

1) X_1, \dots, X_n - выборка с эмпир. ф.у. распр. \hat{F}_n . Найти $\text{cov}(\hat{F}_n(x), \hat{F}_n(y))$

Лемма: $F(x) = \int_{-\infty}^x f(t) dt = \int_R \mathbb{1}\{X: \leq x\} f(t) dt = E \mathbb{1}\{X: \leq x\}$; эмпирическая функция

распределения: $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i: \leq x\} \Rightarrow E \hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n E \mathbb{1}\{X_i: \leq x\} = \frac{1}{n} \cdot n F(x) = F(x)$.

Также посчитаем $E(\hat{F}_n(x) \hat{F}_n(y)) = \frac{1}{n^2} E\left(\sum_{i=1}^n \mathbb{1}\{X_i: \leq x\} \cdot \sum_{j=1}^n \mathbb{1}\{X_j: \leq y\}\right) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n E(\mathbb{1}\{X_i: \leq x\} \mathbb{1}\{X_j: \leq y\})$,

где $E = \begin{cases} E(\mathbb{1}\{X_i: \leq x \text{ и } X_j: \leq y\}) & , i=j \\ E \mathbb{1}\{X_i: \leq x\} \cdot E \mathbb{1}\{X_j: \leq y\} & , i \neq j \end{cases} = \begin{cases} F(\min(x, y)) & , i=j \\ F(x) \cdot F(y) & , i \neq j \end{cases} \Rightarrow E(\hat{F}_n(x) \hat{F}_n(y)) =$

$= \frac{1}{n^2} (n F(\min(x, y)) + n(n-1) F(x) F(y)) = \frac{1}{n} F(\min(x, y)) + \frac{n-1}{n} F(x) F(y)$. Тогда ковариация

$\text{cov}(\hat{F}_n(x), \hat{F}_n(y)) = E(\hat{F}_n(x) \hat{F}_n(y)) - E \hat{F}_n(x) \cdot E \hat{F}_n(y) = \left(\frac{1}{n} F(\min(x, y)) + \frac{n-1}{n} F(x) F(y)\right) - F(x) F(y) =$

$= \frac{1}{n} F(\min(x, y)) - \frac{1}{n} F(x) F(y)$

2) $X_1, \dots, X_n \sim P$. Найти оценку коэф.-та энтальпии методом подстановки. Описать схему метода Бутстрэпа для оценки дисперсии полученной оценки.

коэф. энтальпии: $\gamma = \frac{1}{\sigma^4} E(X - \alpha)^4 - 3$, оценка коэф.-та энтальпии методом подстановки:

$\hat{\gamma} = \frac{1}{\hat{\sigma}^4} \int (x - \hat{\alpha})^4 d\hat{F}_n(x) - 3 = \frac{1}{n \hat{\sigma}^4} \sum_{i=1}^n (X_i - \hat{\alpha})^4 - 3$ (аналогично материалу слайдов).

Схема метода бутстрэпа для оценки дисперсии оценки коэф. энтальпии

1) сгенерируем B бутстрепных выборок из эмпир. распределения \hat{F}_n , т.е. $X_b^* = (X_{b1}^*, \dots, X_{bn}^*)$,

$b \in \overline{1, B}$, каждая выборка - по n случайных (индексов из $\overline{1, \dots, n}$) элементов из $\{X_1, \dots, X_n\}$ с возвращением;

2) по каждой бутстрепной выборке посчитаем $\hat{\gamma}$, получим $\hat{\gamma}_1^* = \hat{\gamma}(X_1^*), \dots, \hat{\gamma}_B^* = \hat{\gamma}(X_B^*)$;

3) по выборке $\hat{\gamma}_b^*, b \in \overline{1, B}$ используем для асимптотич. оценки дисперсии оценки $\hat{\gamma}$,

т.е. $\hat{V}_{boot} = \hat{D} \hat{\gamma} = \frac{1}{B} \sum_{b=1}^B (\hat{\gamma}_b^*)^2 - \left(\frac{1}{B} \sum_{b=1}^B \hat{\gamma}_b^*\right)^2$. \hat{V}_{boot} - искомый ответ

3) формулировки

а) оценка параметра - статистика (ф.у. от выборки), применяемая для получения приближенного значения некоторого параметра

эмпир. распределение - $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i: \leq x\}$, грубое приближение ф.у. распределения исходя из реальной выборки X_1, \dots, X_n

метод подстановки - метод получения оценки, при котором используется эмпир. плотность $\hat{F}_n'(x)$, как, например, в задаче 2, где E посчитано с помощью $d\hat{F}_n$

метод Монте-Карло - метод численного интегрирования, состояющий из интегрирования n раз функции от равномерно распределенной $(b: \int_b^a f(x) dx = (a-b) \int_b^a f(x) dx) \sim U(b, a)$ функции $f(x)$ по формуле $\int_b^a f(x) dx \approx (a-b) \cdot \frac{1}{n} \sum_{i=1}^n f(X_i)$. Метод работает для \mathbb{R}^n , $n \gg 1$.

ОМП - оценка, максимизирующая ф.у. $L_X(\theta) = \prod_{i=1}^n p_\theta(x_i)$ - ф.у. правдоподобия

3 (продолжение)

а) асимптотика метода подстановки: \sqrt{n}

асимптотика метода Монте Карло: \sqrt{n}

Аппроксимация методом М-К оценки методом подстановки: \sqrt{B}

4) x_1, \dots, x_n — реализация выборки, x_1^*, \dots, x_n^* — построенная по ней бутстрепная выборка. С какой вероятностью этот x_i попадет в бутстрепную выборку? посчитать среднее # уникал. эл-тов в бутстрепной выборке, если в иск. выборке все наблюдения различны.

$$1) P(x_i \in X^*) = 1 - P(x_i \notin X^*) = 1 - P(x_i \neq x_1^*) \cdot \dots \cdot P(x_i \neq x_n^*) = 1 - \left(\frac{n-1}{n}\right)^n = 1 - \left(1 - \frac{1}{n}\right)^n$$

Интересно, что $\left(1 - \frac{1}{n}\right)^n \xrightarrow{n \rightarrow \infty} e^{-1}$, т.е. для достаточно больших n : $P(x_i \in X^*) \approx 1 - e^{-1} \xrightarrow{n \rightarrow \infty} 1 - \frac{1}{e}$

$$2) \# \text{ уникальных} = \sum_{i=1}^n \mathbb{1}\{x_i \in X^*\} \Rightarrow E(\# \text{ уникальных}) = \sum_{i=1}^n P(x_i \in X^*) = n \left(1 - \left(1 - \frac{1}{n}\right)^n\right) \approx n(1 - e^{-1})$$

↑
(так работает, если в иск. выборке все наблюдения различны)

↑
из линейности математического ожидания

↓
 n

Вывод: при достаточно больших размерах выборки n каждый эл-т почти наверняка попадет в бутстрепную выборку, и почти наверняка эл-ты бутстрепной выборки будут уникальными

5) $X_1, \dots, X_{(n)} \sim U[0, \theta]$, $\hat{\theta} = X_{(n)}$. Почему бутстреп для оценки распределения $\hat{\theta}$ работает плохо?

Обозначим бутстреп. выборку $X_b^* = (X_{b1}^*, \dots, X_{bn}^*)$, $b \in \overline{1, B}$, и макс $X_b^* := X_{b(n)}^*$. Тогда

$$\hat{E} \hat{X}_{(n)} = \frac{1}{B} \sum_{b=1}^B X_{b(n)}^* \leq \frac{1}{B} \cdot B X_{(n)} = X_{(n)}, \text{ т.к. максимум выборки может не попасть в бутстреп. выборку.}$$

↑
(оценка макс. оценки $X_{(n)}$ для θ методом бутстреп)

Более того, $\hat{E} X_{(n)} = X_{(n)}$ лишь если эл-т $X_{(n)}$ попадет в каждую из B бутстреп. выборок,

что согласно замеч. 4) происходит с вероятностью $\approx 1 - e^{-1}$. То есть $P(\hat{E} X_{(n)} = X_{(n)}) =$

$$= (1 - e^{-1})^B \sim (1 - e^{-1})^n < 1$$

т.е. с ненулевой вероятностью оценка макс. ожидаемая оценка параметра будет занижена, и бутстреп работает плохо, особенно — для малых выборок

6 X_1, \dots, X_n - выборка, X_1^*, \dots, X_n^* - постр. по ней бутстр. выборка, $\bar{X}^* := \frac{1}{n} \sum_{i=1}^n X_i^*$. Найти $D(\bar{X}^* | X_1, \dots, X_n) - ?$ $D \bar{X}^* - ?$

$$D(\bar{X}^* | X_1, \dots, X_n) = \hat{D} \bar{X}^* = \frac{1}{B} \sum_{b=1}^B (\bar{X}_b^*)^2 - \left(\frac{1}{B} \sum_{b=1}^B \bar{X}_b^* \right)^2 \stackrel{B=1}{=} (\bar{X}^*)^2 - (\bar{X}^*)^2 = 0$$

Аналогично, при заданной реализации выборки случайность исчезает и дисперсия однозначно = 0.

$$D \bar{X}^* = \frac{1}{n^2} \sum_{i=1}^n D X_i^* = \frac{1}{n^2} \sum_{i=1}^n D X_i = \frac{1}{n^2} \cdot n D X_1 = \frac{D X_1}{n} \neq 0$$

В этом случае некоторая случайная погрешность сохраняется и дисперсия $\neq 0$

7 Регрессия Y по X , оценка совместной плотн. (X, Y) с ядром $q(x) \cdot q(y)$ $E(Y|X)$. Найти оценку и сравнить с оценкой в модели ядерной регрессии

(ядро: $q_h(x) = \frac{1}{h} q\left(\frac{x-X_i}{h}\right)$, обобщим: $f_h(x, y) = \frac{1}{h_x h_y} q\left(\frac{x-X_i}{h_x}\right) q\left(\frac{y-Y_i}{h_y}\right)$. При $h=1$.

~~Решение:~~ Тогда $E(Y|X) = \int y f(y|x) dy = \int \frac{y f(x, y)}{f(x)} dy =$

$$= \int \frac{y \sum q(x-X_i) q(y-Y_i)}{\sum q(x-X_i)} dy = \frac{\sum (q(x-X_i) \int y q(y-Y_i) dy)}{\sum q(x-X_i)} = \frac{\sum q(x-X_i) Y_i}{\sum q(x-X_i)}, \text{ что и есть оценка в модели ЯЯ. регрессии}$$