)}} = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$$



抽样分布

定理四： 设 X_1, X_2, \dots, X_{n_1} 与 Y_1, Y_2, \dots, Y_{n_2} 分别是来自正态总体 $X \sim N(\mu_1, \sigma_1^2)$ 和 $Y \sim N(\mu_2, \sigma_2^2)$ 的样本，且这两个样本相互独立。设 \bar{X} ， \bar{Y} 分别是这两个样本的样本均值， S_1^2 ， S_2^2 分别是这两个样本的样本方差，则有

1.

$$\frac{S_1^2/S_2^2}{\sigma_1^2/\sigma_2^2} \sim F(n_1 - 1, n_2 - 1)$$

2. 当 $\sigma_1^2 = \sigma_2^2 = \sigma^2$ 时，

$$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{S_w \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(n_1 + n_2 - 2)$$

$$\text{其中 } S_w^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}, S_w = \sqrt{S_w^2}$$



抽样分布

证1：由定理二

$$\frac{(n_1 - 1)S_1^2}{\sigma_1^2} \sim \chi^2(n_1 - 1), \frac{(n_2 - 1)S_2^2}{\sigma_2^2} \sim \chi^2(n_2 - 1)$$

由假设 S_1^2, S_2^2 相互独立，则由 F 分布的定义知

$$\frac{(n_1 - 1)S_1^2}{(n_1 - 1)\sigma_1^2} / \frac{(n_2 - 1)S_2^2}{(n_2 - 1)\sigma_2^2} \sim F(n_1 - 1, n_2 - 1)$$

即

$$\frac{S_1^2/S_2^2}{\sigma_1^2/\sigma_2^2} \sim F(n_1 - 1, n_2 - 1)$$



抽样分布

证2: 易知 $\bar{X} - \bar{Y} \sim N(\mu_1 - \mu_2, \frac{\sigma^2}{n_1} + \frac{\sigma^2}{n_2})$, 即有

$$U = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim N(0,1).$$

又由给定条件知

$$\frac{(n_1 - 1)S_1^2}{\sigma^2} \sim \chi^2(n_1 - 1), \frac{(n_2 - 1)S_2^2}{\sigma^2} \sim \chi^2(n_2 - 1)$$

且它们相互独立, 故由 χ^2 分布的可加性知

$$V = \frac{(n_1 - 1)S_1^2}{\sigma^2} + \frac{(n_2 - 1)S_2^2}{\sigma^2} \sim \chi^2(n_1 + n_2 - 2).$$



抽样分布

证2(续):

由本章附录2知 U 与 V 相互独立, 从而按 t 分布的定义知

$$\frac{U}{\sqrt{V/(n_1 + n_2 - 2)}} = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{S_w \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(n_1 + n_2 - 2)$$

本节所介绍的几个分布以及四个定理, 在下面各章都起着重要的作用。应注意, 它们都是在总体为正态这一基本假定下得到的。



抽样分布

例： $X \sim N(72, 100)$, 为使样本均值大于70的概率不小于90%, 则样本容量至少取多少?

解： 设样本容量为 n , 则 $\bar{X} \sim N(72, \frac{100}{n})$

$$\begin{aligned} \text{故 } P(\bar{X} > 70) &= 1 - P(\bar{X} \leq 70) \\ &= 1 - \Phi\left(\frac{70 - 72}{\frac{10}{\sqrt{n}}}\right) = \Phi(0.2\sqrt{n}) \end{aligned}$$

令 $\Phi(0.2\sqrt{n}) \geq 0.9$ 得 $0.2\sqrt{n} \geq 1.29$

即 $n \geq 41.6025$

所以取 $n = 42$

抽样分布

例：从正态总体 $X \sim N(\mu, \sigma^2)$ 中，抽取了 $n = 20$ 的样本 $(X_1, X_2, \dots, X_{20})$

(1) 求 $P\left(0.37\sigma^2 \leq \frac{1}{20} \sum_{i=1}^{20} (X_i - \bar{X})^2 \leq 1.76\sigma^2\right)$

(2) 求 $P\left(0.37\sigma^2 \leq \frac{1}{20} \sum_{i=1}^{20} (X_i - \mu)^2 \leq 1.76\sigma^2\right)$

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

即 $\frac{19S^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^{20} (X_i - \bar{X})^2 \sim \chi^2(19)$

故 $P\left(0.37\sigma^2 \leq \frac{1}{20} \sum_{i=1}^{20} (X_i - \bar{X})^2 \leq 1.76\sigma^2\right)$

$$= P\left(7.4 \leq \frac{1}{\sigma^2} \sum_{i=1}^{20} (X_i - \bar{X})^2 \leq 35.2\right)$$

$$= P\left(\frac{1}{\sigma^2} \sum_{i=1}^{20} (X_i - \bar{X})^2 \geq 7.4\right) - P\left(\frac{1}{\sigma^2} \sum_{i=1}^{20} (X_i - \bar{X})^2 \geq 35.2\right)$$

查表 $= 0.99 - 0.01 = 0.98$

抽样分布

解：(2) $\sum_{i=1}^{20} \left(\frac{X_i - \mu}{\sigma} \right)^2 \sim \chi^2(20)$

故

$$\begin{aligned} & P \left(0.37\sigma^2 \leq \frac{1}{20} \sum_{i=1}^{20} (X_i - \mu)^2 \leq 1.76\sigma^2 \right) \\ &= P \left(7.4 \leq \sum_{i=1}^{20} \left(\frac{X_i - \mu}{\sigma} \right)^2 \leq 35.2 \right) \\ &= P \left(\sum_{i=1}^{20} \left(\frac{X_i - \mu}{\sigma} \right)^2 \geq 7.4 \right) - P \left(\sum_{i=1}^{20} \left(\frac{X_i - \mu}{\sigma} \right)^2 \geq 35.2 \right) \\ &= 0.995 - 0.025 = 0.97 \end{aligned}$$

抽样分布

例： $(X_1, X_2, \dots, X_{10})$ 是来自总体 $N(0, 0.3^2)$ 的样本，求 $P(\sum_{i=1}^{10} X_i^2 > 1.44)$.

解： 将 $X_i, i = 1, 2, \dots, 10$ 标准化得 $\frac{X_i - 0}{0.3} \sim N(0, 1), i = 1, 2, \dots, 10$.

由 χ^2 分布的构造得 $\sum_{i=1}^{10} \left(\frac{X_i}{0.3}\right)^2 \sim \chi^2(10)$. 因此有

$$\begin{aligned} P\left(\sum_{i=1}^{10} X_i^2 > 1.44\right) &= P\left\{\sum_{i=1}^{10} \left(\frac{X_i}{0.3}\right)^2 > \frac{1.44}{0.09}\right\} \\ &= P\left\{\sum_{i=1}^{10} \left(\frac{X_i}{0.3}\right)^2 > 16\right\} = 0.1 \end{aligned}$$

抽样分布

例： 设 X 与 Y 相互独立， $X \sim N(0,16)$, $Y \sim N(0,9)$ ， X_1, X_2, \dots, X_9 与 Y_1, Y_2, \dots, Y_{16} 分别是取自 X 与 Y 的简单随机样本，求统计量 $Z = \frac{X_1 + X_2 + \dots + X_9}{\sqrt{Y_1^2 + Y_2^2 + \dots + Y_{16}^2}}$ 所服从的分布。

解： $X_1 + X_2 + \dots + X_9 \sim N(0, 9 \times 16)$

$$\frac{1}{3 \times 4} (X_1 + X_2 + \dots + X_9) \sim N(0, 1)$$

$$\frac{1}{3} Y_i \sim N(0, 1), i = 1, 2, \dots, 16$$

$$\sum_{i=1}^{16} \left(\frac{1}{3} Y_i \right)^2 \sim \chi^2(16)$$

$$\text{从而 } \frac{X_1 + X_2 + \dots + X_9}{\sqrt{Y_1^2 + Y_2^2 + \dots + Y_{16}^2}} = \frac{\frac{1}{3 \times 4} (X_1 + X_2 + \dots + X_9)}{\sqrt{\frac{\sum_{i=1}^{16} \left(\frac{1}{3} Y_i \right)^2}{16}}} \sim t(16)$$



谢谢!

