

Бутстреп

$X = (X_1, \dots, X_n) \sim P$ независ.

$T(X)$ — статистика

Пример

$$T(X) = \frac{\overline{\cos X}}{\bar{X}}$$

Хотим оценить дисперсию $T(X)$

Применим бутстреп

Схема :

1) Возьмём $j_1^i - j_n^i \sim U\{1, \dots, n\}$

2) Рассмотрим бутстрепную выборку $X_1^* = (X_{j_1^i} - X_{j_n^i})$
 $X_B^* = (X_{j_1^B} - X_{j_n^B})$

3) Считаем $T(X_1^*) - T(X_B^*)$

4) Оцениваем искомую величину на $T(X_1^*) - T(X_B^*)$

дисперсия :
$$\hat{\sigma}(T(X)) = \frac{1}{n} \sum T_i^2 - \left(\frac{1}{n} \sum T_i \right)^2$$

Пример

Пусть

$$X = (5, 3, 2, 8, 0, 2, 7)$$

$$T(X) = \frac{\cos X}{\bar{X}}$$

$$X_1^* = (2, 2, 7, 5, 3, 5, 0)$$

$$T_1 = 0.02$$

$$X_2^* = (0, 3, 3, 5, 7, 8, 0)$$

$$T_2 = 0.19$$

$$X_3^* = (2, 2, 8, 2, 2, 8, 0)$$

$$T_3 = -0.04$$

$$X_4^* = (3, 2, 2, 8, 3, 2, 2)$$

$$T_4 = 0.318$$

Почему размер бутстрепной выборки = n ?

$$T(X) = \bar{X}$$

Хотим дисперсию.

$$D\bar{X} = D \frac{1}{n} \sum x_i = \frac{1}{n^2} D \sum x_i = \frac{1}{\underline{n}} D x_i$$

Пусть $T_1 \dots T_B$ — статистики по \hat{X}_i размера k

$$DT^* = D\bar{X}^* = \frac{1}{\underline{k}} D x_i$$

Если $T(X)$ «не зависит от n », то можно брать любого размера бутстрепную выборку.

Задача

x_1, \dots, x_n — выборка

x_1^*, \dots, x_n^* — бутстрепная выборка

Найти $E \bar{x}^*$

Бутстреп: $x_i^* = x_1 I\{j_1=1\} + \dots + x_n I\{j_1=n\}$

$$\begin{aligned} E(x_i^* | x_1, \dots, x_n) &= E(x_1 I\{j_1=1\} + \dots + x_n I\{j_1=n\} | x_1, \dots, x_n) = \\ &= \sum_{i=1}^n E(x_i I\{j_1=i\} | x_i) = \sum_{j_1 \in (x_1, \dots, x_n)} x_i E I\{j_1=i\} = \\ &= \frac{1}{n} \sum x_i = \bar{x} \end{aligned}$$

$$E(\bar{x}^* | x_1, \dots, x_n) = E\left(\frac{1}{n} \sum_{i=1}^n (x_i^* | x_1, \dots, x_n)\right) = \frac{1}{n} \sum_{i=1}^n E(x_i^* | x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n \bar{x} = \bar{x}$$

$$E(\bar{X}^*) = E(E(\bar{X}^* | X_1, \dots, X_n)) = \frac{1}{n} \sum_{i=1}^n EX_i = EX_1$$

$$E(\bar{X}) = EX_1$$

Каким брать B ? [кол-во бутстреппных выборок]

$$\lambda = (x_1, \dots, x_n) \sim P$$

1) Хотим сгенерировать $x' \sim P$

Для этого используем метод подстановки и будем генерировать из \hat{P} — эмпирическое распределение P по x_1, \dots, x_n

2) Когда оцениваем дисперсию \Rightarrow используем Монте-Карло

оцениваемая величина

$$\theta \approx \hat{\theta} \approx \tilde{\theta}$$

↑
метод
подстановки

Монте-Карло ск. ск. $\frac{1}{\sqrt{B}}$

ск. ск. $\frac{1}{\sqrt{n}}$ по т. Колмогорова - Смирнова: $\sup_A \sqrt{n} |P(A) - P_n(A)| \rightarrow \zeta \} \Rightarrow \frac{1}{\sqrt{B}} = o\left(\frac{1}{\sqrt{n}}\right)$
т.е. $B \geq n^{1+\epsilon}$

Метод подстановки и метод моментов

Вспомним метод моментов

$$X = (X_1, \dots, X_n) \sim P \in \{P_\theta \mid \theta \in \Theta\}$$

Пусть целевые функции g_1, \dots, g_d

Составляем систему

$$m(\theta) = \begin{pmatrix} E_\theta g_1(X) \\ \vdots \\ E_\theta g_d(X) \end{pmatrix} = \begin{pmatrix} \overline{g_1(X)} \\ \vdots \\ \overline{g_d(X)} \end{pmatrix}$$

Пусть P — непарам.

Введём $G(P) = m^{-1}(E_P g_1(X), \dots, E_P g_d(X))$

$$\forall \theta \in \Theta \quad G(P_\theta) = m^{-1}(m(\theta)) = \theta$$

Тогда $\hat{\theta}$ — оценка по методу подстановки

$$\hat{\theta} = G(\hat{P}_{\theta}) = m^{-1}(\overline{g_1(x)}, \dots, \overline{g_d(x)})$$

Критерий хи-квадрат

Пусть $X = (X_1, \dots, X_n) \sim P$

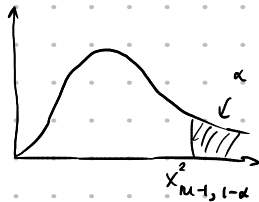
$$H_0: P = P_0 \quad \text{vs.} \quad H_1: P \neq P_0$$

Пусть $\mathcal{X} = \bigcup_{j=1}^m B_j$, тогда $\mu_j = \# \{ i \mid X_i \in B_j \}$

$$p_j^0 = P_0(X_i \in B_j)$$

Рассмотрим $\chi(X) = \sum_{j=1}^m \frac{(\mu_j - np_j^0)^2}{np_j^0} \xrightarrow{d} \chi_{m-1}^2$

Критерий $\{ \chi(X) > \chi_{m-1, 1-\alpha}^2 \}$



Пример $X = (X_1, \dots, X_n)$ - выборка, сгенерированная псевдослучайным генератором чисел

$$H_0: P = U\{1, 2, 3, 4\} \text{ vs. } H_1: P \neq U\{1, 2, 3, 4\}$$

Проверить генератор.

Эксперимент :

#1	=	249
#2	=	254
#3	=	246
#4	=	251
Σ		1000

$$B_j = \{j\} \quad j \in \overline{1, 4}$$

$$P_j^0 = 1/4$$

$$\Rightarrow \chi(X) = \frac{1 + 4^2 + 4^2 + 1}{250} = \frac{34}{250} \approx 0,136$$

Где отвергать? Если генератор даёт 1 2 3 4 по кругу, то
или его не отвергнем

\Rightarrow отвергаем справа и около 0, т.к. и генератор слишком
идеальный

Критерий: $S = \{ \chi(x) < \chi^2_{3, \alpha/2} \} \cup \{ \chi(x) > \chi^2_{3, 1-\alpha/2} \}$

Численные значения:

$$\chi(x) = 0.136$$

$$\chi^2_{3, 1-\alpha} = 7.81$$

\Rightarrow не отвергаем

$$\chi^2_{3, 1-\alpha/2} = 9.35$$

$$\chi^2_{3, \alpha/2} \approx 0.216 \Rightarrow \text{отвергается}$$

Вывод: генератор подозрительный.

Обобщённый критерий хи-квадрат

$$X = (X_1, \dots, X_n) \sim P \in \{P_\theta \mid \theta \in \Theta \subseteq \mathbb{R}^d\}$$

$$H_0: P \in \mathcal{P}_\theta^0 \quad \text{vs.} \quad H_1: P \in \mathcal{P}_\theta'$$

$$\mathcal{X} = \bigcup_{j=1}^m B_j \quad \mu_j = \# \{i \mid x_i \in B_j\} \quad p_j^0(\theta) = P_\theta^0(x_i \in B_j)$$

$$\chi(x) = \sum_{j=1}^m \frac{(\mu_j - n p_j^0(\hat{\theta}))^2}{n p_j^0(\hat{\theta})} \xrightarrow{d_\theta} \chi_{m-d-1}^2$$

$$\hat{\theta} - \text{ОМП где} \quad L_x(\theta) = \log \left[\prod_{i=1}^n \prod_{j=1}^m p_j^0(\theta)^{\mathbb{I}\{x_i \in B_j\}} \right] =$$

категориальное распр.

$$= \sum_{i=1}^n \sum_{j=1}^m \mu_j \log p_j^0(\theta) \rightarrow \max_{\theta}$$

Критерий : $\{x(x) > x_{m-d-1, 1-d}^2\}$

Задача

Данные по бомбёжкам Лондона.

Поделили Лондон на сетку 24×24 ($= 576$)

Посчитали во сколько кварталов с какой частотой попадали бомбы

0 1 2 3
229 211 93 35

4	5	6	7	≥ 4
7	0	0	1	8

Используя обобщённый хи-квадрат, проанализировать случай но ли прилетают бомбы.

Если H_0 верна (бомбы случайны), то число событий (прилётов) имеет пуассоновское распределение.

Залем-то

хотим

$$np_j^o(\hat{\theta}) \geq 5$$

$$\mu_j \geq 5$$

$$p_0^o(\theta) = e^{-\theta}$$

$$p_1^o(\theta) = \theta e^{-\theta}$$

$$p_2^o(\theta) = \frac{\theta^2}{2} e^{-\theta}$$

$$p_3^o(\theta) = \frac{\theta^3}{6} e^{-\theta}$$

$$p_4^o(\theta) = 1 - e^{-\theta} \left(1 + \theta + \frac{\theta^2}{2} + \frac{\theta^3}{6} \right)$$

Pois $p(x) = \frac{\lambda^x}{x!} e^{-\lambda}$

Считаем ОМП :

$$L_x(\theta) = \sum_{j=0}^4 \mu_j \log p_j^0(\theta) =$$

$$= 229(-\theta) + 211(-\theta + \log \theta) + 93(-\theta + 2 \log \theta - \log 2) + 35(-\theta + 3 \log \theta - \log 6) + 8 \log p_4^0$$

$$\frac{\partial L_x(\theta)}{\partial \theta} = (-229 - 211 - 93 - 35) + \frac{211 + 186 + 105}{\theta} + \frac{e^{-\theta} \left(1 + \theta + \frac{\theta^2}{2} + \frac{\theta^3}{6}\right) - e^{-\theta} \left(1 + \theta + \frac{\theta^2}{2}\right)}{1 - e^{-\theta} \left(1 + \theta + \frac{\theta^2}{2} + \frac{\theta^3}{6}\right)} =$$

$$= -568 + \frac{502}{\theta} + \frac{8 e^{-\theta} \theta^3/6}{1 - e^{-\theta} \left(1 + \theta + \frac{\theta^2}{2} + \frac{\theta^3}{6}\right)} = \dots$$

$$\hat{\theta} = 0.93$$

$$p_0^o(\hat{\theta}) = 0.384$$

$$p_1^o(\hat{\theta}) = 0.367$$

$$p_2^o(\hat{\theta}) = 0.171$$

$$p_3^o(\hat{\theta}) = 0.053$$

$$p_4^o(\hat{\theta}) = 0.015$$

$$\chi_{3, 1-\alpha}^2 = 7.81$$

$$K(x) = 1.17$$

\Rightarrow не отвергается

$$p\text{-value} = 0.759$$