

第 1 章 预备知识

1.1 背景介绍

统计推断的一个基本任务是由样本观察值去了解总体. 若根据经验或某些理论, 能在进行统计推断之前对总体作一些假设, 然后基于这些假定进行统计推断, 这种统计方法称为参数统计方法. 如果所知甚少, 在进行统计推断之前不能对总体作任何假设, 或仅能作一些非常一般性 (例如分布是连续的、是对称的等) 的假设, 这时如果仍然使用参数统计方法, 其统计推断的结果显然是不可信的, 甚至有可能是错的. 在对总体的分布不作假设或仅作非常一般性的假设的条件下发展的统计方法称为非参数统计方法.

非参数统计方法是 19 世纪 40 年代以后兴起的. 1942 年, J.Wolfowitz 首先使用非参数统计一词, 早期的非参数统计主要是扩充参数检验的内容, 以使得传统的检验过程可以应用于小样本以及不同分布类型的数据. 比如常用的非参数检验有符号秩检验、双样本 Wilcoxon 检验、多样本 Kruskal-Wallis 检验等.

近年来, 由于统计理论的进一步发展与计算机收集和处理数据能力的提高, 使得发展随数据结构不同而灵活变化的模型的统计推断方法成为可能. 非参数密度估计、非参数回归等内容也成为新的研究和应用主题. 统计研究人员利用统计渐近理论突破了参数回归和模型估计的原有理论框架, 利用各种算法改进模型的计算过程, 通过调整预测偏差和方差的比例来发展适应性更强、解释更为精练、拟合优度更适中和计算更为有效的模型.

本书的主要内容之一是介绍这种研究数据结构的非参数方法, 包括非参数密度估计、非参数回归及其相关问题, 比如分布函数的估计、密度函数的导数的估计、条件密度函数的估计以及和密度函数有关的检验, 等等. 鉴于生存数据普遍存在于很多研究领域, 本书也给出了随机右删失模型下, 几种生存时间的函数的非参数估计方法.

非参数统计问题对总体分布的假定所要求的条件很宽. 因而针对这种问题而构造的非参数统计方法, 不致因为对总体分布的假定不当而导致重大错误, 所以它往往有较好的稳健性. 这是非参数统计方法的一个非常重要的特点. 但它也有以下缺点: 首先因为非参数方法基于更少的信息作出推断, 在模型假定正确的前提下, 非参数统计方法就会比参数统计方法的效果差一些. 例如, 在处理估计问题时, 估计的方差要大一些, 收敛速度要慢一些

(参数估计的速度一般为 $n^{-1/2}$, 但非参数方法的收敛速度比 $n^{-1/2}$ 慢). 又例如, 在给定的显著性水平下进行检验时, 基于非参数估计方法构建的检验方法的第 II 类错误相比基于参数估计方法构建的检验方法要大些. 其次, 非参数方法受数据维数的影响, 存在维数祸根的问题. 具体表现为随着模型变量的维数增加, 所需样本量成指数级增加. 这就导致数据的维数高于三维时, 很多非参数方法的效果并不好. 发展克服或者部分克服维数问题的非参数和半参数方法是当前研究的热点之一.

为了克服非参数方法的缺陷而发展起来的是所谓的半参数统计方法. 半参数统计方法是 20 世纪 70 年代以后发展起来的重要的统计方法. 它在参数模型的基础上引入非参数分量, 从而使这种模型既含有参数分量又含有非参数分量, 兼顾了参数模型的准确和非参数模型的稳健的优点, 相比单纯的参数模型或非参数模型有更大的适应性, 具有更强的解释能力, 并且部分地克服了维数祸根的问题. 半参数模型吸引了很多理论研究领域和应用领域的关注. 本书的主要内容之一是介绍非常典型的并且应用非常广泛的几类半参数模型, 包括部分线性模型、单指标模型等, 以及研究生存数据时使用非常广泛的 Cox 模型.

1.2 收敛方式和极限分布

在介绍非参数和半参数方法之前, 我们给出在后面的内容中经常需要用到的概率论基础知识, 主要包括随机序列的几种收敛方式以及包括弱大数定律、强大数定律和中心极限定理在内的一些统计渐近理论.

1.2.1 依概率收敛

依概率收敛是用概率的方法刻画随机变量的极限.

定义 1.2.1 (依概率收敛) 对随机变量序列 $\{X_n, n = 1, 2, \dots\}$ 和随机变量 X , 若满足: $\forall \varepsilon > 0, \lim_{n \rightarrow \infty} P(|X_n - X| \geq \varepsilon) = 0$, 则称随机变量序列 $\{X_n, n = 1, 2, \dots\}$ 依概率收敛于随机变量 X , 记为 $X_n \xrightarrow{P} X$.

举例: 假设 X_1, X_2, \dots, X_n 是均值为 μ 、方差为 σ^2 的独立同分布序列. \bar{X}_n 为样本均值. 显然 $E(\bar{X}_n) = \mu$ 和 $\text{var}(\bar{X}_n) = \sigma^2/n$. 由切比雪夫不等式, 对于 $\forall \varepsilon > 0$,

$$P(|\bar{X}_n - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{n\varepsilon^2}.$$

所以 $\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| \geq \varepsilon) \leq \lim_{n \rightarrow \infty} \frac{\sigma^2}{n\varepsilon^2} = 0$, 即 $\bar{X}_n \xrightarrow{P} \mu$.

定理 1.2.1 (弱大数定律) 假设 X_1, X_2, \dots, X_n 是独立同分布随机变量, 且 $E|X_1| < \infty$, 则当 $n \rightarrow \infty$ 时有

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} E(X_1).$$

注: (1) 更一般的情况下, $\{X_n, n = 1, 2, \dots\}$ 是独立随机变量序列, 并且 $E(X_i) = \mu_i$, 有

$$\frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n \mu_i \xrightarrow{P} 0.$$

(2) 设 a_n 是 a 的估计, 若 $a_n \xrightarrow{P} a$, 则称 a_n 是 a 的弱相合估计.

因此, 定理 1.2.1 中, \bar{X}_n 是 $E(X_1)$ 的弱相合估计.

(3) 大数定律 (law of large numbers, LLN) 说明当样本量足够大时, 样本均值的随机性消失. 也就是说, 从更多的数据, 可以得到更多样本空间的信息.

下面给出需要经常使用的一个定理. 注意 $\{X_n, n = 1, 2, \dots\}$ 和 $\{Y_n, n = 1, 2, \dots\}$ 为随机变量序列.

定理 1.2.2 若 $X_n \xrightarrow{P} X$, $Y_n \xrightarrow{P} Y$, 则有:

(1) $cX_n \xrightarrow{P} cX$, 其中 c 为常数;

(2) $X_n + Y_n \xrightarrow{P} X + Y$;

(3) $X_n Y_n \xrightarrow{P} XY$;

(4) 若 $Y \neq \mathbf{0}$, 则有 $X_n/Y_n \xrightarrow{P} X/Y$.

定理 1.2.3 若 $X_n \xrightarrow{P} X$, 且 $f(\cdot)$ 是连续函数, 则有 $f(X_n) \xrightarrow{P} f(X)$.

定理 1.2.3 经常被称为 Slutsky 定理.

1.2.2 几乎必然收敛

几乎必然收敛又称为以概率 1 收敛.

定义 1.2.2 (几乎必然收敛) 随机变量序列 $\{X_n, n = 1, 2, \dots\}$, 当 $P(\lim_{n \rightarrow \infty} X_n = X) = 1$ 时, 说它几乎必然 (以概率为 1) 收敛于一个随机变量 X , 记为: $X_n \xrightarrow{\text{a.s.}} X$.

注: 等价地, 若对 $\forall \epsilon > 0$, 有 $P(\lim_{n \rightarrow \infty} |X_n - X| < \epsilon) = 1$, 则 $X_n \xrightarrow{\text{a.s.}} X$.

下面介绍另一个 a.s. 收敛的定义.

定理 1.2.4 $X_n \xrightarrow{\text{a.s.}} X$ 当且仅当对 $\forall \epsilon > 0$, $\lim_{m \rightarrow \infty} P(\sup_{n \geq m} |X_n - X| \leq \epsilon) = 1$.

注: 若 $\forall \epsilon > 0$, $\lim_{n \rightarrow \infty} P(|X_n - X| \leq \epsilon) = 1$, 则 $X_n \xrightarrow{P} X$. 由上面定理知几乎必然收敛强于依概率收敛.

定理 1.2.5 (强大数定律) 假设 X_1, X_2, \dots, X_n 是独立同分布的随机变量序列, 且有 $E|X_1| < \infty$, 则当 $n \rightarrow \infty$ 时, 有

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\text{a.s.}} E(X_1).$$

注: 设 a_n 是 a 的估计, 若 $a_n \xrightarrow{\text{a.s.}} a$, 则称 a_n 是 a 的强相合估计.

因此, 定理 1.2.5 中, \bar{X}_n 是 $E(X_1)$ 的强相合估计.

1.2.3 r 阶收敛

定义 1.2.3 (r 阶中心矩收敛) 对随机变量序列 $\{X_n, n = 1, 2, \dots\}$, 存在 $r > 0$ 有 $E|X_n|^r < \infty$. 若存在一个随机变量 X , 使得 $\lim_{n \rightarrow \infty} E(|X_n - X|^r) = 0$, 则称 X_n 依 r 阶中心矩收敛于 (在 L^r 空间) X , 记为 $X_n \xrightarrow{r.m.} X$ 或 $X_n \xrightarrow{L^r} X$.

注: 一般在 $r = 2$ 情况下讨论. 此时称其为均方收敛.

定义 1.2.4 (r 阶矩收敛) 对随机变量序列 $\{X_n, n = 1, 2, \dots\}$, 存在 $r > 0$ 有 $E|X_n|^r < \infty$. 若 $\lim_{n \rightarrow \infty} E|X_n|^r = E|X|^r$, 则 X_n 依 r 阶矩收敛于 X .

1.2.4 依分布收敛

定义 1.2.5 (依分布收敛) 设 $F_n(x)$ 和 $F(x)$ 分别是随机变量序列 X_n 和随机变量 X 的分布函数. 若 $\lim_{n \rightarrow \infty} F_n(x) = F(x)$ 对 $F(\cdot)$ 的定义域中的任意连续点都成立, 则称随机变量序列 $\{X_n, n = 1, 2, \dots\}$ 依分布收敛于分布函数为 $F(x)$ 的随机变量 X , 记为 $X_n \xrightarrow{d} X$.

注: (1) 对于依分布收敛, $\{X_n, n = 1, 2, \dots\}$ 不需要定义在相同的概率空间. 它不是 $\{X_i\}$ 的收敛, 而是由 $\{X_n, n = 1, 2, \dots\}$ 导出的概率分布 $\{F_n, n = 1, 2, \dots\}$ 的收敛. 可以将其视为在一些概率测度下的弱拓扑的集合的收敛. 因此, 文献中经常称依分布收敛为弱收敛. 它“弱”是因为它是可以由其他的一些收敛得到, 例如依概率收敛和几乎必然收敛.

(2) 此外, $X_n \xrightarrow{d} X$ 当且仅当对任意在紧集上有界连续的函数 f 有 $Ef(X_n) \rightarrow Ef(X)$. 进一步, φ_n 和 φ 分别为 X_n 和 X 的特征函数. $X_n \xrightarrow{d} X$ 当且仅当 $\varphi_n(t) \rightarrow \varphi(t)$. 这些结论的证明以及更多的弱收敛的等价定义可以参考文献 (Pollard, 1984).

1.2.5 收敛方式间的关系

下面讨论随机变量序列的几种收敛方式之间的关系.

定理 1.2.6 (1) 若 $X_n \xrightarrow{a.s.} X$, 则 $X_n \xrightarrow{P} X$.

(2) 若 $X_n \xrightarrow{r.m.} X$, 则 $X_n \xrightarrow{P} X$.

(3) 若 $X_n \xrightarrow{P} X$, 则 $X_n \xrightarrow{d} X$.

证明: (1) 注意到 $\lim_{n \rightarrow \infty} P\{|X_n - X| > \epsilon\} \leq \lim_{n \rightarrow \infty} P\{\bigcup_{k \geq n} \{|X_k - X| > \epsilon\}\} = P\{\bigcap_{n=1}^{\infty} \bigcup_{k \geq n} \{|X_k - X| > \epsilon\}\} = 0$.

(2) 因为对于充分大的 n 有 $\lim_{n \rightarrow \infty} E(|X_n - X|^r) = 0$, $E(|X_n - X|^r) < \infty$, 则由切比雪夫不等式得, 对任意 $\epsilon > 0$, 有 $P(|X_n - X| \geq \epsilon) \leq E(|X_n - X|^r)/\epsilon^r \rightarrow 0, n \rightarrow \infty$. 因此 $\lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1$.

(3) f 是任意有界且一致连续的函数. 令 $M = \sup_x |f(x)|$. 对任意 $\epsilon > 0$, 选择 δ 使得 $|X_n - X| \leq \delta$, 有 $|f(X_n) - f(X)| \leq \epsilon$. 可以得到 $E|f(X_n) - f(X)| \leq \epsilon + 2M \times P\{|X_n - X| >$

$\delta\}$. 这样就有 $|Ef(X_n) - Ef(X)| \leq E|f(X_n) - f(X)| \leq \epsilon + 2MP\{|X_n - X| > \delta\}$. 因为 $X_n \xrightarrow{P} X$, 因此可证得 $Ef(X_n) \rightarrow Ef(X)$. 从而定理 (c) 结论得证.

1.3 中心极限定理和几个常用的定理

1.3.1 中心极限定理

下面介绍几个关于中心极限定理的著名结果.

定理 1.3.1 (Lindeberg-Levy 中心极限定理) 设 $\{X_i\}_{i=1}^n$ 是均值向量为有限向量 μ 、协方差阵为正定阵 Σ 的独立同分布随机向量, 则

$$Z_n \equiv \sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(\mathbf{0}, \Sigma),$$

其中 $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

下面给出另一个常用的中心极限定理.

定理 1.3.2 (Liapounov 中心极限定理) 设 $\{X_{n,i}\}_{i=1}^n$ 是独立随机变量序列, $E(X_{n,i}) = \mu_{n,i}$ 且 $\text{var}(X_{n,i}) = \sigma_{n,i}^2$. 假设存在 $\delta > 0$, 有 $E|X_{n,i}|^{2+\delta} < \infty$. 令 $S_n = \sum_{i=1}^n (X_{n,i} - \mu_{n,i})$, $L_{n,i} = (X_{n,i} - \mu_{n,i})/\sigma_n$, 其中 $\sigma_n^2 = \sum_{i=1}^n \sigma_{n,i}^2$. 若存在 $\delta > 0$, 使得 $\lim_{n \rightarrow \infty} \sum_{i=1}^n E|L_{n,i}|^{2+\delta} = 0$, 则有

$$\frac{S_n}{\sigma_n} = \sum_{i=1}^n L_{n,i} \xrightarrow{d} N(0, 1).$$

注: 上述定理的一个特殊情况是当 $\mu_{n,i} = \mathbf{0}$ 且 $\sigma_{n,i}^2$ 满足 $\lim_{n \rightarrow \infty} \frac{\sigma_n^2}{n} = \sigma^2$. 若存在 $\delta > 0$, 使得 $\lim_{n \rightarrow \infty} \sum_{i=1}^n E \left| \frac{L_{n,i}}{\sqrt{n}} \right|^{2+\delta} = 0$, 则有

$$\frac{S_n}{\sqrt{n}\sigma} \xrightarrow{d} N(0, 1).$$

实际上, 对非参数和半参数模型的估计和检验问题, 统计量经常表现为双求和的形式, 这时经常要用到 U 统计量的中心极限定理, 更多的关于 U 统计量的中心极限定理的内容可以参考文献 (Lee, 1990).

1.3.2 几个常用的定理

下面列举一些在渐近分析中常用的定理. 下面定理中 $\{X_i\}_{i=1}^n, \{Y_i\}_{i=1}^n$ 为独立同分布随机变量序列.

定理 1.3.3 设 f 是一个连续函数, 则有:

(1) 若 $X_n \xrightarrow{a.s.} X$, 则 $f(X_n) \xrightarrow{a.s.} f(X)$.

(2) 若 $X_n \xrightarrow{P} X$, 则 $f(X_n) \xrightarrow{P} f(X)$.

(3) 若 $X_n \xrightarrow{d} X$, 则 $f(X_n) \xrightarrow{d} f(X)$.

特别地, (3) 被称为连续映射定理 (CMT).

定理 1.3.4 若 $X_n \xrightarrow{d} X$ 且 $Y_n \xrightarrow{P} a$, a 是常数, 则

(1) $Y_n X_n \xrightarrow{d} aX$;

(2) $X_n + Y_n \xrightarrow{d} X + a$.

1.3.3 Delta 方法

下面介绍 Delta 方法.

定理 1.3.5 (Delta 方法) 映射 $\phi: \mathbf{R}^k \rightarrow \mathbf{R}$ 关于 θ 连续可微并且 $\phi'(\theta) \neq 0$. 若 $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \Sigma)$, 则 $\sqrt{n}(\phi(\hat{\theta}_n) - \phi(\theta)) \xrightarrow{d} N(0, [\phi'(\theta)]^T \Sigma \phi'(\theta))$.

证明: 由泰勒展开,

$$\phi(\hat{\theta}_n) = \phi(\theta) + [\phi'(\theta)]^T (\hat{\theta}_n - \theta) + o(\|\hat{\theta}_n - \theta\|) = \phi(\theta) + [\phi'(\theta)]^T (\hat{\theta}_n - \theta) + o_p(n^{-\frac{1}{2}}),$$

由此得到:

$$\sqrt{n}(\phi(\hat{\theta}_n) - \phi(\theta)) = [\phi'(\theta)]^T \sqrt{n}(\hat{\theta}_n - \theta) + o_p(1) \xrightarrow{d} N(0, [\phi'(\theta)]^T \Sigma \phi'(\theta)).$$

1.4 记号 $o_p(1)$ 和 $O_p(1)$

在这一节, 介绍以后章节中经常要用到的两个十分重要的记号: $o_p(1)$ 和 $O_p(1)$. 这两个记号与数学分析中的 $o(1)$ 和 $O(1)$ 十分相似, 它们的运算规律也十分相似. 需要注意的是 $o_p(1)$ 或 $O_p(1)$ 是具有某种大样本性质的随机变量序列.

令 ξ_n 为一个随机变量序列, 又设 $n \rightarrow \infty$ 表示一个过程. 说 $\xi_n = o_p(1)$ 是指在 n 趋于无穷时 $\xi_n \xrightarrow{P} 0$, 或对任意 $\varepsilon > 0$, 当 $n \rightarrow \infty$ 时, $P\{|\xi_n| \geq \varepsilon\} \rightarrow 0$.

同时用 $\xi_n = O_p(1)$ 表示 ξ_n 是依概率有界的量, 即对任意 $\varepsilon > 0$, 存在 $M > 0$, 使得

$$P\{|\xi_n| \geq M\} < \varepsilon.$$

例: 设 $\{X_i\}_{i=1}^n$ 为均值为 μ 、方差为 σ^2 的独立随机变量, 则:

(1) 若 $\mu = 0$, 则 $\sum_{i=1}^n X_i$ 是 $O_p(n^{1/2})$; 若 $\mu \neq 0$, 则 $\sum_{i=1}^n X_i$ 是 $O_p(n)$.

(2) 若 $\mu = \sigma^2 = 0$ 不成立, 则 $\sum_{i=1}^n X_i^2$ 是 $O_p(n)$.

注意: (1) 经常出现在渐近理论中. 若 $\mu = 0$, 则由 $\sum_{i=1}^n X_i$ 的均值为 0 且方差为 $n\sigma^2$ 可以得到 $\sum_{i=1}^n \frac{X_i}{n^{1/2}}$ 的均值为 0 且方差为 σ^2 . 所以若 $\mu = 0$, 则 $\sum_{i=1}^n X_i$ 是 $O_p(n^{1/2})$. 然而, 若 $\mu \neq 0$, 则由 $\sum_{i=1}^n X_i$ 的均值为 $n\mu$, 可以得到 $\sum_{i=1}^n X_i$ 是 $O_p(n)$, 而不是 $O_p(n^{1/2})$.

另外, (2) 由以下事实得到: X_i^2 均值 $m_1 = \mu^2 + \sigma^2$ 有限且不为 0, 由强大数定理可得结论.

注: 文献中经常出现“ \sqrt{n} 相合”的概念, 其定义为: 若 $\hat{\theta} - \theta = O_p(n^{-1/2})$, 则称 $\hat{\theta}$ 为 θ 的 \sqrt{n} 相合估计. 若 $\hat{\theta}$ 有分解: $\hat{\theta} = \theta + n^{-1/2} \sum_{i=1}^n \xi_i + o_p(n^{-1/2})$, 其中 $\xi_i, i = 1, 2, \dots, n$ 独立同分布且 $E(\xi_1) = 0$. 上式第二个部分是 $O_p(n^{-1/2})$. 因此有 $\hat{\theta} - \theta = O_p(n^{-1/2})$, 这时 $\hat{\theta}$ 为 θ 的 \sqrt{n} 相合估计.

关于 $o_p(1)$ 和 $O_p(1)$, 它们具有下列性质:

$$\begin{aligned} o_p(1) + o_p(1) &= o_p(1), \\ o_p(1) + O_p(1) &= O_p(1), \\ o_p(1) \cdot o_p(1) &= o_p(1), \\ o_p(1) \cdot O_p(1) &= o_p(1), \\ o_p(1) + c &= O_p(1) \quad (c \neq 0). \end{aligned}$$

另外, 有两个类似的符号 $o(1)$ 和 $O(1)$. 其中 $o(1)$ 表示一个关于 1 的高阶无穷小量, 也就是说表示一个极限为 0 的量; $O(1)$ 表示一个有界量.

举例:

- (1) -4 是 $O(1)$.
- (2) $6n^3$ 是 $O(n^3)$, $o(n^4)$, $o(n^5)$.
- (3) $\frac{7}{n}$ 是 $O(n^{-1})$, 同时也是 $o(1)$.
- (4) $\frac{5}{n} - \frac{3}{n^{3/2}}$ 是 $O(n^{-1})$.

在运算中, $o_p(1)$ 和 $O_p(1)$ 与通常的 $o(1)$ 和 $O(1)$ 一起使用, 具有下面的一些结果:

$$\begin{aligned} o_p(1) + o(1) &= o_p(1), \\ o_p(1) + O(1) &= O_p(1), \\ O_p(1) + o(1) &= O_p(1), \end{aligned}$$

等等. 但注意

$$o_p(1) + o(1) \neq o(1).$$

上面的等式之所以不成立, 其原因是左边为随机变量序列, 而右边为数列, 两者性质不同.

$o_p(1)$ 和 $O_p(1)$ 还有下面的性质:

设 $\xi_n \xrightarrow{d} F$, 则 $\xi_n = O_p(1)$, $\xi_n + o_p(1) \xrightarrow{d} F$, $\xi_n \cdot o_p(1) = o_p(1)$.

除了 $o_p(1)$ 和 $O_p(1)$ 外, 有时还用记号 $o_p(\xi_n)$ 和 $O_p(\xi_n)$. $o_p(\xi_n)$ 表示随机变量序列 $\xi_n o_p(1)$, $O_p(\xi_n)$ 表示随机变量序列 $\xi_n O_p(1)$.

第 2 章 非参数核密度估计

2.1 介 绍

当分析样本数据时,经常希望通过密度函数或分布函数来了解数据的特点.而一般情况下,数据样本对应的总体的密度函数和分布函数是不知道的,这时就有必要对它们进行估计从而获取总体的信息.对密度函数,当已有的经验给出足够信息说明数据所在总体的密度函数形式是已知的,就可以运用参数方法,比如极大似然估计方法.如果密度函数的形式的假定是错误的,那么运用参数推断方法就会得出错误的结论.这种情况下,发展不依赖于密度函数的形式的方法就非常有必要了.非参数密度估计方法正是这样不依赖于密度函数形式的假定而对密度函数进行估计的方法.

因为非参数方法不需要假定密度函数的形式,因此适合很多类型的数据,比如非正态数据、重尾数据等.估计密度函数的非参数方法有核密度估计方法、近邻估计方法、序列估计方法、罚似然估计方法以及局部似然估计方法等,其中使用最广、理论最完善的方法是核密度估计方法.

本章除考虑密度函数的核估计及其渐近性质和带宽选择之外,还考虑了分布函数、密度函数的导数、条件密度函数基于核方法的估计及其渐近性质.

2.2 单元密度函数的估计

2.2.1 核密度估计的提出

这一节,我们考虑一维随机变量的密度函数的非参数估计.假设总体为 X ,其密度函数为 $f(x)$,分布函数为 $F(x)$.有来自总体 X 的独立同分布的样本 $\{X_i, i = 1, 2, \dots, n\}$.如果没有关于 $f(x)$ 和 $F(x)$ 的函数形式的信息,参数估计方法就不再可用.下面发展不需要假定密度函数或分布函数的函数形式的非参数方法来估计密度函数 $f(x)$.

注意到密度函数是分布函数的导数,即有

$$f(x) = F'(x) = \lim_{h \rightarrow 0} \frac{F(x+h) - F(x-h)}{2h}.$$

另外注意到分布函数的一个常用的估计是经验分布函数, 其定义为 $F_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$, 这里 $I(\cdot)$ 是示性函数. 因此将分布函数的估计, 经验分布函数 $F_n(x)$, 代入上式, 可以给出密度函数的估计:

$$f_n(x) = \frac{F_n(x+h) - F_n(x-h)}{2h}, \quad (2.2.1)$$

这里 h 为比较小的数. 注意到 $2nhf_n(x) = \sum_{i=1}^n I(x-h < X_i \leq x+h)$. 易见其服从二项分布 $B(n, F(x+h) - F(x-h))$.

下面计算 $f_n(x)$ 的均值和方差. 当 $h \rightarrow 0$ 时, 有

$$E[f_n(x)] = \frac{1}{2h}[F(x+h) - F(x-h)] \rightarrow f(x).$$

又因为

$$\text{var}[2nhf_n(x)] = n[1 - (F(x+h) - F(x-h))][F(x+h) - F(x-h)],$$

故当 $h \rightarrow 0$ 且 $nh \rightarrow \infty$ 时, 有

$$\text{var}[f_n(x)] = \frac{1}{4nh^2}[1 - (F(x+h) - F(x-h))][F(x+h) - F(x-h)] \rightarrow 0.$$

从上可以看到, 估计量 $f_n(x)$ 是 $f(x)$ 的相合估计, 并且方差趋于 0. 这样看来 $f_n(x)$ 是 $f(x)$ 的一个比较好的估计.

设 $k(u) = \frac{1}{2}I(|u| \leq 1)$. 利用经验分布函数的定义, 对式 (2.2.1) 做适当变形后, 可得到:

$$\begin{aligned} f_n(x) &= \frac{1}{n} \sum_{i=1}^n \frac{1}{h} \frac{1}{2} I\left(-1 \leq \frac{X_i - x}{h} \leq 1\right) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{1}{h} k\left(\frac{X_i - x}{h}\right). \end{aligned} \quad (2.2.2)$$

从上面可见, 密度函数的估计 $f_n(x)$ 实际上是一个加权和, 并且在进行加权求和时, 对处于区间 $[x-h, x+h]$ 的样本点赋予同样的权值, 也即是同等对待的. 这样的处理实际上是不太合理的. 因为直观来说, 估计 X 在 x 点处的密度 $f(x)$, 离 x 较近的样本点应该能提供更多的关于 $f(x)$ 的信息. 因此在定义 $f(x)$ 的估计时, 就要对离 x 较近的样本点赋比较大的权, 距离 x 较远的点赋比较小的权或者赋权为 0. 这可通过选取不同于恒等函数 $k(u) = \frac{1}{2}I(|u| \leq 1)$ 的权函数来实现. 这样的权函数就称为核函数. 前面所用的权函数 $k(u) = \frac{1}{2}I(|u| \leq 1)$ 也是核函数, 称为均匀核函数, 易见它是 $[-1, 1]$ 上均匀分布的密度函数. 上面定义的 $f_n(x)$ 又称为均匀核密度估计或者朴素 (Naive) 估计.

2.2.2 常用的核函数及其性质

除均匀核函数外,常用的核函数有:

Triangle 核函数: $k(u) = (1 - |u|)I(|u| \leq 1)$;

Epanechnikov 核函数: $k(u) = \frac{3}{4}(1 - u^2)I(|u| \leq 1)$;

Quartic 核函数: $k(u) = \frac{15}{16}(1 - u^2)^2I(|u| \leq 1)$;

Triweight 核函数: $k(u) = \frac{35}{32}(1 - u^2)^3I(|u| \leq 1)$;

Gaussian 核函数: $k(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}u^2\right)$;

Cosine 核函数: $k(u) = \frac{\pi}{4} \cos\left(\frac{\pi u}{2}\right)I(|u| \leq 1)$, 等等.

其中均匀核函数、Epanechnikov 核函数、Quartic 核函数和 Triweight 核函数可以看做对称 Beta 族 $K_r(u) = 1/\text{Beta}(0.5, r + 1)(1 - u^2)^r I(|u| \leq 1)$ 对应 $r = 0, 1, 2, 3$ 的情形.

图 2.2.1 画出了常用的核函数的图像,可以看到上面的核函数在零点取最大值.然后两边递减.这与之前考虑估计 $f(x)$ 时,需要对距离 x 较近的点赋予比较大的权的考虑是一致的.因为对给定的带宽 h ,距离 x 较近的点也就是使得 $\frac{X_i - x}{h}$ 的绝对值比较小的点.

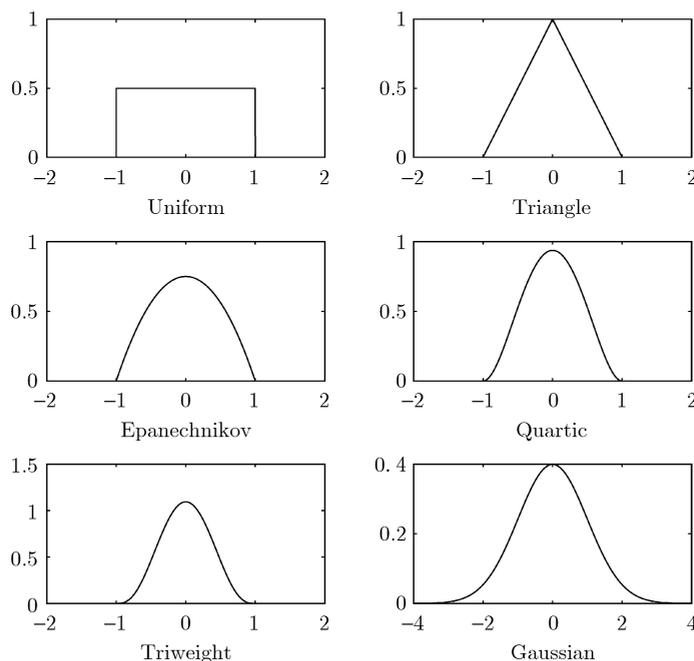


图 2.2.1 常用的核函数