

Семинар - 9.

Линейная регрессия. Доверительные интервалы

Гауссовская линейная регрессия.

$$y = \sum_{j=1}^p \psi_j(x) \theta_j + \varepsilon_i = \psi(x)^T \theta + \varepsilon_i, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

Пусть имеются данные $\{(X_i, Y_i) : 1 \leq i \leq n\}$, причем наблюдения $Y_1, \dots, Y_n \in \mathbb{R}$ получены из модели

$$\mathbf{Y} = \Psi^T \theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_n),$$

где $\theta \in \mathbb{R}^p$, $\Psi = (\psi(X_1), \dots, \psi(X_n)) \in \mathbb{R}^{p \times n}$ – матрица плана полного ранга, то есть $\text{rank}(\Psi) = p$, X_1, \dots, X_n детерминированные, а $\varepsilon_1, \dots, \varepsilon_n$ – независимые центрированные случайные величины с дисперсией σ^2 . Мы предполагаем, что матрица $\Psi \in \mathbb{R}^{p \times n}$ имеет полный ранг,

Оценка в методе наименьших квадратов (ОМНК):

$$\hat{\theta}_{LSE} = \arg \min_{\theta \in \Theta} \sum_{i=1}^n (Y_i - \psi(X_i)^T \theta)^2 = \arg \min_{\theta \in \Theta} \|\mathbf{Y} - \Psi^T \theta\|^2 = (\Psi \Psi^T)^{-1} \Psi \mathbf{Y}.$$

Введем ортопроектор на линейное подпространство размерности p , натянутое на строки матрицы Ψ : $(\psi_j(X_1), \dots, \psi_j(X_n))^T \in \mathbb{R}^n$, $1 \leq j \leq n$

$$\Pi = \Psi^T (\Psi \Psi^T)^{-1} \Psi \in \mathbb{R}^{n \times n}.$$

Тогда оценка регрессионной функции $\psi(x)^T \theta$ в точках X_1, \dots, X_n :

$$\hat{f} = \Psi^T \hat{\theta}_{LSE} \equiv \Pi \mathbf{Y}.$$

Утверждение 1 (Геометрическая интерпретация \hat{f}). \hat{f} является проекцией вектора наблюдений на линейную оболочку строк матрицы плана Ψ , то есть на $\text{Im}(\Psi^T)$.

Лемма (Свойство гауссовых векторов). Пусть вектор X имеет нормальное распределение $\mathcal{N}(\mu, \Sigma_n)$. Тогда для произвольной матрицы $B \in \mathbb{R}^{m \times n}$ вектор BX имеет нормальное распределение $\mathcal{N}(B\mu, B\Sigma_n B^T)$.

Лемма. Пусть вектор X имеет стандартное нормальное распределение $\mathcal{N}(0, I_n)$ и пусть B – ортопроектор. Тогда $Q = X^T BX$ имеет распределение хи-квадрат $\chi^2(\text{tr } B)$.

Задачи

- Найти оценку неизвестного параметра θ методом максимального правдоподобия. Будет ли она несмешенной? Будет ли она оптимальна? Сравнить с оценкой наименьших квадратов в модели линейной регрессии (то есть когда вектор ошибок не обязательно гауссовский).
- Найти оценку максимального правдоподобия для σ^2 в гауссовой линейной регрессии. Будет ли она несмешенной?
- Доказать, что $\frac{1}{\sigma^2} \|\Pi \varepsilon\|^2$ имеет распределение $\chi^2(p)$. - не успели
- Доказать, что $\hat{\theta}_{LSE}$ и $\mathbf{Y} - \Psi^T \hat{\theta}_{LSE}$ независимы. Доказать, что $\frac{1}{\sigma^2} \|\mathbf{Y} - \Psi^T \hat{\theta}\|^2$ имеет распределение хи-квадрат $\chi^2(n-p)$. - не успели.
- Постройте точный $1 - \alpha$ -доверительный интервал для θ_j , где $j = \{1, \dots, p\}$, если
 - σ^2 известен;
 - σ^2 неизвестна.
- Постройте точное $1 - \alpha$ -доверительное множество для θ , если σ^2 неизвестна.
- Постройте точный $1 - \alpha$ -доверительный интервал для σ^2 .

1. Найти оценку неизвестного параметра θ методом максимального правдоподобия. Будет ли она несмещенной? Будет ли она оптимальна? Сравнить с оценкой наименьших квадратов в модели линейной регрессии (то есть когда вектор ошибок не обязательно гауссовский).

$$Y = \Psi^T \theta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

$\uparrow X$ вектор

$$Y_1 = \Psi(X_1)^T \theta + \varepsilon_1$$

$$\begin{matrix} \uparrow & \uparrow \\ R^P & R^P \end{matrix}$$

$$= \sum_{j=1}^P \Psi_j \theta_j$$

$$\Psi = \begin{pmatrix} \Psi_1 \\ \Psi_2 \\ \vdots \\ \Psi_P \end{pmatrix} = \begin{pmatrix} 1 \\ X \\ X^2 \\ X^3 \\ X^4 \end{pmatrix}$$

$$\Psi \in \mathbb{R}^{P \times n}$$

$$= \mathbb{R}^{P \times L}$$

$$X_1, \dots, X_n \sim P_\theta$$

$$(Y_1, \dots, Y_n) \sim N(\Psi^T \theta, \sigma^2)$$

$$L(\theta) = \prod_{i=1}^n p_\theta(Y_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(Y_i - \Psi(X_i)^T \theta)^2}{2\sigma^2}}$$

$$Y_i \sim \Psi(X_i)^T \theta + N(0, \sigma^2) \sim N(\Psi(X_i)^T \theta, \sigma^2)$$

$$\ln L(\theta) = \sum_{i=1}^n \left(\underbrace{\frac{1}{\sqrt{2\pi\sigma^2}}}_{C} - \frac{(Y_i - \Psi(X_i)^T \theta)^2}{2\sigma^2} \right)$$

$$\hat{\theta} = ?$$

$$= nC - \frac{1}{2\sigma^2} \sum (Y_i - \Psi(X_i)^T \theta)^2$$

$$= nC - \frac{1}{2\sigma^2} \|Y - \Psi^T \theta\|_2^2 \quad \leftarrow \text{переведем в векторной форме}$$

$$df = \langle \nabla f(\theta), d\theta \rangle$$

$$L(\theta) = \text{Const} \cdot \frac{1}{\sqrt{2\sigma^2}} \langle Y - \Psi^T \theta, Y - \Psi^T \theta \rangle$$

$$dL(\theta) = \frac{1}{\sigma^2} \langle \underbrace{\Psi^T \theta}_1, Y - \Psi^T \theta \rangle$$

$$= \underbrace{\langle \Psi(Y - \Psi^T \theta), d\theta \rangle}_{\sigma^2} \quad i = \nabla f(\theta)$$

$$\langle Ax, y \rangle = \langle x, A^T y \rangle$$

$$\left. \frac{\partial f(\theta)}{\partial \theta} \right|_{\theta=\hat{\theta}} = \Psi(Y - \Psi^T \hat{\theta}) = 0$$

$$\Psi Y = (\Psi \Psi^T)^{-1} \Psi Y = \hat{\theta}$$

Несмешенная оцк.

$$(\Psi \Psi^T) \Psi E Y = \theta$$

$$\Psi^T \theta$$

2. Найти оценку максимального правдоподобия для σ^2 в гауссовой линейной регрессии. Будет ли она несмешенной?

$$\ln L(\theta) = \sum_{i=1}^n \left(\ln \left(\frac{1}{\sqrt{2\pi}\sigma^2} \right) - \frac{(Y_i - \Psi(X_i)^T \theta)^2}{2\sigma^2} \right) = \frac{n}{2} \left(n \cdot \ln(2\pi) + \ln \sigma^2 \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \Psi(X_i)^T \theta)^2$$

$$\frac{\partial}{\partial \sigma^2} \ln L(\theta) = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \|Y - \Psi^T \theta\|_2^2 = 0 \quad \left| \frac{\sigma^2}{n} \right.$$

$$\sigma^2 = \frac{1}{n} \|Y - \Psi^T \theta\|_2^2$$

т.ч. несмешн. оценк.

$$\hat{\sigma}^2 = \frac{1}{n} \|Y - \Psi^T \theta\|_2^2 = \frac{1}{n} \|Y - \Pi Y\|_2^2 = \frac{1}{n} \|(I - \Pi)Y\|_2^2$$

$$\Pi = \Psi^T (\Psi \Psi^T)^{-1} \Psi$$

5. Постройте точный $1 - \alpha$ -доверительный интервал для θ_j , где $j = \{1, \dots, p\}$, если

- (a) σ^2 известен;
- (b) σ^2 неизвестна.

$$\hat{\theta} \sim \mathcal{N}\left(\underbrace{(\Psi\Psi^\top)^{-1}\Psi\Psi^\top\theta}_{\text{см. теор. с лекции}}, \underbrace{(\Psi\Psi^\top)^{-1}\Psi\Psi^\top(\Psi\Psi^\top)^{-1}\sigma^2}_{\sim N(\theta, (\Psi\Psi^\top)^{-1}\sigma^2)}\right)$$

$$Y \sim \mathcal{N}(\Psi^\top\theta, \sigma^2 I_n)$$

$$G = \frac{(\hat{\theta}_j - \theta_j)}{\sqrt{(\Psi\Psi^\top)_{jj}} \hat{\sigma}} \sim N(0, 1)$$

$$\hat{\theta}_j \sim \mathcal{N}(\theta_j, (\Psi\Psi^\top)_{jj} \hat{\sigma}^2)$$

$\delta)$ $\hat{\sigma}^2$ не известна.

$$\hat{\sigma}^2 = \frac{1}{n-p} \|(\mathbf{I} - \mathbf{P})Y\|_2^2$$

$$\frac{\hat{\theta}_j - \theta_j}{\sqrt{(\Psi\Psi^\top)_{jj}} \hat{\sigma}} \sim t_{n-p} \quad \text{расп.}$$

$$\frac{1}{\hat{\sigma}^2} \|(\mathbf{I} - \mathbf{P})Y\|_2^2 \sim \chi^2(n-p) \quad \text{расп.}$$

trace = n trace P . \uparrow загадка "