




# 数理统计




数理统计是具有广泛应用的一个数学分支，它以概率论为理论基础，根据试验或观察得到的数据，来研究随机现象，对研究对象的客观规律性作出种种合理的估计和判断。




数理统计的内容包括：

如何收集、整理数据资料；


如何对所得的数据资料进行分析和研究，从而对所研究的对象的性质、特点作出推断。



假定某市成年男性的身高服从正态分布，  
希望得到平均身高 $\mu$ ：



在数理统计中，不是对所研究的对象全体（称为**总体**）进行观察，而是抽取其中的部分（称为**样本**）进行观察获得数据（**抽样**），并通过这些数据对总体进行推断。



假定某市成年男性的身高服从正态分布，  
希望得到平均身高 $\mu$ ：

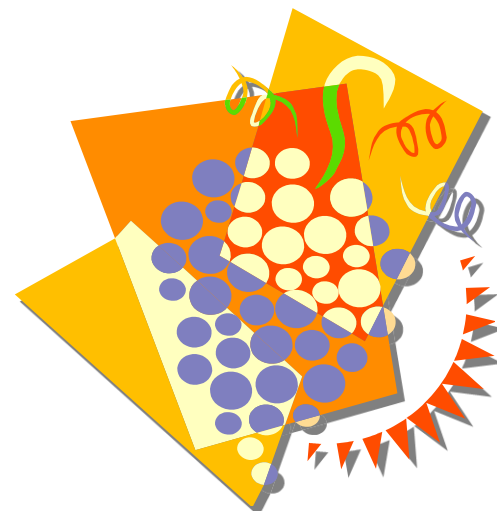
- $\mu$ 的大小如何；
- $\mu$ 大概落在什么范围内；
- 能否认为某一说法成立（如  $\mu \leq 1.68$ ）。

# 第六章 样本及抽样分布

第一节 随机样本

第二节 直方图和箱线图

第三节 抽样分布







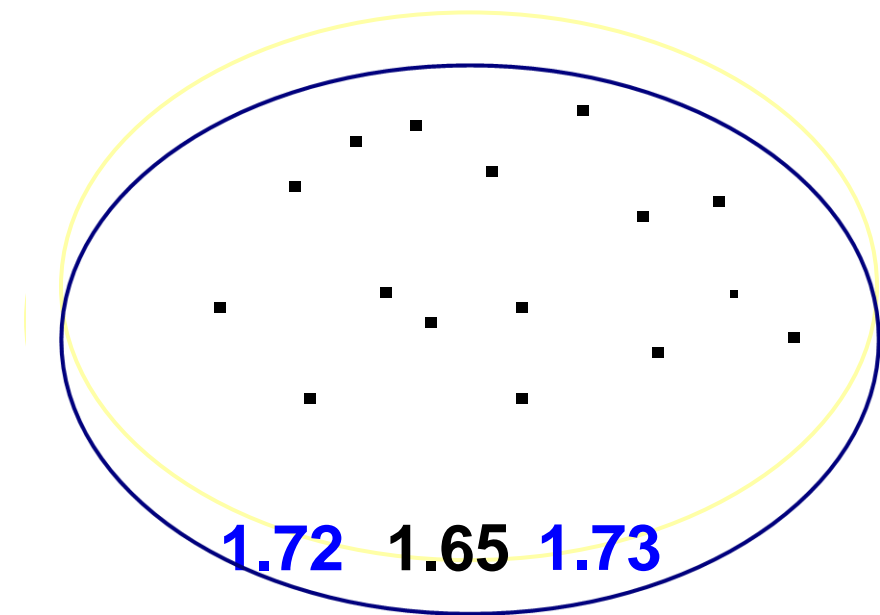
**总体：研究对象的全体**

**个体：每个对象**



**总体：研究对象的全体**

**个体：每个对象**



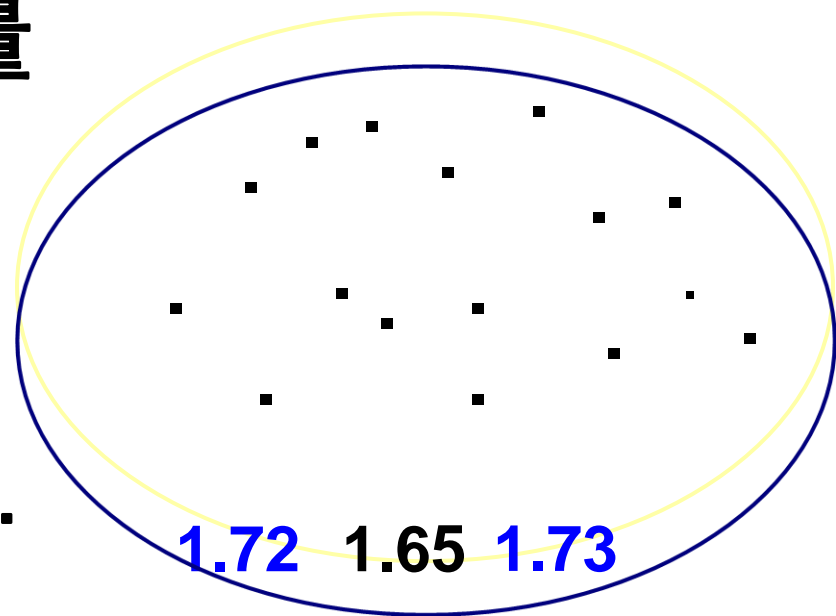
**总体：研究对象的全体**

**个体：每个对象**



**总体：研究对象的某项数量指标的**  
**全部可能的观察值**

**个体：**  
**每一个可能观察值为个体。**

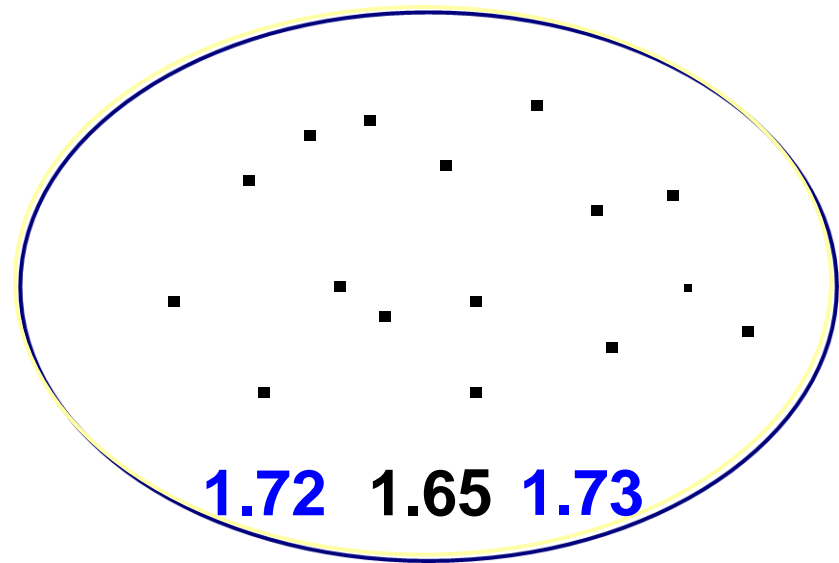




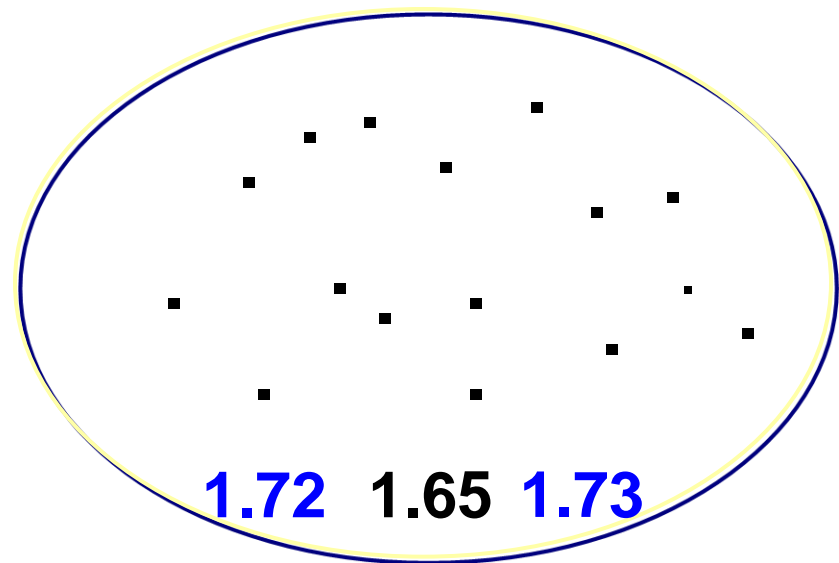
总体中所包含的个体的个数称为总体的**容量**.


容量为有限的称为**有限总体**;  
容量为无限的称为**无限总体**.

◆从总体中任取一个个体，以 $X$ 表示结果.



- ◆从总体中任取一个个体，以 $X$ 表示结果.
- ◆随机变量 $X$ 的所有可能取值就是总体中的数值；
- ◆ $X$ 的取值规律就是总体中数值的规律.





例 从0,1,2中随机抽取一个数, 用 $X$ 表示抽取结果,可以得到 $X$ 的分布律

例 从0,1,2中随机抽取一个数, 用 $X$ 表示抽取结果,可以得到 $X$ 的分布律

$X$	0	1	2
$P$	$1/3$	$1/3$	$1/3$



例 从0,1,2中随机抽取一个数, 用 $X$ 表示抽取结果,可以得到 $X$ 的分布律

$X$	0	1	2
$P$	$1/3$	$1/3$	$1/3$

- ◆  $X$ 的所有可能取值就是总体中的数值;
- ◆  $X$ 的取值规律就是总体中数值的规律.

例 从0,1,2,2中随机抽取一个数，用 $X$ 表示抽取结果,可以得到 $X$ 的分布律

例 从0,1,2,2中随机抽取一个数，用 $X$ 表示抽取结果,可以得到 $X$ 的分布律

$X$	0	1	2
$P$	1/4	1/4	1/2

例 从0,1,2,2中随机抽取一个数，用 $X$ 表示抽取结果,可以得到 $X$ 的分布律

$X$	0	1	2
$P$	1/4	1/4	1/2

- ◆  $X$ 的所有可能取值就是总体中的数值;
- ◆  $X$ 的取值规律就是总体中数值的规律.



例 考察某厂的产品质量,

总体 = {该厂生产的全部合格品与不合格品}

**例** 考察某厂的产品质量，以**0**记合格品，以**1**记不合格品，则

总体 = {该厂生产的全部合格品与不合格品}

**例** 考察某厂的产品质量，以**0**记合格品，以**1**记不合格品，则


总体 = {该厂生产的全部合格品与不合格品}  
= {由**0**或**1**组成的一堆数}

**例** 考察某厂的产品质量，以**0**记合格品，以**1**记不合格品，则

总体 = {该厂生产的全部合格品与不合格品}  
= {由**0**或**1**组成的一堆数}

若以  $p$  表示这堆数中**1**的比例（不合格品率），





若从该批产品中随机抽取一件，用  $X$  表示  
这一件产品的不合格数，

若从该批产品中随机抽取一件，用  $X$  表示这一件产品的不合格数，不难看出  $X$  服从一个二点分布  $b(1, p)$  .

$X$	0	1
$P$	$1 - p$	$p$

若从该批产品中随机抽取一件，用  $X$  表示这一件产品的不合格数，不难看出  $X$  服从一个二点分布  $b(1, p)$  .

$X$	0	1
$P$	$1 - p$	$p$

- ◆  $X$ 的所有可能取值就是总体中的数值;
- ◆  $X$ 的取值规律就是总体中数值的规律.



**从总体中抽取的部分个体称为一个样本.**

**从总体中抽取的部分个体称为一个样本.**



为调查大学生的阅读情况，某同学  
在图书馆抽取了部分同学进行调查。

从总体中抽取的部分个体称为一个样本.

◆从总体 $X$ 中随机抽取一个个体,  
以 $X_1$ 表示其结果,  $X_1$ 和 $X$ 有相同的分布.

从总体中抽取的部分个体称为一个样本.

- ◆从总体 $X$ 中随机抽取一个个体,  
以 $X_1$ 表示其结果,  $X_1$ 和 $X$ 有相同的分布.
- ◆放回, 从总体 $X$ 中再随机抽取一个个体,  
以 $X_2$ 表示其结果,  $X_2$ 和 $X$ 有相同的分布.  
.....
- ◆放回, 从总体 $X$ 中再随机抽取一个个体,  
以 $X_n$ 表示其结果,  $X_n$ 和 $X$ 有相同的分布.

从总体中抽取的部分个体称为一个样本.

◆从总体 $X$ 中随机抽取一个个体,  
以 $X_1$ 表示其结果,  $X_1$ 和 $X$ 有相同的分布.

◆放回, 从总体 $X$ 中再随机抽取一个个体,  
以 $X_2$ 表示其结果,  $X_2$ 和 $X$ 有相同的分布.  
.....

◆放回, 从总体 $X$ 中再随机抽取一个个体,  
以 $X_n$ 表示其结果,  $X_n$ 和 $X$ 有相同的分布.

◆ $X_1 \dots X_n$ 为来自总体 $X$ 的简单随机样本.



- 对于有限总体，采用放回抽样就能得到简单随机样本，但放回抽样使用起来不方便，当个体的总数 $N$ 比要得到的样本容量 $n$ 大很多时，在实际中可将不放回抽样近似地当做放回抽样来处理.

- 对于有限总体，采用放回抽样就能得到简单随机样本，但放回抽样使用起来不方便，当个体的总数 $N$ 比要得到的样本容量 $n$ 大很多时，在实际中可将不放回抽样近似地当做放回抽样来处理.
- 至于无限总体，因抽取一个个体不影响它的分布，所以总是用不放回抽样.

## 样本的两重性

- 抽取前无法预知它们的数值，因此，样本是随机变量，用大写字母  $X_1, X_2, \dots, X_n$  表示；

## 样本的两重性

- 抽取前无法预知它们的数值，因此，样本是随机变量，用大写字母  $X_1, X_2, \dots, X_n$  表示；
- 抽取后经观测就有确定的观测值，因此，样本又是一组数值。此时用小写字母  $x_1, x_2, \dots, x_n$  。



P133 综合上述，给出定义

若 $X_1, X_2, \dots, X_n$ 为 $F$ 的一个样本,  
则 $X_1, X_2, \dots, X_n$ 相互独立, 且它们的分布函数  
都是 $F$ , 所以 $(X_1, X_2, \dots, X_n)$ 的分布函数为

$$F^*(x_1, x_2, \dots, x_n) = \prod_{i=1}^n F(x_i).$$

若 $X_1, X_2, \dots, X_n$ 为 $F$ 的一个样本,  
则 $X_1, X_2, \dots, X_n$ 相互独立, 且它们的分布函数  
都是 $F$ , 所以 $(X_1, X_2, \dots, X_n)$ 的分布函数为

$$F^*(x_1, x_2, \dots, x_n) = \prod_{i=1}^n F(x_i).$$

若 $X$ 具有概率密度 $f$ , 则 $(X_1, X_2, \dots, X_n)$ 的  
概率密度为

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

## 第三节 抽样分布

由样本去推断总体情况，需要对样本进行“加工”，这就要构造一些样本的函数，它把样本中所含的（某一方面）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.

**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.

$$X_1, X_2, \dots, X_n \xrightarrow{\text{观察值}} x_1, x_2, \dots, x_n$$

$$g(X_1, X_2, \dots, X_n) \xrightarrow{\text{观察值}} g(x_1, x_2, \dots, x_n)$$

## 思考

设  $X_1, \dots, X_n$  为来自总体  $X \sim N(\mu, \sigma^2)$  的一个样本, 其中  $\mu$  未知,  $\sigma^2$  已知, 问下列哪些是统计量?

$$\frac{X_1 + X_n}{2};$$

$$\frac{X_1 + \dots + X_n}{n} - \mu ;$$

$$\frac{(X_1 + X_n)^2}{\sigma^2};$$

设  $X_1, \dots, X_n$  为来自总体  $X$  的一个样本,

数理统计中最常用的统计量及其观察值有:

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

观察值记为 
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (2)$$

设  $X_1, \dots, X_n$  为来自总体  $X$  的一个样本,

数理统计中最常用的统计量及其观察值有:

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

观察值记为 
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (2)$$

2. 样本方差 
$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad (3)$$

观察值记为 
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (4)$$

设  $X_1, \dots, X_n$  为来自总体  $X$  的一个样本,

数理统计中最常用的统计量及其观察值有:

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

观察值记为 
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (2)$$

2. 样本方差 
$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left( \sum_{i=1}^n X_i^2 - n\bar{X}^2 \right) \quad (3)$$

观察值记为 
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) \quad (4)$$

### 3. 样本标准差

$$S = \sqrt{S^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2} \quad (5)$$

它的观察值记为

$$s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (6)$$

### 3. 样本标准差

$$S = \sqrt{S^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2} \quad (5)$$

它的观察值记为

$$s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (6)$$

### 4. 样本 $k$ 阶原点矩

$$A_k = \frac{1}{n} \sum_{i=1}^n X_i^k, \quad k = 1, 2, \dots \quad (7)$$

它的观察值记为

$$a_k = \frac{1}{n} \sum_{i=1}^n x_i^k, \quad k = 1, 2, \dots \quad (8)$$



## 5. 样本 $k$ 阶中心矩

$$B_k = \frac{1}{n} \sum_{i=1}^n \left( X_i - \overline{X} \right)^k, \quad k = 1, 2, \dots \quad (9)$$

它的观察值记为

$$b_k = \frac{1}{n} \sum_{i=1}^n \left( x_i - \bar{x} \right)^k, \quad k = 1, 2, \dots \quad (10)$$

设总体 $X$ 的均值为 $\mu$ ,方差为 $\sigma^2$ ,

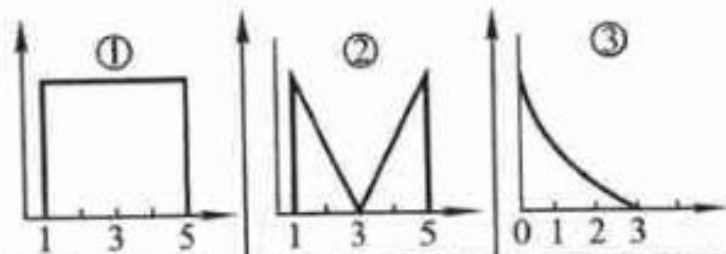
$X_1, X_2, \dots, X_n$ 是 $X$ 的一个样本.

$$E(\bar{X}) = \mu,$$

$$D(\bar{X}) = \frac{\sigma^2}{n}$$

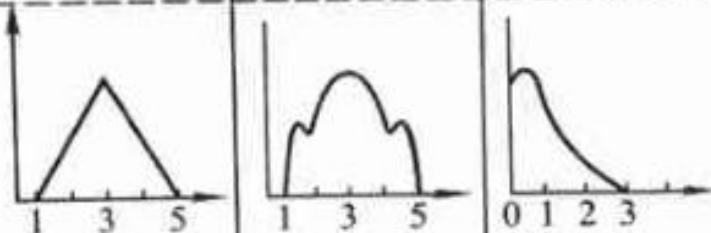
$$E(S^2) = \sigma^2$$

总体分布



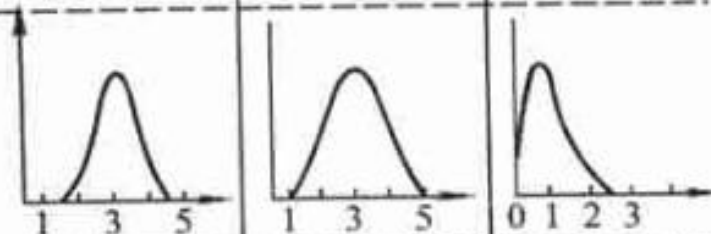
$\bar{X}$  的分布

( $n=2$ )



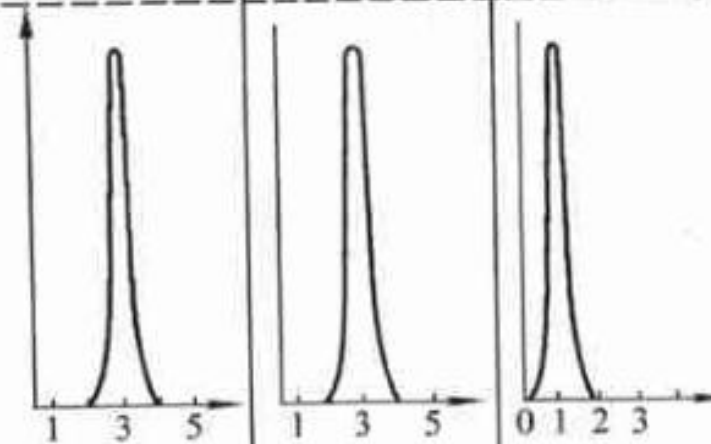
$\bar{X}$  的分布

( $n=5$ )



$\bar{X}$  的分布

( $n=30$ )



# 来自正态总体的几个常用统计量的分布

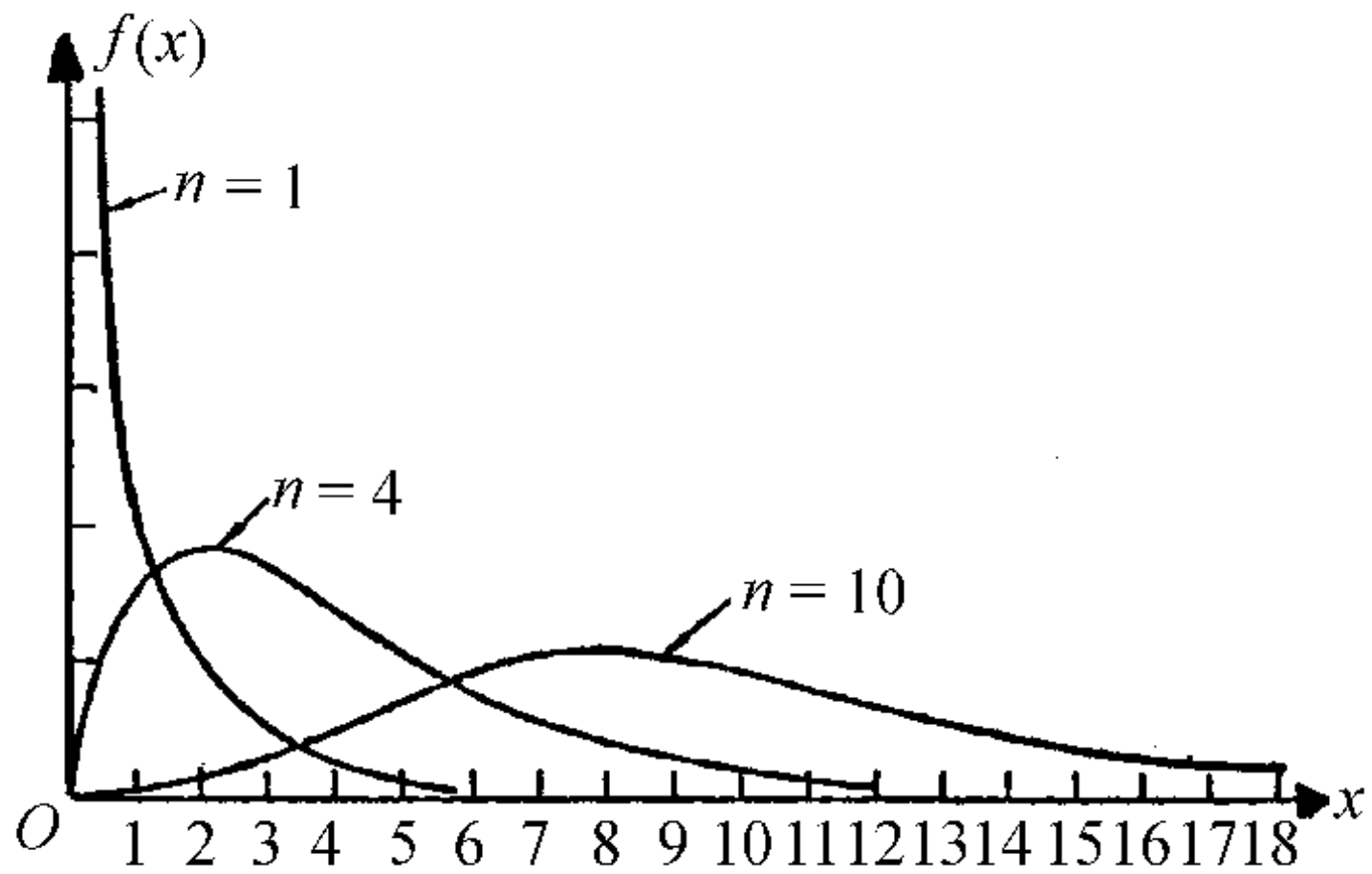
# 来自正态总体的几个常用统计量的分布


## (一) $\chi^2$ 分布

$X_1, X_2, \dots, X_n$  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  $\chi^2$  分布. 记为  $\chi^2 \sim \chi^2(n)$ .





若总体  $X \sim N(0,1)$ , 从此总体中取一个容量为 3 的样本  $X_1, X_2, X_3$ , 设

$$Y = (X_1 + X_2 + X_3)^2$$

试决定常数  $C$ , 使随机变量  $CY$  服从  $\chi^2$  分布.

## $\chi^2$ 分布的可加性

设 $\chi_1^2 \sim \chi^2(n_1)$ ,  $\chi_2^2 \sim \chi^2(n_2)$ , 并且 $\chi_1^2$ ,  $\chi_2^2$ 独立,  
则有

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



3. 若  $X \sim \chi^2(n)$ , 则

$$E(X)=n, D(X)=2n.$$

事实上, 由  $X_i \sim N(0,1)$ , 故  $E(X_i^2) = D(X_i) = 1$

$$D(X_i^2) = E(X_i^4) - [E(X_i^2)]^2 = 3 - 1 = 2$$

$$E(\chi^2) = \sum_{i=1}^n E(X_i^2) = n, D(\chi^2) = \sum_{i=1}^n D(X_i^2) = 2n.$$

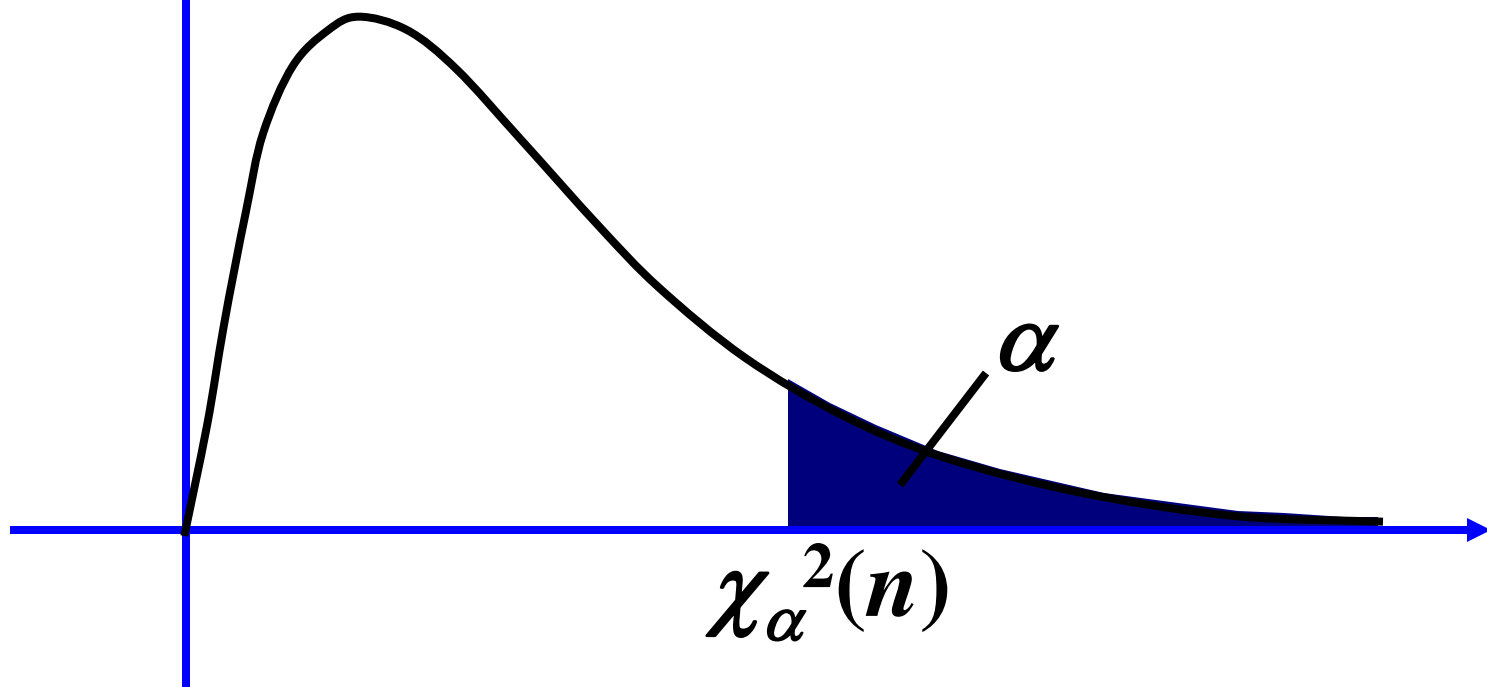
# $\chi^2$ 分布的分位数

## $\chi^2$ 分布的分位数

对于给定的正数 $\alpha$ ,  $0 < \alpha < 1$ , 称满足

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

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



## (二) $t$ 分布 (学生氏分布)

## (二) $t$ 分布

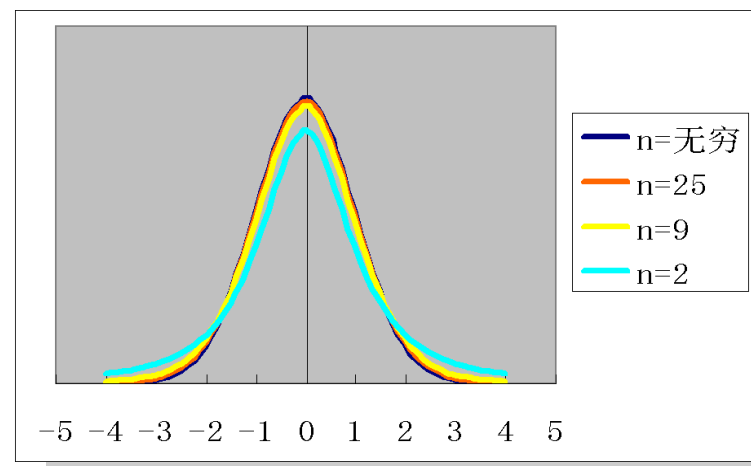
设  $X \sim N(0,1)$ ,  $Y \sim \chi^2(n)$ , 且  $X, Y$  相互独立, 称

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

服从自由度为  $n$  的  $t$  分布. 记为  $t \sim t(n)$ .

## $t$ 分布的性质：

$t$ 分布的密度函数关于 $t = 0$ 对称.当 $n$ 充分大时,  
其图形近似于标准正态分布概率密度的图形,



$t$  分布的概率密度曲线

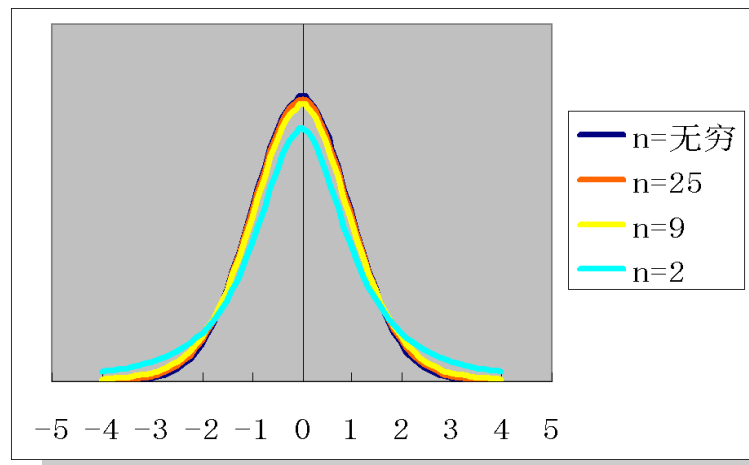
## $t$ 分布的性质:

$t$ 分布的密度函数关于 $t = 0$ 对称.当 $n$ 充分大时,  
其图形近似于标准正态分布概率密度的图形,

再由 $\Gamma$ 函数的性质有

$$\lim_{n \rightarrow \infty} h(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}.$$

即当 $n$ 足够大时,  $t \overset{\text{近似}}{\sim} N(0,1)$ .

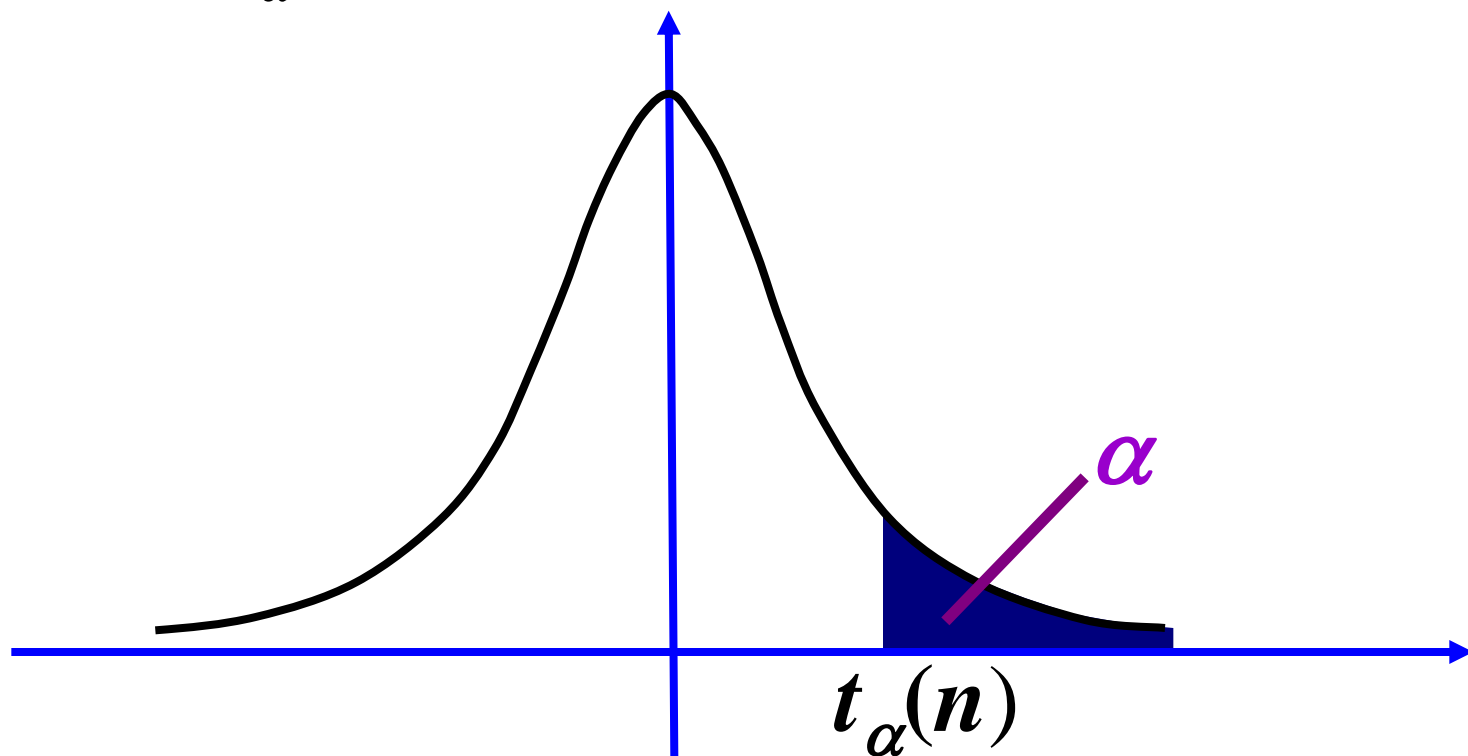


$t$ 分布的概率密度曲线

**$t$ 分布的分位数** 对于给定的 $\alpha$ ,  $0 < \alpha < 1$ ,

称满足条件  $P\{t > t_{\alpha}(n)\} = \int_{t_{\alpha}(n)}^{\infty} h(t)dt = \alpha$

的点 $t_{\alpha}(n)$ 为 $t(n)$ 分布的上 $\alpha$ 分位数.





### (三) $F$ 分布

设  $U \sim \chi^2(n_1)$ ,  $V \sim \chi^2(n_2)$ , 且  $U, V$  相互独立,

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  $F \sim F(n_1, n_2)$ .

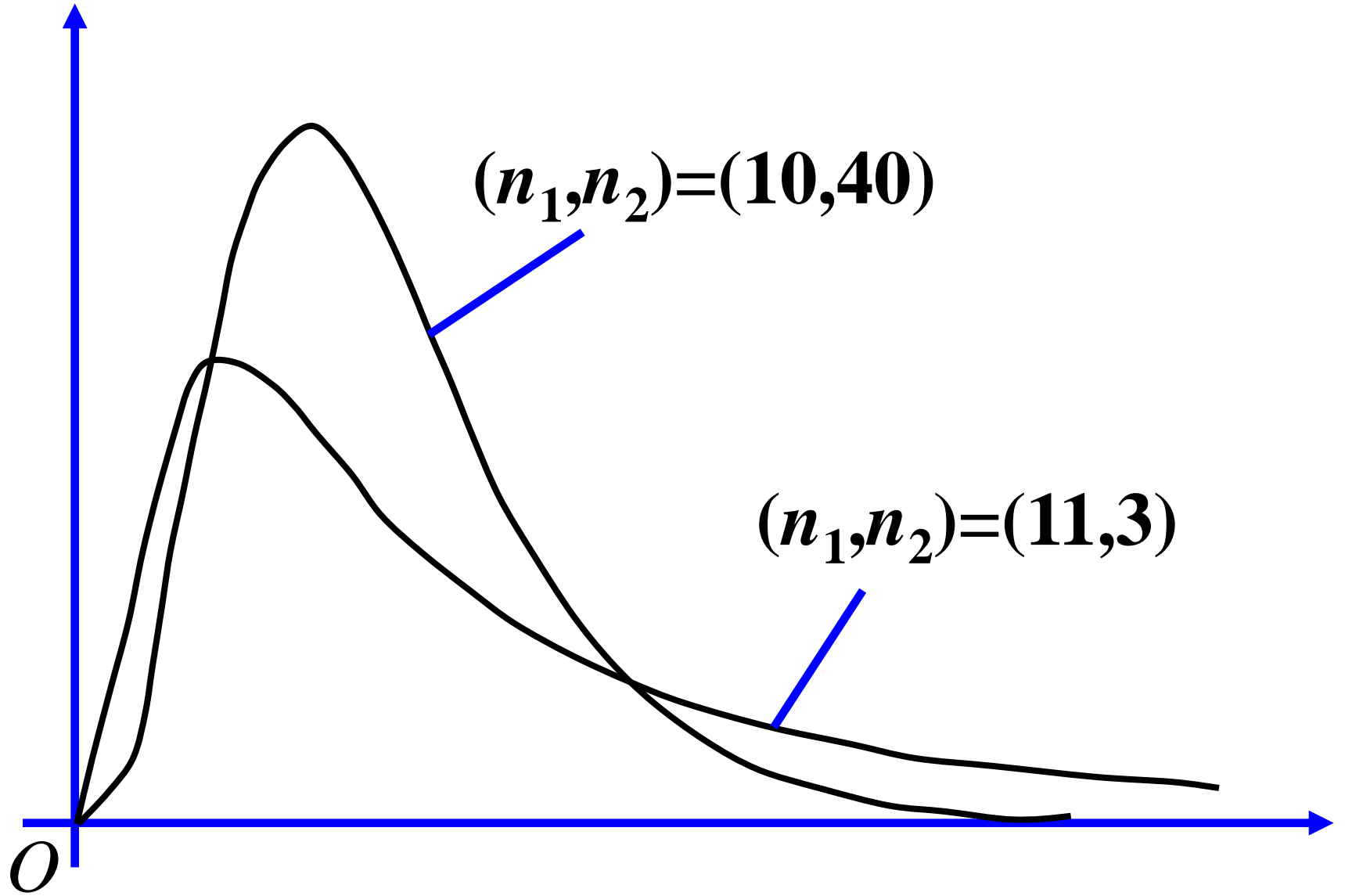
### (三) $F$ 分布

设  $U \sim \chi^2(n_1)$ ,  $V \sim \chi^2(n_2)$ , 且  $U, V$  相互独立,

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  $F \sim F(n_1, n_2)$ .

若  $F \sim F(n_1, n_2)$ , 则  $1 / F \sim F(n_2, n_1)$ .

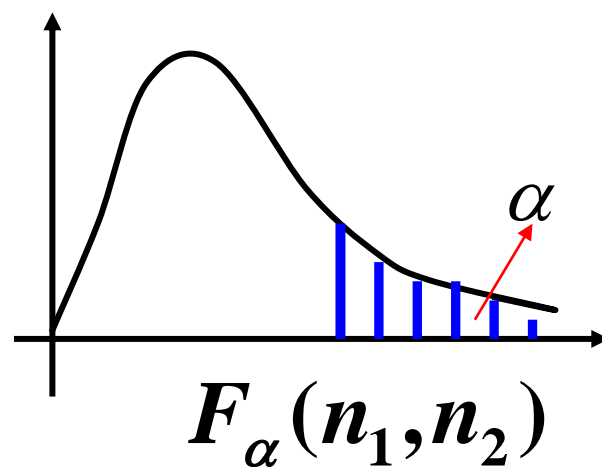


# $F$ 分布的分位数

对于给定的 $\alpha(0 < \alpha < 1)$ , 称满足条件

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

的点 $F_{\alpha}(n_1, n_2)$ 为 $F$ 分布的上 $\alpha$ 分位数。



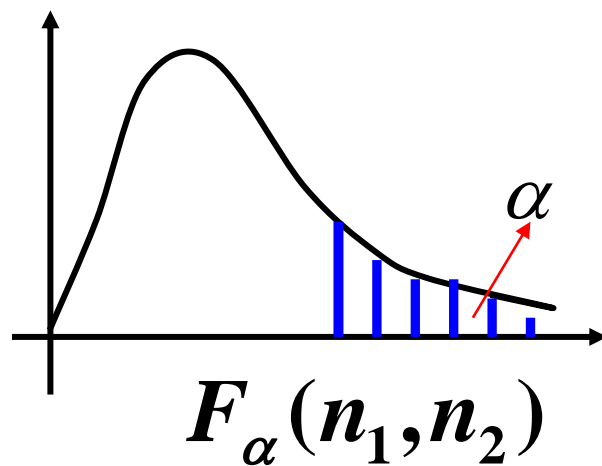
# $F$ 分布的分位数

对于给定的 $\alpha(0 < \alpha < 1)$ , 称满足条件

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

的点 $F_{\alpha}(n_1, n_2)$ 为 $F$ 分布的上 $\alpha$ 分位数。

$$F_{1-\alpha}(n_1, n_2) = 1 / F_{\alpha}(n_2, n_1)$$



$$Z_{0.025} = \underline{\hspace{2cm}}; \quad Z_{0.95} = \underline{\hspace{2cm}};$$

$$\chi^2_{0.05}(3) = \underline{\hspace{2cm}}; \quad \chi^2_{0.95}(3) = \underline{\hspace{2cm}};$$

$$t_{0.05}(5) = \underline{\hspace{2cm}}; \quad t_{0.975}(5) = \underline{\hspace{2cm}};$$

$$F_{0.025}(10,10) = \underline{\hspace{2cm}}; \quad F_{0.95}(8,10) = \underline{\hspace{2cm}};$$

# 来自正态总体的几个常用统计量的分布

(一)  $\chi^2$ 分布

(二)  $t$  分布

(三)  $F$  分布

## (四) 正态总体的样本均值与样本方差的分布

设产品的某质量指标  $X \sim N(\mu, \sigma^2)$

考虑总体均值  $\mu$ 、总体方差  $\sigma^2$  的统计推断问题



设总体 $X$ 的均值为 $\mu$ ,方差为 $\sigma^2$ ,

$X_1, X_2, \dots, X_n$ 是 $X$ 的一个样本.

$$E(\bar{X}) = \mu,$$

$$D(\bar{X}) = \frac{\sigma^2}{n}$$

$$E(S^2) = \sigma^2$$

**定理一** 设 $X_1, X_2, \dots, X_n$ 是来自总体 $N(\mu, \sigma^2)$ 的样本,  $\bar{X}$ 是样本均值, 则有

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

**定理二** 设 $X_1, X_2, \dots, X_n$ 是来自总体 $N(\mu, \sigma^2)$ 的样本,  $\bar{X}$ 是样本均值, 则有

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

**定理三** 设 $X_1, X_2, \dots, X_n$ 是来自总体 $N(\mu, \sigma^2)$ 的样本,  $\bar{X}$ 和 $S^2$ 是样本均值和样本方差, 则有

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

$$(2) \bar{X} \text{ 与 } S^2 \text{ 独立}$$

## 定理四

设 $X_1, X_2, \dots, X_n$ 是取自正态总体  $N(\mu, \sigma^2)$  的样本,  $\bar{X}$ 和 $S^2$ 分别为样本均值和样本方差, 则有

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

设甲厂产品的某质量指标  $X \sim N(\mu_1, \sigma_1^2)$

设乙厂产品的某质量指标  $Y \sim N(\mu_2, \sigma_2^2)$

为了比较产品质量指标，需要考虑

$\mu_1 - \mu_2, \sigma_1^2 / \sigma_2^2$  的统计推断问题。

设产品的某质量指标  $X \sim N(\mu_1, \sigma_1^2)$ ,

由于原材料的改变、或设备条件发生变化、或技术革新等因素的影响, 使得产品质量指标可能发生变化, 此时产品的质量指标为  $Y \sim N(\mu_2, \sigma_2^2)$ .

为了了解产品质量指标有多大的变化, 需要考虑  $\mu_1 - \mu_2, \sigma_1^2 / \sigma_2^2$  的统计推断问题.

## 定理五（两正态总体的样本方差）

设  $X \sim N(\mu_1, \sigma_1^2)$ ,  $Y \sim N(\mu_2, \sigma_2^2)$ , 且  $X$  与  $Y$  独立,  $X_1, X_2, \dots, X_{n_1}$  是取自  $X$  的样本,  $Y_1, Y_2, \dots, Y_{n_2}$  是取自  $Y$  的样本,  $\bar{X}$  和  $\bar{Y}$  分别是这两个样本的样本均值,  $S_1^2$  和  $S_2^2$  分别是这两个样本的样本方差, 则有

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

## 定理五（两正态总体样本均值）

设  $X \sim N(\mu_1, \sigma_1^2)$ ,  $Y \sim N(\mu_2, \sigma_2^2)$ , 且  $X$  与  $Y$  独立,  
 $X_1, X_2, \dots, X_{n_1}$  是取自  $X$  的样本,  $Y_1, Y_2, \dots, Y_{n_2}$  是  
取自  $Y$  的样本,  $\bar{X}$  和  $\bar{Y}$  分别是这两个样本的样本  
均值,  $S_1^2$  和  $S_2^2$  分别是这两个样本的样本方差,  
则有

$$\frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t \quad (n_1 + n_2 - 2)$$



## 第三节 抽样分布

- 统计量
- 三大抽样分布
- 正态总体的样本均值与样本方差的分布

作业：

◆P147 1题 4题