

fit, predict, и всё готово

Линейная регрессия

Чертоги разума

Модель:

$$\mathbf{X} = Z\boldsymbol{\theta} + \boldsymbol{\varepsilon} \quad \text{или, что то же самое,} \quad X_i = \sum_{j=1}^k z_{ij}\theta_j + \varepsilon_i, \quad i \in \{1, \dots, n\}.$$

где:

- $Z \in \mathbb{R}^{n \times k}$ — фиксированная матрица признаков: z_{ij} — j -ый признак для i -ого объекта;
- $\boldsymbol{\theta} \in \mathbb{R}^k$ — вектор неизвестных параметров;
- $\boldsymbol{\varepsilon}$ — случайный n -мерный вектор (шум);
- \mathbf{X} — n -мерный (тоже случайный) вектор целевых величин (таргетов): X_i — таргет i -ого объекта.

Ограничения, накладываемые на модель:

L1 Для всех i выполнено $E\varepsilon_i = 0$ (в среднем ошибки нет);

L2 Дисперсия ε_i одинакова и равна неизвестному параметру σ^2 , причём ε_i попарно нескоррелированы, то есть $D\varepsilon = \sigma^2 E_n$ (наблюдения друг на друга не влияют).

L3 Столбцы $\mathbf{z}_1, \dots, \mathbf{z}_k$ матрицы Z линейно независимы.

Определение. Оценка $\hat{\boldsymbol{\theta}} = (Z^T Z)^{-1} Z^T \mathbf{X}$ называется *оценкой по методу наименьших квадратов* (или *MНK-оценкой*).

Свойства МНК-оценки.

- $E_{\boldsymbol{\theta}, \sigma^2} \hat{\boldsymbol{\theta}} = \boldsymbol{\theta}$, $D_{\boldsymbol{\theta}, \sigma^2} \hat{\boldsymbol{\theta}} = \sigma^2 (Z^T Z)^{-1}$;
- $\hat{\boldsymbol{\theta}}$ является лучшей оценкой в среднеквадратичном подходе в классе несмешённых линейных (по \mathbf{X}) оценок (*теорема Гаусса-Маркова*);
- Статистика $\hat{\sigma}^2 = \frac{1}{n-k} \|\mathbf{X} - Z\hat{\boldsymbol{\theta}}\|^2$ несмешённо оценивает σ^2 .

В модели **гауссовской линейной регрессии**, когда $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, E_n)$, справедлив следующий результат:

$$\hat{\boldsymbol{\theta}} \sim \mathcal{N}(\boldsymbol{\theta}, \sigma^2 (Z^T Z)^{-1}), \quad \frac{1}{\sigma^2} \|\mathbf{X} - Z\hat{\boldsymbol{\theta}}\|^2 \sim \chi^2_{n-k}, \quad \hat{\boldsymbol{\theta}} \perp \hat{\sigma}^2.$$

На основании этих фактов можно получить статистику Стьюдента для линейной комбинации параметров $\mathbf{c}^T \boldsymbol{\theta}$, по которой построить точный ДИ:

$$\frac{\mathbf{c}^T \hat{\boldsymbol{\theta}} - \mathbf{c}^T \boldsymbol{\theta}}{\sqrt{\hat{\sigma}^2 \mathbf{c}^T (Z^T Z)^{-1} \mathbf{c}}} \sim T_{n-k}.$$

Esercizi

1. Пусть в модели линейной регрессии исходный вектор параметров $\boldsymbol{\theta}$ линейно выражается через вектор параметров $\boldsymbol{\psi}$, то есть $\boldsymbol{\theta} = S\boldsymbol{\psi}$, где $S \in GL_k(\mathbb{R})$. Покажите, что $\widehat{\boldsymbol{\theta}} = S\widehat{\boldsymbol{\psi}}$.

2. Пусть X_1, \dots, X_n — независимые случайные величины, представимые в виде

$$X_i = a + bw_i + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

где w_i — фиксированные константы. Постройте МНК-оценки для параметров a и b .

3. **Недоопределение.** В модели линейной регрессии $\mathbf{X} = Z\boldsymbol{\theta} + \boldsymbol{\varepsilon}$, где $\mathbf{X} \in \mathbb{R}^n$, $Z \in \mathbb{R}^{n \times k}$, $\theta_i \neq 0$, МНК-оценка $\widehat{\boldsymbol{\theta}}$ вектора $(\theta_1, \dots, \theta_m)^T$ строится только по первым $m < k$ столбцам матрицы Z (статистик думает, что целевая величина зависит лишь от первых m признаков, когда на самом деле от всех k). Докажите, что полученная оценка может оказаться смещённой.

4. **Переопределение.** В модели линейной регрессии $\mathbf{X} = Z\boldsymbol{\theta} + \boldsymbol{\varepsilon}$, где $\mathbf{X} \in \mathbb{R}^n$, $Z \in \mathbb{R}^{n \times k}$, МНК-оценка $\widehat{\boldsymbol{\theta}}$ вектора $(\theta_1, \dots, \theta_m)^T$ строится сразу по $m > k$ признакам с матрицей

$$Z' = (\mathbf{z}_1, \dots, \mathbf{z}_k, \mathbf{z}_{k+1}, \dots, \mathbf{z}_m) = (Z, \dots)$$

(статистик думает, что целевая величина зависит от всех m признаков, когда на самом деле только от первых k , то есть $\theta_{k+1} = \dots = \theta_m = 0$). Докажите, что для всех $i = 1, \dots, k$ выполнено неравенство $D_{\boldsymbol{\theta}}\widehat{\theta}_i \geq D_{\boldsymbol{\theta}}\widetilde{\theta}_i$, где $\boldsymbol{\theta}$ — МНК-оценка, построенная по матрице Z (иными словами, добавление бесполезных признаков увеличило дисперсию оценок действительно нужных параметров).

5. Пусть $X_1, \dots, X_n \sim \mathcal{N}(a, \sigma^2)$ и $Y_1, \dots, Y_m \sim \mathcal{N}(b, \sigma^2)$ — независимые выборки. Сведя задачу к линейной регрессии, постройте точный доверительный интервал для $a - b$.

6. Постройте точный доверительный интервал уровня доверия γ для целевой величины нового объекта с признаками \mathbf{z}_0 (то есть $\mathbf{z}_0^T \boldsymbol{\theta} + \varepsilon$, где $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ не зависит от $\boldsymbol{\varepsilon}$) для модели гауссовской линейной регрессии $\mathbf{X} = Z\boldsymbol{\theta} + \boldsymbol{\varepsilon}$.

Замечание. ДИ для случайной величины называют *предсказательным интервалом*.

7. Докажите, что статистика $(\widehat{\boldsymbol{\theta}}, \|\mathbf{X} - Z\widehat{\boldsymbol{\theta}}\|^2)$ является **(а)** достаточной; **(б)** полной в модели гауссовской линейной регрессии.

8. В модели линейной регрессии $\mathbf{X} = Z\boldsymbol{\theta} + \boldsymbol{\varepsilon}$, где $E\boldsymbol{\varepsilon} = \mathbf{0}$, $D\boldsymbol{\varepsilon} = \sigma^2 V$ (V — п.о. матрица) найдите лучшую оценку среди линейных несмешённых в среднеквадратичном подходе.