

Эконометрический ликбез: непараметрические и полупараметрические МЕТОДЫ

Непараметрическая регрессия*

Станислав Анатольев[†]

Российская экономическая школа, Москва, Россия

Настоящее эссе повествует о принципах и методологии непараметрического оценивания регрессии среднего. Акцент делается на методах ядерного сглаживания, но дается и обзор неядерных методов.

1 Введение

Пусть имеется случайная выборка $\{(x_i, y_i)\}_{i=1}^n$ из популяции пар (x, y) . Нас интересует оценивание регрессии среднего $g(x) = \mathbb{E}[y|x]$ в предположении, что она существует для всех x носителя и является гладкой. Для этого чаще всего пользуются *параметрическими* методами, когда предполагается, что регрессионная функция имеет известную функциональную форму и конечное число неизвестных параметров. Оценивание этих параметров автоматически дает оценки для $g(x)$. Естественно, неверная спецификация функциональной формы может привести к серьезным искажениям при оценивании и инференции, причем часто непредсказуемым (см. Крил, 2008).

В настоящем эссе мы даем обзор *непараметрического* оценивания регрессии среднего, то есть такого, при котором избегают параметрических предположений о функциональной форме. Как и большая часть соответствующей литературы, мы делаем акцент на методах ядерного сглаживания. В отличие от остальной литературы, однако, мы рассматриваем с самого начала регрессию среднего, избегая предварительного разговора об оценивании плотности. Другие источники информации на данную тему включают обзоры Härdle & Linton (1994) и Расин (2008), а также монографии Härdle (1990), Pagan & Ullah (1999) и Li & Racine (2007).

2 Построение непараметрической оценки

Мы вначале предполагаем, что регрессор x – единственный. Впоследствии мы обсудим и случай многопеременной регрессии.

2.1 Дискретный регрессор

Прежде всего рассмотрим случай дискретного регрессора. Пусть носитель x сосредоточен в $a_{(1)}, \dots, a_{(k)}$, где $a_{(1)} < \dots < a_{(k)}$, и k конечно (если носитель – бесконечное, но счетное множество, мало что меняется в анализе). Зафиксируем $a_{(j)}$, $j = 1, \dots, k$. Заметим, что

$$g(a_{(j)}) = \mathbb{E}[y|x = a_{(j)}] = \frac{\mathbb{E}[y \mathbb{I}\{x = a_{(j)}\}]}{\mathbb{E}[\mathbb{I}\{x = a_{(j)}\}]}$$

*Работа основана на лекциях, читаемых автором в РЭШ. Цитировать как: Анатольев, Станислав (2009) «Непараметрическая регрессия», Квантиль, №7, стр. 37–52. Citation: Anatolyev, Stanislav (2009) “Nonparametric regression,” Quantile, No.7, pp. 37–52.

[†]Адрес: 117418, г. Москва, Нахимовский проспект, 47, офис 1721(3). Электронная почта: sanatoly@nes.ru

из-за справедливости следующих равенств:

$$\begin{aligned}\mathbb{E} [\mathbb{I} \{x = a_{(j)}\}] &= \mathbb{P}\{x_i = a_{(j)}\}, \\ \mathbb{E} [y \mathbb{I} \{x = a_{(j)}\}] &= \mathbb{E} [y|x = a_{(j)}] \mathbb{P}\{x = a_{(j)}\}.\end{aligned}$$

Согласно принципу аналогий можно построить оценку $\hat{g}(a_{(j)})$ как

$$\hat{g}(a_{(j)}) = \frac{\sum_{i=1}^n y_i \mathbb{I} \{x_i = a_{(j)}\}}{\sum_{i=1}^n \mathbb{I} \{x_i = a_{(j)}\}}, \quad (1)$$

что можно интерпретировать как среднее по наблюдениям, попадающим в вертикальное сечение $x = a_{(j)}$.

2.2 Непрерывный регрессор

Если регрессор непрерывно распределен, описанный метод не работает, так как в произвольном сечении $x = a$ нечего усреднять (ибо туда попадет максимум одно наблюдение), хотя и $f(a) \neq 0$, где $f(a)$ – значение плотности регрессора $f(x)$ в a . Поэтому необходимо привлечь информацию откуда-то еще. Если регрессионная кривая непрерывна, наблюдения, попадающие в окрестность a , являются наиболее информативными о значении регрессии в a .

Выберем положительное h и назовем его *шириной окна* (хотя впоследствии необязательно будет задействовано окно в буквальном смысле). Обобщим формулу (1) следующим образом:

$$\hat{g}(a) = \frac{\sum_{i=1}^n y_i \mathbb{I} \{a - h \leq x_i \leq a + h\}}{\sum_{i=1}^n \mathbb{I} \{a - h \leq x_i \leq a + h\}}, \quad (2)$$

т.е. мы усредняем y по наблюдениям, попадающим в окно $[a - h, a + h]$. При изменении a $\hat{g}(a)$ описывает оцененную регрессионную кривую. Отметим, что последняя имеет разрывы из-за попадания в окно новых наблюдений и выпадения из него старых.

Информация от наблюдений, попавших в окно $[a - h, a + h]$, используется одинаково. То есть те наблюдения, которые попали в окно, участвуют в оценивании с равным весом, в то время как наблюдения, не попавшие в окно, вообще в нем не участвуют. Разумной представляется идея сделать веса зависимыми от расстояния от x_i до a , а также, возможно, использовать информацию во всех наблюдениях. Введем с этой целью симметричную *ядерную функцию* $K(u)$, интегрирующуюся в единицу, т.е. $\int K(u) du = 1$, где интеграл \int берется по всей области определения, которой может быть отрезок, обычно $[-1, 1]$, или вся числовая ось. Популярными ядрами являются:

$$\text{Равномерное:} \quad K(u) = \frac{1}{2} \mathbb{I} \{|u| \leq 1\},$$

$$\text{Треугольное:} \quad K(u) = (1 - |u|) \mathbb{I} \{|u| \leq 1\},$$

$$\text{Епанечниково:} \quad K(u) = \frac{3}{4} (1 - u^2) \mathbb{I} \{|u| \leq 1\},$$

$$\text{Гауссово (нормальное):} \quad K(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right).$$

Область определения первых трех ядер – отрезок $[-1, 1]$, в то время как последнее имеет бесконечный носитель. Следовательно, при использовании равномерного, треугольного или Епанечникова ядра оценка будет использовать информацию в ограниченном окне в окрестности a , а оценка, использующая гауссово ядро, будет использовать информацию из всех наблюдений. Заметим также, что в принципе ядро не обязано быть всюду неотрицательным.

Далее, введем обозначение

$$K_h(u) = \frac{1}{h} K\left(\frac{u}{h}\right).$$

Теперь обобщим формулу (2) и получим

$$\hat{g}(a) = \frac{\sum_{i=1}^n y_i K_h(x_i - a)}{\sum_{i=1}^n K_h(x_i - a)}. \quad (3)$$

Оценка (3) называется *оценкой Надарайа–Уотсона* для регрессии среднего, в честь Nadaraya (1965) и Watson (1964). Заметим, что нормализация делением на h в определении $K_h(u)$ не влияет на численное значение оценки и сделано лишь для удобства, в чем мы убедимся позднее.

Заметим, что при использовании трех последних ядер (треугольного, Епанечникова и гауссова), так же как и многих других, оцененная регрессионная кривая непрерывна, так как новые и старые наблюдения вводятся в формулу и выводятся из нее непрерывным образом по мере того как a меняется. При этом велика роль параметра ширины окна. Если h слишком велика, оценка задействует слишком много нерелевантной информации, что увеличивает смещение и приводит к явлению *сверхсглаживания*. Сверхсглаженная кривая слишком «линейная», в то время как ей положено быть более извилистой и более пристально отслеживать рисунок наблюдений. Если же h слишком мала, оценка задействует маловато точек, что увеличивает дисперсию и приводит к явлению *недосглаживания*. Недосглаженная кривая слишком извилистая, так как она слишком пристально отслеживает отдельные наблюдения.

3 Асимптотические свойства

В случае дискретного регрессора легко получить, используя ЗБЧ и ЦПТ, что

$$\sqrt{n} (\hat{g}(a_{(j)}) - g(a_{(j)})) \xrightarrow{d} \mathcal{N} \left(0, \frac{\mathbb{V}[y|x = a_{(j)}]}{\mathbb{P}\{x = a_{(j)}\}} \right). \quad (4)$$

Интерпретация выражения для асимптотической дисперсии следующая. Качество оценивания положительно связано с частотой попадания точек в вертикальное сечение $x = a_{(j)}$ и отрицательно связано со степенью разброса точек вдоль него. Заметим, что скорость сходимости оценки – параметрическая, \sqrt{n} . Действительно, задачу можно трактовать как параметрическую, ибо вектор параметров $(a_{(1)}, \dots, a_{(k)})'$ конечномерен.

В случае непрерывно распределенного регрессора необходимо, чтобы ширина окна асимптотически падала до нуля, иначе смещение из-за нерелевантности используемой информации из соседних к a наблюдений испортит состоятельность оценки. Таким образом, необходимо установить правило $h \rightarrow 0$ по мере того как $n \rightarrow \infty$. С другой стороны, ширина окна не должна падать слишком быстро, иначе дисперсия из-за малого количества участвующих в оценивании точек не будет падать, что также испортит состоятельность оценки. Точнее, поскольку дисперсия обратно пропорциональна эффективному количеству участвующих наблюдений, которое в свою очередь прямо пропорционально nh , необходимо установить правило $nh \rightarrow \infty$ по мере того как $h \rightarrow 0$ и $n \rightarrow \infty$.

Обозначим

$$\sigma_K^2 = \int u^2 K(u) du$$

и

$$R_K = \int K(u)^2 du.$$

Эти две константы зависят только от выбранного ядра. Подразумевается, что обе величины конечны. Далее, установим дополнительное требование к скорости падения ширины окна:

$$\lambda = \lim_{n \rightarrow \infty} \sqrt{nh^5},$$

предполагая $\lambda < \infty$. Заметим, что λ может равняться, а может и не равняться нулю. Мы также предполагаем непрерывность и ограниченность $g(x)$, $g'(x)$, $g''(x)$, $f(x)$ и $f'(x)$ всюду на области определения.

Рассмотрим разницу между оценкой и оцениваемой величиной:

$$\hat{g}(a) - g(a) = \frac{\hat{q}_1(a) + \hat{q}_2(a)}{\hat{f}(a)},$$

где

$$\begin{aligned}\hat{q}_1(a) &= \frac{1}{n} \sum_{i=1}^n e_i K_h(x_i - a), \\ \hat{q}_2(a) &= \frac{1}{n} \sum_{i=1}^n (g(x_i) - g(a)) K_h(x_i - a), \\ \hat{f}(a) &= \frac{1}{n} \sum_{i=1}^n K_h(x_i - a),\end{aligned}$$

и за e_i обозначены, как обычно, регрессионные ошибки в точках выборки: $e_i = y_i - g(x_i)$.

Начнем со знаменателя

$$\hat{f}(a) = \frac{1}{n} \sum_{i=1}^n K_h(x_i - a).$$

Величина $\hat{f}(a)$, называемая *оценкой плотности Надарайа–Уотсона*, и правда оценивает плотность регрессора $f(x)$ в точке a , отсюда и обозначения. Действительно, рассмотрим

$$\begin{aligned}\mathbb{E}[\hat{f}(a) - f(a)] &= \mathbb{E}[K_h(x - a)] - f(a) \\ &= \frac{1}{h} \int K\left(\frac{x - a}{h}\right) f(x) dx - f(a) \\ &= \int K(u) f(a + hu) du - f(a) \\ &= \int K(u) f(a + hu) du - f(a).\end{aligned}$$

Разложим $f(a + hu)$ до первого порядка:

$$\begin{aligned}\mathbb{E}[\hat{f}(a) - f(a)] &= \int K(u) (f(a) + f'(a) hu + o(h)) du - f(a) \\ &= f(a) \int K(u) du + f'(a) h \int u K(u) du + o(h) - f(a) \\ &= o(h),\end{aligned}$$

так как ядро интегрируется в единицу и симметрично. Использую ту же технологию,

$$\begin{aligned}\mathbb{V}[\hat{f}(a)] &= \frac{1}{n} \mathbb{V}[K_h(x - a)] \\ &= \frac{1}{n} \mathbb{E}[K_h(x - a)^2] - \frac{1}{n} \mathbb{E}[K_h(x - a)]^2 \\ &= \frac{1}{n} \int \left(\frac{1}{h} K\left(\frac{x - a}{h}\right)\right)^2 f(x) dx - \frac{1}{n} \left(\int \frac{1}{h} K\left(\frac{x - a}{h}\right) f(x) dx\right)^2 \\ &= \frac{1}{n h} \int K(u)^2 f(a + hu) du - \frac{1}{n} \left(\int K(u) f(a + hu) du\right)^2 \\ &= \frac{1}{n h} \int K(u)^2 (f(a) + o(1)) du - \frac{1}{n} O(1) \\ &= O\left(\frac{1}{n h}\right).\end{aligned}$$

Поскольку $h \rightarrow 0$ и $nh \rightarrow \infty$, имеем $\mathbb{E}[\hat{f}(a) - f(a)] \rightarrow 0$ и $\mathbb{V}[\hat{f}(a)] \rightarrow 0$, откуда следует, что действительно $\hat{f}(a) \xrightarrow{P} f(a)$.

Теперь рассмотрим первую часть числителя

$$\hat{q}_1(a) = \frac{1}{n} \sum_{i=1}^n e_i K_h(x_i - a).$$

Это среднее регрессионных ошибок, взвешенных ядром. Как и в параметрическом анализе, такое среднее должно дать асимптотическую нормальность. Разница в том, что дисперсия отдельного слагаемого здесь не постоянна, а зависит от асимптотически падающего h . Следовательно, необходимо рассчитать предел этой дисперсии. Представим $\sqrt{nh}\hat{q}_1(a)$ как

$$\sqrt{nh}\hat{q}_1(a) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{e_i}{\sqrt{h}} K\left(\frac{x_i - a}{h}\right).$$

Дисперсия отдельного слагаемого равна

$$\begin{aligned} \mathbb{V}\left[\frac{e}{\sqrt{h}} K\left(\frac{x-a}{h}\right)\right] &= \mathbb{E}\left[\frac{\sigma^2(x)}{h} K\left(\frac{x-a}{h}\right)^2\right] \\ &= \frac{1}{h} \int \sigma^2(x) K\left(\frac{x-a}{h}\right)^2 f(x) dx \\ &= \int \sigma^2(a+hu) K(u)^2 f(a+hu) du \\ &= \int \sigma^2(a) K(u)^2 f(a) du + o(1) \\ &= \sigma^2(a) f(a) R_K + o(1). \end{aligned}$$

По ЦПТ Линдберга–Леви, $\sqrt{nh}\hat{q}_1(a) \xrightarrow{d} \mathcal{N}(0, \sigma^2(a) f(a) R_K)$.

Наконец, рассмотрим вторую часть числителя

$$\hat{q}_2(a) = \frac{1}{n} \sum_{i=1}^n (g(x_i) - g(a)) K_h(x_i - a).$$

Это среднее взвешенных ядром отклонений значений регрессионной функции в точках наблюдений от ее значения в точке a , где она оценивается. Эти отклонения рождают смещение. Конечно же, $\hat{q}_2(a)$ обладает и дисперсией, но она мала по сравнению с дисперсией, рождаемой в $\hat{q}_1(a)$ регрессионными ошибками.

Рассмотрим математическое ожидание $\hat{q}_2(a)$:

$$\begin{aligned} \mathbb{E}[\hat{q}_2(a)] &= \mathbb{E}[(g(x) - g(a)) K_h(x - a)] \\ &= \frac{1}{h} \int (g(x) - g(a)) K\left(\frac{x-a}{h}\right) f(x) dx \\ &= \int (g(a+hu) - g(a)) K(u) f(a+hu) du \\ &= \int \left(g'(a) hu + \frac{g''(a)}{2} (hu)^2 + o(h^2)\right) K(u) (f(a) + f'(a) hu + o(h)) du \\ &= g'(a) f(a) h \int u K(u) du \\ &\quad + \left(g'(a) f'(a) + \frac{g''(a)}{2} f(a)\right) h^2 \int u^2 K(u) du + o(h^2) \\ &= h^2 f(a) B(a) \sigma_K^2 + o(h^2), \end{aligned}$$

где

$$B(a) = \frac{g'(a)f'(a)}{f(a)} + \frac{g''(a)}{2}.$$

Легко также получить, что

$$\mathbb{V}[\hat{q}_2(a)] = o\left(\frac{1}{nh}\right).$$

Эти два результата приводят к тому, что

$$\sqrt{nh}\hat{q}_2(a) \xrightarrow{p} \lambda f(a) B(a) \sigma_K^2.$$

Все выведенное выше в совокупности дает

$$\begin{aligned} \sqrt{nh}(\hat{g}(a) - g(a)) &= \frac{\sqrt{nh}\hat{q}_1(a) + \sqrt{nh}\hat{q}_2(a)}{\hat{f}(a)} \\ &\xrightarrow{d} \frac{\mathcal{N}(0, \sigma^2(a) f(a) R_K) + \lambda f(a) B(a) \sigma_K^2}{f(a)} \\ &\sim \mathcal{N}\left(\lambda B(a) \sigma_K^2, \frac{\sigma^2(a)}{f(a)} R_K\right). \end{aligned}$$

Отметим две интересные черты этого асимптотического результата. Во-первых, непараметрическая скорость сходимости \sqrt{nh} меньше, чем параметрическая \sqrt{n} , поскольку асимптотически $h \rightarrow 0$. Этот факт отражает меньшую точность оценивания бесконечномерных объектов, чем точность оценивания конечномерных. Во-вторых, хотя асимптотическое распределение и нормальное, оно нецентрированное. Асимптотическое смещение отражает тот факт, что информация, используемая при оценивании, не является всецело релевантной.

Формула для асимптотической дисперсии похожа на свой аналог в случае дискретного регрессора: скедастическая функция в числителе и «вероятностная масса» в знаменателе. Асимптотическое смещение зависит от множества характеристик форм регрессионной и плотностной функций, зашитых в формуле для $B(a)$. Одна часть асимптотического смещения пропорциональна $g'(a)f'(a)$ и отражает смещение, возникающее, когда регрессионная кривая имеет наклон, а наблюдения падают несимметрично слева и справа от a , в результате чего точки слева и точки справа создают неодинаковое и взаимно не компенсирующееся смещение вверх и вниз. Вторая часть асимптотического смещения пропорциональна $g''(a)$ и отражает смещение, возникающее, когда регрессионная кривая локально нелинейна, в результате чего ординаты точек слева и точек справа от a несимметрично распределены выше и ниже $g(a)$, даже если их абсциссы распределены симметрично слева и справа от a . Отметим, что и дисперсия, и смещение в общем случае обратно пропорциональны плотности $f(a)$, что отражает тот факт, что точность оценивания, и в смысле дисперсии, и в смысле смещения, низка при оценивании $g(a)$ около тех границ носителя x , где плотность убывает в ноль.

Полученный асимптотический результат также означает, что оптимальная скорость падения ширины окна $h \propto n^{-1/5}$, так как в этом случае $\lambda > 0$, и асимптотические смещение и дисперсия уравновешены. С другой стороны, если положить $h = o(n^{-1/5})$, можно добиться того, что λ будет равно нулю, и асимптотическое смещение исчезнет. Конечно, это удобно с точки зрения реализации, т.к. в таком случае нет необходимости оценивать компоненты $B(a)$, но подобные действия скрывают истинную картину, связанную со смещением оценки, и могут привести к плохому качеству асимптотического приближения.

Выведенный выше асимптотический результат означает, что приближенно

$$\hat{g}(a) \sim \mathcal{N}\left(g(a) + \frac{\sqrt{nh^5} B(a) \sigma_K^2}{\sqrt{nh}}, \frac{\sigma^2(a) R_K}{f(a) nh}\right).$$

Как обычно, асимптотическое распределение можно использовать для тестирования статистических гипотез о $g(a)$ и построения доверительных интервалов для этой величины. Доверительный интервал, например, выглядит так:

$$\hat{g}(a) - h^2 \hat{B}(a) \sigma_K^2 \mp z_{1-\alpha/2} \sqrt{\frac{\hat{\sigma}^2(a) R_K}{\hat{f}(a) nh}},$$

где $z_{1-\alpha/2}$ – $(1 - \alpha/2)$ -квантиль стандартного нормального распределения, а $\hat{f}(a)$, $\hat{\sigma}^2(a)$ и $\hat{B}(a)$ – непараметрические оценки соответствующих функций в точке a . Заметим, что если построить доверительные интервалы для $g(a)$ при всех значениях a (на решетке) с вероятностью покрытия $1 - \alpha$, возникнет *поточечный доверительный коридор* для $g(x)$. Естественно, неверно говорить, что истинная регрессионная кривая находится внутри этого коридора с вероятностью $1 - \alpha$.

4 Выбор ширины окна

4.1 Правило подстановки

Из выведенного выше результата следует, что асимптотическая среднеквадратическая ошибка оценивания равна

$$AMSE(a) = h^4 B(a)^2 \sigma_K^4 + \frac{\sigma^2(a) R_K}{f(a) nh}.$$

Если минимизировать эту величину по h , то возникнет *правило подстановки* для (локально) оптимальной ширины окна

$$h^*(a) = \left(\frac{\sigma^2(a) R_K}{4f(a) B(a)^2 \sigma_K^4} \right)^{1/5} n^{-1/5}.$$

Заметим, что скорость падения равна оптимальной, выведенной выше.

Практическое применение правила подстановки заключается в (непараметрическом!) оценивании $f(a)$, $f'(a)$, $g'(a)$, $g''(a)$, $\sigma^2(a)$, вычислении R_K , σ_K^4 и подстановке полученных результатов в формулу для $h^*(a)$. Эта трудоемкая процедура дает численное значение оптимальной ширины окна всего для одного значения a , так что ее приходится повторить для всех a . Конечно, это очень нелегко, и исследователи чувствовали бы себя комфортнее с одним-единственным значением *глобально оптимальной ширины окна* h^* , общей для всех a .

Глобально оптимальную ширину окна легко вывести, минимизируя критерий интегрированной асимптотической среднеквадратической ошибки оценивания

$$IAMSE(a) = \int AMSE(a) da = h^4 \sigma_K^4 \int B(a)^2 da + \frac{R_K}{nh} \int \frac{\sigma^2(a)}{f(a)} da.$$

Естественно, можно ввести какую-либо взвешивающую схему, если есть причины не использовать равномерное взвешивание для разных a . При равномерном взвешивании

$$h^* = \left(\frac{\int \sigma^2(a) / f(a) da R_K}{4 \int B(a)^2 da \sigma_K^4} \right)^{1/5} n^{-1/5}.$$

Данная стратегия, однако, не лишает исследователя необходимости оценивать $f(a)$, $f'(a)$, $g'(a)$, $g''(a)$, $\sigma^2(a)$, то есть процедура настолько же трудоемкая, как и прежде. Статистик Бернард Сильверман предложил универсальную формулу, основанную на вышеописанной процедуре, но специфичной для определенного частного случая, например, в предположении

нормальной плотности f .¹ Вот эта универсальная формула, называемая *правилом Сильвермана*:

$$h^S = 1.364 \left(\frac{R_K}{\sigma_K^4} \right)^{1/5} \hat{\sigma}_x n^{-1/5},$$

где $\hat{\sigma}_x^2$ – выборочная дисперсия регрессора. В частности, для гауссова ядра формула выглядит как $h^S = 1.06 \hat{\sigma}_x n^{-1/5}$.

На практике правило Сильвермана обычно обеспечивает приемлемые результаты, кроме, возможно, тех случаев, когда оценивание происходит вблизи краев носителя x . Иногда, впрочем, вид оцененной регрессионной кривой неудовлетворителен, и подобную ширину окна используют лишь как стартовое значение при поиске более приемлемой.

4.2 Кросс-валидация

Радикально иное правило выбора глобальной ширины окна – это *кросс-валидация*. Она основана на качестве подгонки, а не на асимптотических свойствах оценки. Можно обычную среднеквадратическую ошибку, оцененную в точках выборки,

$$MSE(h) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{g}(x_i))^2,$$

устремить к нулю, что приводит к идеальной подгонке, если положить h равным своему наименьшему значению, когда каждое наблюдение объясняется только им же самим. Естественно, это не что иное как экстремальная степень недосглаживания, и оно нас не устраивает. От источника проблемы можно легко избавиться, если запретить объяснять наблюдение самим собой. Исходя из этого, разумным критерием будет *функция кросс-валидации*

$$CV(h) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{g}^{-i}(x_i))^2,$$

где $\hat{g}^{-i}(x_i)$ – это оценка Надарайа–Уотсона в точке x_i , которая (оценка) использует все наблюдения, за исключением i -го, т.е.

$$\hat{g}^{-i}(x_i) = \frac{\sum_{j=1, j \neq i}^n y_j K_h(x_j - x_i)}{\sum_{j=1, j \neq i}^n K_h(x_j - x_i)}.$$

Оптимальная в смысле кросс-валидации ширина окна h^{CV} минимизирует $CV(h)$. К сожалению, при практическом применении данная ширина окна часто приводит к сильному недосглаживанию. Следовательно, в этих случаях ее можно использовать в качестве начального приближения для приемлемой ширины окна, и окончательное решение опять же остается за визуальным анализом.

5 Многопеременная ядерная регрессия

До сих пор регрессор x был единственным. Как обобщить оценку Надарайа–Уотсона на постановку с множественными регрессорами?

Обозначим через d размерность x . У ядра теперь будет d -мерный аргумент, отображающий расстояние между a и каждым x_i в d -мерном пространстве в скалярный вес: $K : \mathbb{R}^d \rightarrow \mathbb{R}$.

¹На самом деле Сильверман использовал оптимальную ширину окна для оценки плотности Надарайа–Уотсона. См. одну из задач в конце эссе.

Определим $d \times d$ -мерную симметричную и положительно определенную *матрицу ширины окна* H . Далее, положим

$$K_H(u) = \frac{1}{\det H} K(H^{-1}u).$$

Оценка Надарайа–Уотсона выглядит по-прежнему как

$$\hat{g}(a) = \frac{\sum_{i=1}^n y_i K_H(x_i - a)}{\sum_{i=1}^n K_H(x_i - a)}.$$

Есть много способов организовать структуру K и H . На практике широко используются два способа.

Первый использует разложение d -мерного пространства на произведение d одномерных:

$$K_H(u) = \prod_{\ell=1}^d K_{h_{\ell}, \ell}(u_{\ell}) = \prod_{\ell=1}^d \frac{1}{h_{\ell}} K_{\ell}\left(\frac{u_{\ell}}{h_{\ell}}\right),$$

где K_{ℓ} и h_{ℓ} – ядро и ширина окна, соответственно, по ℓ -ой размерности. Такое $K_H(u)$ называется *ядром-произведением*. K_{ℓ} потенциально могут быть разными по разным координатам, но обычно они одинаковы. Матрица ширины окна равна $H = \text{diag}\{h_{\ell}\}_{\ell=1}^d$, и каждая h_{ℓ} выбирается отдельно (например, в простейшем случае – с помощью правила Сильвермана).

Ядро-произведение игнорирует зависимость между регрессорами. Кроме того, выбор d значений ширины окна не особенно привлекателен. У второго метода нет этих недостатков. Матрица ширины окна имеет следующую структуру:

$$H = hS^{1/2},$$

где h – единая ширина окна, S – выборочная дисперсионная матрица регрессоров, а $S^{1/2}$ – квадратный корень из нее (например, определяемый через разложение Холецки).

К сожалению, точность оценивания при больших значениях d мала. Это явление называют *проклятием размерности*. Причина в том, что при прочих равных в d -мерное гиперокно попадает намного меньше наблюдений, чем в его одномерный аналог, и дисперсия оценки быстро возрастает с увеличением d . В частности, скорость сходимости (при второй схеме) равна $\sqrt{nh^d}$, а оптимальная ширина окна – $O(n^{-1/(d+4)})$. На практике непараметрическое оценивание обычно возможно только при очень малом количестве регрессоров d вроде 1, 2, 3, редко выше, и требует больших выборок для достижения приличной надежности оценивания.

6 Локальная полиномиальная регрессия

Вернемся к скалярному x . Заметим, что оценку Надарайа–Уотсона можно представить в виде решения задачи минимизации взвешенной суммы квадратов:

$$\hat{g}(a) = \arg \min_{\beta_0} \sum_{i=1}^n (y_i - \beta_0)^2 K_h(x_i - a).$$

Это означает, что оценивание Надарайа–Уотсона – это локальная (в том смысле, что взвешивание основано на локальности наблюдений к a) регрессия на константе. Нет причин останавливаться на регрессии на константе. Естественной является идея расширить ее на линейную регрессию не только на константе, но также и на x . Результатом будет *локальная линейная регрессия*:

$$\hat{g}_1(a) = (1, 0) \arg \min_{(\beta_0, \beta_1)} \sum_{i=1}^n (y_i - \beta_0 - \beta_1(x_i - a))^2 K_h(x_i - a).$$

Вектор $(1, 0)$ отбирает первый элемент вектора 2×1 . В качестве сопутствующего результата второй элемент этого вектора дает локально линейную оценку наклона регрессионной кривой в a , то есть ее первой производной $g'(a)$.

Оценка локальной линейной регрессии имеет то преимущество перед оценкой Надарайа–Уотсона, что она учитывает наклонность регрессионной кривой, что имеет значение, если наблюдения распределены неравномерно вокруг a , и таким образом уменьшает смещение. Это проявляется в асимптотических свойствах оценки, идентичных свойствам оценки Надарайа–Уотсона, кроме того момента, что теперь

$$B(a) = \frac{g''(a)}{2}.$$

Локальную полиномиальную регрессию можно обобщить далее и получить *локальную полиномиальную регрессию порядка p* :

$$\hat{g}_p(a) = (1, 0, \dots, 0) \arg \min_{(\beta_0, \dots, \beta_p)} \sum_{i=1}^n (y_i - \beta_0 - \dots - \beta_p (x_i - a)^p)^2 K_h(x_i - a).$$

Данная оценка учитывает не только наклонность, но и свойства кривизны регрессионной кривой в точке a , и в качестве сопутствующего результата дает оценки производных $g(x)$ до p -го порядка в точке a . Удобно, что оценку можно записать в явном виде как оценку взвешенного метода наименьших квадратов:

$$\hat{g}_p(a) = (1, 0, \dots, 0) (X'WX)^{-1} X'WY,$$

где $Y = (y_1, \dots, y_n)'$, $X = (X_1, \dots, X_n)'$, $X_i = (1, x_i - a, \dots, (x_i - a)^p)'$, и, наконец, $W = \text{diag} \{K_h(x_i - a)\}_{i=1}^n$. Если $p > 1$, асимптотическое смещение $B(a)$ равно нулю. Это означает, что оптимальная ширина окна и результирующая скорость сходимости оценки меняются.

При практическом применении существуют серьезные недостатки использования локальной полиномиальной регрессии. Использование информации в локальных непараметрических методах очень ограничено, а увеличение p означает уменьшение степеней свободы. В качестве экстремального примера рассмотрим равномерное ядро и узкое окно, настолько узкое, что в него попадает лишь два наблюдения. Оценка Надарайа–Уотсона усредняет ординаты этих двух наблюдений, локальная линейная регрессия соединяет их прямой и берет ординату пересечения с вертикальной прямой $x = a$, а локальная полиномиальная регрессия при $p > 1$ попросту не существует.

На практике стоит ограничиваться малыми p , вроде 0, 1 или 2, не больше.

7 Временные ряды

Если вместо случайной выборки у нас стационарный временный ряд, основные методы, приведенные выше, в целом работают. Интересно, что, в отличие от параметрических задач, в асимптотической дисперсии не возникают автоковариации y_t в $x_t = a$. Причина в том, что эффект последних асимптотически исчезает на фоне дисперсии y_t в $x_t = a$ (см. Bierens, 1994). Таким образом, можно использовать те же самые формулы.

Типичное приложение непараметрических методов во временных рядах – непараметрическая авторегрессия

$$y_t = g(y_{t-1}, y_{t-2}, \dots, y_{t-k}) + \sigma(y_{t-1}, y_{t-2}, \dots, y_{t-k}) \eta_t$$

где $\mathbb{E}[\eta_t | y_{t-1}, y_{t-2}, \dots] = 0$ и $\mathbb{V}[\eta_t | y_{t-1}, y_{t-2}, \dots] = 1$.

Оценку Надарайа–Уотсона $\hat{g}(a_1, \dots, a_k)$ авторегрессионной функции $g(y_{t-1}, \dots, y_{t-k})$ в $(y_{t-1}, \dots, y_{t-k}) = (a_1, \dots, a_k)$ можно построить по знакомой схеме, а оценку автоскедастичной функции $\sigma^2(y_{t-1}, \dots, y_{t-k})$ как

$$\hat{\sigma}^2(y_{t-1}, \dots, y_{t-k}) = \hat{\delta}(a_1, \dots, a_k) - \hat{g}(a_1, \dots, a_k)^2,$$

где $\hat{\delta}(a_1, \dots, a_k)$ – оценка Надарайа–Уотсона непараметрической регрессии y_t^2 на $(y_{t-1}, \dots, y_{t-k})$ в точке $(y_{t-1}, \dots, y_{t-k}) = (a_1, \dots, a_k)$. Если порядок авторегрессии k неизвестен, его можно оценить совместно с регрессионной функцией (см., например, Tschernig & Yang, 2000). Обзор непараметрических методов в контексте временных рядов содержится в Heiler (2001).

8 Другие непараметрические методы

Непараметрический метод может быть одного из двух типов: локальный или глобальный. Оценка Надарайа–Уотсона, локальная линейная и локальная полиномиальная регрессии принадлежат классу локальных методов, так как оценивание $g(x)$ в точке a использует информацию в наблюдениях, находящихся вблизи a . Глобальные непараметрические методы вместо этого пытаются подогнать всю кривую ко всем точкам выборки одновременно. При этом влияние одного наблюдения не так ограничено, и значения его абсциссы и ординаты влияют не только на оцененную регрессию поблизости этой точки, но и на положение и форму всей оцененной регрессионной кривой.

Независимо от того, локальный метод или глобальный, всегда наличествует *параметр сглаживания*, контролирующий степень последнего. Выбор этого параметра осуществляется исследователем. Параметр сглаживания в уже рассмотренных методах – это ширина окна.

8.1 Метод ближайших соседей

Другим локальным непараметрическим методом является *оценка k ближайших соседей*

$$\hat{g}_{NN}(a) = \frac{1}{k} \sum_{i=1}^n y_i 1_i,$$

где $1_i = \mathbb{I}\{x_i \text{ является одним из } k \text{ ближайших соседей к } a\}$, и о соседстве судится по тому, насколько близки абсциссы точек выборки к a . Здесь k является параметром сглаживания, и для состоятельности необходимы условия $k \rightarrow \infty$ и $k/n \rightarrow 0$ по мере того как $n \rightarrow \infty$. Другой версией является *симметризованная оценка k ближайших соседей*. Предполагая, что k четно,

$$\hat{g}_{SNN}(a) = \frac{1}{k} \sum_{i=1}^n y_i 1_i,$$

где $1_i = \mathbb{I}\{x_i \text{ является одним из левых } k/2 \text{ или правых } k/2 \text{ ближайших соседей к } a\}$. Асимптотические свойства этих оценок схожи со свойствами оценки Надарайа–Уотсона.

Одним из преимуществ методов ближайших соседей перед ядерными методами является существование оценки при любом раскладе, так как у любой точки есть соседи в любой выборке, в то время как в окно могут не попасть наблюдения совсем. Кроме того, количество усредняемых величин контролируется напрямую. Конечно же, обе идеи можно скомбинировать.

8.2 Разложение в серии (решето)

Оценивание сериями, или решето, – один из глобальных методов, основанный на разложении гладкой функции по базису в функциональном пространстве. Положим, существует, предпочтительно ортогональный, упорядоченный базис $\{\psi_j(x)\}_{j=0}^{\infty}$, такой, что

$$g(x) = \sum_{j=0}^{\infty} \gamma_j \psi_j(x).$$

Термин «упорядоченный» означает, что $\psi_0(x)$ наиболее важный член, $\psi_1(x)$ менее важный, $\psi_2(x)$ еще менее важный и т.д. Например, базисом могут быть полиномы Эрмита или синусы и косинусы Фурье. Выберем параметр отсечения J , такой, что асимптотически $J \rightarrow \infty$ и $J/n \rightarrow 0$ по мере того как $n \rightarrow \infty$. Тогда

$$\hat{g}_S(x) = \sum_{j=0}^J \hat{\gamma}_j \psi_j(x),$$

где $(\hat{\gamma}_0, \dots, \hat{\gamma}_J)'$ – вектор МНК-оценок в «линейной регрессии» y на $(\psi_0(x), \dots, \psi_J(x))'$. Этот метод также иногда называют полунепараметрическим, так как он непараметрический по сути, но параметрический по технической реализации. См. обзор Chen (2007), а также статью Кристенсен (2009) в настоящем номере журнала.

8.3 Искусственные нейронные сети

Идея искусственных нейронных сетей похожа на идею решета, т.е. аппроксимацию нелинейной функции линейной комбинацией неких базовых функций, в таком количестве, в каком необходимо для достижения хорошего приближения. В качестве базисных обычно используются логистические функции:

$$\hat{g}_{ANN}(x) = \hat{\phi}_{0,0} + \hat{\phi}_{1,0}x + \sum_{j=1}^J \frac{\hat{\gamma}_j}{1 + \exp(\hat{\phi}_{0,j} + \hat{\phi}_{1,j}x)},$$

где коэффициенты оцениваются нелинейным методом наименьших квадратов вместо МНК, используемого в решете. Расширенные версии задействуют интересные иерархические структуры, см. Franses & van Dijk (2000, глава 5).

8.4 Сплаины

Другим глобальным методом является оценивание сплайнами. Положим, мы хотим подогнать всю регрессионную кривую, используя обычный критерий подгонки, сумму квадратов ошибок. Это, конечно же, немедленно приведет к интерполяции, т.е. недосглаживанию максимальной степени. Однако мы можем добавить штрафной член, наказывающий за недосглаживание. Интерполирующая кривая слишком извилистая, то есть обладает большой второй производной по абсолютной величине. Эти идеи приводят к следующей оценке сплайнами:

$$\hat{g}_{CS}(x) = \arg \min_{\hat{g}(x)} \sum_{i=1}^n (y_i - \hat{g}(x_i))^2 + \lambda \int (\hat{g}''(u))^2 du,$$

где параметром сглаживания является λ . Если λ мало, решение близко к интерполирующей кривой, в то время как если λ очень большое, решение близко к линейному предиктору с точки зрения критерия наименьших квадратов. Индекс CS расшифровывается как «кубический сплайн», что означает решение в виде кусочно кубического полинома с непрерывно дифференцируемыми переходами в абсциссах точек наблюдений. Серьезное пособие по сплайнам – Wahba (1990).

9 Задачи

9.1 Оценка плотности Надарайа–Уотсона

Вывести асимптотическое распределение оценки Надарайа–Уотсона плотности скалярной случайной величины x , имеющей непрерывное распределение, аналогично тому, как выведено асимптотическое распределение оценки Надарайа–Уотсона регрессионной функции, при аналогичных предположениях. Дать интерпретацию зависимости выражений для асимптотического смещения и асимптотической дисперсии от формы плотности.

9.2 Несмещенность ядерных оценок

Рассмотрим оценку Надарайа–Уотсона $\hat{g}(x)$ условного среднего $g(x)$ для случайной выборки. Показать, что если $g(x) = c$, где c – некоторая константа, то оценка $\hat{g}(x)$ несмещена. Какова интуиция за этим результатом? Выяснить, при каких обстоятельствах оценка локальной линейной регрессии $g(x)$ будет несмещена. Будет ли оценка плотности $f(x)$ несмещена?

9.3 Оценивание при ограничении на форму

Фирмы производят продукт, используя технологию $f(l, k)$. Функциональная форма f неизвестна, но известно, что она обладает свойством постоянного эффекта от масштаба. Для фирмы i наблюдается труд l_i , капитал k_i и выпуск y_i , а порождающий данные процесс принимает форму $y_i = f(l_i, k_i) + \varepsilon_i$, где $\mathbb{E}[\varepsilon_i] = 0$, и ошибка ε_i независима от (l_i, k_i) . Для случайной выборки $\{y_i, l_i, k_i\}_{i=1}^n$ предложить непараметрическую оценку $f(l, k)$, которая бы тоже обладала свойством постоянного эффекта от масштаба.

9.4 Непараметрическая функция риска

Пусть z_1, \dots, z_n – скалярные независимые одинаково распределенные случайные величины с неизвестной плотностью $f(\cdot)$ и функцией распределения $F(\cdot)$. Предположим, что распределение z имеет носитель \mathbb{R} . Возьмем $t \in \mathbb{R}$, такое что $0 < F(t) < 1$. Целью является оценивание функции риска

$$H(t) = \frac{f(t)}{1 - F(t)}.$$

Предложить непараметрическую оценку $\hat{F}(t)$ для $F(t)$. Обозначим за $\hat{f}(t)$ оценку Надарайа–Уотсона для $f(t)$, и выберем ширину окна h так, что $nh^5 \rightarrow 0$. Предложить оценку $\hat{H}(t)$ для $H(t)$, использующую $\hat{F}(t)$ и $\hat{f}(t)$, и найти ее асимптотическое распределение.

10 Решения задач

10.1 Оценка плотности Надарайа–Уотсона

Применим знакомую технологию:

$$\begin{aligned} \mathbb{E}[\hat{f}(a) - f(a)] &= \int K(u) \left(f(a) + hu f'(a) + \frac{1}{2}(hu)^2 f''(a) + O(h^3) \right) du - f(a) \\ &= \frac{h^2}{2} f''(a) \sigma_K^2 + O(h^3) \end{aligned}$$

и

$$\begin{aligned} \mathbb{V}[\hat{f}(a)] &= \frac{1}{nh} \int K(u)^2 (f(a) + o(1)) du - \frac{1}{n} O(1) \\ &= \frac{1}{nh} f(a) R_K + o\left(\frac{1}{nh}\right). \end{aligned}$$

Используя ЦПТ Линдберга–Леви, по мере того как $n \rightarrow \infty$, $h \rightarrow 0$ и $nh \rightarrow \infty$, имеем

$$\sqrt{nh} \left(\hat{f}(a) - f(a) \right) \xrightarrow{d} \mathcal{N} \left(\frac{1}{2} \lambda f''(a) \sigma_K^2, f(a) R_K \right),$$

при условии, что $\lambda \equiv \lim_{n \rightarrow \infty} \sqrt{nh^5}$ существует и конечно.

Асимптотическое смещение пропорционально $f''(a)$, значение которой говорит о том, насколько плотность в окрестности a отличается от оцениваемой плотности в a . Отметим, что асимптотическое смещение не зависит от $f(a)$, т.е. как часто наблюдения попадают в данную область, и от $f'(a)$, т.е. отличаются ли плотности слева и справа от a . Асимптотическая дисперсия пропорциональна $f(a)$, плотности в a , что может показаться странным (большая частота выпадения наблюдений приводит к худшему качеству оценивания). Однако мы оцениваем $f(a)$, так что ее большее значение также означает больший разброс оценки вокруг истинной величины, и этот эффект превалирует (эффект частоты дает $\propto f(a)^{-1}$, эффект размера дает $\propto f(a)^2$).

10.2 Несмещенность ядерных оценок

Математическое ожидание равно

$$\begin{aligned} \mathbb{E}[\hat{g}(a)] &= \mathbb{E} \left[\mathbb{E} \left[\frac{\sum_{i=1}^n y_i K_h(x_i - a)}{\sum_{i=1}^n K_h(x_i - a)} \middle| x_1, \dots, x_n \right] \right] \\ &= \mathbb{E} \left[\frac{\sum_{i=1}^n \mathbb{E}[y_i | x_i] K_h(x_i - a)}{\sum_{i=1}^n K_h(x_i - a)} \right] = \mathbb{E} \left[\frac{\sum_{i=1}^n c K_h(x_i - a)}{\sum_{i=1}^n K_h(x_i - a)} \right] = c, \end{aligned}$$

т.е. оценка $\hat{g}(a)$ несмещена для $c = g(a)$. Причина проста: все элементы выборки одинаково релевантны при оценивании тривиального условного среднего, так что от участия точек вдали от a смещение не возникает.

Оценка локальной линейной регрессии будет несмещена, если $g(x) = c + bx$. Тогда все элементы выборки одинаково релевантны при оценивании, так как оценивается линейная, хоть и локальная, регрессия. Действительно,

$$\hat{g}_1(a) = \bar{y} + \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) K_h(x_i - a)}{\sum_{i=1}^n (x_i - \bar{x})^2 K_h(x_i - a)} (a - \bar{x}),$$

так что

$$\begin{aligned} \mathbb{E}[\hat{g}_1(a)] &= \mathbb{E}[\mathbb{E}[\bar{y} | x_1, \dots, x_n]] \\ &\quad + \mathbb{E} \left[\mathbb{E} \left[\frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) K_h(x_i - a)}{\sum_{i=1}^n (x_i - \bar{x})^2 K_h(x_i - a)} \middle| x_1, \dots, x_n \right] (a - \bar{x}) \right] \\ &= \mathbb{E}[c + b\bar{x}] \\ &\quad + \mathbb{E} \left[\mathbb{E} \left[\frac{\sum_{i=1}^n (c + bx_i - c - b\bar{x})(x_i - \bar{x}) K_h(x_i - a)}{\sum_{i=1}^n (x_i - \bar{x})^2 K_h(x_i - a)} \middle| x_1, \dots, x_n \right] (a - \bar{x}) \right] \\ &= \mathbb{E}[c + b\bar{x} + b(a - \bar{x})] = c + bx. \end{aligned}$$

Что же касается плотности, вряд ли стоит ожидать несмещенности. Действительно,

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n K_h(x_i - a),$$

так что

$$\mathbb{E}[\hat{f}(a)] = \mathbb{E}[K_h(x_i - a)] = \frac{1}{h} \int K \left(\frac{x_i - a}{h} \right) f(x) dx.$$

Это матожидание сильно зависит от ширины окна и ядра, и вряд ли будет равно $f(x)$, кроме как в особых условиях (например, равномерная $f(x)$, a далеко от границ и т.д.).

10.3 Оценивание при ограничении на форму

Технология с постоянным эффектом от масштаба обладает свойством

$$f(l, k) = kf\left(\frac{l}{k}, 1\right).$$

Регрессия для нормированных переменных выглядит как

$$\frac{y_i}{k_i} = f\left(\frac{l_i}{k_i}, 1\right) + \frac{\varepsilon_i}{k_i}.$$

Поэтому можно построить (одномерную!) ядерную оценку для $f(l, k)$ как

$$\hat{f}(l, k) = k \times \frac{\sum_{i=1}^n \frac{y_i}{k_i} K_h\left(\frac{l_i}{k_i} - \frac{l}{k}\right)}{\sum_{i=1}^n K_h\left(\frac{l_i}{k_i} - \frac{l}{k}\right)}.$$

По сути, мы присваиваем больший вес тем наблюдениям, которые ближе к лучу l/k .

10.4 Непараметрическая функция риска

Простая непараметрическая оценка для $F(t) \equiv \Pr\{z \leq t\}$ – это выборочная частотность

$$\hat{F}(t) = \frac{1}{n} \sum_{j=1}^n \mathbb{I}\{z_j \leq t\}.$$

Согласно ЗБЧ, она состоятельна для $F(t)$. Согласно ЦПТ, ее скорость сходимости равна \sqrt{n} .

Мы знаем из разделов 9.1 и 10.1, что

$$\sqrt{nh}(\hat{f}(t) - f(t)) \xrightarrow{d} \mathcal{N}(0, R_K f(t)),$$

используя $\lambda = 0$. По принципу аналогий,

$$\hat{H}(t) = \frac{\hat{f}(t)}{1 - \hat{F}(t)}.$$

По теореме Слущкого, эта оценка состоятельна для $H(t)$, а также

$$\begin{aligned} \sqrt{nh}(\hat{H}(t) - H(t)) &= \frac{\sqrt{nh}(\hat{f}(t) - f(t))}{1 - \hat{F}(t)} + \sqrt{h}f(t) \frac{\sqrt{n}(\hat{F}(t) - F(t))}{(1 - \hat{F}(t))(1 - F(t))} \\ &\xrightarrow{d} \frac{1}{1 - F(t)} \mathcal{N}(0, R_K f(t)) + 0 \sim \mathcal{N}\left(0, R_K \frac{f(t)}{(1 - F(t))^2}\right). \end{aligned}$$

Причина того, что неопределенность в $\hat{F}(t)$ не влияет на асимптотическое распределение $\hat{H}(t)$, в том, что $\hat{F}(t)$ сходится быстрее, чем $\hat{f}(t)$.

К сожалению, $\hat{H}(t)$ выглядит не очень аппетитно из-за скачков в $\hat{F}(t)$.

Список литературы

- Крил, М. (2008). Некоторые ловушки параметрической инференции. *Квантиль* 4, 1–6.
- Кристенсен, Д. (2009). Полупараметрическая эконометрика: вводный курс. *Квантиль* 7, 53–83.
- Расин, Дж. (2008). Непараметрическая эконометрика: вводный курс. *Квантиль* 4, 7–56.
- Bierens, H.J. (1994). *Topics in Advanced Econometrics: Estimation, Testing, and Specification of Cross-Section and Time Series Models*. New York: Cambridge University Press.
- Chen, X. (2007). Large sample sieve estimation of semi-nonparametric models. Глава 76 в *Handbook of Econometrics* (под редакцией J.J. Heckman & E.E. Leamer), том 6/2. Elsevier Science.
- Franses, P. & D. van Dijk (2000). *Nonlinear Time Series Models in Empirical Finance*. New York: Cambridge University Press.
- Härdle, W. (1990). *Applied Nonparametric Regression*. New York: Cambridge University Press.
- Härdle, W. & O. Linton (1994). Applied nonparametric methods. Глава 38 в *Handbook of Econometrics* (под редакцией R. Engle & D. McFadden), том 4. Elsevier Science.
- Heiler, S. (2001). Nonparametric time series analysis: nonparametric regression, locally weighted regression, autoregression, and quantile regression. Глава в *A Course in Time Series Analysis* под редакцией D. Peña, G. Tiao & R. Tsay. Wiley.
- Li, Q. & J.S. Racine (2007). *Nonparametric Econometrics: Theory and Practice*. Princeton University Press.
- Nadaraya, E.A. (1965). On nonparametric estimates of density functions and regression curves. *Theory of Applied Probability* 10, 186–190.
- Pagan, A. & A. Ullah (1999). *Nonparametric Econometrics*. New York: Cambridge University Press.
- Tschernig, R. & L. Yang (2000). Nonparametric lag selection for time series. *Journal of Time Series Analysis* 21, 457–487.
- Wahba, G. (1990). *Spline Models for Observational Data*. Philadelphia: SIAM.
- Watson, G.S. (1964). Smooth regression analysis. *Sankhya* 26, 359–372.

Nonparametric regression

Stanislav Anatolyev

New Economic School, Moscow, Russia

This essay covers the principles and methodology of nonparametric estimation of a mean regression. The emphasis is put on kernel smoothing, but non-kernel methods are also reviewed.

Полупараметрическое моделирование и оценивание*

Деннис Кристенсен[†]

Колумбийский Университет, Нью-Йорк, США

Центр эконометрического анализа временных рядов, Орхус, Дания

Полупараметрические модели характеризуются тем, что включают в себя конечномерные и бесконечномерные (функциональные) параметры. Из-за этого они обладают дополнительной гибкостью по сравнению с параметрическими моделями, и в то же время можно строить оценки для параметрических компонент, сходящиеся со скоростью, обычной для параметрических оценок. Эти две особенности делают полупараметрические модели и оценки все более и более популярными в прикладной экономике. В настоящем эссе содержится выборочный обзор литературы по полупараметрическому моделированию и оцениванию, с уклоном в полупараметрические регрессионные модели. Особое внимание уделяется построению двухшаговых полупараметрических оценок и выводу их асимптотических свойств. Также кратко обсуждаются оценивание «решетом» и полупараметрическая эффективность.

1 Введение

В течение последних тридцати лет полупараметрическое моделирование и оценивание экономических процессов привлекает много внимания. Главная причина такой популярности в том, что подобная модель является компромиссным решением между двумя крайностями, полностью параметрическим и полностью непараметрическим моделированием. В первом случае для объяснения выборочных данных используются параметрическая модель, и естественная оценка в этом случае – оценка максимального правдоподобия (ММП-оценка). При правильной спецификации ММП-оценка обладает обычными приятными свойствами, такими как максимальная эффективность. Но если некоторые компоненты модели определены неверно, ММП-оценка асимптотически смещена, и выводы, полученные на основе оцененной модели, могут быть крайне обманчивыми. В противоположность этому полностью непараметрические модели дают максимальную гибкость, сводя к минимуму вероятность неправильно специфицировать модель. С другой стороны, непараметрическое оценивание требует много исходных данных, и в малых выборках мы получим довольно неточные оценки. Особенно это проявляется в моделях большой размерности, где точность оценок падает по мере добавления новых переменных. Это явление называется «проклятие размерности».

Полупараметрические модели являются компромиссным решением между непараметрическим и параметрическим подходами. Они включают в себя как непараметрические, так и параметрические компоненты. Поэтому полупараметрическая модель сохраняет до некоторой степени гибкость полностью непараметрической модели и гораздо менее подвержена неправильной спецификации по сравнению с полностью параметрической моделью. В то же время параметрическую компоненту полупараметрической модели вообще можно оценить с точностью, сравнимой с той, которую мы бы получили, используя (правильно специфицированную) полностью параметрическую модель.

*Перевод О. Еремина и С. Анатольева. Цитировать как: Кристенсен, Деннис (2009). «Полупараметрическое моделирование и оценивание», Квантиль, №7, стр. 53–83. Citation: Kristensen, Dennis (2009). “Semiparametric modelling and estimation,” *Quantile*, No.7, pp. 53–83.

[†]Адрес: Economics Department, Columbia University, 1018 International Affairs Building, 420 West 118th Street, New York, NY 10027, USA. Электронная почта: dk2313@columbia.edu

В настоящем эссе мы даем краткое введение и обзор методов по полупараметрическому моделированию и оцениванию, уделяя особое внимание регрессионным моделям. Мы вводим основные понятия полупараметрического моделирования и оценивания в рамках регрессионных моделей по трем причинам. Во-первых, эти модели широко используются в экономике и поэтому должны быть хорошо знакомы рядовому читателю. Во-вторых, с регрессионными моделями довольно просто работать, что позволяет легко ввести основные полупараметрические понятия и методики. В-третьих, большинство способов, которые мы рассматриваем в рамках регрессий, можно перенести на многие другие виды моделей. Чтобы продемонстрировать это, мы кратко коснемся полупараметрических копул и покажем, как представленные методы можно применить в этих условиях.

Сперва мы рассмотрим оценивание некоторых ведущих полупараметрических моделей. Затем мы определим основной подход, с помощью которого можно проанализировать асимптотические свойства этих оценок. Основной класс оценок, который мы будем рассматривать, – это так называемые двухшаговые полупараметрические оценки, в которых на первом шаге оцениваются непараметрические компоненты модели, которые затем используются для оценки параметрических компонент. Мы выведем ряд условий общего типа, при которых полупараметрическая оценка состоятельна и асимптотически нормально распределена, и далее обсудим подробнее, как эти условия можно проверить для конкретной модели.

В качестве альтернативной стратегии оценивания мы кратко рассмотрим класс полупараметрических оценок, основанных на так называемом методе «решето». В то же время мы не детализируем основную теорию этих оценок. В заключение мы немного обсудим полупараметрическую эффективность и как она используется для создания новых оценок. Опять же, эта часть статьи не требует специальных знаний, мы лишь попробуем объяснить различные концепции на интуитивном уровне.

Мы не ставим перед собой цель охватить в данном обзоре всю литературу по данной теме. Следует иметь в виду, что есть много замечательных обзоров литературы по полупараметрическому моделированию и оцениванию. Среди них Ichimura & Todd (2007), Härdle, Müller, Sperlich & Werwatz (2004), Horowitz (2009), Li & Racine (2007), Pagan & Ullah (1999), Powell (1994) и Robinson (1988), дополняющие и расширяющие наш обзор по ряду направлений.

Эссе организовано следующим образом. В разделах 2–4 мы рассмотрим несколько примеров полупараметрических моделей и обсудим, как осуществляется их оценивание. В разделе 5 мы проанализируем свойства основного класса двухшаговых полупараметрических оценок, который включает некоторые из оценок, рассмотренных в предыдущем разделе. Мы сосредоточимся на оценках, основанных на ядерном сглаживании, так как они просты для анализа и популярны в прикладных исследованиях. В разделе 6 мы кратко рассмотрим метод одновременного оценивания обоих компонент, используя, так называемый, метод «решето» для непараметрической компоненты. Полупараметрическую эффективность мы обсудим в разделе 7. В заключение в разделе 8 мы укажем список работ, где более детально рассмотрены упомянутые темы.

Разделы 2–4 и 6–7 не требуют глубоких знаний эконометрической теории, но раздел 5 потребует некоторых усилий от менее подготовленного в техническом плане читателя. Чтобы сохранить сложность изложения на приемлемом уровне, мы не приводим доказательств, и большинство математических выводов даны в общих чертах. Заинтересовавшийся читатель может обратиться к автору за доказательством всех теорем. Кроме того, в разделе 8 приведены статьи, содержащие более точные результаты и строгие доказательства.

2 Полупараметрические регрессии

В самой общей форме регрессионную модель можно представить в виде

$$Y = m(X) + \varepsilon, \quad \mathbb{E}[\varepsilon|X] = 0, \quad (1)$$

где $Y \in \mathbb{R}$ – отклик (зависимая переменная), $X \in \mathbb{R}^d$ – набор из $d \geq 1$ регрессоров (независимых переменных), и $\varepsilon \in \mathbb{R}$ – ошибка. Регрессионная функция $m : \mathbb{R}^d \mapsto \mathbb{R}$ объясняет, как условное математическое ожидание Y меняется с X :

$$\mathbb{E}[Y|X = x] = m(x).$$

Также, пусть $f_{\varepsilon|X}(e|x)$ – условная плотность распределения ε при заданном $X = x$.¹ Предположим, у нас имеется случайная выборка (Y_i, X_i) из модели, где $i = 1, \dots, n$. Нам интересно осуществить инференцию относительно функций $m(x)$ и $f_{\varepsilon|X}(e|x)$.

В полностью параметрическом случае мы предполагаем, что как регрессионная функция, m , так и (условная) функция распределения ошибок, $f_{\varepsilon|X}$, известны вплоть до некоторого параметра конечной размерности. То есть, у нас есть конкретные параметрические функции $m(x; \beta)$ и $f_{\varepsilon|X}(e|x; \sigma)$, где $\beta \in \mathcal{B}$ включает коэффициенты регрессии, характеризующие вид функции m , а $\sigma \in \Sigma$ – параметр, задающий вид (условной) функции распределения ошибок. Предполагая верную спецификацию модели, то есть $m(x) = m(x; \beta_0)$ и $f_{\varepsilon|X}(e|x) = f_{\varepsilon|X}(e|x; \sigma_0)$ для некоторого $\theta_0 = (\beta_0, \sigma_0)$, естественной оценкой модели будет ММП-оценка

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log f_{\varepsilon|X}(Y_i - m(X_i; \beta) | X_i; \sigma).$$

Часто используют гауссовскую регрессионную модель, где вектор ошибок не зависит от X и распределен нормально: $N(0, \sigma^2)$. В этом случае ММП-оценка $\theta = (\beta, \sigma^2)$ совпадает с оценкой наименьших квадратов: $\hat{\beta}_{\text{MLE}} = \hat{\beta}_{\text{LS}}$ и $\hat{\sigma}_{\text{MLE}}^2 = \hat{\sigma}_{\text{LS}}^2$ где

$$\hat{\beta}_{\text{LS}} = \arg \min_{\beta \in \mathcal{B}} \sum_{i=1}^n (Y_i - m(X_i; \beta))^2, \quad \hat{\sigma}_{\text{LS}}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - m(X_i; \hat{\beta}))^2.$$

Что касается спецификации регрессионной функции, широко используется линейная функция: $m(x; \beta) = \beta_1 x_1 + \dots + \beta_d x_d$, и ММП-оценка превращается в оценку метода наименьших квадратов (МНК),

$$\hat{\beta}_{\text{OLS}} = \left(\sum_{i=1}^n X_i X_i' \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right). \quad (2)$$

При выполнении условий регулярности оценка $\hat{\beta}_{\text{MLE}}$ \sqrt{n} -состоятельна и асимптотически нормально распределена. Например, при гауссовых ошибках ММП-оценка удовлетворяет

$$\sqrt{n}(\hat{\beta}_{\text{MLE}} - \beta_0) \rightarrow^d N(0, V), \quad V = \sigma^2 \mathbb{E} [\dot{m}(x; \beta) \dot{m}(x; \beta)']^{-1},$$

где $\dot{m}(x; \beta) = \partial m(x; \beta) / \partial \beta$ (см., например, Amemiya, 1985). В свою очередь, это значит, что регрессионную функцию можно оценить как $\hat{m}_{\text{MLE}}(x) = m(x; \hat{\beta}_{\text{MLE}})$.

Однако параметрическая модель может быть неверно специфицирована, то есть $m(x; \beta) \neq m(x)$ для всех $\beta \in \mathcal{B}$ и/или $f_{\varepsilon|X}(e|x; \sigma) \neq f_{\varepsilon|X}(e|x)$ для всех значений $\sigma \in \Sigma$. В этом случае оценка регрессионной функции $\hat{m}_{\text{MLE}}(x)$, вообще говоря, несостоятельна и ведет к неправильному представлению о том, как X влияет на Y . Чтобы избежать этого, можно использовать полностью параметрические оценки m , например, ядерные оценки или оценки «решетом». Мы сосредоточимся на ядерных оценках и вкратце дадим общее представление о решетчатых оценках. Подробнее смотрите статьи Härdle (1992), Silverman (1986), Расин (2008) и Анатольев (2009). Оценки «решетом» мы кратко обсудим в разделе 6. Ядерные оценки

¹Мы предполагаем, как это обычно принято в литературе по не- и полупараметрическому моделированию, что все переменные имеют непрерывное распределение.

образует особый класс непараметрических оценок, которые используют локальные данные для инференции о характеристиках распределения. Предположим, X непрерывно распределена с плотностью распределения $f(x)$. Тогда плотность можно оценить непараметрически с помощью ядерной оценки плотности: для каждого значения $x \in \mathbb{R}^d$ она вычисляется по формуле

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n K_h(X_i - x), \quad (3)$$

где $K_h(x) = K(x/h)/h^d$, $K: \mathbb{R}^d \mapsto \mathbb{R}$ – ядерная функция, а $h > 0$ – ширина окна. Исследователь сам выбирает K и h . Ядерная оценка плотности похожа на гистограмму плотности, где ширина окна задает ширину ячейки гистограммы, а ядро – вес наблюдения для каждой ячейки. Наибольший вес имеют наблюдения, близкие к x , а расположенные далеко от x играют малую роль или вовсе не учитываются. Аналогичным образом ядерная оценка регрессии $m(x) = \mathbb{E}[Y|X=x]$ в данной точке $x \in \mathbb{R}^d$ принимает форму взвешенного среднего:

$$\hat{m}(x) = \frac{\sum_{i=1}^n Y_i K_h(X_i - x)}{\sum_{i=1}^n K_h(X_i - x)}. \quad (4)$$

Повторим, что это локальная оценка, которая использует наблюдения X_i , близкие к x , для извлечения информации о значении $m(\cdot)$ в точке x .

Ядерные оценки регрессии очень робастны. Оценка $\hat{m}(x)$ состоятельна при $h \rightarrow 0$ и $nh^d \rightarrow \infty$ независимо от истинного значения регрессионной функции m . Но при этом страдает точность оценивания: в конечной выборке дисперсия оценки выше по сравнению с дисперсией параметрических оценок. На теоретическом уровне это следует из того, что скорость сходимости ядерной оценки $\sqrt{n^{4/(4+d)}}$ ниже, чем \sqrt{n} – скорость сходимости параметрических оценок при $h \rightarrow 0$.² Заметим, что точность непараметрических оценок зависит от размерности X , $d \geq 1$: с ростом d снижается скорость сходимости непараметрической оценки (упоминавшееся ранее проклятье размерности). Вдобавок к этому, даже если удастся получить точные непараметрические оценки, ядерные оценки $\hat{f}(x)$ и $\hat{m}(x)$ сложно представить и интерпретировать при большом d .

Таким образом, выбор между разными моделями и способами оценивания – это выбор между риском неправильной спецификации модели и степенью точности оценок. ММП-оценка полностью параметрической модели обладает наибольшей точностью, но очень велика вероятность получить неверные выводы из-за неправильной спецификации. Напротив, полностью непараметрическая модель гарантирует правильную спецификацию, но дает очень низкую точность оценки. Это приводит нас к использованию полупараметрических моделей и оценок, которые обладают относительно высокой гибкостью и в то же время лучшей скоростью сходимости для некоторых параметров модели.

2.1 Одноиндексная модель

Широко распространенная полупараметрическая регрессионная модель, так называемая одноиндексная, имеет следующий вид:

$$Y = g(\beta'X) + \varepsilon, \quad \mathbb{E}[\varepsilon|X] = 0, \quad (5)$$

где функция $g: \mathbb{R} \mapsto \mathbb{R}$ и параметр $\beta \in \mathbb{R}^d$ неизвестны. Мы не делаем никаких предположений относительно (условного) распределения ε , и считаем g и $f_{\varepsilon|X}$ непараметрическими функциями. В этом случае $\gamma = (g, f_{\varepsilon|X})$ – бесконечномерная (непараметрическая) компонента, а β – параметрическая компонента.

²Мы предполагаем, что интересующая нас функция дважды непрерывно дифференцируема.

Название «одноиндексная» происходит от того, что g является функцией *индекса* $\beta'X \in \mathbb{R}$, а не всего вектора $X \in \mathbb{R}^d$. Таким образом, мы полагаем, что X влияет на Y только посредством индекса $\beta'X$, что является ограничением относительно полной регрессионной функции m из уравнения (1). Поэтому по сравнению с полностью непараметрической моделью есть риск неправильной спецификации модели.

С другой стороны, легко интерпретировать влияние X на Y , которое описывается конечномерным параметром β и одномерной функцией g . При этом g имеет область определения \mathbb{R} , в то время как область определения функции m из уравнения (1) – \mathbb{R}^d . Поэтому, независимо от размерности X , оценивание g остается одномерной задачей, что по существу избавляет нас от проклятия размерности.

Данный подход применим для некоторых видов трансформационных моделей. Предположим, случайная величина Y^* удовлетворяет

$$Y^* = \beta'_0 X + \eta,$$

где $\eta \sim F_\eta$ не зависит от X . Однако мы наблюдаем не Y^* , а только

$$Y = t(Y^*)$$

для некоторого преобразования t , которое может быть как известным, так и неизвестным. Видно, что

$$\mathbb{E}[Y|X = x] = \mathbb{E}[t(\beta'_0 X + \eta) | X = x] = \int t(\beta'_0 x + v) dF_\eta(v).$$

Определив g и ε как

$$g(z) := \int t(z + v) dF_\eta(v), \quad \varepsilon := Y - \mathbb{E}[Y|X = x],$$

этот класс моделей можно представить в форме уравнения (5). Трансформационные модели включают в себя модели с ограниченной зависимой переменной, такие как цензурированные регрессионные модели и модели длительности. Например, при $t(y) = 1\{y > 0\}$ трансформационная модель превращается в модель бинарного выбора, такую как логит и пробит. Если преобразование $t(y) = y \cdot 1\{y > 0\}$, мы получим Тобит-модель. Преимущество полупараметрического подхода в том, что мы можем оценить β без необходимости настаивать на точной форме t и F_η .

Теперь мы хотим построить оценку одноиндексной модели. С этой целью сперва нам надо обсудить вопрос идентификации интересующих нас параметров. То есть можем ли мы однозначно определить параметры β и функции g и $f_{\varepsilon|X}$ из данных? Заметим, что параметр β нельзя идентифицировать, если $\Pr\{\delta'X = c\} = 1$ для некоторых констант $c \in \mathbb{R}$ и $\delta \in \mathbb{R}^d$. Более того, нам нужно нормировать β , чтобы идентифицировать g . Чтобы понять, почему это так, положим $\tilde{g}(z) = g(a + bz)$ для некоторых констант $(a, b) \in \mathbb{R}^2$, что эквивалентно $\tilde{g}(-a + 1/bz) = g(z)$. Обе спецификации неразличимы:

$$g(\beta'x) = \tilde{g}(-a + (1/b)\beta'x) = \tilde{g}(\tilde{\beta}'x),$$

где $\tilde{\beta} = (1/b)\beta$, то есть по исходной выборке мы не сможем различить \tilde{g} и g . Поэтому потребуем, чтобы X не содержал констант, а один из коэффициентов β был равен единице; положим $\beta_1 = 1$ (всегда можно переобозначить порядок компонент X). Наконец, заметим, что если g – линейная функция, то мы не можем идентифицировать β (если g неизвестна).

При сделанных предположениях мы теперь можем вывести оценки для β и g . Сперва предположим, что функция g известна. Тогда естественной оценкой β будет МНК-оценка

$$\hat{\beta}_g = \arg \min_{\beta \in \mathcal{B}} \frac{1}{n} \sum_{i=1}^n [Y_i - g(\beta'X_i)]^2. \quad (6)$$

И наоборот, предположим, что $\beta \in \mathbb{R}^d$ известна, тогда естественной непараметрической оценкой g будет стандартная ядерная регрессионная оценка

$$\hat{g}(z; \beta) = \frac{\sum_{i=1}^n Y_i K_h(\beta' X_i - z)}{\sum_{i=1}^n K_h(\beta' X_i - z)}.$$

Однако, если β и g неизвестны, обе эти оценки будут недостижимыми. Вместо этого мы предлагаем использовать комбинированный подход: подставить непараметрическую оценку $\hat{g}(z; \beta)$ в функцию наименьших квадратов уравнения (6) и получить достижимую оценку β :

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n [Y_i - \hat{g}(\beta' X_i; \beta)]^2.$$

Когда мы получили $\hat{\beta}$, очевидно, оценкой $g(z)$ будет $\hat{g}(z; \hat{\beta})$.

Альтернативный способ – метод оценивания средних производных, предложенный Powell, Stock & Stoker (1989). Пусть g дифференцируемая функция, тогда выполняется следующее равенство:

$$\frac{\partial \mathbb{E}[Y|X=x]}{\partial x} = \beta g'(\beta' x),$$

где $g'(x) = \partial g(z) / \partial z$. Следовательно, для любой ограниченной функции w

$$\mathbb{E} \left[\frac{\partial \mathbb{E}[Y|X]}{\partial x} w(X) \right] = \beta \mathbb{E} [w(X) g'(\theta' X)].$$

Это доказывает, что параметр δ , заданный формулой

$$\delta := \mathbb{E} \left[\frac{\partial \mathbb{E}[Y|X]}{\partial x} w(X) \right]$$

неотличим с точностью до коэффициента нормирования ($\mathbb{E} [w(X) g'(\theta' X)]$) от β . Теперь построим оценку δ с использованием весовой функции w , такой что $w(x) = f(x)$, где f – плотность распределения X . Сперва рассмотрим

$$\begin{aligned} \mathbb{E} \left[\frac{\partial \mathbb{E}[Y|X]}{\partial x} f(X) \right] &= \int_{\mathbb{R}^d} \frac{\partial \mathbb{E}[Y|X=x]}{\partial x} f^2(x) dx \\ &= -2 \int_{\mathbb{R}^d} \mathbb{E}[Y|X=x] f(x) \frac{\partial f(x)}{\partial x} dx \\ &= -2 \mathbb{E} \left[\mathbb{E}[Y|X] \frac{\partial f(X)}{\partial x} \right] \\ &= -2 \mathbb{E} \left[Y \frac{\partial f(X)}{\partial x} \right]. \end{aligned}$$

Последнее выражение в правой части является основой для оценки δ : заменив популяционное математическое ожидание выборочным, а плотность f ядерной оценкой \hat{f} согласно равенству (3), получим:

$$\hat{\delta} = \frac{1}{n} \sum_{i=1}^n Y_i \frac{\partial \hat{f}(X_i)}{\partial x}.$$

Преимущество $\hat{\delta}$ над $\hat{\beta}$ в том, что первая оценка имеет явный вид и не требует численной оптимизации.

Можно распространить одноиндексную модель на более общий класс моделей:

$$Y = g(v(X; \beta_0)) + \varepsilon, \quad \mathbb{E}[\varepsilon|X] = 0,$$

для некоторой функции $v: \mathbb{R} \times \mathcal{B} \mapsto \mathbb{R}$, известной с точностью до β_0 . Методика оценивания, изложенная выше в общих чертах, применима и для этого более общего случая.

2.2 Частично линейная модель

Если предположить, что m из (1) линейна по некоторым своим аргументам, то можно получить другую спецификацию. Положим $X = (X_1, X_2)$, где $X_i \in \mathbb{R}^{d_i}$, $i = 1, 2$, и $d = d_1 + d_2$, так что

$$Y = \beta'_0 X_1 + g(X_2) + \varepsilon, \quad \mathbb{E}[\varepsilon|X] = 0, \quad (7)$$

для некоторых $g : \mathbb{R}^{d_2} \mapsto \mathbb{R}$ и $\beta \in \mathbb{R}^{d_1}$. Как и раньше, мы не специфицируем $\varepsilon|X$.

По сравнению с общей регрессионной моделью из уравнения (1), мы налагаем следующее ограничение на вид регрессионной функции: $m(x) = \beta'_0 x_1 + g(x_2)$. То есть Y аддитивна по X_1 и X_2 , причем Y зависит от X_1 линейно. Наша модель включает параметрическую компоненту β_0 и две непараметрические компоненты g и $f_{\varepsilon|X}$, то есть является полупараметрической моделью.

Опять же, мы должны наложить ограничения на модель, чтобы g и β были идентифицируемы. Среди компонент X не должно быть константы, так как при $\tilde{g}(x_2) = g(x_2) - a$, $a \in \mathbb{R}$, мы не сможем различить $\beta'x_1 + g(x_2)$ и $(a + \beta'x_1) + \tilde{g}(x_2)$. В действительности мы должны предположить, что матрица

$$\Omega = \mathbb{E}[(X_1 - \mathbb{E}[X_1|X_2])(X_1 - \mathbb{E}[X_1|X_2])']$$

невырождена. Если это условие не выполняется, мы не сможем различать линейный и нелинейный члены. Чтобы показать это, рассмотрим

$$\mathbb{E}[Y|X_2] = \beta'_0 \mathbb{E}[X_1|X_2] + g(X_2) + \mathbb{E}[\varepsilon|X_2],$$

что означает

$$Y - \mathbb{E}[Y|X_2] = \beta'_0 (X_1 - \mathbb{E}[X_1|X_2]) + \eta, \quad (8)$$

где $\eta = \varepsilon - \mathbb{E}[\varepsilon|X_2]$ удовлетворяет условию $\mathbb{E}[\eta|X] = 0$. Таким образом, чтобы идентифицировать β_0 , нужно, чтобы Ω была невырождена.

Уравнение (8) является базовым для оценки, основанной на остатках: строим ядерные оценки для $m_Y(x_2) = \mathbb{E}[Y|X_2 = x_2]$ и $m_{X_1}(x_2) = \mathbb{E}[X_1|X_2 = x_2]$,

$$\hat{m}_Y(x_2) = \frac{\sum_{i=1}^n Y_i K_h(X_{2,i} - x_2)}{\sum_{i=1}^n K_h(X_{2,i} - x_2)}, \quad \hat{m}_{X_1}(x_2) = \frac{\sum_{i=1}^n X_{1,i} K_h(X_{2,i} - x_2)}{\sum_{i=1}^n K_h(X_{2,i} - x_2)},$$

и подставляем их в (8). После этого мы можем оценить β с помощью МНК,

$$\hat{\beta} = \left(\sum_{i=1}^n \hat{Z}_i \hat{Z}_i' \right)^{-1} \sum_{i=1}^n \hat{Z}_i (Y_i - \hat{m}_Y(X_{2,i}))', \quad (9)$$

где $\hat{Z}_i = X_{1,i} - \hat{m}_{X_1}(X_{2,i})$.

Этот метод оценивания можно расширить на следующую более общую модель:

$$Y = v(X_1; \beta) + g(X_2) + \varepsilon, \quad \mathbb{E}[\varepsilon|X] = 0, \quad (10)$$

где $v : \mathbb{R}^{d_1} \times \mathcal{B} \mapsto \mathbb{R}$ известна с точностью до $\theta \in \Theta$. Однако полученная оценка имеет неявный вид, и приходится применять численные методы оптимизации.

3 Спецификация распределения ошибок

До этого момента мы обсуждали только как можно смоделировать и оценить функциональный вид m в общей регрессионной модели с помощью полупараметрических методов. В этом разделе мы сосредоточимся на ошибке ε и обсудим, как разные предположения относительно ошибки приводят к разным (полупараметрическим) методикам оценивания регрессионной функции. В некоторых случаях можно вывести оценку интересного нам параметра, не оценивая бесконечномерные компоненты. Однако эти оценки неэффективны, и для улучшения их эффективности можно использовать методы полупараметрического оценивания.

3.1 Линейная регрессионная модель

Рассмотрим стандартную линейную регрессионную модель:

$$Y = \beta'X + \varepsilon, \quad (11)$$

где $\mathbb{E}[\varepsilon|X] = 0$. Обычно она считается полностью параметрической моделью, но в нашей терминологии это полупараметрическая модель, так как распределения $\varepsilon|X$, $f_{\varepsilon|X}$, определены не полностью. Если $f_{\varepsilon|X}$ не задана, мы имеем параметрическую компоненту, θ , и непараметрическую компоненту, $f_{\varepsilon|X}$.

Если предположить, что ошибка имеет нормальное распределение, то, как было показано в предыдущем разделе, ММП-оценка превращается в обычную МНК-оценку как в уравнении (2). Однако МНК-оценку можно рассматривать как полупараметрическую оценку θ , так как она остается \sqrt{n} -состоятельной независимо от точности спецификации $f_{\varepsilon|X}$. Более того, привлекательное свойство МНК состоит в том, что для его применения не нужно оценивать $f_{\varepsilon|X}$ в отличие от полупараметрических оценок, рассмотренных в предыдущем разделе, где нужно было получить предварительную оценку непараметрической компоненты, чтобы затем оценить параметрический компонент.

Однако читателю может быть интересно, существуют ли другие оценки, получше. Очевидно, если мы наложим (правильную) параметрическую структуру на $f_{\varepsilon|X}$, мы можем использовать ММП, который вообще более эффективен, чем МНК. Но даже если мы не налагаем никакой формы на распределение, то, как мы увидим далее, МНК-оценка в общем случае неэффективна в классе полупараметрических оценок.

3.2 Гетероскедастичность неизвестной формы

Мы рассматриваем линейную модель (11), но теперь предполагаем, что ошибки гетероскедастичны,

$$\mathbb{E}[\varepsilon^2|X = x] = \sigma^2(x), \quad (12)$$

где условная функция дисперсии, $\sigma^2(\cdot)$, *неизвестна*.

Обычная МНК-оценка (2) по-прежнему состоятельна и асимптотически нормально распределена, но теперь ее асимптотическое распределение следующее:

$$\sqrt{n}(\hat{\theta}_{\text{OLS}} - \theta) \rightarrow^d N\left(0, \mathbb{E}[XX']^{-1} \mathbb{E}[\sigma^2(X) XX']^{-1} \mathbb{E}[XX']^{-1}\right).$$

В частности, она уже не эффективна, как мы далее увидим. Сперва рассмотрим случай, когда условная дисперсионная функция $\sigma^2(x)$ *известна*. Тогда мы можем применить взвешенный метод наименьших квадратов (ВМНК),

$$\tilde{\theta}_{\text{WLS}} = \left(\sum_{i=1}^n \sigma^{-2}(X_i) X_i X_i'\right)^{-1} \left(\sum_{i=1}^n \sigma^{-2}(X_i) X_i Y_i\right), \quad (13)$$

который дает меньшую асимптотическую дисперсию оценки по сравнению с МНК:

$$\sqrt{n}(\tilde{\theta}_{\text{WLS}} - \theta) \rightarrow^d N\left(0, \mathbb{E}[\sigma^{-2}(X) XX']^{-1}\right),$$

где

$$\mathbb{E}[\sigma^{-2}(X) XX']^{-1} \leq \mathbb{E}[XX']^{-1} \mathbb{E}[\sigma^2(X) XX']^{-1} \mathbb{E}[XX']^{-1}$$

которое превращается в равенство тогда и только тогда, когда $\sigma^2(X) = \sigma^2 = \mathbb{E}[\varepsilon^2]$ константа почти наверное.

Если функция условной дисперсии $\sigma^2(x)$ неизвестна, то $\tilde{\theta}_{\text{WLS}}$ недостижима. В этом случае можно положить, что $\sigma^2(x)$ зависит от параметров, и оценить эти параметры, используя стандартные методы. Однако в этом случае нужно правильно специфицировать функциональный вид условной дисперсии. Чтобы избежать неправильной спецификации, следует использовать непараметрическую оценку $\sigma^2(x)$. Чтобы мотивировать использование оценки, заметим сперва, что $\sigma^2(x)$ по определению является обычным условным средним ε^2 , ср. (12). Естественной оценкой условного среднего является ядерная регрессионная оценка, представленная в (4). В идеале хотелось бы вычислить $\hat{\sigma}^2(x) = \sum_{i=1}^n \varepsilon_i^2 K_h(X_i - x) / \sum_{i=1}^n K_h(X_i - x)$. Однако, так как ε_i , $i = 1, \dots, n$, ненаблюдаемы, мы заменяем их остатками. Это приводит к следующей трехшаговой процедуре:

1. Рассчитать МНК-оценку, $\hat{\theta}_{\text{OLS}}$, по формуле (2).
2. Рассчитать остатки, $\hat{\varepsilon}_i = Y_i - \hat{\theta}'_{\text{OLS}} X_i$, $i = 1, \dots, n$, и с их помощью непараметрически оценить условную дисперсию:

$$\hat{\sigma}^2(x) = \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2 K_h(X_i - x)}{\sum_{i=1}^n K_h(X_i - x)}.$$

3. Получить ВМНК-оценку по формуле (13), но вместо $\sigma^2(x)$ подставить $\hat{\sigma}^2(x)$:

$$\hat{\theta}_{\text{WLS}} = \left(\sum_{i=1}^n \hat{\sigma}^{-2}(X_i) X_i X_i' \right)^{-1} \left(\sum_{i=1}^n \hat{\sigma}^{-2}(X_i) X_i Y_i \right), \quad (14)$$

И вновь этот метод оценивания можно обобщить для более сложной параметрической модели

$$Y = g(X; \theta) + \varepsilon, \quad \mathbb{E}[\varepsilon|X] = 0,$$

где $g : \mathbb{R}^d \times \Theta \mapsto \mathbb{R}$ известна с точностью до $\theta \in \Theta$.

3.3 Предположение о независимости

Можно использовать изложенную выше идею, чтобы получить ММП-оценки в случае, когда распределение ошибок имеет неизвестный вид. Мы сохраняем линейную спецификацию (11), но теперь предполагаем, что

ε и X независимы,

так что $f_{\varepsilon|X}(\varepsilon|X) = f_{\varepsilon}(\varepsilon)$, где

$$\mathbb{E}[\varepsilon] = \int_{\mathbb{R}} z f_{\varepsilon}(z) dz = 0, \quad \sigma^2 = \int_{\mathbb{R}} z^2 f_{\varepsilon}(z) dz < \infty.$$

По сравнению с предыдущим разделом мы ввели дополнительное предположение о независимости регрессоров и ошибок. Однако мы не предполагаем, что известно распределение ε , так что и в этом случае модель остается полупараметрической.

Предположение о независимости позволяет оценить параметрическую компоненту с помощью полупараметрической ММП-оценки: допустим, плотность f_{ε} известна, тогда можно использовать ММП-оценку

$$\tilde{\theta}_{\text{MLE}} = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_i \log f_{\varepsilon}(Y_i - \theta' X_i), \quad (15)$$

которая при условии регулярности будет удовлетворять

$$\sqrt{n}(\tilde{\theta}_{\text{MLE}} - \theta) \rightarrow^d N(0, H_0^{-1}),$$

где

$$H_0 = \mathbb{E} \left[\frac{\partial \log f_\varepsilon(Y_i - \theta' X_i)}{\partial \theta} \frac{\partial \log f_\varepsilon(Y_i - \theta' X_i)}{\partial \theta'} \right] = \int \frac{f'_\varepsilon(z)^2}{f_\varepsilon(z)} dz \mathbb{E}[X X'].$$

Однако плотность f_ε неизвестна, и, следовательно, $\tilde{\theta}_{\text{MLE}}$ недоступна. С другой стороны, заметим, что МНК по-прежнему применим и дает состоятельную оценку. Однако МНК-оценка будет неэффективна по сравнению с ММП-оценкой, так как $\int f'_\varepsilon(z)^2 / f_\varepsilon(z) dz \leq \sigma^2$, что превращается в равенство тогда и только тогда, когда f_ε – плотность нормального распределения $N(0, \sigma^2)$.

Чтобы улучшить эффективность ММП, мы собираемся получить полупараметрическую версию ММП-оценки с помощью трехшаговой процедуры:

1. Рассчитать МНК-оценку, $\hat{\theta}_{\text{OLS}}$, по формуле (2).
2. Рассчитать остатки, $\hat{\varepsilon}_i = Y_i - \hat{\theta}'_{\text{OLS}} X_i$, $i = 1, \dots, n$, и с их помощью непараметрически оценить плотность безусловного распределения f_ε , например, как

$$\hat{f}_\varepsilon(x) = \frac{1}{n} \sum_{i=1}^n K_h(\hat{\varepsilon}_i - x). \quad (16)$$

3. Получить ММП-оценку по формуле (15), но вместо f_ε подставить \hat{f}_ε ,

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log \hat{f}_\varepsilon(Y_i - \theta' X_i). \quad (17)$$

И вновь данный метод оценивания можно обобщить для более сложной параметрической модели. Предположим, например, что

$$Y = g(X; \theta) + \sigma(X; \theta) \varepsilon,$$

где $g, \sigma : \mathbb{R}^d \times \Theta \mapsto \mathbb{R}$ известны с точностью до $\theta \in \Theta$, ε и X независимы, и вдобавок

$$\mathbb{E}[\varepsilon] = \int_{\mathbb{R}} z f_\varepsilon(z) dz = 0, \quad \mathbb{E}[\varepsilon^2] = \int_{\mathbb{R}} z^2 f_\varepsilon(z) dz = 1.$$

Предположим, мы получили предварительную оценку θ , например ММП-оценку при нормальных ошибках, $\hat{\theta}_{\text{QMLE}}$, которая остается состоятельной даже если ошибки не нормально распределены. Тогда мы можем рассчитать остатки

$$\hat{\varepsilon}_i = \frac{Y_i - g(X_i; \hat{\theta}_{\text{QMLE}})}{\sigma(X_i; \hat{\theta}_{\text{QMLE}})}, \quad i = 1, \dots, n,$$

и затем непараметрически оценить плотность по формуле (16). На заключительном этапе мы определим

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \left\{ \log \hat{f}_\varepsilon \left(\frac{Y_i - g(X_i; \theta)}{\sigma(X_i; \theta)} \right) + \log(\sigma(X_i; \theta)) \right\}.$$

Мы ожидаем, что $\hat{\theta}$ будет наиболее эффективной оценкой, в отличие от $\hat{\theta}_{\text{QMLE}}$.

4 Копулы

Чтобы показать, что полупараметрическое моделирование имеет приложения за пределами регрессионного подхода, рассмотрим последний пример, касающийся копул. Копулы оказались удобным инструментом при моделировании ситуаций с многомерной зависимостью. В частности, их применяют в финансах, см., например, Genest, Gendron & Bourdeau-Brien (2009). В этом разделе мы рассмотрим полупараметрическое семейство копул и соответствующие оценки.

Пусть $Z = (Z_1, Z_2) \in \mathbb{R}^2$ – двумерная непрерывная случайная величина. Обозначим плотность функции совместного распределения (ФПР) и кумулятивную функцию распределения (КФР), соответственно, за f и F :

$$\Pr \{Z_1 \leq z_1, Z_2 \leq z_2\} = F(z_1, z_2) = \int_{-\infty}^{z_2} \int_{-\infty}^{z_1} f(v_1, v_2) dv_1 dv_2,$$

и пусть f_k и F_k – маргинальные ФПР и КФР Z_k соответственно:

$$\Pr \{Z_k \leq z\} = F_k(z) = \int_{-\infty}^z f_k(v) dv, \quad k = 1, 2.$$

Так называемые копулы используют, чтобы смоделировать структурную зависимость между Z_1 и Z_2 . При этом используется следующий известный факт: существует единственная функция $C : [0, 1]^2 \mapsto [0, 1]$, такая что

$$F(z_1, z_2) = C(F_1(z_1), F_2(z_2)),$$

(см. Joe, 1997). Функция C называется *копулой* для Z . Легко видеть, что C – КФР равномерно распределенной случайной величины $U := (F_1(Z_1), F_2(Z_2))$:

$$C(u_1, u_2) = \Pr \{F_1(Z_1) \leq u_1, F_2(Z_2) \leq u_2\}.$$

Кроме того, плотность совместного распределения Z можно записать в виде

$$f(z_1, z_2) = c(F_1(z_1), F_2(z_2)) f_1(z_1) f_2(z_2),$$

где $c : [0, 1] \times [0, 1] \mapsto \mathbb{R}_+$ – функция плотности распределения U .

Теперь можно смоделировать совместное распределение Z , задав два маргинальных распределения и копулу. В полностью параметрическом подходе это можно сделать, например, по формуле

$$f(z_1, z_2; \xi) = c(F_1(z_1; \alpha_1), F_2(z_2; \alpha_1); \theta) f_1(z_1; \alpha_1) f_2(z_2; \alpha_1),$$

где $\xi = (\theta', \alpha'_1, \alpha'_2)' \in \Theta \times \mathcal{A}_1 \times \mathcal{A}_1$ – конечномерный параметр. После чего можно оценить ξ с помощью ММП:

$$\hat{\xi} = \arg \max_{\xi \in \Xi} \frac{1}{n} \sum_{i=1}^n \{\log c(F_1(Z_{1,i}; \alpha_1), F_2(Z_{2,i}; \alpha_1); \theta) + \log f_1(Z_{1,i}; \alpha_1) + \log f_2(Z_{2,i}; \alpha_1)\}.$$

При большой размерности ξ это может оказаться сложной численной задачей. Вместо этого можно оценить параметры с помощью двухшаговой процедуры оценивания: сперва оценим $(\alpha'_1, \alpha'_2)'$ как

$$\hat{\alpha}_k = \arg \max_{\alpha_k \in \mathcal{A}_k} \frac{1}{n} \sum_{i=1}^n \log f_k(Z_{k,i}; \alpha_k), \quad k = 1, 2,$$

а затем –

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log c(F_1(Z_{1,i}; \hat{\alpha}_1), F_2(Z_{2,i}; \hat{\alpha}_2); \theta).$$

Эта двухшаговая оценка менее эффективна по сравнению с полной ММП-оценкой, но ее проще реализовать.

Тривиальная полупараметрическая модель с копулой имеет следующий вид. Мы, как и раньше, задаем параметрическое семейство копул, $c(u_1, u_2; \theta)$, но теперь оставляем оба маргинальных распределения незадаанными. Мы хотим непараметрически оценить маргинальные распределения и на их основе сделать выводы относительно θ . Пусть

$$\hat{F}_k(z_k) = \frac{1}{n} \sum_{i=1}^n 1\{Z_{k,i} \leq z_k\}, \quad k = 1, 2,$$

– выборочная кумулятивная функция распределения. Тогда естественной оценкой θ будет

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log c(\hat{F}_1(Z_{1,i}), \hat{F}_2(Z_{2,i}); \theta).$$

5 Класс двухшаговых оценок

В предыдущих двух разделах мы рассмотрели несколько примеров полупараметрических моделей и вывели оценки интересующих нас параметров. В этом разделе мы исследуем методы, с помощью которых можно проанализировать асимптотические свойства этих оценок. В частности, мы выясним, при каких условиях оценки будут \sqrt{n} -состоятельны и асимптотически нормально распределены.

Мы начнем с представления основного класса полупараметрических двухшаговых оценок. На первом шаге вычисляется предварительная непараметрическая оценка. На втором шаге эта непараметрическая оценка подставляется в целевую функцию, которая минимизируется с целью получить оценку параметрической компоненты. Этот достаточно большой класс оценок среди прочих включает в себя оценки, рассмотренные в предыдущих разделах. Мы дадим общие условия для состоятельности и асимптотической нормальности оценок параметрических компонент при соответствующих условиях регулярности.

Наша задача оценивания имеет много общего со стандартной параметрической двухшаговой задачей оценивания, где сперва оценивается шумовой параметр, а затем эту оценку используют для получения оценки другого, интересующего нас, параметра. Единственное отличие заключается в том, что в нашем случае предварительная оценка является функцией, бесконечномерным параметром. Тем не менее, с небольшими видоизменениями можно использовать идею доказательства для параметрической двухшаговой оценки.

5.1 Структура

Мы хотим оценить конечномерный параметр $\theta \in \Theta \subseteq \mathbb{R}^k$ с помощью случайной целевой функции $Q_n(\theta, \gamma)$, где $\gamma \in \Gamma$ – некоторый бесконечномерный параметр, чаще всего функция. Целевая функция как правило будет являться функцией наблюдаемых данных, (Y_i, X_i) , где $i = 1, \dots, n$, однако мы не будем явно задавать эту зависимость и обозначим ее лишь с помощью индекса n . Мы предполагаем, что пространство параметров Γ – линейное пространство с нормой $\|\cdot\|$. Это может быть, например, супремум-норма, $\|\gamma\| = \sup_x |\gamma(x)|$, или норма L_q , $\|\gamma\| = (\int |\gamma(x)|^q w(x) dx)^{1/q}$ для некоторой функции весов $w(x) \geq 0$.

Если известно истинное значение параметра γ , которое мы обозначим за γ_0 , то можно оценить θ как

$$\tilde{\theta} = \arg \min_{\theta \in \Theta} Q_n(\theta, \gamma_0). \quad (18)$$

В этом случае, чтобы вывести асимптотические свойства $\tilde{\theta}$, можно использовать стандартный результат для параметрических оценок (см., например, Newey & McFadden, 1994).

Мы рассмотрим случай, когда γ_0 неизвестна, и, следовательно, $\tilde{\theta}$ недостижима. Однако предположим, что доступна предварительная оценка, $\hat{\gamma}$. Тогда, заменяя γ_0 на $\hat{\gamma}$, мы можем записать

$$\hat{\theta} = \arg \min_{\theta \in \Theta} Q_n(\theta, \hat{\gamma}). \quad (19)$$

Назовем $\hat{\theta}$ полупараметрической двухшаговой оценкой.

Прежде всего мы сделаем минимальные предположения относительно вида $Q_n(\theta, \gamma)$ и $\hat{\gamma}$ и потребуем лишь, чтобы последняя была состоятельной оценкой γ_0 и достаточно быстро сходилась. В большинстве случаев целевая функция имеет вид

$$Q_n(\theta, \gamma) = \frac{1}{n} \sum_{i=1}^n q(Z_i; \theta, \gamma), \quad (20)$$

но мы не будем ограничиваться только такими функциями.

Прежде чем перейти к анализу общей двухшаговой оценки, сперва покажем, как с помощью подходящего выбора $q(z; \theta, \gamma)$ можно представить оценки из предыдущих разделов в виде (19)–(20):

Пример 1: Одноиндексная модель. При $\gamma = g$ оценку для этой модели можно записать в виде (19)–(20), где q имеет вид

$$q(z; \theta, \gamma) = [y - \gamma(\theta'x)]^2,$$

а оценку $\hat{\gamma}_\theta$ можно выбрать как

$$\hat{\gamma}_\theta(z) = \frac{\sum_{i=1}^n Y_i K_h(\theta'X_i - z)}{\sum_{i=1}^n K_h(\theta'X_i - z)}.$$

Пример 2: Частично линейная модель. В этом случае, чтобы оценка приняла желаемый вид, надо взять

$$q(z; \theta, \gamma) = [y - \gamma_1(x_2) - \theta'(x_1 - \gamma_2(x_2))]^2,$$

где $\gamma_1(x_2) = \mathbb{E}[Y|X_2 = x_2]$, $\gamma_2(x_2) = \mathbb{E}[X_1|X_2 = x_2]$. Предварительными оценками будут

$$\hat{\gamma}_1(x_2) = \frac{\sum_{i=1}^n Y_i K_h(X_{2,i} - x_2)}{\sum_{i=1}^n K_h(X_{2,i} - x_2)}, \quad \hat{\gamma}_2(x_2) = \frac{\sum_{i=1}^n X_{1,i} K_h(X_{2,i} - x_2)}{\sum_{i=1}^n K_h(X_{2,i} - x_2)}.$$

Пример 3: Эффективное оценивание при гетероскедастичности. Определим функцию q как

$$q(z; \theta, \gamma) = \gamma^{-1}(x) [y - \theta'x]^2, \quad (21)$$

где $\gamma(x) = \sigma^2(x)$. ВМНК-оценка является частным случаем двухшаговой оценки. Здесь предварительными оценками будут

$$\hat{\gamma}(x_2) = \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2 K_h(X_i - x)}{\sum_{i=1}^n K_h(X_i - x)}. \quad (22)$$

Пример 4: Полупараметрические копулы. Функция q , определяющая оценку параметра копулы θ , задается формулой

$$q(z; \theta, \gamma) = \log c(\gamma_1(z_1), \gamma_2(z_2); \theta),$$

где $\gamma_k(z) = F_k(z)$ – маргинальная кумулятивная функция распределения Z_k , $k = 1, 2$, которая оценивается как

$$\hat{\gamma}_k(z) = \frac{1}{n} \sum_{i=1}^n 1\{Z_{k,i} \leq z\}, \quad k = 1, 2.$$

В двух из этих примеров, а именно в частично линейной модели и регрессионной модели с неизвестной гетероскедастичностью, можно вывести в явном виде выражение для $\hat{\theta}$. То есть, можно провести прямой анализ этих конкретных оценок, и он будет гораздо проще по сравнению с косвенным анализом, предложенным ниже. Но в общем случае явные выражения для оценок недоступны, и все сводится к анализу свойств целевой функции $Q_n(\theta, \gamma)$.

В рамках этого подхода мы сперва установим общие условия, при которых $\hat{\theta}$ состоятельна и сходится к нормальному распределению. Налагая ограничения на вид целевой функции $Q_n(\theta, \gamma)$ и оценку $\hat{\gamma}$, мы в общих чертах покажем, как эти условия можно проверить в том конкретном случае, когда $\hat{\gamma}$ – ядерная оценка. Наконец, мы выясним асимптотические свойства первого порядка ВМНК-оценки, определенной в разделе 3.2, проверив выполнение общих условий для этой оценки.

5.2 Состоятельность

Доказательство состоятельности более или менее идентично доказательству в случае параметрических двухшаговых оценок. Единственное концептуальное отличие заключается в том, что теперь мы работаем с бесконечномерным параметром. Мы наложим следующие ограничения на целевую функцию:

С.1 Существует функция $Q(\theta, \gamma)$, такая что $\sup_{\theta \in \Theta} |Q_n(\theta, \gamma_0) - Q(\theta, \gamma_0)| \xrightarrow{P} 0$.

С.2 Для всех $\varepsilon > 0$, $\inf_{\|\theta - \theta_0\| > \varepsilon} Q(\theta, \gamma_0) > Q(\theta_0, \gamma_0)$.

С.3 Для некоторого $\lambda > 0$ и $B_n = O_P(1)$,

$$\sup_{\theta \in \Theta} |Q_n(\theta, \gamma) - Q_n(\theta, \gamma_0)| \leq B_n \|\gamma - \gamma_0\|^\lambda$$

для γ из окрестности γ_0 .

Условие (С.1) означает, что недостижимая конечномерная целевая функция, $Q_n(\theta, \gamma_0)$, имеет определенный предел, $Q(\theta, \gamma_0)$. (С.2) – условие идентификации, то есть того, что целевая функция имеет единственный минимум $\theta_0 = \arg \min_{\theta \in \Theta} Q(\theta, \gamma_0)$. Легко показать, что условие (С.2) следует из следующих трех условий: Θ компакт, $\theta \mapsto Q(\theta, \gamma_0)$ непрерывна, и $Q(\theta, \gamma_0) > Q(\theta_0, \gamma_0)$ для всех $\theta \neq \theta_0$. Более примитивные условия для (С.1) можно найти в Newey (1991). Условия (С.1)–(С.2) означают, что недостижимая оценка $\tilde{\theta}$, заданная уравнением (18), состоятельна, $\tilde{\theta} \xrightarrow{P} \theta_0$; см., например, Newey & McFadden (1994, Теорема 2.1).

Последнее условие (С.3) утверждает, что разница между двумя целевыми функциями, $Q_n(\theta, \hat{\gamma})$ и $Q_n(\theta, \gamma_0)$, асимптотически пренебрежима: $Q_n(\theta, \hat{\gamma}) \xrightarrow{P} Q_n(\theta, \gamma_0)$ при $\hat{\gamma} \xrightarrow{P} \gamma_0$. Заметим, что норма $\|\gamma - \gamma_0\|$ функциональная, что обсуждалось в начале этого раздела. При этих предположениях достижимая оценка сходится к недостижимой, $\hat{\theta} = \tilde{\theta} + o_P(1)$.

Условия (С.1) и (С.3) можно заменить на следующие два условия:

С.1' $\sup_{\theta \in \Theta, \|\gamma - \gamma_0\| < \delta} |Q_n(\theta, \gamma) - Q(\theta, \gamma)| \xrightarrow{P} 0$, $\delta > 0$.

С.3' $\sup_{\theta \in \Theta} |Q(\theta, \gamma) - Q(\theta, \gamma')| \rightarrow 0$ при $\gamma \rightarrow \gamma'$.

Для проверки условий (С.1') и (С.3') можно воспользоваться теорией эмпирических процессов, см. Andrews (1994a,b), Chen, Linton & van Keilegom (2003), van der Vaart & Wellner (1996). Эта проверка обычно включает в себя проверку условия Липшица, выраженное в (С.3).

Итоговые результаты о состоятельности сформулированы в следующих теоремах:

Теорема 1 Пусть $Q_n(\theta, \gamma)$ удовлетворяет (С.1)–(С.3). Если $\hat{\gamma} \in \Gamma$ с некоторого момента по мере $\hat{\gamma} \xrightarrow{P} \gamma_0$, тогда $\hat{\theta} \xrightarrow{P} \theta_0$.

Доказательство Доступно по запросу к автору.

Замечание 2 В случае, когда $\hat{\gamma}$ зависит от θ , нужно усилить условие состоятельности до $\sup_{\theta \in \Theta} \|\hat{\gamma}_\theta - \gamma_\theta\| \xrightarrow{P} 0$.

Теперь проверим условия (С.1)–(С.3) для ВМНК-оценки:

Пример 3 (продолжение). Обозначим как σ_0^2 и θ_0 истинные значения параметров. Мы предполагаем, что $X \in \mathcal{X}$, где $\mathcal{X} \subset \mathbb{R}^d$ компакт, и $\mathbb{E}[Y^2] < \infty$. Для идентификации нужно, чтобы $\mathbb{E}[X X' \sigma_0^{-2}(X)]$ была невырожденной. Предположение о компактности носителя \mathcal{X} можно снять, но тогда нужно будет ввести подравнивание оценки, см. Robinson (1987).

Также предположим, что $\sigma_0^2(x), f(x) > 0$ дважды непрерывно дифференцируемы. В частности, $\underline{\sigma}^2 := \inf_{x \in \mathcal{X}} \sigma_0^2(x) > 0$. Введем ограничение, что Θ компакт, то есть существует константа c , такая что $|\theta'x| \leq c$ для любых $\theta \in \Theta$ и $x \in \mathcal{X}$. Определим норму σ^2 как $\|\sigma^2\|_\infty = \sup_{x \in \mathcal{X}} |\sigma^2(x)|$ и предположим, что выполнено

$$\|\hat{\sigma}^2 - \sigma_0^2\|_\infty \xrightarrow{P} 0, \quad (23)$$

где $\hat{\sigma}^2$ – ядерная оценка из формулы (22); это, например, можно сделать, используя результаты в Kristensen (2009).

Сперва мы покажем (С.3). Целевая функция принимает вид (21). Разложение в ряд Тейлора до первого порядка малости функции $a \mapsto 1/a$ дает

$$\begin{aligned} q(z; \theta, \hat{\sigma}^2) - q(z; \theta, \sigma_0^2) &= [y - \theta'x]^2 \left(\frac{1}{\hat{\sigma}^2(x)} - \frac{1}{\sigma_0^2(x)} \right) \\ &= [y - \theta'x]^2 \frac{-1}{[\lambda_x \hat{\sigma}^2(x) + (1 - \lambda_x) \sigma_0^2(x)]^2} (\hat{\sigma}^2(x) - \sigma_0^2(x)) \end{aligned}$$

для некоторого $\lambda_x \in [0, 1]$. В силу (23) $\inf_{x \in \mathcal{X}} \hat{\sigma}^2(x) \geq \underline{\sigma}^2/2$ почти наверное с некоторого момента по мере $n \rightarrow \infty$. Следовательно,

$$\frac{1}{[\lambda_x \hat{\sigma}^2(x) + (1 - \lambda_x) \sigma_0^2(x)]^2} \leq \frac{1}{[\lambda_x \underline{\sigma}^2/2 + (1 - \lambda_x) \underline{\sigma}^2/2]^2} = \frac{4}{\underline{\sigma}^4} < \infty.$$

Также, так как $\theta \mapsto (Y - \theta'X)^2$ непрерывна и $(Y - \theta'X)^2 \leq 3Y^2 + 3c^2$, где $\mathbb{E}[Y^2] < \infty$, из известных фактов о равномерной сходимости (см., например, Newey, 1991) следует, что

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n (Y_i - \theta'X_i)^2 - \mathbb{E}[(Y_i - \theta'X_i)^2] \right| \xrightarrow{P} 0$$

и $\sup_{\theta \in \Theta} E[(Y - \theta'X)^2] < \infty$. Следовательно, $\sup_{\theta \in \Theta} n^{-1} \sum_{i=1}^n (Y_i - \theta'X_i)^2 = O_P(1)$. Далее перепишем

$$\begin{aligned} \sup_{\theta \in \Theta} |Q_n(\theta, \hat{\sigma}^2) - Q_n(\theta, \sigma_0^2)| &\leq \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n |q(Z_i; \theta, \hat{\sigma}^2) - q(Z_i; \theta, \sigma_0^2)| \\ &\leq \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n [Y_i - \theta'X_i]^2 \times \sup_x \left| \frac{1}{\hat{\sigma}^2(x)} - \frac{1}{\sigma_0^2(x)} \right|, \end{aligned}$$

где

$$\sup_x \left| \frac{1}{\hat{\sigma}^2(x)} - \frac{1}{\sigma_0^2(x)} \right| \leq \frac{4 \|\hat{\sigma}^2 - \sigma_0^2\|_\infty}{\sigma^4} = o_P(1).$$

Следовательно, (С.3) выполняется для

$$B_n = \frac{4}{\sigma^2} \times \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n [Y_i - \theta'X_i]^2,$$

и $\lambda = 1$.

Теперь проверим (С.1). С помощью еще одного приложения равномерного закона больших чисел легко показать, что

$$\sup_{\theta \in \Theta} |Q_n(\theta, \sigma_0^2) - Q(\theta, \sigma_0^2)| \xrightarrow{P} 0,$$

где $\theta \mapsto Q(\theta, \sigma_0^2) = \mathbb{E} \left[[Y - \theta'X]^2 \sigma_0^{-2}(X) \right]$ непрерывна.

Наконец, чтобы проверить (С.2), заметим, что для любых $\theta \neq \theta_0$,

$$\begin{aligned} Q(\theta, \sigma_0^2) &= \mathbb{E} \left[[(\theta_0 - \theta)'X + \varepsilon]^2 \sigma_0^{-2}(X) \right] \\ &= (\theta_0 - \theta)' \mathbb{E} [X X' \sigma_0^{-2}(X)] (\theta_0 - \theta) + \mathbb{E} [\varepsilon^2 \sigma_0^{-2}(X)] \\ &> \mathbb{E} [\varepsilon^2 \sigma_0^{-2}(X)] \\ &= Q(\theta_0, \sigma_0^2). \end{aligned}$$

Согласно замечаниям после условий (С.1)–(С.3), отсюда следует (С.2). Мы проверили все условия и согласно теореме 1 $\hat{\theta}$ состоятельна.

5.3 Асимптотическая нормальность

Для доказательства асимптотической нормальности мы воспользуемся тем же принципом, что и в случае параметрических двухшаговых оценок. Можно было бы разложить оценку, полученную на первом шаге, в ряд Тейлора, таким образом принимая во внимание дополнительную выборочную ошибку из-за оценивания на первом шаге. Однако, в наших условиях оценка первого шага является функцией, то есть бесконечномерным объектом. Поэтому, чтобы воспользоваться тем же методом, что и для параметрических двухшаговых оценок, нам сперва нужно обобщить понятие производных с конечномерного случая на бесконечномерный.

Пусть $T : \Gamma \mapsto \mathbb{R}^d$ – функционал, который ставит в соответствие $\gamma \in \Gamma$ некоторый евклидовый вектор. Например, $T(\gamma) = \int \gamma(x) dx$, $T_x(\gamma) = \partial \gamma(x) / \partial x$ и $T_x(\gamma_1, \gamma_2) = \gamma_1(x) \gamma_2(x)$.

Определение 3 *Говорят, что T дифференцируема по направлению в точке $\gamma \in \Gamma$, если существует линейный непрерывный функционал $\dot{T}(\gamma) [\cdot] : \Gamma \mapsto \mathbb{R}^d$ такой что*

$$\dot{T}(\gamma) [h] = \lim_{t \rightarrow 0} \frac{T(\gamma + th) - T(\gamma)}{t},$$

для всех $h \in \Gamma$.

\dot{T} называется *производной по направлению* от T . В конечномерном случае, если T дифференцируема с производной $\partial T(\gamma)/\partial\gamma$, тогда $\dot{T}(\gamma)[h]$ – дифференциал T :

$$\dot{T}(\gamma)[h] = \frac{\partial T(\gamma)}{\partial\gamma} h.$$

Как правило, можно перенести результаты с конечномерного случая, используя производную по направлению. В частности, по-прежнему работает правило дифференцирования сложной функции.

Примеры функционалов. (i) $\Gamma = \{\gamma : \int |\gamma(x)| dx < \infty\}$ и $T(\gamma) = \int \gamma(x) dx$. Тогда $\dot{T}(\gamma)[h] = \int h(x) dx$ – линейная, непрерывная в пространстве L_1 , и

$$T(\gamma + th) - T(\gamma) = \int (\gamma + th)(x) dx - \int \gamma(x) dx = t \int h(x) dx = t\dot{T}(\gamma)[h].$$

(ii) $\Gamma = \{\gamma | \exists \partial\gamma(x)/\partial x\}$ и $T(\gamma) = \partial\gamma(x)/\partial x$. Тогда $\dot{T}(\gamma)[h] = \partial h(x)/\partial x$:

$$T(\gamma + th) - T(\gamma) = \frac{\partial(\gamma + h)}{\partial x} - \frac{\partial\gamma}{\partial x} = t \frac{\partial h}{\partial x} = t\dot{T}(\gamma)[h].$$

Эти два случая просты, так как T – линейный функционал.

(iii) $T_x(\gamma) = F(\gamma(x))$, тогда $\dot{T}_x(\gamma)[h] = F'(\gamma(x))h(x)$.

(iv) $T(\gamma) = \int F(\gamma(x)) dx$. Тогда $\dot{T}(\gamma)[h] = \int F'(\gamma(x))h(x) dx$ при определенных ограничениях на F и Γ .

Теперь мы хотим использовать производные по направлению для вычисления дополнительной выборочной дисперсии нашей оценки $\hat{\theta}$ из-за присутствия $\hat{\gamma}$. Сперва рассмотрим следующие функционалы: скор и Гессиян целевой функции

$$S_n(\theta, \gamma) = \frac{\partial Q_n(\theta, \gamma)}{\partial\theta}, \quad H_n(\theta, \gamma) = \frac{\partial^2 Q_n(\theta, \gamma)}{\partial\theta\partial\theta'}.$$

Теперь предположим, что существует производная по направлению $h \in \Gamma$ функции $S_n(\theta, \gamma)$ по γ в $(\theta, \gamma) = (\theta_0, \gamma_0)$, которую обозначим $\dot{S}_n(\theta_0; \gamma_0)[h]$. Тогда мы можем сформулировать условия, при которых имеет место асимптотическая нормальность:

N.1 $\|\hat{\gamma} - \gamma_0\| = o_P(n^{-1/4})$ и $\hat{\theta} \xrightarrow{P} \theta_0$.

N.2 $\theta_0 \in \text{int}(\Theta)$.

N.3 $Q_n(\theta, \gamma)$ дважды непрерывно дифференцируемая по θ в окрестности \mathcal{N} точки θ_0 .

N.4 Производная по направлению $\dot{S}_n(\theta_0; \gamma_0)[h]$ существует и удовлетворяет условию

$$\left\| S_n(\theta_0, \gamma) - S_n(\theta_0, \gamma_0) - \dot{S}_n(\theta_0; \gamma_0)[\gamma - \gamma_0] \right\| \leq B_n \|\gamma - \gamma_0\|^2,$$

где $B_n = O_P(1)$.

N.5 $\sqrt{n} \left\{ S_n(\theta_0, \gamma_0) + \dot{S}_n(\theta_0; \gamma_0)[\hat{\gamma} - \gamma_0] \right\} \xrightarrow{d} N(0, \Omega_0)$.

N.6 $\|H_n(\theta, \gamma) - H_n(\theta, \gamma_0)\| \leq B_n \|\gamma - \gamma_0\|^\lambda$, где $B_n = O_P(1)$.

N.7 $\sup_{\theta \in \mathcal{N}} \|H_n(\theta, \gamma_0) - H(\theta, \gamma_0)\| \xrightarrow{P} 0$, где $H_0 = H(\theta_0, \gamma_0)$ невырожденна.

Как и в случае состоятельности, эти условия состоят из двух частей. (N.2), (N.3), (N.5) (положив $\dot{S}_n(\theta_0; \hat{\gamma} - \gamma_0) = 0$) и (N.7) означают, что недостижимая оценка, $\tilde{\theta}$, при известном γ_0 \sqrt{n} -асимптотически нормально распределена (см., например, Newey & McFadden, 1994, Теорема 3.1). Оставшиеся условия, (N.1), (N.4) и (N.6), позволяют показать, что достижимая оценка также \sqrt{n} -асимптотически нормально распределена.

Теорема 4 Пусть $\hat{\theta} \rightarrow^P \theta_0$, и выполняются (N.1)–(N.7). Тогда

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N(0, H_0^{-1} \Omega_0 H_0^{-1}),$$

где Ω_0 и H_0 определены как в (N.5) и (N.7).

Доказательство Доступно по запросу к автору.

Как мы упоминали до теоремы 4, недостижимая оценка, $\tilde{\theta}$, также \sqrt{n} -асимптотически нормально распределена, если выполнены (N.1)–(N.7). Однако, $\tilde{\theta}$ в общем случае имеет меньшую асимптотическую дисперсию и потому эффективнее, чем $\hat{\theta}$. Асимптотические дисперсии этих двух оценок равны только если дополнительный член асимптотически исчезает, то есть $\sqrt{n}\dot{S}_n(\theta_0; \gamma_0) [\hat{\gamma} - \gamma_0] = o_P(1)$. В этом случае $\tilde{\theta}$ и $\hat{\theta}$ эквивалентны в первом приближении. Однако, в большинстве случаев из-за того, что мы не знаем γ_0 , эффективность страдает, то есть $\mathbb{V}[\hat{\theta}] > \mathbb{V}[\tilde{\theta}]$.

Иногда вместо условий (N.3)–(N.7) проще проверить альтернативные условия:

N.3' $Q_n(\theta, \gamma)$ непрерывно дифференцируема по θ в окрестности \mathcal{N} точки θ_0 .

N.4' Существует функционал $S(\theta, \gamma)$, такой что $\nu_n(\theta, \gamma) := S_n(\theta, \gamma) - S(\theta, \gamma)$ удовлетворяет

$$\sup_{\|\theta - \theta_0\| < \delta, \|\gamma - \gamma_0\| < \delta} \|\nu_n(\theta, \gamma) - \nu_n(\theta_0, \gamma_0)\| = o_P(1/\sqrt{n}),$$

и $S(\theta_0, \gamma_0) = 0$.

N.5' Производная по направлению $\dot{S}(\theta, \gamma)[h]$ от функции $S(\theta, \gamma)$ существует и удовлетворяет

$$\left\| S(\theta_0, \gamma) - S(\theta_0, \gamma_0) - \dot{S}(\theta_0, \gamma_0)[\gamma - \gamma_0] \right\| \leq B \|\gamma - \gamma_0\|^2$$

для некоторой константы $B < \infty$.

N.6' $\sqrt{n}(S_n(\theta_0, \gamma_0) + (\theta_0; \gamma_0) [\hat{\gamma} - \gamma_0]) \rightarrow^d N(0, \Omega_0)$.

N.7' Функция $S(\theta, \gamma)$ непрерывно дифференцируема по θ в окрестности \mathcal{N} точки θ_0 с непрерывной производной $H(\theta, \gamma)$, которая удовлетворяет

$$\sup_{\theta \in \mathcal{N}} \|H(\theta, \gamma) - H(\theta, \gamma_0)\| \leq B \|\gamma - \gamma_0\|^\lambda,$$

где $H_0 = H(\theta_0, \gamma_0)$ невырожденна.

Заметим, что мы ослабили условие (N.3), так что теперь требуется только, чтобы $Q_n(\theta, \gamma)$ имела одну производную. Условие (N.4') более общее, но его можно проверить эмпирическими методами (см. необходимые условия, например, в Chen, Linton & van Keilegom, 2003). В большинстве случаев $S(\theta, \gamma)$ можно выбрать как $S(\theta, \gamma) = \partial Q(\theta, \gamma) / \partial \theta$, при этом условие идентификации в (C.2) обычно гарантирует, что $S(\theta_0, \gamma_0) = 0$.

Также заметим, что условия (N.5') и (N.6') включают в себя предельную скор-функцию $S(\theta, \gamma)$, а не $S_n(\theta, \gamma)$, как в (N.5) и (N.6).

Теорема 5 Пусть $\hat{\theta} \xrightarrow{P} \theta_0$, и выполняются условия (N.1)–(N.2) и (N.3')–(N.7'). Тогда заключение теоремы 4 остается верным.

Доказательство Доступно по запросу к автору.

Если в теореме 4 требуется выполнение ЦПТ для $\sqrt{n}(S_n(\theta_0, \gamma_0) + \dot{S}_n(\theta_0; \gamma_0)[\hat{\gamma} - \gamma_0])$, то в теореме 5 – для $\sqrt{n}(S_n(\theta_0, \gamma_0) + \dot{S}(\theta_0; \gamma_0)[\hat{\gamma} - \gamma_0])$. При применении этих теорем основная проблема заключается в том, чтобы проверить выполнение ЦПТ для соответствующего слагаемого. На первый взгляд может показаться, что это невозможно из-за наличия $\dot{S}_n(\theta_0; \gamma_0)[\hat{\gamma} - \gamma_0]$ и $\dot{S}(\theta_0; \gamma_0)[\hat{\gamma} - \gamma_0]$, так как оба слагаемых включают в себя непараметрическую оценку, $\hat{\gamma}$, которая вообще сходится медленнее, чем \sqrt{n} . Однако, если подставить $\hat{\gamma}$ в производную по направлению, то происходит дополнительное полное сглаживание непараметрической оценки. Как мы увидим далее, это сглаживание в общем случае увеличивает скорость сходимости и делает возможным проверить условия (N.5) или (N.6').

Чтобы продемонстрировать, как можно проверить \sqrt{n} -сходимость, мы ограничимся случаем, когда скор-функция имеет вид

$$S_n(\theta, \gamma) = \frac{1}{n} \sum_{i=1}^n s(Z_i; \theta, \gamma) + o_P(n^{-1/2}),$$

где $Z_i \in \mathbb{R}^d$, $i = 1, \dots, n$, независимо и одинаково распределены. Это ограничение выполняется, когда, например, $Q_n(\theta, \gamma)$ задано как в (20), при этом $s(z; \theta, \gamma) = \partial q(z; \theta, \gamma) / \partial \theta$. При этих ограничениях производные по направлению $S_n(\theta, \gamma)$ и $S(\theta, \gamma)$ задаются как

$$\dot{S}_n(\theta, \gamma)[h] = \frac{1}{n} \sum_{i=1}^n \dot{s}(Z_i; \theta, \gamma)[h], \quad \dot{S}(\theta, \gamma)[h] = \mathbb{E}[\dot{s}(Z; \theta, \gamma)[h]],$$

где $\dot{s}(z; \theta, \gamma)[h]$ – производная по направлению h функции $s(z; \theta, \gamma)[h]$ по γ .

Сперва мы сформулируем достаточные условия, при которых выполняются (N.4)–(N.5):

N.4.i $\|s(z; \theta_0, \gamma) - s(z; \theta_0, \gamma_0) - \dot{s}(z; \theta_0, \gamma_0)[\gamma - \gamma_0]\| \leq b(z) \|\gamma - \gamma_0\|^2$, причем $\mathbb{E}[b(Z)] < \infty$.

N.5.i $n^{-1} \sum_{i=1}^n \dot{s}(Z_i; \theta_0, \gamma_0)[\hat{\gamma} - \gamma_0] = \dot{S}(\theta_0, \gamma_0)[\hat{\gamma} - \gamma_0] + o_P(1/\sqrt{n})$.

N.5.ii Существует функция $\delta : \mathbb{R}^d \mapsto \mathbb{R}^k$ с $\mathbb{E}[\delta(Z)] = 0$ и $E[\|\delta(Z)\|^2] < \infty$, такая что

$$\dot{S}(\theta_0, \gamma_0)[\hat{\gamma} - \gamma_0] = \frac{1}{n} \sum_{i=1}^n \delta(Z_i) + o_P(1/\sqrt{n}).$$

N.5.iii $\mathbb{E}[s(Z; \theta_0, \gamma_0)] = 0$ и $E[\|s(Z; \theta_0, \gamma_0)\|^2] < \infty$.

Лемма 6 Пусть $\dot{s}(z; \theta, \gamma)[h]$ существует и удовлетворяет (N.4.i)–(N.5.iii). Тогда условия (N.4) и (N.5) выполняются, причем

$$\Omega = \mathbb{E}[[s(Z; \theta_0, \gamma_0) + \delta(Z)][s(Z; \theta_0, \gamma_0) + \delta(Z)]'].$$

Доказательство Доступно по запросу к автору.

Хотя (N.4.i)–(N.5.iii) более простые условия, во многих случаях по-прежнему не очевидно, как их проверить. В частности, не так очевидно доказать, что выполняются предположения (N.5.i)–(N.5.ii). Поэтому в дальнейшем мы предположим, что $\hat{\gamma}$ может быть записана в виде

$$\hat{\gamma}(x) = \frac{1}{n} \sum_{i=1}^n w_n(x, Z_i) + o_P(n^{-1/4}),$$

для некоторой функции w_n , которая может зависеть от размера выборки n . Это ограничение выполняется, например, для ядерной оценки, если положить $w_n(x, Z_i) = Y_i K_h(X_i - x)$. Можно легко убедиться, что оценки решетом тоже сюда подпадают (см., например, Newey, 1997). Далее, мы опустим зависимость от (θ_0, γ_0) , и будем писать, например, просто $\dot{S}[\hat{\gamma} - \gamma_0]$ вместо $\dot{S}(\theta_0, \gamma_0)[\hat{\gamma} - \gamma_0]$.

Проверка (N.5.i). Сперва заметим, что так как \dot{s} линейная функция, можно записать

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \dot{s}(Z_i) [\hat{\gamma} - \gamma_0] - \dot{S}[\hat{\gamma} - \gamma_0] &= \frac{1}{n} \sum_{i=1}^n \dot{s}(Z_i) [\hat{\gamma} - \bar{\gamma}] - \dot{S}[\hat{\gamma} - \bar{\gamma}] \\ &+ \frac{1}{n} \sum_{i=1}^n \dot{s}(Z_i) [\bar{\gamma} - \gamma_0] - \dot{S}[\bar{\gamma} - \gamma_0] \\ &=: I_{n,1} + I_{n,2}, \end{aligned}$$

где $\bar{\gamma}(x) = \mathbb{E}[\hat{\gamma}(x)]$. Определив

$$V_n(x, x') = \dot{s}(z)[w_n(\cdot, x')], \quad V_n = \mathbb{E}[V_n(Z_1, Z_2)],$$

$$V_{n,1}(x) = \mathbb{E}[V_n(x, Z)], \quad V_{n,2}(x) = \mathbb{E}[V_n(Z, x)],$$

и опять пользуясь тем, что \dot{s} – линейный функционал, запишем:

$$\begin{aligned} I_{n,1} &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \dot{s}(Z_i) [w_n(\cdot, Z_j)] - \frac{1}{n} \sum_{i=1}^n \dot{s}(Z_i) [E[w_n(\cdot, Z)]] \\ &\quad - \frac{1}{n} \sum_{i=1}^n \dot{S}[w_n(\cdot, Z_j)] + \dot{S}[E[w_n(\cdot, Z)]] \\ &= \frac{1}{n} \sum_{i,j=1}^n V_n(Z_i, Z_j) - \frac{1}{n} \sum_{i=1}^n V_{n,1}(Z_i) - \frac{1}{n} \sum_{j=1}^n V_{n,2}(Z_j) - V_n. \end{aligned}$$

Теперь можно использовать результаты для так называемых U-статистик (см., например, Lee, 1990 в качестве введения), чтобы показать, что правая часть есть $o_P(n^{-1/2})$ в самом общем случае. Итак, с первым членом порядок. Что касается второго, обычно можно показать, что

$$\|\dot{s}(z; \theta_0, \gamma_0)[h]\| \leq b(z) \|h\|,$$

в случае чего

$$\mathbb{E}[\|I_{n,2}\|] \leq \mathbb{E}[b(Z)] \|E[\hat{\gamma}] - \gamma_0\|,$$

и теперь нужно показать, что смещение достаточно быстро стремится к нулю: $\|E[\hat{\gamma}] - \gamma_0\| = o_P(n^{-1/2})$. В случае, когда $\hat{\gamma}$ – ядерная оценка, это условие можно проверить в самом общем случае, комбинируя так называемые ядра высокого порядка с недостаточным сглаживанием.

Проверка (N.5.ii). Это условие обычно проверяют, сначала показав, что существует функция d , такая что

$$\dot{S}[h] = \int d(x) h(x) dx.$$

Как правило, это выводится непосредственно, если известно явное выражение для \dot{S} , см. Newey (1994b). В качестве альтернативы можно воспользоваться теоремой представления Рисса, как это сделано в Ait-Sahalia (1993).

При заданном представлении, как правило, можно найти функцию δ . Предположим, например, что $\gamma_0(x) = f_X(x) \mathbb{E}[Y|X=x]$, $\hat{\gamma}(x) = 1/n \sum_{i=1}^n Y_i K((X_i - x)/h)/h^d$, и $\dot{S}[h] = \int d(x) h(x) dx$. Тогда сперва запишем

$$\dot{S}[\hat{\gamma} - \gamma_0] = \dot{S}[\hat{\gamma}] - \dot{S}[\gamma_0].$$

Первое слагаемое удовлетворяет

$$\dot{S}[\hat{\gamma}] = \int d(x) \hat{\gamma}(x) dx = \frac{1}{n} \sum_{i=1}^n Y_i \frac{1}{h^d} \int d(x) K\left(\frac{X_i - x}{h}\right) dx = \frac{1}{n} \sum_{i=1}^n Y_i d(X_i) + o_P\left(\frac{1}{\sqrt{n}}\right),$$

где последнее равенство следует из соответствующего условия регулярности, так как

$$h^{-d} \int d(x) K\left(\frac{X_i - x}{h}\right) dx \rightarrow d(X_i)$$

при $h \rightarrow 0$. Второе слагаемое можно записать как

$$\dot{S}[\gamma_0] = \int d(x) \gamma_0(x) dx = \int d(x) f_X(x) \mathbb{E}[Y|X=x] dx = \mathbb{E}[Yd(X)].$$

Полагая $\delta(z)$ равным

$$\delta(z) := yd(x) - \mathbb{E}[Yd(X)],$$

мы получим желаемое:

$$\dot{S}[\hat{\gamma} - \gamma_0] = \frac{1}{n} \sum_{i=1}^n \delta(Z_i) + o_P(1/\sqrt{n}).$$

Дальнейшие шаги по проверке (N.5.ii) для ядерной оценки можно найти в Newey (1994b).

Пример 3 (продолжение). Допустим, мы уже проверили, что

$$\|\hat{\sigma}^2 - \sigma_0^2\|_\infty = o_P(n^{-1/4}),$$

для некоторой последовательности ширины окна (см., например, Kristensen, 2009).

Чтобы вывести асимптотическое распределение $\hat{\theta}$, сперва найдем скор-функцию и Гессиян:

$$\begin{aligned} S_n(\theta, \sigma^2) &= \frac{\partial Q_n(\theta, \sigma^2)}{\partial \theta} = -\frac{2}{n} \sum_{i=1}^n \sigma^{-2}(X_i) [Y_i - \theta' X_i] X_i, \\ H_n(\theta, \sigma^2) &= \frac{\partial^2 Q_n(\theta, \sigma^2)}{\partial \theta \partial \theta'} = \frac{2}{n} \sum_{i=1}^n \sigma^{-2}(X_i) X_i X_i'. \end{aligned}$$

Сделаем следующую догадку относительно производной по направлению скор-функции:

$$\dot{S}_n[h] = \frac{1}{n} \sum_{i=1}^n \dot{s}(Z_i)[h],$$

где

$$\begin{aligned} \dot{s}(z)[h] &= [y - \theta_0' x] x' \sigma_0^{-4}(x) h(x) = d(y, x) h(x), \\ d(y, x) &= [y - \theta_0' x] x' \sigma_0^{-4}(x), \end{aligned}$$

и h – направление производной по σ^2 . Здесь мы опустили зависимость $\dot{S}_n[h]$ от θ_0 и $\sigma_0^2(x)$. Заметим, что

$$d(Y, X) = [Y - \theta_0' X] X \sigma_0^{-4}(X) h(X) = \varepsilon X \sigma_0^{-4}(X) h(X).$$

Проверяем, что необходимые условия удовлетворены. Во-первых, $h \mapsto \dot{S}_n[h]$ линейна; во-вторых, разложение в ряд Тейлора до второго порядка малости функции $a \mapsto 1/a$ дает

$$\frac{1}{a} - \frac{1}{a_0} = -\frac{1}{a_0^2} (a - a_0) + \frac{2}{(\lambda a + (1 - \lambda) a_0)^3} (a - a_0)^2$$

для некоторого $\lambda \in [0, 1]$. Полагая $a = \hat{\sigma}^2(x)$, $a_0 = \sigma_0^2(x)$, получим

$$\hat{\sigma}^{-2}(x) - \sigma_0^{-2}(x) = -\sigma_0^{-4}(x) (\hat{\sigma}^2(x) - \sigma_0^2(x)) + \frac{2 (\hat{\sigma}^2(x) - \sigma_0^2(x))^2}{(\lambda_x \hat{\sigma}^2(x) + (1 - \lambda_x) \sigma_0^2(x))^3},$$

где

$$\left| \frac{2}{(\lambda_x \hat{\sigma}^2(x) + (1 - \lambda_x) \sigma_0^2(x))^3} \right| \leq \frac{16}{\underline{\sigma}^6},$$

с той же самой аргументацией, что и при доказательстве состоятельности. Следовательно,

$$\begin{aligned} & \left\| S_n(\theta_0, \hat{\sigma}^2) - S_n(\theta_0, \sigma_0^2) - \dot{S}_n[\hat{\sigma}^2 - \sigma_0^2] \right\| \\ & \leq \frac{2}{n} \sum_{i=1}^n \|\varepsilon_i X'_i\| |\hat{\sigma}^{-2}(X_i) - \sigma_0^{-2}(X_i) + \sigma_0^{-4}(X_i) (\hat{\sigma}^2(X_i) - \sigma_0^2(X_i))| \\ & \leq \left(\frac{2}{n} \sum_{i=1}^n \|\varepsilon_i X'_i\| \right) \sup_{x \in \mathcal{X}} |\hat{\sigma}^{-2}(x) - \sigma_0^{-2}(x) + \sigma_0^{-4}(x) (\hat{\sigma}^2(x) - \sigma_0^2(x))|, \end{aligned}$$

где $n^{-1} \sum_{i=1}^n \|\varepsilon_i X'_i\| = \mathbb{E}[\|\varepsilon X'\|] + o_P(1)$ по ЗБЧ, в то время как

$$\begin{aligned} \sup_{x \in \mathcal{X}} |\hat{\sigma}^{-2}(x) - \sigma_0^{-2}(x) + \sigma_0^{-4}(x) (\hat{\sigma}^2(x) - \sigma_0^2(x))| & \leq \frac{16}{\underline{\sigma}^6} \sup_{x \in \mathcal{X}} |\hat{\sigma}^2(x) - \sigma_0^2(x)|^2 \\ & = \frac{16}{\underline{\sigma}^6} \|\hat{\sigma}^2 - \sigma_0^2\|_\infty^2 \\ & = o_P\left(\frac{1}{\sqrt{n}}\right). \end{aligned}$$

Делаем вывод, что

$$\sqrt{n} S_n(\theta_0, \hat{\sigma}^2) = \sqrt{n} S_n(\theta_0, \sigma_0^2) + \sqrt{n} \dot{S}_n[\hat{\sigma}^2 - \sigma_0^2] + o_P(1).$$

Далее, для любого h ,

$$\begin{aligned} \dot{S}[h] & = \mathbb{E}[\dot{s}(Z)[h]] = \mathbb{E}[d(Y, X)h(X)] \\ & = \mathbb{E}[\varepsilon \sigma_0^{-4}(X)h(X)X'] = \mathbb{E}[\mathbb{E}[\varepsilon|X]\sigma_0^{-4}(X)h(X)X'] \\ & = 0, \end{aligned}$$

так как $\mathbb{E}[\varepsilon|X] = 0$. Это означает, что добавочное слагаемое $\delta(Z) \equiv 0$. Таким образом, если мы сможем проверить, что

$$\sqrt{n} \left(\dot{S}_n[\hat{\sigma}^2 - \sigma_0^2] - \dot{S}[\hat{\sigma}^2 - \sigma_0^2] \right) \rightarrow^P 0,$$

то сможем заключить, что

$$\sqrt{n} \dot{S}_n[\hat{\sigma}^2 - \sigma_0^2] = o_P(1).$$

Гессиан удовлетворяет

$$\begin{aligned} \|H_n(\hat{\sigma}^2) - H_n(\sigma_0^2)\| &\leq \frac{2}{n} \sum_{i=1}^n \|X_i\|^2 |\hat{\sigma}^{-2}(X_i) - \sigma_0^{-2}(X_i)| \\ &\leq \left(\frac{2}{n} \sum_{i=1}^n \|X_i\|^2 \right) \sup_{x \in \mathcal{X}} |\hat{\sigma}^{-2}(x) - \sigma_0^{-2}(x)| \\ &= O_P(1) \times o_P(1), \end{aligned}$$

и $H_n(\sigma_0^2) \xrightarrow{P} H_0 = \mathbb{E}[\sigma_0^{-2}(X) X X']$ по ЗБЧ.

Собрав полученные результаты, получим

$$\begin{aligned} \sqrt{n}(\hat{\theta} - \theta_0) &= H_n^{-1}(\bar{\theta}, \hat{\sigma}^2) \sqrt{n} S_n(\theta_0, \hat{\sigma}^2) \\ &= H_n^{-1}(\bar{\theta}, \hat{\sigma}^2) \left\{ \sqrt{n} S_n(\theta_0, \sigma_0^2) + \sqrt{n} \dot{S}_n[\hat{\sigma}^2 - \sigma_0^2] + o_P(1) \right\} \\ &= H_n^{-1}(\bar{\theta}, \hat{\sigma}^2) \left\{ \sqrt{n} S_n(\theta_0, \sigma_0^2) + o_P(1) \right\} \\ &\xrightarrow{d} N(0, H_0^{-1} \Omega_0 H_0^{-1}), \end{aligned}$$

где

$$\Omega_0 = \mathbb{E}[\sigma^{-4}(X) \varepsilon^2 X X'] = \mathbb{E}[\sigma^{-2}(X) X X'],$$

так что

$$H_0^{-1} \Omega_0 H_0^{-1} = \mathbb{E}[\sigma^{-2}(X) X X']^{-1}.$$

Заметим, что недостижимая ВМНК-оценка

$$\tilde{\theta} = \arg \min_{\theta \in \Theta} Q_n(\theta, \sigma_0^2)$$

имеет такое же асимптотическое распределение, что и $\hat{\theta}$. Следовательно, достижимая ВМНК-оценка, основанная на непараметрической оценке $\hat{\sigma}^2(x)$, асимптотически эквивалентна недостижимой ВМНК-оценке. Так что при замене σ_0^2 на $\hat{\sigma}^2$ в оценке θ мы не теряем в эффективности. Однако заметим, что это не означает, что $\hat{\theta}$ обязательно наиболее эффективная из доступных оценок, так как не выполняется равенство Рао–Крамера (за исключением случая, когда нормированные ошибки $\sigma^{-1}(X)\varepsilon$ имеют независимое и одинаковое стандартное нормальное распределение).

5.4 Оценивание дисперсий

Положив $\hat{H} = H_n(\hat{\theta}, \hat{\gamma})$, можно оценить H_0 , при этом оценка будет состоятельна при выполнении условий теоремы 4. Если нам известна оценка $\hat{\delta}$ параметра δ , то при условиях регулярности

$$\hat{\Omega} = \frac{1}{n} \sum_{i=1}^n \left[s(Z_i; \hat{\theta}, \hat{\gamma}) + \hat{\delta}(Z_i) \right] \left[s(Z_i; \hat{\theta}, \hat{\gamma}) + \hat{\delta}(Z_i) \right]'$$

будет состоятельной оценкой для Ω_0 . Обычно можно вывести явное выражение для δ , так как $\delta(z) = \delta(z; \theta_0, \gamma_0)$, при этом естественно взять $\hat{\delta}(z) = \delta(z; \hat{\theta}, \hat{\gamma})$. Стандартными способами можно показать, что оценка дисперсии будет состоятельной, если s и δ удовлетворяют условию Липшица по (θ, γ) , см., например, Newey & McFadden (1994, Теорема 8.13).

Однако в сложных моделях нельзя получить выражение для δ в явном виде (см., например, Kristensen, 2008). В этом случае можно воспользоваться бутстрапом (Chen, Linton & van Keilegom, 2003) или численными методами (Newey, 1994a).

6 Оценивание «решетом»

Несмотря на то, что во многих случаях полупараметрические модели можно оценивать двухшаговой процедурой, существует альтернативный подход, который состоит в одновременном оценивании как параметрической, так и непараметрической компонент. Обсудим, как это реализуется в контексте метода «решето».

Как и в предыдущем разделе, мы хотим оценить параметр $\theta \in \Theta \subseteq \mathbb{R}^k$, используя целевую функцию $Q_n(\theta, \gamma)$, где $\gamma \in \Gamma$ – бесконечномерный параметр. Вместо того, чтобы использовать предварительную оценку (если такая вообще доступна) и воспользоваться двухшаговой процедурой, можно оценить θ и γ одновременно, используя так называемое решето. Метод «решето» – это общий непараметрический метод, в котором бесконечномерные функциональные пространства заменяются на приближающие их конечномерные (так называемое решето) в конечных выборках. Ошибка приближения, возникающая из-за перехода к конечномерным пространствам, асимптотически исчезает с ростом объема выборки ввиду увеличения размерности решета.

Для того, чтобы дать определение полупараметрической оценке методом решето, сначала введем дополнительные обозначения. Пусть выбрана последовательность аппроксимирующих конечномерных пространств $\{\Gamma_J\}$, такая что $\Gamma_J \subseteq \Gamma$, $J \geq 1$, и $\bigcup_{J=1}^{\infty} \Gamma_J = \Gamma$. Определим

$$(\hat{\theta}, \hat{\gamma}) = \arg \min_{\theta \in \Theta, \gamma \in \Gamma_{J_n}} Q_n(\theta, \gamma) \quad (24)$$

для некоторой последовательности $J_n \rightarrow \infty$ при $n \rightarrow \infty$. Здесь мы оцениваем θ и γ одновременно, используя ту же целевую функцию, Q_n . При использовании же двухшаговых оценок, рассмотренных в предыдущем разделе, использовались две разные целевые функции для получения оценок θ и γ соответственно.

Общие результаты относительно состоятельности и скорости сходимости для случая, где $Q_n(\theta, \gamma)$ принимает вид выборочного среднего как в уравнении (20), можно найти в работах Shen & Wong (1994) и Shen (1997). Более того, в этих статьях выводятся условия, при которых оценка $\hat{\theta}$ \sqrt{n} -асимптотически нормальна. ОММ-подобные решетчатые оценки для моделей, определяемых через условия на условные моменты, разработаны и проанализированы в Ai & Chen (2003); см. также Blundell, Chen & Kristensen (2007). Поскольку условия, необходимые для получения указанных результатов, являются сугубо техническими (как и доказательства), мы не станем здесь вдаваться в детали.

Один из недостатков подхода, рассмотренного выше, связан с его практической реализацией. При двухшаговом оценивании $\hat{\gamma}$ используется в качестве предварительной оценки, так что достаточно решить оптимизационную задачу малой размерности $\hat{\theta} = \arg \min_{\theta \in \Theta} Q_n(\theta, \hat{\gamma})$. С другой стороны, решетчатые оценки требуют одновременной оптимизации как по θ , так и по γ . В частности, размерность γ может оказаться достаточно большой в стандартных задачах и будет расти экспоненциально с ростом числа переменных, от которых берется функция. Таким образом, численная задача решения $(\hat{\theta}, \hat{\gamma})$ в уравнении (24) является задачей «высокой» размерности и может оказаться расчетно невыполнимой. Тем не менее, во многих случаях можно найти решение в явном виде, что упрощает вычисления.

Рассмотрим два примера, в которых показано, как можно получить полупараметрические решетчатые оценки.

Пример 1 (продолжение). Когда $\gamma = g$, целевая функция для одноиндексной модели имеет форму (20), где q задается как

$$q(z; \theta, \gamma) = [y - \gamma(\beta'x)]^2.$$

Пусть Γ – некоторое функциональное пространство, для которого существует решетка в форме

$$\Gamma_J = \left\{ \gamma_J(z) = \sum_{j=1}^J \alpha_j \varphi_j(z) : \alpha_j \in \mathbb{R}, j = 1, \dots, J \right\} \quad (25)$$

где $\varphi_1(z), \varphi_2(z), \dots$ – известные базисные функции. Решетчатая оценка, определенная в уравнении (24), принимает вид

$$(\hat{\beta}, \hat{\gamma}) = \arg \min_{\beta, A_{J_n}} \sum_{i=1}^n [Y_i - A'_{J_n} \Phi_{J_n}(\beta' X_i)]^2,$$

где $A_{J_n} = (\alpha_1, \dots, \alpha_{J_n})'$ и $\Phi_{J_n}(z) = (\varphi_1(z), \dots, \varphi_{J_n}(z))'$. Для любого заданного значения β условия первого порядка по A_{J_n} имеют вид

$$\sum_{i=1}^n [Y_i - A'_{J_n} \varphi_{J_n}(\theta' X_i)] \Phi_{J_n}(\beta' X_i) = 0,$$

откуда получаем решение:

$$\hat{A}_{J_n}(\beta) = \left(\sum_{i=1}^n \Phi_{J_n}(\beta' X_i) \Phi_{J_n}(\beta' X_i)' \right)^{-1} \sum_{i=1}^n \Phi_{J_n}(\beta' X_i) Y_i.$$

Подставив, получаем «профильную» оценку:

$$\hat{\beta} = \arg \min_{\theta} \sum_{i=1}^n \left[Y_i - \hat{A}_{J_n}(\beta)' \Phi_{J_n}(\beta' X_i) \right]^2.$$

Здесь сложность вычислений ограничивается лишь численной оптимизацией по θ . Заметим, что в данном случае одновременная оценка совпадает с оценкой по двухшаговому методу, где оценка сериями используется в качестве предварительной оценки γ .

Пример 4 (продолжение). В случае полупараметрической модели с копулой $\gamma = (f_1, f_2)$, и целевую функцию снова можно записать в виде (20), где q задается:

$$q(z; \theta, \gamma) = \{ \log c(F_1(z_1), F_2(z_2); \theta) + \log f_1(z_1) + \log f_2(z_2) \}.$$

Эта оценка была предложена в Chen, Fan & Tsyrennikov (2006), где также указан способ, которым можно построить решетчатое пространство для двух плотностей. Решетчатую оценку в общем случае нельзя записать в явном виде, и необходимо решать задачу численной оптимизации, что делает этот метод не слишком привлекательным. Решетчатая оценка, составленная таким образом, в общем случае будет более эффективной, чем двухшаговая оценка, построенная в разделе 4.

Хотя в первом из этих примеров решетчатая и двухшаговая ядерная оценки очень похожи, в общем случае решетчатая оценка будет отличаться. В частности, решетчатые оценки будут более эффективными по сравнению с двухшаговыми ядерными оценками в силу их построения, поскольку параметрическая и непараметрическая компоненты оцениваются одновременно. Это приводит нас к вопросу об эффективности полупараметрических оценок.

7 Полупараметрическая эффективность

Итак, любая полупараметрическая модель полностью характеризуется параметрической компонентой, θ_0 , и непараметрической, $\gamma_0(\cdot)$. Нас интересуют параметр θ_0 и насколько эффективно можно его оценить, ничего не зная априори о непараметрической компоненте, $\gamma_0(\cdot)$. В общем случае ответить на этот вопрос нелегко, но существует конструктивный подход, с помощью которого можно рассчитать границы уровня эффективности для θ_0 .

Содержательный смысл описываемых границ следующий. Рассмотрим оценивание двух статистических моделей, в которых вторая модель содержится в первой (встроена в нее). Естественно ожидать, что оценивание второй модели будет проще по сравнению с первой. В частности, если обе модели содержат общий параметр, например θ , то следует ожидать, что этот параметр во второй модели будет оценен точнее. Таким образом, если мы сможем рассчитать эффективность оценивания θ во второй модели, то мы получим границы эффективности θ для первой модели.

Stein (1956) использовал эту идею для построения границ эффективности в задачах полупараметрического оценивания. В качестве первой, более сложной модели он выбрал полупараметрическую модель, которая характеризуется набором $(\theta_0, \gamma_0(\cdot))$. А в качестве второй, более простой модели он использовал полностью параметрическую подмодель. Выберем некоторое параметрическое семейство функций, $\gamma(\cdot; \alpha)$, где $\alpha \in \mathcal{A} \subseteq \mathbb{R}^l$ – параметр, и предположим, что параметрическая подмодель содержит истинную функцию $\gamma_0(\cdot)$, $\gamma_0(\cdot) = \gamma(\cdot; \alpha_0)$ для некоторого $\alpha_0 \in \mathcal{A}$. Таким образом, вторая модель характеризуется набором (θ_0, α_0) .

Оценивание полупараметрической модели должно быть как минимум не легче, чем полностью параметрической подмодели. Таким образом, следует ожидать, что в полупараметрической модели мы не сможем оценить θ_0 точнее. Поскольку параметрическая модель полностью специфицирована в терминах (θ, α) , мы можем записать плотность модели как функцию (θ, α) , $(\theta, \alpha) \mapsto p(z; \theta, \gamma(\cdot; \alpha))$. Естественной в таком случае будет ММП-оценка, а ее точность определяется соответствующей информационной матрицей Фишера,

$$\mathcal{I} = \mathbb{E} \left[\frac{\partial^2 \log p(Z; \theta, \gamma(\cdot; \alpha))}{\partial (\theta, \alpha) \partial (\theta, \alpha)'} \right] = \begin{bmatrix} \mathcal{I}_{\theta\theta} & \mathcal{I}_{\theta\alpha} \\ \mathcal{I}_{\theta\alpha} & \mathcal{I}_{\alpha\alpha} \end{bmatrix}.$$

Для любой заданной параметрической спецификации, $\gamma(\cdot; \alpha)$, уровень эффективности для θ задается неравенством Рао–Крамера,

$$\mathcal{I}_p = \mathcal{I}_{\theta\theta} - \mathcal{I}_{\theta\alpha} \mathcal{I}_{\alpha\alpha}^{-1} \mathcal{I}_{\theta\alpha}.$$

То есть, асимптотическая дисперсия ММП-оценки равна \mathcal{I}_p^{-1} . Это значение дисперсии выражает цену, которую мы вынуждены платить за незнание $\gamma(\cdot)$ (что соответствует α в параметрической модели). Если α известно, θ можно оценить с асимптотической дисперсией $\mathcal{I}_{\theta\theta}^{-1}$. Если α неизвестно и его нужно оценивать, дисперсия будет $\mathcal{I}_p^{-1} \geq \mathcal{I}_{\theta\theta}^{-1}$, где равенство возможно тогда и только тогда, когда $\mathcal{I}_{\theta\alpha} = 0$.

Рассмотрим теперь оценку, не полагающуюся на параметрическую информацию о $\gamma(\cdot)$, и пусть \mathcal{I}_{sp}^{-1} – ее асимптотическая дисперсия. Тогда должно выполняться $\mathcal{I}_p^{-1} \leq \mathcal{I}_{sp}^{-1}$. Утверждение верно независимо от выбора параметрической подмодели, так что $\sup_{\gamma(\cdot; \alpha)} \mathcal{I}_p^{-1} \leq \mathcal{I}_{sp}^{-1}$. Из этого следует следующее определение границ эффективности полупараметрической оценки как асимптотической дисперсии «наименее благоприятной» параметрической подмодели:

$$\text{semiparametric efficiency bound (SEB)} = \sup_{\gamma(\cdot; \alpha)} \mathcal{I}_p^{-1}.$$

Одним из привлекательных является класс полупараметрических оценок, который хорош настолько же, как если бы наличествовала полная информация о непараметрической компоненте. Назовем полупараметрическую оценку $\hat{\theta}$ адаптивной, если

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N(0, \mathcal{I}_{\theta\theta}^{-1}).$$

Необходимым условием для этого является $\mathcal{I}_{\theta\alpha} = 0$ для всех параметрических подмоделей, что в общем случае не выполняется. Например, ни одна из полупараметрических оценок, рассмотренных в разделе 2, не является адаптивной. С другой стороны, очевидно, что ММП-подобные оценки, рассмотренные в разделе 3.3, адаптивны при выполнении определенных условий регулярности на распределение ошибок.

Чтобы понять, как связаны границы эффективности с асимптотическими результатами предыдущего раздела, напомним, что дисперсия полупараметрической оценки имеет вид $H_0^{-1}\Omega_0H_0^{-1}$, где $H_0 = \mathcal{I}_{\theta\theta}$ и

$$\Omega_0 = \mathbb{E} \left[\{s(Z; \theta_0, \gamma_0) + \delta(Z)\} \{s(Z; \theta_0, \gamma_0) + \delta(Z)\}' \right].$$

Здесь $s(Z; \theta_0, \gamma_0) = \partial \log p(Z; \theta, \gamma) / \partial \theta$, где $\delta(Z)$ – добавочный член, связанный с использованием оценки γ_0 вместо ее истинного значения. Вопрос о границе эффективности полупараметрической оценки, грубо говоря, – это вопрос о поиске оценки с минимально возможным δ в терминах дисперсии. В частности, если $\delta = 0$, оценка *адаптивная*, ср. определение выше, поскольку дисперсия полностью совпадает с дисперсией в ситуации, когда нам известно γ_0 .

Следует отметить, что полупараметрическая оценка, достигающая границы эффективности, может не существовать, поскольку такая граница не обязательно является точной. Примеры такого рода можно найти в Ritov & Bickel (1987), где граница эффективности полупараметрической оценки четко определена, но \sqrt{n} -состоятельной полупараметрической оценки не существует.

Несмотря на то, что граница эффективности имеет интуитивный смысл, в общем случае в полупараметрических задачах для нее сложно получить явное выражение. Поэтому мы не станем далее рассматривать этот вопрос. Вместо этого рассмотрим подробнее интуитивный смысл такой границы и покажем, как она меняется в зависимости от предположений исследователя о виде модели. В качестве простого примера рассмотрим следующую полупараметрическую регрессионную модель:

$$Y = m(X; \theta) + \varepsilon,$$

где условное среднее полностью параметризовано, и единственными условиями на ошибки являются ограничения $\mathbb{E}[\varepsilon|X] = 0$ и $\mathbb{E}[\varepsilon^2] < \infty$. В этом случае $\gamma = F_{\varepsilon|X}(e|x)$ – непараметрическая компонента. В зависимости от дополнительных предположений исследователя касательно $F_{\varepsilon|X}$ возникают различные границы эффективности. Например, если учитывать только ограничение на условное среднее $\mathbb{E}[\varepsilon|X] = 0$, то SEB принимает вид

$$\text{SEB} = \mathbb{E} \left[\sigma^{-2}(X) \dot{m}(X; \theta, \gamma) \dot{m}(X; \theta, \gamma)' \right]^{-1},$$

где $\dot{m}(X; \theta, \gamma) = \partial m(X; \theta, \gamma) / \partial \theta$, а $\sigma^2(X) = \mathbb{E}[\varepsilon^2|X]$ – условная дисперсия. Этого можно достичь, например, используя решетчатую оценку Ai & Chen (2003).

С другой стороны, если предположить, что ε и X независимы и что ε имеет симметричное распределение, граница эффективности приобретает вид

$$\text{SEB} = \mathbb{E} \left[\frac{\partial \log f_{\varepsilon}(Y - m(X; \theta))}{\partial \theta} \frac{\partial \log f_{\varepsilon}(Y - m(X; \theta))}{\partial \theta} \right]^{-1}.$$

В этом случае граница эффективности совпадает со значением в неравенстве Рао–Крамера. Адаптивная оценка может быть получена способом, аналогичным описанному в разделе 3.3.

8 Примечания

Оценка одноиндексной модели была предложена в работе Ichimura (1993), где также выведены ее теоретические свойства. Асимптотическая теория для оценки средних производных

была получена в Powell, Stock & Stoker (1989); см. также Hristache, Juditsky & Spokoiny (2001). В случае модели бинарного выбора в качестве альтернативы можно использовать ММП-оценки β_0 в одноиндексной модели, см. Klein & Spady (1993).

Robinson (1988b) и Speckman (1988) предложили оценку, основанную на остатках частично-линейной модели, которая задается уравнением (7), и получили ее асимптотическое распределение. Andrews (1994a) получил результат для расширенной версии (10).

Robinson (1987) вывел асимптотику для ВМНК-оценки для случая гетероскедастичности неизвестного вида и показал, что эта оценка достигает границы полупараметрической эффективности. Ai & Chen (2003) предложили полупараметрические решетчатые оценки для класса полупараметрических моделей, задаваемых ограничениями на условные моменты.

Введение и описание общих результатов для копул содержится в Joe (1997). Свойства полупараметрических оценок с копулами из раздела 4 были получены в работе Genest, Ghoudi & Rivest (1995), в то время как свойства решетчатых оценок из раздела 6 проанализированы в Chen, Fan & Tsyrennikov (2006).

Для дальнейшего чтения по теме функциональных производных см. Luenberger (1969) Kantorovich & Akilov (1982).

Наши асимптотические результаты для двухшаговых оценок похожи на полученные в работах Andrews (1994a), Chen, Linton & van Keilegom (2003), Newey & McFadden (1994), Newey (1994b), Pakes & Olley (1995). В этих статьях рассматриваются общие условия состоятельности и асимптотической нормальности двухшаговых полупараметрических оценок. Что касается свойств более высоких порядков для полупараметрических оценок, сошлемся на Linton (1995, 1996).

Результаты касательно непараметрических решетчатых оценок можно найти в Andrews (1991), Fenton & Gallant (1996), Gallant & Nychka (1987), Newey (1997), Shen & Wong (1994). Про их использование в полупараметрическом оценивании см. Ai & Chen (2003), Shen (1997). Chen (2007) приводит обзор как непараметрических, так и полупараметрических методов оценивания с использованием решетчатых оценок.

Newey (1990) является хорошим введением в теорию границ эффективности полупараметрических оценок, а так же методов их получения; общий подход к расчету границ эффективности можно найти в Bickel, Klaassen, Ritov & Wellner (1993), а также в Severini & Tripathi (2001). Chamberlain (1987, 1992) выводит границы эффективности для случая ограничений на условные моменты и для полупараметрических регрессий. Manski (1984) рассматривает границы эффективности для адаптивных оценок нелинейных регрессионных моделей при наличии предположения о независимости; см. Drost & Klaassen (1997) относительно схожих результатов для случая моделей гетероскедастичных временных рядов.

Мы не обсуждаем вопросы практического применения полупараметрических оценок и отсылаем читателя к обзору Ichimura & Todd (2007). Cattaneo, Crump & Jansson (2009) детально обсуждают выбор границ для оценок средних производных, тогда как Härdle, Hall & Ichimura (1993) предлагают специальный метод для одноиндексных моделей.

Всюду в нашей работе предполагалась независимость и одинаковая распределенность данных. Большинство результатов для предложенных оценок подходят для стационарных последовательностей и последовательностей с перемешиванием; см., например, Ang & Kristensen (2009), Chen, Wu & Yu (2009), Hidalgo (1992), Kristensen (2008), а также Li & Wooldridge (2002). Вопрос границ полупараметрической эффективности для временных рядов тем не менее недостаточно хорошо исследован, если отбросить условие марковости и допустить произвольную форму зависимости; см. Bickel & Kwon (2001), а также Schick & Wefelmeyer (2005), где приводятся некоторые результаты и обсуждения.

Благодарности

Автор хотел бы поблагодарить Бруно Джиованети и Шин Канайа за полезные комментарии и предложения.

Список литературы

- Анатольев, С. (2009). Непараметрическая регрессия. *Квантиль* 7, 37–52.
- Расин, Дж. (2008). Непараметрическая эконометрика: вводный курс. *Квантиль* 4, 7–56.
- Ai, C. & X. Chen (2003). Efficient estimation of models with conditional moment restrictions containing unknown functions. *Econometrica* 71, 1795–1844.
- Aït-Sahalia, Y. (1993). The Delta method for nonparametric kernel functionals. Manuscript, University of Chicago.
- Amemiya, T. (1985). Non-linear regression models. Глава 6 в *Handbook of Econometrics* (под редакцией M.D. Intriligator & Z. Griliches), том 1, 333–389. Elsevier Science.
- Amemiya, T. (1985). *Advanced Econometrics*. Harvard University Press.
- Andrews, D.W.K. (1991). Asymptotic normality of series estimators for nonparametric and semiparametric regression models. *Econometrica* 59, 307–45.
- Andrews, D.W.K. (1994a). Asymptotics for semiparametric econometric models via stochastic equicontinuity. *Econometrica* 62, 43–72.
- Andrews, D.W.K. (1994b). Empirical process methods in econometrics. Глава 37 в *Handbook of Econometrics* (под редакцией R. Engle & D. McFadden), том 4, 2246–2294. Elsevier Science.
- Ang, A. & D. Kristensen (2009). Testing conditional factor models. CREATES Research Papers 2009–09, University of Aarhus.
- Bickel, P.J., C.A.J. Klaassen, Y. Ritov & J.A. Wellner (1993). *Efficient and Adaptive Estimation for Semiparametric Models*. The John Hopkins University Press.
- Bickel, P.J. & J. Kwon (2001). Inference for semiparametric models: Some questions and an answer. *Statistica Sinica* 11, 863–960.
- Blundell, R., X. Chen & D. Kristensen (2007). Semi-nonparametric IV estimation of shape invariant Engel curves. *Econometrica* 75, 1613–1670.
- Cattaneo, M.D., R.K. Crump & M. Jansson (2009). Small bandwidth asymptotics for density-weighted average derivatives. Manuscript, University of California–Berkeley.
- Chamberlain, G. (1987). Asymptotic efficiency in estimation of conditional moment restrictions. *Journal of Econometrics* 34, 305–334.
- Chamberlain, G. (1992). Efficiency bounds for semiparametric regression. *Econometrica* 60, 567–596.
- Chen, X. (2007). Large sample sieve estimation of semi-nonparametric models. Глава 76 в *Handbook of Econometrics* (под редакцией J.J. Heckman & E.E. Leamer), том 6/2. Elsevier Science.
- Chen, X., Y. Fan & V. Tsyrennikov (2006). Efficient estimation of semiparametric multivariate copula models. *Journal of American Statistical Association* 101, 1228–1240.
- Chen, X., O. Linton & I. van Keilegom (2003). Estimation of semiparametric models when the criterion function is not smooth. *Econometrica* 71, 1591–1608.
- Chen, X., W.B. Wu & Y. Yu (2009). Efficient estimation of copula-based semiparametric Markov models. Cowles Foundation Discussion Papers, No. 1691.
- Drost, F.C. & C.A.J. Klaassen (1997). Efficient estimation in semiparametric GARCH models. *Journal of Econometrics* 81, 193–221.
- Fenton, V.M. & A.R. Gallant (1996). Convergence rates of SNP density estimators. *Econometrica* 64, 719–727.
- Gallant, A.R. & D.W. Nychka (1987). Semi-nonparametric maximum likelihood estimation. *Econometrica* 55, 363–390.
- Genest, C., K. Ghoudi & L.-P. Rivest (1995). A semiparametric estimation procedure of dependence parameters in multivariate families of distributions. *Biometrika* 82, 543–552.

- Genest, C., M. Gendron & M. Bourdeau-Brien (2009). The advent of copulas in finance. *European Journal of Finance* 15, в печати.
- Härdle, W. (1990). *Applied Nonparametric Regression*. Cambridge: Cambridge University Press.
- Härdle, W., P. Hall & H. Ichimura (1993). Optimal smoothing in single-index models. *Annals of Statistics* 21, 157–178.
- Härdle, W. & O. Linton (1994). Applied nonparametric methods. Глава 38 в *Handbook of Econometrics* (под редакцией R. Engle & D. McFadden), том 4. Elsevier Science.
- Härdle, W., M. Müller, S. Sperlich & A. Werwatz (2004). *Nonparametric and Semiparametric Models*. New York: Springer-Verlag.
- Hidalgo, J. (1992). Adaptive estimation in time series regression models with heteroskedasticity of unknown form. *Econometric Theory* 8, 161–187.
- Horowitz, J. (2009). *Semiparametric and Nonparametric Methods in Econometrics*. Springer-Verlag.
- Hristache, M., A. Juditsky & V. Spokoiny (2001). Direct estimation of the index coefficients in a single-index model. *Annals of Statistics* 29, 595–623.
- Ichimura, H. (1993). Semiparametric least squares (SLS) and weighted SLS estimation of single-index models. *Journal of Econometrics* 58, 71–120.
- Ichimura, H. & P.E. Todd (2007). Implementing nonparametric and semiparametric estimators. Глава в *Handbook of Econometrics* (под редакцией J.J. Heckman & E.E. Leamer), том 6/2, 5369–5468. Elsevier Science.
- Joe, H. (1997). *Multivariate Models and Dependence Concepts*. Chapman and Hall.
- Kantorovich, L.V. & G.P. Akilov (1982). *Functional Analysis*. Pergamon Press, Oxford.
- Luenberger, D. G. (1969). *Optimization by Vector Space Methods*. John Wiley.
- Klein, R. & R. Spady (1993). An efficient semiparametric estimator for binary response models. *Econometrica* 61, 387–421.
- Kristensen, D. (2008). Pseudo-maximum-likelihood estimation in two classes of semiparametric diffusion models. Manuscript, Columbia University.
- Kristensen, D. (2009). Uniform convergence rates of kernel estimators with heterogeneous, dependent data. *Econometric Theory* 25, в печати.
- Lee, A.J. (1990). *U-Statistics, Theory and Practice*. Marcel Dekker.
- Li, Q. & J.S. Racine (2007). *Nonparametric Econometrics: Theory and Practice*. Princeton University Press.
- Li, Q. & J.M. Wooldridge (2002). Semiparametric estimation of partially linear models for dependent data with generated regressors. *Econometric Theory* 18, 625–645.
- Linton, O.B. (1995). Second order approximation in the partially linear regression model. *Econometrica* 63, 1079–1112.
- Linton, O.B. (1996). Edgeworth approximation for MINPIN estimators in semiparametric regression models. *Econometric Theory* 12, 30–60.
- Manski, C. (1984). Adaptive estimation of non-linear regression. *Econometric Reviews* 3, 145–194.
- Newey, W.K. (1990). Semiparametric efficiency bounds. *Journal of Applied Econometrics* 5, 99–135.
- Newey, W.K. (1991). Uniform convergence in probability and stochastic equicontinuity. *Econometrica* 59, 1161–1167.
- Newey, W.K. (1994a). Kernel estimation of partial means and a general variance estimator. *Econometric Theory* 10, 233–253.
- Newey, W.K. (1994b). The asymptotic variance of semiparametric estimators. *Econometrica* 62, 1349–1362.
- Newey, W.K. (1997). Convergence rates and asymptotic normality for series estimators. *Journal of Econometrics* 79, 147–168.
- Newey, W.K. & D.L. McFadden (1994). Large Sample Estimation and Hypothesis Testing. Глава 36 в *Handbook of Econometrics* (под редакцией R. Engle & D. McFadden), том 4, 2111–2245. Elsevier Science.
- Pagan, A. & A. Ullah (1999). *Nonparametric Econometrics*. Cambridge University Press.

- Pakes, A. & S. Olley (1995). A limit theorem for a smooth class of semiparametric estimators. *Journal of Econometrics* 65, 295–332.
- Powell, J.L. (1994). Estimation of Semiparametric Models. Глава 41 в *Handbook of Econometrics* (под редакцией R. Engle & D. McFadden), том 4, 2443–2521. Elsevier Science.
- Powell, J.L., J.H. Stock & T.M. Stoker (1989). Semiparametric estimation of index coefficients. *Econometrica* 51, 1403–1430.
- Robinson, P.M. (1987). Asymptotically efficient estimation in the presence of heteroskedasticity of unknown form. *Econometrica* 55, 875–891.
- Robinson, P.M. (1988a). Semiparametric econometrics: A survey. *Journal of Applied Econometrics* 3, 35–51.
- Robinson, P.M. (1988b). Root-N-consistent semiparametric regression. *Econometrica* 56, 931–954.
- Severini, T.A. & G. Tripathi (2001). A simplified approach to computing efficiency bounds in semiparametric models. *Journal of Econometrics* 102, 23–66.
- Shen, X. (1997). On methods of sieves and penalization. *Annals of Statistics* 25, 2555–2591.
- Shen, X. & W.H. Wong (1994). Convergence rate of sieve estimates. *Annals of Statistics* 22, 580–615.
- Schick, A. & W. Wefelmeyer (2006). Efficient estimators for time series. Глава в *Frontiers in Statistics* (под редакцией J. Fan & H. L. Koul), 45–62. Imperial College Press.
- Silverman, B.W. (1986). *Density Estimation for Statistics and Data Analysis*. Chapman and Hall.
- Speckman, P. (1988). Kernel smoothing in partial linear models. *Journal of Royal Statistical Society B* 50, 413–436.
- van der Vaart, A & J. Wellner (1996). *Weak Convergence and Empirical Processes*. Springer-Verlag.

Semiparametric modelling and estimation

Dennis Kristensen

Columbia University, New York, USA

Center for Research in Econometric Analysis of Time Series, Aarhus, Denmark

Semiparametric models are characterized by a finite- and infinite-dimensional (functional) component. As such they allow for added flexibility over fully parametric models, and at the same time estimators of parametric components can be developed that exhibit standard parametric convergence rates. These two features have made semiparametric models and estimators increasingly popular in applied economics. We give a partial overview over the literature on semiparametric modelling and estimation with particular emphasis on semiparametric regression models. The main focus is on developing two-step semiparametric estimators and deriving their asymptotic properties. We do however also briefly discuss sieve-based estimators and semiparametric efficiency.

