

## §2.7 充分统计量

### 充分统计量

统计量(statistics)是对样本的加工. “好”的统计量应该能够将样本中关于总体分布的未知信息尽可能地集中起来.

若要研究某参数分布族中的某个未知参数,为此抽取了一组样本,样本中所包含的信息可分成两部分.

其一是关于该未知参数的信息;

其二关于样本结构的信息等其它信息.

### Example

某厂要了解其半成品的不合格率 $p$ , 检验员检查了10件产品, 检查的结果是, 除前两件是不合格品(记为 $X_1 = 1, X_2 = 1$ )外, 其它都是合格品(记为 $X_i = 0, i = 3, 4, \dots, 10$ ).

这时样本 $(X_1, X_2, \dots, X_{10})$ 提供了两种信息:

- (1) 10 次检验中, 不合格品出现了几次,
- (2) 不合格品出现在哪几次试验上.

第二种信息(试验编号的信息)对了解不合格率 $p$ 是没有什么帮助的. 例如, 设在另一次试验中, 试验结果为 $(0, 0, \dots, 0, 1, 1)$ . 这两次试验结果所含的关于 $p$ 的信息应该是一样的.

在上例中, 当厂长问及检查结果时, 看看检验员的如下两种回答:

1. 10件中有两件不合格——所用统计量为 $T_1 = \sum_{i=1}^{10} X_i$ , 其值

为 $t_1 = \sum_{i=1}^{10} x_i = 2$ ;

2. 前两件不合格( $X_1 = 1, X_2 = 1$ )——所用统计量为 $T_2 = X_1 + X_2$ , 其值

为 $t_2 = x_1 + x_2 = 2$ .

显然, 第二种回答是不能令人满意的, 因为统计量 $T_2$ 不包含样本中有关 $p$ 的全部信息. 而第一种回答综合了样本中有关 $p$ 的全部信息.

一个“好”的统计量, 应该能够将样本中所包含的关于未知参数的信息全部集中起来. 这样的统计量就是该未知参数的充分统计量.

如何将这样一个直观的想法用严格的数学形式来表示呢?

$$\begin{aligned} & \text{样本 } \tilde{X} \text{ 中的信息} \\ &= \text{统计量 } T(\tilde{X}) \text{ 中所含样本 } \tilde{X} \text{ 的信息} \\ &+ \text{在知道 } T(\tilde{X}) \text{ 后样本 } \tilde{X} \text{ 中还含有的剩余信息.} \end{aligned}$$

设有参数分布族  $\mathcal{F} = \{F_\theta : \theta \in \Theta\}$ ,  $\tilde{X} = (X_1, X_2, \dots, X_n)$  是从某总体  $F_\theta \in \mathcal{F}$  中抽取的样本,  $T = T(X_1, X_2, \dots, X_n)$  是一个统计量. 样本  $\tilde{X}$  有一个样本分布(即样本的联合分布)

$$F_\theta(\tilde{x}) = \prod_{i=1}^n F_\theta(x_i),$$

统计量  $T$  也有一抽样分布  $F_\theta^T(t)$ .

充分性要求统计量 $T(\tilde{X})$ 包含样本 $\tilde{X}$ 中有关 $\theta$ 的全部信息. 也就是说样本分布 $F_{\theta}(\tilde{x})$ 所含的关于 $\theta$ 的全部信息都包含在抽样分布 $F_{\theta}^T(t)$ 中.

即除了 $F_{\theta}^T(t)$ 所含的关于 $\theta$ 的信息外,  $F_{\theta}(\tilde{x})$ 不再含的关于 $\theta$ 的信息.

换言之,即在统计量 $T$ 的取值给定后, 譬如 $T = t$ 后, 样本的条件分布 $F_{\theta}(\tilde{x}|T = t)$ 已不再依赖于参数 $\theta$ .

**定义** 设有一个分布族  $\mathcal{F} = \{F\}$ ,  $(X_1, X_2, \dots, X_n)$  是从某总体  $F \in \mathcal{F}$  中抽取的一个样本.  $T = T(X_1, X_2, \dots, X_n)$  为一个 (一维或多维的) 统计量. 如果当给定  $T = t$  下, 样本  $(X_1, X_2, \dots, X_n)$  的条件分布与总体分布  $F$  无关, 则称  $T$  为此分布族的 **充分统计量** (*sufficient statistics*).

如果  $\mathcal{F}$  是参数型分布族  $\{F_\theta : \theta \in \Theta\}$ , 当给定  $T = t$  下, 样本  $(X_1, X_2, \dots, X_n)$  的条件分布与参数  $\theta$  无关, 则亦称  $T$  为是参数  $\theta$  的充分统计量.

**说明:** 在实际应用时, 对于离散型或连续型总体而言, 条件分布常用条件概率分布列或条件概率密度函数来代替.

### Theorem

充分统计量的一一变换仍是充分统计量.

设 $s = \psi(t)$ 是一一变换,  $S = \psi(T)$ , 则事件 $\{S = s\}$  与 $\{T = \psi^{-1}(s)\}$ 等价. 一般地,  $\{S \in A\}$  与 $\{T \in \psi^{-1}(A)\}$ 等价. 所以

$$\begin{aligned} & P(X_1 < x_1, X_2 < x_2 \cdots, X_n < x_n | S = s) \\ &= P(X_1 < x_1, X_2 < x_2 \cdots, X_n < x_n | T = \psi^{-1}(s)). \end{aligned}$$



### Theorem

设 $\tilde{Y}$ 是样本 $\tilde{X}$ 的一一变换, 则统计量 $T = T(\tilde{X})$ 对 $\tilde{X}$ 的充分性等价于 $T$ 对 $\tilde{Y}$ 的充分性.

设 $\tilde{Y} = \tilde{\psi}(\tilde{X})$ 是一一变换, 则事件 $\{\tilde{X} \in A\}$  与 $\{\tilde{Y} \in \tilde{\psi}(A)\}$ 等价. 所以

$$P(\tilde{X} \in A | T = t) = P(\tilde{Y} \in \tilde{\psi}(A) | T = t).$$

### Example

**例** 设 $X_1, X_2, \dots, X_n$ 是来自二点分布族 $\{B(1, p), 0 < p < 1\}$ 中的某总体的一个样本, 其中 $0 < p < 1, n > 2$ . 考察统计量

$$T_1 = \sum_{i=1}^n X_i, \quad T_2 = X_1 + X_2.$$

这时, 样本  $\tilde{X} = (X_1, X_2, \cdots, X_n)$  的联合分布列是

$$P(X_1 = x_1, X_2 = x_2, \cdots, X_n = x_n) = p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i},$$

其中  $x_i = 0$  或  $1, i = 1, 2, \cdots, n$ . 而  $T_1$  的分布列为

$$P(T_1 = t) = \binom{n}{t} p^t (1-p)^{n-t}, \quad t = 0, 1, \cdots, n.$$

在给定  $T_1 = t (t = 0, 1, \cdots, n)$  的条件下, 若  $x_1 + \cdots + x_n \neq t$ , 则

$$\begin{aligned} & P(X_1 = x_1, X_2 = x_2, \cdots, X_n = x_n | T_1 = t) \\ &= \frac{P(X_1 = x_1, X_2 = x_2, \cdots, X_n = x_n, T_1 = t)}{P(T_1 = t)} = 0; \end{aligned}$$

若  $x_1 + \cdots + x_n = t$ , 则

$$\begin{aligned} & P(X_1 = x_1, X_2 = x_2, \cdots, X_n = x_n | T_1 = t) \\ &= \frac{P(X_1 = x_1, X_2 = x_2, \cdots, X_n = t - \sum_{i=1}^{n-1} x_i)}{P(T_1 = t)} \\ &= \frac{p^t (1-p)^{n-t}}{\binom{n}{t} p^t (1-p)^{n-t}} = \binom{n}{t}^{-1}, \quad x_i = 0 \text{ 或 } 1, \quad i = 1, 2, \cdots, n. \end{aligned}$$

此条件分布列与参数  $p$  无关, 故  $T_1$  为  $p$  的充分统计量.

对于 $T_2$ , 在给定 $T_2 = t(t = 0, 1, 2)$ 下, 当 $x_1 + x_2 \neq t$ 时, 条件概率是0;  
当 $x_1 + x_2 = t$ 时,

$$\begin{aligned} & P(X_1 = x_1, X_2 = x_2, \cdots, X_n = x_n | T_2 = t) \\ &= \frac{P(X_1 = x_1, X_2 = t - x_1, \cdots, X_n = x_n)}{P(T_2 = t)} \\ &= \frac{p^{t + \sum_{i=3}^n x_i} (1 - p)^{n - t - \sum_{i=3}^n x_i}}{\binom{2}{t} p^t (1 - p)^{2 - t}} \\ &= \binom{2}{t}^{-1} p^{\sum_{i=3}^n x_i} (1 - p)^{n - 2 - \sum_{i=3}^n x_i}, \quad x_i = 0 \text{ 或 } 1, \quad i = 1, 2, \cdots, n. \end{aligned}$$

这个条件分布仍与参数 $p$ 有关,  $T_2$ 不是 $p$ 的充分统计量.

对充分统计量的理解:

设想有两个试验员, 试验员A可观察到样本 $\tilde{X}$ , 试验员B只能观察到统计量 $T(\tilde{X})$ . 如果 $T$ 是充分统计量, 他们得到关于 $\theta$ 的信息应该是等价的.

### Example

例 设 $X_1, X_2, \dots, X_n$ 是来自几何分布

$$P(X = x) = \theta(1 - \theta)^x, \quad x = 0, 1, 2, \dots$$

的一个样本, 其中 $0 < \theta < 1$ , 样本容量 $n \geq 2$ . 则 $T = \sum_{i=1}^n X_i$ 是参数 $\theta$ 的充分统计量.

证明:我们先求 $T$ 的分布.设想有一系列独立重复试验,每次试验的成功率为 $\theta$ ,记 $T_1$ 为第一次成功时前面试验的总次数(即失败的次数), $T_2$ 为第一次成功后到第二次成功之间的失败次数, ...,  $T_n$ 为第 $n-1$ 次成功后到第 $n$ 次成功之间的失败次数,这时 $T_1, T_2, \dots, T_n$ 独立同分布且与 $X$ 的分布相同. 因而 $T^* = \sum_{i=1}^n T_i$ 与 $T$ 同分布. 因为 $T^*$ 是第 $n$ 次成功时,失败试验的总次数, $T^* = t$ 意味着共进行了 $t+n$ 次试验,其中第 $t+n$ 次成功,前 $t+n-1$ 次试验中有 $n-1$ 次成功, $t$ 次失败. 因此

$$P(T = t) = P(T^* = t) = \binom{t+n-1}{n-1} \theta^n (1-\theta)^t, \quad t = 0, 1, 2, \dots.$$



所以在 $T = t(t = 0, 1, \dots)$ 时, 样本的条件分布为: 若 $x_1 + \dots + x_n \neq t$ , 则

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T = t) = 0;$$

若 $x_1 + \dots + x_n = t$ , 则对 $x_i = 0, 1, 2, \dots, i = 1, 2, \dots, n$ ,

$$\begin{aligned} & P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T = t) \\ &= \frac{P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n, T = t)}{P(T = t)} \\ &= \frac{P(X_1 = x_1)P(X_2 = x_2) \cdots P(X_n = t - \sum_{i=1}^{n-1} x_i)}{P(T = t)} \\ &= \frac{\theta(1 - \theta)^{x_1} \cdots \theta(1 - \theta)^{x_{n-1}} \theta(1 - \theta)^{t - \sum_{i=1}^{n-1} x_i}}{\binom{t+n-1}{n-1} \theta^n (1 - \theta)^t} = \binom{t+n-1}{n-1}^{-1}. \end{aligned}$$

这个条件分布与参数 $\theta$ 无关, 所以 $T = \sum_{i=1}^n X_i$ 是参数 $\theta$ 的充分统计量.

### Example

**例** 设总体 $X$ 来自指数分布族 $\{E(\lambda) : \lambda > 0\}$ , 即其概率密度函数为

$$p(x, \lambda) = \lambda e^{-\lambda x}, \quad x > 0,$$

其中 $\lambda > 0$ 为未知参数.  $X_1, X_2, \dots, X_n$ 为来自该总体的简单随机样本. 证明: 样本均值 $\bar{X}$ 为 $\lambda$ 的充分统计量.

证: 令  $T = \sum_{j=1}^n X_j$ . 因为指数分布  $E(\lambda) = \Gamma(1, \lambda)$ . 由gamma分布的可加性知  $T \sim \Gamma(n, \lambda)$ , 密度函数为

$$p_T(t; \theta) = [(n-1)!]^{-1} \lambda^n t^{n-1} e^{-\lambda t}, \quad t > 0.$$

另一方面, 线性变换

$$\begin{cases} x_1 = x_1 \\ \vdots \\ x_{n-1} = x_{n-1} \\ t = \sum_{i=1}^n x_i \end{cases}$$

的Jacobian行列式为

$$\frac{D(x_1, \cdots, x_{n-1}, x_n)}{D(x_1, \cdots, x_{n-1}, t)} = 1.$$

所以  $X_1, \dots, X_{n-1}, T$  的联合密度函数为

$$p(x_1, \dots, x_{n-1}, t; \theta) = \lambda e^{-\lambda x_1} \dots \lambda e^{-\lambda x_{n-1}} \lambda e^{-\lambda(t - \sum_{i=1}^{n-1} x_i)} = \lambda^n e^{-\lambda t},$$
$$0 < x_i, i = 1, 2, \dots, n-1, t \geq \sum_{i=1}^{n-1} x_i.$$

因此当给定  $T = t$  时,  $X_1, X_2, \dots, X_{n-1}$  的条件密度为

$$p_{X_1, \dots, X_{n-1}|T}(x_1, \dots, x_{n-1}|t) = (n-1)! t^{-(n-1)},$$
$$x_i > 0, i = 1, 2, \dots, n-1, t \geq \sum_{i=1}^{n-1} x_i.$$

与  $\lambda$  无关. 因此  $T$  是参数  $\lambda$  的充分统计量, 从而  $\bar{X} = T/n$  与  $T$  具有一一对应关系, 因此也是参数  $\lambda$  的充分统计量.

## 因子分解定理

根据充分统计量的定义及其解释, 在对总体未知参数进行推断时, 应在可能的情况下尽量找出关于未知参数的充分统计量. 虽然可以直接根据定义来验证一个统计量是否充分(这种做法通常比较繁琐), 但是无法提供寻找充分统计量的途径.

用 $\tilde{X}$ 记样本 $(X_1, X_2, \dots, X_n)$ .

**定理** (因子分解定理) 设样本 $\tilde{X}$ 的联合pdf或pmf为 $p(\tilde{x}; \theta)$ , 其中 $\theta$ 为未知参数. 则 $T = T(\tilde{X})$ 为(关于 $\theta$ 的)充分统计量, 当且仅当

$$p(\tilde{x}; \theta) = g(T(\tilde{x}); \theta)h(\tilde{x}),$$

其中 $g(t; \theta)$ 是定义在统计量 $T(\tilde{X})$ 取值空间 $\mathcal{T}$ 上的函数,  $h(\tilde{x})$ 与 $\theta$ 无关.

证明: 此定理的严格证明需要测度论的知识, 超出了本课程的范围. 下面只给出离散场合下的证明. 这时

$$p(\tilde{x}; \theta) = \mathbf{P}_\theta(\tilde{X} = \tilde{x}).$$

对  $t \in \mathcal{T}$ , 令集合

$$A(t) = \{\tilde{x} : T(\tilde{x}) = t\}.$$



**充分性:** 设 $p(\tilde{x}; \theta)$ 有定理中的因子分解形式. 在给定 $T = t$ 下, 当 $\tilde{x} \notin A(t)$ 时,

$$P_{\theta}(\tilde{X} = \tilde{x} | T = t) = 0;$$

当 $\tilde{x} \in A(t)$ 时,

$$\begin{aligned} P_{\theta}(\tilde{X} = \tilde{x} | T = t) &= \frac{P_{\theta}(\tilde{X} = \tilde{x}, T = t)}{P_{\theta}(T = t)} \\ &= \frac{P_{\theta}(\tilde{X} = \tilde{x})}{P_{\theta}(T = t)} = \frac{p(\tilde{x}; \theta)}{\sum_{\tilde{y} \in A(t)} p(\tilde{y}; \theta)} \\ &= \frac{g(t; \theta)h(\tilde{x})}{\sum_{\tilde{y} \in A(t)} g(t; \theta)h(\tilde{y})} = \frac{h(\tilde{x})}{\sum_{\tilde{y} \in A(t)} h(\tilde{y})}. \end{aligned}$$

此条件分布与参数 $\theta$ 无关. 所以 $T(\tilde{X})$ 是充分统计量.

**必要性:** 设 $T(\tilde{X})$ 是参数 $\theta$ 的充分统计量, 由定义知,  $P_{\theta}(\tilde{X} = \tilde{x}|T = t)$ 与参数 $\theta$ 无关, 它只能是 $\tilde{x}$ 的函数, 记之为 $h(\tilde{x})$ .

对任意的 $\tilde{x}$ , 记 $t = T(\tilde{x})$ . 则

$$\begin{aligned} p(\tilde{x}; \theta) &= P_{\theta}(\tilde{X} = \tilde{x}) \\ &= P_{\theta}(\tilde{X} = \tilde{x}, T = t) \\ &= P_{\theta}(\tilde{X} = \tilde{x}|T = t)P_{\theta}(T = t) \\ &= h(\tilde{x})g(t; \theta) = h(\tilde{x})g(T(\tilde{x}); \theta). \end{aligned}$$

显然满足因子分解条件. 这样就证明了定理的必要性.

### Example

例 设  $\tilde{X} = (X_1, \dots, X_n)$  为取自在  $(0, \theta)$  区间上均匀分布总体的样本, 其中  $\theta > 0$  为未知参数. 样本的联合概率密度函数为

$$p(\tilde{x}; \theta) = \theta^{-n} I_{\{0 < x_{(1)} \leq x_{(n)} < \theta\}} = \theta^{-n} I_{\{x_{(n)} < \theta\}} \cdot I_{\{0 < x_{(1)} \leq x_{(n)}\}},$$

其中  $I_A$  为集合  $A$  的示性函数,  $x_{(1)} = \min\{x_i, i = 1, \dots, n\}$ ,

$x_{(n)} = \max\{x_i, i = 1, \dots, n\}$ . 令

$$g(T(\tilde{x}); \theta) = \theta^{-n} I_{\{x_{(n)} < \theta\}}, \quad h(\tilde{x}) = I_{\{0 < x_{(1)} \leq x_{(n)}\}},$$

其中  $T(\tilde{x}) = x_{(n)}$ . 则  $p(\tilde{x}; \theta)$  满足因子分解条件, 因而统计量  $T(\tilde{X}) = X_{(n)} = \max\{X_i, i = 1, \dots, n\}$  是充分统计量.

### Example

例 次序统计量是充分统计量.

设 $(X_{(1)}, X_{(2)}, \dots, X_{(n)})$ 为样本 $(X_1, X_2, \dots, X_n)$ 的次序统计量, 总体有密度函数(或分布列函数) $p_\theta(x)$ . 则样本的联合密度函数(或分布列函数)为

$$p(\tilde{x}; \theta) = p_\theta(x_1)p_\theta(x_2) \cdots p_\theta(x_n) = p_\theta(x_{(1)})p_\theta(x_{(2)}) \cdots p_\theta(x_{(n)}).$$

### Example

例 设样本 $X_1, \dots, X_n$ 取自正态总体 $\sim N(\mu, \sigma^2)$ , 其中 $\mu$ 和 $\sigma^2$ 均未知. 样本的联合密度函数为

$$\begin{aligned} p(\tilde{x}; \mu, \sigma^2) &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\} \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \left( \sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2 \right) \right\}. \end{aligned}$$

取 $h(\tilde{x}) = 1$ , 由因子分解定理知,  $(\sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i)$ 是 $(\mu, \sigma^2)$ 的充分统计量.

### Example

**例** 设样本 $X_1, \dots, X_n$ 取自正态总体 $\sim N(\mu, \sigma^2)$ , 其中 $\mu$ 和 $\sigma^2$ 均未知. 则由前一页知 $\tilde{T} = (T_1, T_2)$ 为 $(\mu, \sigma^2)$ 的充分统计量, 其中 $T_1 = \sum_i X_i^2$ ,  $T_2 = \sum_i X_i$ . 又 $(\bar{X}, S^2)$ 为 $\tilde{T}$ 的一一变换, 所以 $(\bar{X}, S^2)$ 为 $(\mu, \sigma^2)$ 的充分统计量.

利用因子分解定理可以对指数型分布族找到充分统计量.

假定样本的pdf或pmf所在的分布族是指数型分布族, 即

$$p(\tilde{x}; \theta) = c^*(\theta) \exp \left\{ \sum_{j=1}^k Q_j(\theta) T_j^*(\tilde{x}) \right\} h^*(\tilde{x}),$$

则  $\tilde{T} = (T_1^*, \dots, T_k^*)$  为充分统计量.

注:

- 在给出样本的联合概率密度函数或者联合分布列时, 需注意支撑;
- 充分统计量的维数不一定等于未知参数的维数;
- 充分统计量往往不唯一, 通常会根据统计推断的目的, 构造维数最小的(希望在不损失所需信息下, 压缩性强, 简化程度高).



## 极小充分统计量

### Example

例 设  $\tilde{X} = (X_1, X_2, \dots, X_n)$  是取自  $B(1, p)$  的一个样本,  $0 < p < 1$ . 样本的联合分布列(pmf)为

$$\begin{aligned} &P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) \\ &= p^{\sum_{i=1}^n x_i} (1 - p)^{n - \sum_{i=1}^n x_i} \cdot I\{x_i = 0, 1, i = 1, 2, \dots, n\}. \end{aligned}$$

取  $h(\tilde{x}) = I\{x_i = 0, 1, i = 1, 2, \dots, n\}$ , 则由因子分解定理, 知

$$T_1 = (X_1, X_2, \dots, X_n)$$

$$T_2 = (X_1 + X_2, X_3, \dots, X_n)$$

...

$$T_k = (X_1 + X_2 + \dots + X_k, X_{k+1}, \dots, X_n)$$

...

$$T_n = X_1 + X_2 + \dots + X_n$$

都是  $p$  的充分统计量.

**定义** 设 $S$ 是分布族 $\mathcal{F}$ 的充分统计量, 假如对 $\mathcal{F}$ 的任意一个充分统计量 $T$ , 存在一个函数 $f(\cdot)$ , 使得

$$S = f(T),$$

则称 $S$ 是此分布族 $\mathcal{F}$ 的极小充分统计量.

## Theorem

**定理** 设总体是指数型分布族, 则样本  $\tilde{X}$  的联合 *pdf* 或 *pmf* 可以写为

$$p(\tilde{x}; \theta) = c^*(\theta) \exp \left\{ \sum_{j=1}^k Q_j(\theta) T_j^*(\tilde{x}) \right\} h^*(\tilde{x}),$$

若  $\tilde{Q}(\theta) = (Q_1(\theta), \dots, Q_k(\theta))$  的值域有非空的内部, 则  $\tilde{T} = (T_1^*, \dots, T_k^*)$  为极小充分统计量.