

Лекция 6 (от 7.10)

2.8. Приближенный поиск ОМП

Метод Ньютона: Пусть $f : \mathbb{R} \rightarrow \mathbb{R}$ — функция. Нужно решить уравнение $f(x) = 0$.

x_0 — начальное приближение

Формула касательной в точке x_k : $y = f(x_k) + f'(x_k)(x - x_k)$. Получим соотношение

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}.$$

Пусть $X = (X_1, \dots, X_n)$ — выборка из неизвестного распределения $P \in \{P_\theta \mid \theta \in \Theta\}$, $\Theta \subset \mathbb{R}^d$. Пусть θ^* — ОМП. Хотим приблизить оценку θ^* .

Уравнение правдоподобия: $\frac{\partial l_X(\theta)}{\partial \theta} = 0$. Применим метод Ньютона для функции $l'_X(\theta)$.

$\hat{\theta}_0$ — начальное приближение. Шаг метода:

$$\hat{\theta}_{k+1} = \hat{\theta}_k - \underbrace{(l''_X(\hat{\theta}_k))^{-1}}_{\text{матрица}} \cdot \underbrace{l'_X(\hat{\theta}_k)}_{\text{вектор}}.$$

Теорема: В условиях регулярности L1 — L9, если $\hat{\theta}_0$ — а.н.о, то

1. $\hat{\theta}_1$ — а.н.о с асимт. дисперсией $(i(\theta))^{-1}$.
2. $\hat{\theta}_1$ асимптотически эквивалентна ОМП θ^* , т.е

$$\sqrt{n}(\hat{\theta}_1 - \theta^*) \xrightarrow{P_\theta} 0.$$

Доказательство: (для $d = 1$, идея)

УТВ. (б/д): $\hat{\theta}_1 - \theta^* = (\hat{\theta}_0 - \theta^*)\varepsilon_n(\theta)$, где $\varepsilon_n(\theta) \xrightarrow{P_\theta} 0$.

$$(2). \sqrt{n}(\hat{\theta}_1 - \theta^*) = \sqrt{n}(\hat{\theta}_0 - \theta^*)\varepsilon_n(\theta) =$$

$$= \underbrace{\sqrt{n}(\hat{\theta}_0 - \theta)\varepsilon_n(\theta)}_{\xrightarrow{d_\theta} \mathcal{N}(0, \dots)} + \underbrace{\sqrt{n}(\theta - \theta^*)\varepsilon_n(\theta)}_{\xrightarrow{d_\theta} 0}.$$

По лемме Слуцкого первое слагаемое $\xrightarrow{d_\theta} 0$, второе слагаемое $\xrightarrow{d_\theta} 0$. Применяя еще раз лемму Слуцкого для их суммы, получим $\sqrt{n}(\hat{\theta}_1 - \theta^*) \xrightarrow{d_\theta (\iff P_\theta, \text{т.к const})} 0$.

$$(1). \sqrt{n}(\hat{\theta}_1 - \theta) = \underbrace{\sqrt{n}(\hat{\theta}_1 - \theta^*)}_{\xrightarrow{P_\theta} 0 \text{ (из (2))}} - \underbrace{\sqrt{n}(\hat{\theta}_0 - \theta)}_{\xrightarrow{d_\theta} \mathcal{N}(0, \frac{1}{i(\theta)}) \text{ (ОМП)}}}. \text{ По лемме Слуцкого}$$

$$\sqrt{n}(\hat{\theta}_1 - \theta) \xrightarrow{d_\theta} \mathcal{N}\left(0, \frac{1}{i(\theta)}\right) \quad \square.$$

Замечание: Утверждение теоремы не изменится, если заменить $l''_X(\theta)$ на $E_\theta l''_X(\theta) = -ni(\theta)$, т.е.

$$\hat{\theta}_{k+1} = \hat{\theta}_k + \frac{i(\theta)^{-1}}{n} l'_X(\theta).$$

Оценка $\hat{\theta}_1$ называется *одношаговой оценкой*.

Смысл:

Отклонение $\hat{\theta}_1$ от θ^* на порядок меньше, чем отклонение θ^* от θ . Значит отклонение $\hat{\theta}_1$ от θ тоже имеет порядок $\sqrt{\frac{1/i(\theta)}{n}}$.

Пример (γ -котики):

$\hat{\mu}$ — а.н.о. с асимпт. дисперсией $\pi^2/4 \approx 2.47$. При этом $i(\theta) = 1/2$, т.е. наименьшая возможная асимпт. дисперсия равна 2. Запишем одношаговую оценку:

$$\hat{\theta}_1 = \hat{\mu} + \frac{\sum_{i=1}^n \frac{X_i - \hat{\mu}}{1 + (X_i - \hat{\mu})^2}}{\sum_{i=1}^n \frac{1 - (X_i - \hat{\mu})^2}{(1 + (X_i - \hat{\mu})^2)^2}}.$$

$\hat{\theta}_1$ — наиболее асимптотически эффективная оценка.

2.9. Робастность и симметричные распределения

Пусть $X = (X_1, \dots, X_n)$ — выборка из $\mathcal{N}(\theta, \sigma^2)$, σ известна.

Оценка $\hat{\theta} = \bar{X}$ обладает всеми хорошими свойствами (сильная состоятельность, асимптотическая нормальность, ОМП и т. д.). Однако если в данных есть выбросы, то все свойства теряются.

Для того, чтобы визуализировать выбросы в данных, можно использовать *ящик с усами* (*box plot*).

Будем рассматривать только одномерный случай.

Определение: *Робастная оценка* — оценка, допускающая отклонение от заданной модели.

Определение: Пусть оценка имеет вид $\hat{\theta} = f(X_{(1)}, \dots, X_{(n)})$.

Пусть k_n^* — наименьшее число k , т. ч. выполнено одно из условий:

1. Если $x_1, \dots, x_{k+1} \rightarrow -\infty$, а x_{k+2}, \dots, x_n фиксированы, то $f(x_1, \dots, x_n) \rightarrow -\infty$.
2. Если $x_{n-k}, \dots, x_n \rightarrow +\infty$, а x_1, \dots, x_{n-k+1} фиксированы, то $f(x_1, \dots, x_n) \rightarrow +\infty$.

Тогда число $\tau_{\hat{\theta}} = \lim_{n \rightarrow \infty} \frac{k_n^*}{n}$ называется *асимптотической толерантностью оценки* $\hat{\theta}$.

Смысл: $\tau(\theta)$ — наибольшая доля выбросов, которые способна выдержать оценка, не смещаясь на $\pm\infty$.

Примеры:

- $\bar{X} : k_n^* = 0, \tau_{\bar{X}} = 0$
- $\hat{\mu} : k_n^* = \lceil n/2 \rceil - 1, \tau_{\hat{\mu}} = 1/2$.

Далее будем рассматривать класс распределений $\mathcal{P} = \{P_\theta | \theta \in \Theta\}$, т. ч.

- P_0 имеет плотность $p_0(x)$ — симметричная, непрерывная, носитель плотности имеет вид $(-c, c)$, $0 < c \leq +\infty$.
- θ — параметр сдвига, т. е. $p_\theta(x) = p_0(x - \theta)$.

Будем искать оценки, которые:

1. Достаточно эффективные в классе \mathcal{P} (в асимптотическом подходе).
2. Робастные — допускают отклонение от \mathcal{P} .

1. Усеченное среднее

Определение: Пусть $\alpha \in (0, 1/2)$, $k = \lceil \alpha n \rceil$. Тогда усеченным средним по выборке X_1, \dots, X_n называется оценка

$$\bar{X}_\alpha = \frac{1}{n - 2k} (X_{(k-1)} + \dots + X_{(n-k)}).$$

- $\alpha = 0$: $\bar{X}_\alpha = \bar{X}$
- $\alpha = 1/2$: $\bar{X}_\alpha = \hat{\mu}$.

Асимптотическая толерантность: $\tau_{\bar{X}_\alpha} = \alpha$.

Теорема (б/д): Пусть $X = (X_1, \dots, X_n)$ — выборка из распределения $P \in \mathcal{P}$. Тогда

$$\sqrt{n}(\bar{X}_\alpha - \theta) \xrightarrow{d_\theta} \mathcal{N}(0, \sigma_\alpha^2), \text{ где}$$

$$\sigma_\alpha^2 = \frac{2}{(1 - 2\alpha)^2} \left(\int_0^\theta x^2 p_0(x) dx + \alpha u_{1-\alpha}^2 \right),$$

$u_{1-\alpha}$ — $(1 - \alpha)$ -квантиль распределения P_0 .

Пример: для $\mathcal{N}(0, 1)$

α	0	1/20	1/8	1/4	3/8	1/2
$\text{ARE}_{\bar{X}_\alpha, \bar{X}}$	1	0.99	0.94	0.84	0.74	0.64

При $\alpha = 1/8$ достигается защита от 12.5% загрязнения выборки, но эффективность теряется на 6%.

Утв: Если $D_\theta X_1 < +\infty$, то $\text{ARE}_{\bar{X}_\alpha, \bar{X}} \geq (1 - 2\alpha)^2$.

\triangle \bar{X}_α — а.н.о θ с асимпт. дисперсией σ_α^2 .

Из ЦПТ: \bar{X} — а.н.о θ с асимпт. дисперсией $D_\theta X_1$. Так как дисперсия не зависит от сдвига, посчитаем дисперсию при $\theta = 0$:

$$\frac{1}{2} D_\theta X_1 = \frac{1}{2} \int_{\mathbb{R}} x^2 p_0(x) dx = \int_0^{+\infty} x^2 p_0(x) dx =$$

$$\begin{aligned}
&= \int_0^{u_{1-\alpha}} x^2 p_0(x) dx + \int_{u_{1-\alpha}}^{+\infty} x^2 p_0(x) dx \geq \\
&\geq \int_0^{u_{1-\alpha}} x^2 p_0(x) dx + \underbrace{u_{1-\alpha}^2 \int_{u_{1-\alpha}}^{+\infty} p_0(x) dx}_{=\alpha} = \int_0^{u_{1-\alpha}} x^2 p_0(x) dx + \alpha u_{1-\alpha}^2 = \frac{\sigma_\alpha^2 (1-2\alpha)^2}{2}.
\end{aligned}$$

Отсюда $\text{ARE}_{\overline{X}_\alpha, \overline{X}} = \frac{D_\theta X_1}{\sigma_\alpha^2} \geq (1-2\alpha)^2 \quad \square$.

α	0	1/20	1/8	1/4	3/8	1/2
$(1-2\alpha)^2$	1	0.81	0.56	0.25	0.06	0

При $\alpha = 1/8$ возможна потеря эффективности до 44%.

2. Медиана средних Уолша

Определение: $Y_{ij} = \frac{X_i + X_j}{2}$ — среднее Уолша.

$W = \text{med}\{Y_{ij}, 1 \leq i \leq j \leq n\}$ — медиана средних Уолша.

Теорема: Пусть $X = (X_1, \dots, X_n)$ — выборка из распределения $P \in \mathcal{P}$. Тогда

$$\sqrt{n}(W - \theta) \xrightarrow{d_\theta} \mathcal{N}(0, \sigma^2), \text{ где}$$

$$\sigma^2 = \frac{1}{12 \left(\int_{\mathbb{R}} p_0^2(x) dx \right)^2}.$$

Пример: $\mathcal{N}(0, 1) : \text{ARE}_{W, \overline{X}} \approx 0.955$ (потеря эффективности на 4.5%).

Утверждение: Для $P_\theta \in \mathcal{P}$ $\text{ARE}_{W, \overline{X}} \geq \frac{108}{125} = 0.864$ (в худшем случае теряем 14% эффективности). Равенство достигается при

$$p_0(x) = \frac{3\sqrt{5}}{100} (5 - x^2) I\{|x| < \sqrt{5}\}.$$

Утверждение: $\tau_W \approx 0.293$ (доказательство см. в ДЗ).

Глава 3. Сложные оценки параметров

3.1. Доверительные интервалы

Определение: Пусть $X = (X_1, \dots, X_n)$ — выборка из неизвестного распределения $P \in \{P_\theta \mid \theta \in \Theta\}$.

- Если $\Theta \subset \mathbb{R}$, то пара статистик $(T_1(X), T_2(X))$ называется *доверительным интервалом для θ уровня доверия α* , если

$$\forall \theta \in \Theta \quad P_\theta(T_1(X) \leq \theta \leq T_2(X)) \geq \alpha.$$

- Если $\Theta \subset \mathbb{R}^d$, то статистика $S(X) \subset \Theta$ называется *доверительной областью для θ уровня доверия α* , если

$$\forall \theta \in \Theta \quad P_{\theta}(\theta \in S(X)) \geq \alpha.$$

- Если равенство точное, то интервал называется *точным*.

Замечание:

1. Если $X = (X_1, \dots, X_n)$ — выборка, то утверждение $P_{\theta}(T_1(X) \leq \theta \leq T_2(X)) = \alpha$ имеет смысл ($(T_1(X), T_2(X))$ — доверительный интервал).
2. Если $x = (x_1, \dots, x_n)$ — реализация выборки, то утверждение $P_{\theta}(T_1(x) \leq \theta \leq T_2(x)) = \alpha$ некорректно.

$(T_1(x), T_2(x))$ — реализация доверительного интервала.

Первая магическая константа статистики: $\alpha = 0.95$ (она же 0.05).