

# Bootstrap

## 1.1 Бутстрапирование совокупности бернуллиевских случайных величин (КР-2020, №4)

### Условие задачи

Исходная выборка  $y$  — вектор из  $n$  независимых случайных величин, равновероятно принимающих значения 0 и 1 (коммент: можно и для любого распределения, у которого есть матожидание и дисперсия). Пусть  $y^*$  — одна из бутстрэп-выборок.

- Найдите  $\mathbb{E}[y_i], \text{Var}(y_i), \mathbb{E}[\bar{y}], \text{Var}(\bar{y})$
- $\mathbb{E}[y_i^*], \text{Var}(y_i^*), \mathbb{E}[\bar{y}^*], \text{Var}(\bar{y}^*)$
- $\text{Cov}(y_i, y_i^*), \text{Cov}(\bar{y}, \bar{y}^*)$

### Решение

- Прежде всего заметим, что по условию  $y_i \sim \text{Bern}(p)$  с  $p = 1/2$ . Отсюда сразу получаем, что  $\mathbb{E}[y_i] = p = 1/2$  и  $\text{Var}(y_i) = p(1-p) = 1/4$ . Более того,

$$\mathbb{E}[\bar{y}] = \frac{1}{n} \mathbb{E}[y_1 + \dots + y_n] = \frac{np}{n} = p = 1/2.$$

Если же принять во внимание независимость случайных величин  $y_1, \dots, y_n$ , получим

$$\text{Var}(\bar{y}) = \text{Var}(y_1 + \dots + y_n) \frac{1}{n^2} = \frac{n}{n^2} \text{Var}(y_1) = \frac{1}{4n}.$$

- Ключевое соображение состоит в том, чтобы прибегнуть к представлению случайной величины  $y_i^*$  в виде линейной комбинации величин  $y_1, \dots, y_n$  со случайными коэффициентами  $d_1, \dots, d_n$ , каждый из которых принимает значения либо 0, либо 1. Более того,  $d_j = 1$  тогда и только тогда, когда  $y_i^* = y_j$ . Обозначим через  $d$  вектор из случайных величин  $d_1, \dots, d_n$ .

Про это нужно думать следующим образом: генерация конкретной бутстрэп-выборки равносильна генерации квадратной матрицы  $D$  размера  $n \times n$ , каждый столбец которой содержит ровно одну единицу, а все остальные её элементы равны нулю. При этом позиция, содержащая единицу, определяется случайно и равновероятно для каждого столбца. Тогда

$$(y_1^*, \dots, y_n^*) = (y_1, \dots, y_n) \cdot D.$$

Теперь несложно заметить, что величина  $y_i^*$  при условии коэффициентов  $d$  имеет распределение  $\text{Bern}(p)$  независимо от значения  $d$ . Действительно, при фиксированном  $d$  мы сразу можем сказать, с какой из случайных величин  $y_1, \dots, y_n$  совпала  $y_i^*$ . Но это означает, что условное распределение для  $y_i^*$  совпадает с безусловным. Таким образом,  $y_i^* \sim \text{Bern}(p)$ . Отсюда сразу получаем, что  $\mathbb{E}[y_i^*] = 1/2$ , а  $\text{Var}(y_i^*) = 1/4$ . Более того,

$$\mathbb{E}[\bar{y}^*] = \frac{1}{n} \mathbb{E}[y_1^* + \dots + y_n^*] = 1/2.$$

Теперь приступим к поиску  $\text{Var}(\bar{y}^*)$ , предварительно отметив, что  $\mathbb{E}[d_j] = 1/n$ , а случайные величины  $y_j$  и  $d_j$  являются независимыми для всех  $j$  (это простое наблюдение, очевидно следующее из самой процедуры бутстрэпа).

Итак,

$$\text{Var}(\bar{y}^*) = \frac{1}{n^2} \text{Var}(y_1^* + \dots + y_n^*) = \frac{1}{n} \text{Var}(y_i^*) + \frac{(n-1)}{n} \text{Cov}(y_1^*, y_2^*).$$

Знаем, что  $\text{Var}(y_i^*) = 1/4$ .

- Таким образом, остаётся вычислить  $\text{Cov}(y_1^*, y_2^*)$ , для чего мы и воспользуемся линейным разложением:

$$\text{Cov}(y_1^*, y_2^*) = \text{Cov}(d_1 y_1 + \dots + d_n y_n, d'_1 y_1 + \dots + d'_n y_n) = \sum_{i=1}^n \text{Cov}(d_i y_i, d'_i y_i) + \sum_{i \neq j} \text{Cov}(d_i y_i, d'_j y_j) = \sum_{i=1}^n \text{Cov}(d_i y_i, d'_i y_i).$$

Вычислим одно из слагаемых получившейся суммы:

$$\text{Cov}(d_i y_i, d'_i y_i) = \mathbb{E}[d_i d'_i y_i^2] - \mathbb{E}[d_i y_i] \mathbb{E}[d'_i y_i] = \mathbb{E}[d_i]^2 \mathbb{E}[y_i^2] - \mathbb{E}[d_i]^2 \mathbb{E}[y_i]^2 = \frac{1}{n^2} \text{Var}(y_i) = \frac{1}{4n^2}.$$

В итоге имеем

$$\text{Cov}(y_1^*, y_2^*) = \frac{1}{4n}.$$

Подставим теперь ответ в выражение для  $\text{Var}(\bar{y}^*)$ :

$$\text{Var}(\bar{y}^*) = \frac{1}{4n} + \frac{(n-1)}{4n^2}.$$

- Остаётся вычислить

$$\text{Cov}(\bar{y}, \bar{y}^*) = \frac{1}{n^2} \left( \sum_{i=1}^n \text{Cov}(y_i, y_i^*) + \sum_{i \neq j} \text{Cov}(y_i, y_j^*) \right).$$

Рассмотрим вначале слагаемое первого типа.

$$\text{Cov}(y_i, y_i^*) = \mathbb{E}[y_i(d_1 y_1 + \dots + d_n y_n)] - \mathbb{E}[y_i] \mathbb{E}[y_i^*] = \frac{n-1+2}{4n} - \frac{1}{4} = \frac{1}{4n}.$$

Второе слагаемое на самом деле равно первому, так как мы снова сравниваем  $y_i$  с величиной, которая совпадёт с ней с вероятностью  $1/n$ ; таким образом,

$$\text{Cov}(y_i, y_j^*) = \text{Cov}(y_i, y_i^*) = \frac{1}{4n}.$$

Подставив полученные выражения в изначальное, получаем ответ:

$$\text{Cov}(\bar{y}, \bar{y}^*) = \frac{1}{n^2} \left( \frac{1}{4} + \frac{(n-1)}{4} \right) = \frac{1}{4n}.$$

## 1.2 Бутстрапирование мёда (КР-2021, №4)

### Условие задачи

У Винни-Пуха 1 000 000 наблюдений в минуту — потоковые данные по мёду от пчёл. А общее количество наблюдений  $n$  необозримо велико.

- К какому распределению стремится количество копий  $i$ -го исходного наблюдения в бутстрэп-выборке с ростом  $n$ ?
- (позже...) Вместо честного наивного бутстрэпа Винни-Пух использует аппроксимацию. В реальном времени каждое поступающее наблюдение  $y_i$  он заменяет на его  $k_i$  копий, где количество  $k_i$  выбирает случайно, независимо от  $y$  и предыдущих  $k_j$ , согласно распределению найденному в предыдущем пункте. Например, если  $k_i = 0$ , то наблюдение  $y_i$  Винни-Пух не запоминает, а если  $k_i = 2$ , то наблюдение  $y_i$  учитывается дважды. Во сколько раз отличается дисперсия обычного среднего выборочного  $\bar{y}_n$  и среднего выборочного сделанных копий, которое посчитает Винни-Пух для  $n$  исходных наблюдений?

### Решение