Св-ва оценок в модели лин. регр.

при несмещ. и гомоскедастичности шума

Предполагаемая зависимость

$$y(x) = \Theta^T X$$

Наблюдаемая зависимость

$$y(x) = X\Theta + \varepsilon$$

$$y \in \mathbb{R}^n \quad X \in \mathbb{R}^{n \times d} \quad \Theta \in \mathbb{R}^d \quad \varepsilon \in \mathbb{R}^n$$

слуг.   известен   неизв.   слуг.
        неслуг.             неизв.

Задача: оценить $\Theta$

$$RSS(\Theta) = \sum_{i=1}^{n} (y_i - \Theta^T x_i)^2 = \| y - X\Theta \|^2$$

residal sum
of squares

$$\hat{\Theta} = \underset{\theta \in \textcircled{H}}{\arg\min} \; RSS(\theta) \; - \; \text{МНК}$$
оценка

Если $X^T X$ — не выр., то $\hat{\Theta} = (X^T X)^{-1} X^T Y$

Свойства:

1) Если $E\varepsilon = 0$, то $E\hat{\theta} = \theta$

$$E\hat{y}(x) = y(x), \text{ где } \hat{y}(x) = \hat{\theta}^T X$$

[несмещённость шума]

2) Если $D\varepsilon = \sigma^2 I_n$, $E\varepsilon = 0$, то

$$D\hat{\theta} = \sigma^2 (X^T X)^{-1}$$

$$D\hat{y}(x) = \sigma^2 x^T (X^{-1} X)^{-1} x$$

$det(X^{-1}X) \sim 0$ ⟹ большая дисперсия

⟹ больно

[гомоскедастичность шума]

▲ $E\hat{\theta} = E((X^{-1}X)^{-1}X^TY) = (X^TX)^{-1}X^TE(X\theta+\varepsilon) = $ $\boxed{\underset{n\times d}{D}\ \underset{d\times1}{A\xi} = \underset{n\times d}{A}\ \underset{d\times d}{D\xi}\ \underset{d\times n}{A^T}}$

$\qquad = (X^TX)^{-1}X^TX\theta = \theta$

$D\hat{\theta} = D((X^{-1}X)^{-1}X^TY) = (X^TX)^{-1}X^T\ D(X\theta+\varepsilon)\ X(X^TX)^{-1} = $

$\qquad = (X^TX)^{-1}\sigma^2 I_n X^TX(X^TX)^{-1} = \sigma^2(X^TX)^{-1}$

**Утв.** Если $E\varepsilon = 0$ $D\varepsilon = \sigma^2 I_n$, то

$$\text{оценка } \hat{\sigma}^2 = \frac{RSS(\hat{\theta})}{n-d} = \frac{\|y - x\hat{\theta}\|^2}{n-d} \quad - \quad \text{несмещ. оценка } \sigma^2$$

▲ $E \, RSS(\hat{\theta}) = E \sum_{i=1}^{n} (y_i - \hat{\theta}^T x_i)^2 = D \sum_{i=1}^{n} (y_i - \hat{\theta}^T x_i) =$

$$\left[ \begin{array}{c} E(y_i - \hat{\theta}^T x_i) = 0 \\ \| \\ \theta^T x_i + \varepsilon \end{array} \right]$$

$$= tr \, D(y - x\hat{\theta})$$

$$D(y - x\hat{\theta}) = D\left[\left(I_n - \underbrace{x(x^\top x)^{-1} x^\top}_{A \, - \, \text{симметр.}}\right) y\right] = D\left[(I_n - A) y\right] =$$

$\left[\begin{array}{l} A \text{ обратима} \implies \varepsilon \equiv 0 \\ n < d \implies rk \, A < n \implies A \text{ необратима} \end{array}\right.$

$$= (I_n - A) \, D \, y \, (I_n - A)^\top = (I_n - A) \, \sigma^2 I_n \, (I_n - A)^\top =$$

$$= \sigma^2 (I_n - A^\top - A + A A^\top) =$$

$$\left[ A A^\top = x(x^\top x)^{-1} \cancel{x^\top x} (x^\top x)^{-1} x^\top = A \right]$$

$$= \sigma^2 I_n (I_n - A)$$

$$\left[ tr \, A = tr\left( x(x^\top x)^{-1} x^\top \right) = tr\left( x^\top x (x^\top x)^{-1} \right) = tr(I_d) = d \right]$$

св-во следа

$$\implies tr \, D(y - x\hat{\theta}) = tr\left( \sigma^2 (I_n - A) \right) = \sigma^2 (n - d)$$

Пример $\quad X_1 — X_n \sim Exp(\theta) \qquad [$ когда-то было $]$

$$\hat{\theta}_1 = 1/\overline{x}$$

$$\hat{\theta}_2 = -\ln \overline{I\{x > 1\}}$$

сильносост.
и ас. норм.
оценки $\theta$

Какая лучше?

Хотим оценить $\tau(\theta) \in \mathbb{R}^d$

**Опр.** Ф-ция $L: \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}_+$, характеризующая степень откло-нения оценки от $\tau(\theta)$ наз-ся **функцией потерь**

Примеры:

1) $d = 1$ $\quad L(x, y) = (x - y)^2$ — квадратичная

2) $d = 1$ $\quad L(x, y) = |x - y|$ — абсолютная

3) $d = 1$ $\quad L(x, y) = \ln(1 + |x - y|)$

4) $d > 1$ $\quad L(x,y) = (x-y)^T A (x-y)$ ,

где $A$ — сим. неотр. полуопр.

$A = I_d$ , тогда $\quad h(x,y) = \sum\limits_{i=1}^{d} (x_i - y_i)^2$

Если $\hat{\theta}$ — оценка $\tau(\theta)$, то при таком оценивании

$$L(\hat{\theta}, \tau(\theta)) - \text{штраф}$$

При таком подходе есть недостаток:

штраф случаен, при разных реализациях получаем разные штрафы

Опр. **Ф-ия риска** оценки $\hat{\theta}$ величины $\tau(\theta)$

$$R_{\hat{\theta}, \tau(\theta)}(\theta) = E_\theta L(\hat{\theta}, \tau(\theta)) \quad \leftarrow \text{ф-ия от } \theta$$

(risk function)

Если $L(x, y) = (x - y)^2$ — кв. ф-ция потерь, то

$$MSE_{\hat{\theta}, \tau(\theta)} = E_\theta (\hat{\theta} - \tau(\theta))^2 \text{ — средняя квадратичная ошибка}$$

$$MAE_{\hat{\theta}, \tau(\theta)} = E_\theta |\hat{\theta} - \tau(\theta)| \text{ — средняя абсолютная ошибка}$$
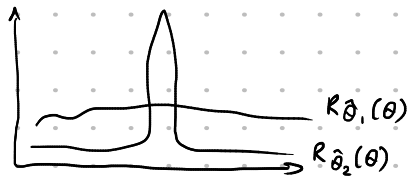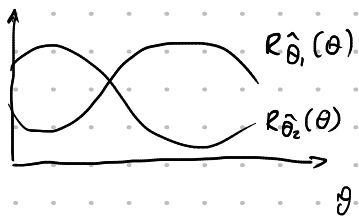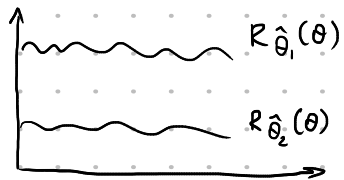
Замечание:

Если $\tau(\theta) = \theta$, то $\tau$ будем опускать

Пример    $X_1 \ldots X_n$ — выборка    $E\,X_1 = \theta$
$D\,X_1 < +\infty$

Посчитаем    MSE    для    $\underset{\overset{\shortparallel}{\hat\theta_1}}{X_1}$,    $\underset{\overset{\shortparallel}{\hat\theta_2}}{\overline{X}}$    — оц. $\theta$

$$MSE_{\hat\theta_1}(\theta) = E_\theta \big( X_1 - \underbrace{\theta}_{E X_1} \big)^2 = D\,X_1$$

$$MSE_{\hat\theta_2}(\theta) = E_\theta \big( \overline{X} - \underbrace{\theta}_{E\overline{X}} \big)^2 = D\,\overline{X} = \frac{D X_1}{n}$$

$R_{\hat{\theta}_1}(\theta)$

$R_{\hat{\theta}_2}(\theta)$

$R_{\hat{\theta}_1}(\theta)$

$R_{\hat{\theta}_2}(\theta)$

$\vartheta$

$R_{\hat{\theta}_1}(\theta)$

$R_{\hat{\theta}_2}(\theta)$

Что лучше ?

# Подходы к сравнению оценок

1) Равномерный подход

- $\hat{\Theta}_1$ не хуже $\hat{\Theta}_2$, если $\forall \theta \in \Theta \quad R_{\hat{\Theta}_1}(\theta) \leq R_{\hat{\Theta}_2}(\theta)$

- $\hat{\Theta}_1$ лучше $\hat{\Theta}_2$, если $\hat{\Theta}_1$ не хуже $\hat{\Theta}_2$

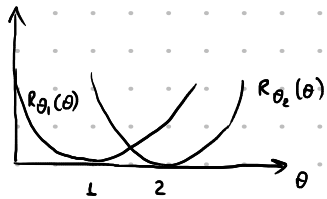$$\text{и} \quad \exists \theta \in \Theta \quad R_{\Theta_1}(\theta) < R_{\Theta_2}(\theta)$$

- $\hat{\Theta} \in \mathcal{K}$ — класс оценок, $\hat{\Theta}$ — наилучшая в $\mathcal{K}$, если

    она лучше любой другой

- Если $L(x,y) = (x-y)^2$ — ф-ция потерь

    то говорят о среднеквадратичном подходе

Утв. В классе $\mathcal{K}$ может не быть наилучшей оценки.

Пример $\quad \textcircled{\text{H}} = \mathbb{R} \qquad \mathcal{K} = \{ \hat{\Theta}_1 = 1, \ \hat{\Theta}_2 = 2 \}$

**Утв.** MSE допускает *bias-variance* разложение

$$MSE_{\hat{\theta}, \tau}(\theta) = E_\theta(\hat{\theta} - \tau(\theta))^2 = \underbrace{D\,\hat{\theta}}_{variance} + \underbrace{(E_\theta\,\hat{\theta} - \tau(\theta))^2}_{bias^2}$$

**Следствие:**

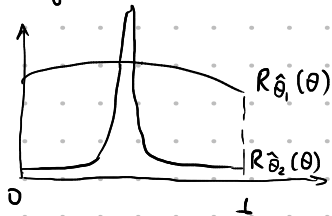наилучший по MSE в классе $\mathcal{K} = \{$ несмещ. оц. $\}$ будет

оценка с наим. дисперсией

$$E_\theta(\hat{\theta} - \tau(\theta))^2 = \underbrace{E_\theta(\hat{\theta} - E_\theta\,\hat{\theta})^2}_{D\hat{\theta}} + E\underbrace{(E_\theta\,\hat{\theta} - \tau(\theta))^2}_{const} + 2\,E_\theta\,\overbrace{(\hat{\theta} - E_\theta\,\hat{\theta})}^{0\ \text{т.к. несмещ.}}\overbrace{(\hat{\theta} - \tau(\theta))}^{const} =$$

$$= D\,\hat{\theta} + (E_\theta\,\hat{\theta} - \tau(\theta))^2$$

2) Байесовский подход

Пусть $Q$ — распределение на $\textcircled{H}$

- $\hat{\Theta}_1$ не хуже $\hat{\Theta}_2$, если $E_Q R_{\hat{\Theta}_1, \tau}(\Theta) \leq E_Q R_{\hat{\Theta}_2, \tau}(\Theta)$



$Q = U[0, 1]$

при таком $Q$ мы по факту считаем интеграл

$\hat{\Theta}_2$ лучше $\hat{\Theta}_1$

3) Минмаксный подход

· $\hat{\theta}_1$ не хуже $\hat{\theta}_2$, если $\sup\limits_{\theta \in \textcircled{H}} R_{\hat{\theta}_1}(\theta) \leq \sup\limits_{\theta \in \textcircled{H}} R_{\hat{\theta}_2}(\theta)$

4) Асимптотический подход (для ас. норм. оц., $\mathbb{R}$)

$\hat{\theta}_1, \hat{\theta}_2$ — а.н.о. с ас.д. $\sigma_1^2, \sigma_2^2$

· $\hat{\theta}_1$ не хуже $\hat{\theta}_2$, если $\forall \theta \in \textcircled{H} \quad \sigma_1^2(\theta) \leq \sigma_2^2(\theta)$

· $\hat{\theta}_1$ лучше $\hat{\theta}_2$, если $\qquad \downarrow \quad + \quad \exists \theta \in \textcircled{H} \quad \sigma_1^2(\theta) < \sigma_2^2(\theta)$

· $ARE_{\hat{\theta}_1, \hat{\theta}_2} = \dfrac{\sigma_2^2(\theta)}{\sigma_1^2(\theta)}$    относительная асимптотическая эффективность

показывает насколько $\hat{\theta}_1$ лучше $\hat{\theta}_2$

Опр $\hat{\theta}$ наз-ся *асимптотически эффективной*, если она наилучшая

в классе всех а.н.о. с непр. а.д. в асимпт. подходе

вроде без непр не будет наил оц. в классе

[ т.е. наим. ас диcп. среди всех а.н.о.]

[L1 — L8]  ОМП — ас. эф. оценка, т.е. $i(\theta)^{-1}$ — наим ас. д.

матрички сравнивать страшно

этого мы делать не будем

Пример $\quad X_1 - X_n \sim \mathcal{N}(\theta, 1)$

ЦПТ: $\overline{X}$ — а.н.о. $\theta$ с а.д. $\sigma^2(\theta)$

т. о выборочной медиане $\hat{\mu}$ — а.н.о. $\theta$ с ас.д. $\sigma_2^2(\theta) = \dfrac{\pi}{2}$

$ARE_{\overline{X}, \mu} = \dfrac{\pi}{2} \approx 1.57$, т.е. $\widetilde{x}$ в $\sim 1.57$ раза лучше $\mu$

и/б здесь $\dfrac{\pi^2}{4}$

↓

Пример $X_1 \dots X_n \sim Exp(\theta)$

[раньше было]

$\hat{\theta}_1 = \dfrac{1}{\overline{x}}$   а.н.о. с   а.д.    $\sigma_1^2(\theta) = \theta^2$

$\hat{\theta}_2 = -\ln \overline{I\{X > 1\}}$   а.н.о. с   а.д.    $\sigma_2^2(\theta) = e^\theta - 1$

$ARE_{\hat{\theta}_1, \hat{\theta}_2} = \dfrac{e^\theta - 1}{\theta^2}$

# Достаточные статистики

Опр   Пусть  $X = (X_1 \dots X_n)$ — выборка   из  $P \in \mathcal{P}$

Статистика  $S(x)$  наз-ся  _достаточной_  для  семейства  $\mathcal{P}$,

если  условное  распр.  $P_\theta(X \in B \mid S(x))$  не зависит  от  $\theta$  $\forall B$

Смысл:   вся  информация  о  $\theta$,  содержащаяся  в  выборке

содержится  в  достаточной  статистике

Тривиальный  пример :  $S(x) = (X_1 \dots X_n)$

Следствие :  Если  данные  поступают  последовательно,  то  достаточно

хранить  только  $S(x)$

Важным является случай, когда размерность достаточной статистики меньше размера выборки. $|S(x)| \ll n$

$$\left[\begin{array}{c} \text{на самом деле такое } S(x) \text{ не всегда сущ.} \\ \text{например для распр. Коши} \end{array}\right.$$

Пример     $X_1 \dots X_n \sim Bern(\theta)$

Какая информация есть в выборке?

Кол-во успехов     $S(x) = \sum\limits_{i=1}^{n} X_i$

Порядок успех-неудач (какой-то уёбищный порядок без количеств)

$S(x)$ — достаточная

△ $P_\theta\left(X = \bar{x} \mid S(X) = s\right) = \dfrac{P_\theta\left(X_1 = x_1 \dots X_n = x_n, \ \sum X_i = s\right)}{P_\theta\left(\sum X_i = s\right)} = \begin{bmatrix} \text{если } \sum x_i \neq s \\ \text{то } 0 \end{bmatrix}$

$= \dfrac{\theta^{\sum x_i}(1-\theta)^{n-\sum x_i} \ \mathbb{I}\{\sum x_i = s\}}{\theta^s(1-\theta)^{n-s} \ C_n^s} = \dfrac{\mathbb{I}\{\sum x_i = s\}}{C_n^s} \leftarrow$ не зависит от $\theta$

Теор. ( Критерий факторизации Неймана - Фишера)

Пусть $X = (X_1 \ldots X_n)$ — выборка из распр $P \in \mathcal{P}$ — домин.
                                                         семейство

Тогда $S(X)$ — дос. стат. $\iff$ справедлива факторизация

$$\rho_\theta(X) = \psi(S(X), \theta) \cdot \underbrace{h(x)}_{\substack{\text{не зависит} \\ \text{от } \theta}}$$

Пример $\quad X_1 \dots X_n \sim \Gamma(\alpha, \beta) \qquad \theta = (\alpha, \beta)$

$$\rho_\theta(x) = \frac{\alpha^\beta}{\Gamma(\beta)} x^{\beta-1} e^{-\alpha x}$$

$$\rho_\theta(x_1 \dots x_n) = \frac{\alpha^{\beta n}}{\Gamma(\beta)^n} \left(\prod_{i=1}^n x_i\right)^{\beta-1} e^{-\alpha \sum x_i} = \psi\left(\underset{u}{S(x)}, \underset{\prime\prime}{\theta}\right) \cdot h(x)$$

$$(\Sigma x_i, \Pi x_i) \quad (\alpha, \beta)$$

$$\psi(x, y, \alpha, \beta) = \frac{\alpha^{\beta n}}{\Gamma(\beta)^n} y^{\beta-1} e^{-\alpha x}$$

$$h(x) = 1$$

$$S(x) = \left(\Sigma x_i, \Pi x_i\right)$$

$$S(x) = \left(\Sigma x_i, \Sigma \ln x_i\right)$$

$\Downarrow$

достаточные
статистики