

Бутстрэп

Бутстрэп — это набор практических методов, который основан на многократной генерации выборок на базе одной имеющейся выборки.

Бутстрэп используется для оценки каких-то параметров распределений, построения доверительных интервалов и т.д.

Мы рассмотрим **параметрический** и **непараметрический** бутстрэп. Начнем с параметрического.

Бутстрэп

Параметрический бутстрэп

Идея заключается в том, что если оценка $\hat{\theta}$ близка к настоящему параметру θ_0 , то распределение $F_{\hat{\theta}}$ будет похоже на F_{θ_0} . Поэтому можно генерировать новые выборки из $F_{\hat{\theta}}$.

Мы здесь предполагаем, что семейство распределений F_{θ} непрерывно зависит от параметра.

Бутстрэп

Пример

Допустим, мы построили какую-то оценку $\hat{\theta}$ неизвестного параметра θ . Ни один из методов построения оценок, которые мы изучали, не гарантирует несмещенность.

Попытаемся исправить смещенность оценки с помощью параметрического бутстрэпа.

Бутстрэп

Это можно сделать следующим образом:

- ▶ сгенерировать выборку Y_1, \dots, Y_n из $F_{\hat{\theta}}$, и подсчитать по ней $\hat{\theta}(Y_1, \dots, Y_n)$;
- ▶ «оценить» смещение $\mathbb{E}[\hat{\theta}(X_1, \dots, X_n)] - \theta_0$ с помощью $\hat{\theta}(X_1, \dots, X_n) - \hat{\theta}(Y_1, \dots, Y_n)$;
- ▶ посчитать «поправленную» оценку

$$2\hat{\theta}(X_1, \dots, X_n) - \hat{\theta}(Y_1, \dots, Y_n).$$

Бутстрэп

Бутстрэп имеет несколько неоспоримых плюсов — он прост в использовании, не требует сложных вычислений и применим даже к весьма громоздким моделям.

С другой стороны, мы не можем явным образом оценить его погрешность, а в случае, если оценка $\hat{\theta}$ значимо промахнулась мимо θ_0 , рискуем неправильно изменить оценку.

Бутстрэп

Как строить доверительные интервалы с помощью бутстрэпа?

Существует и несколько методов построения доверительных интервалов. Наиболее простой из них — **pivotal** интервал.

Бутстрэп

Идея: рассмотрим оценку $\hat{\theta}$ параметра θ_0 .

- ▶ возьмем несколько выборок из $F_{\hat{\theta}}$ и построим на их основе другие оценки $\hat{\theta}_1, \dots, \hat{\theta}_m$;
- ▶ упорядочим $\hat{\theta}_i$ и выберем те из них, $\hat{\theta}_-$ и $\hat{\theta}_+$, которые стоят на местах $[(\alpha/2)m]$ и $[(1 - \alpha/2)m]$ по возрастанию;
- ▶ тогда нашим интервалом будет

$$(\hat{\theta}_-, \hat{\theta}_+).$$

Бутстрэп

Непараметрический бутстрэп

Очень часто бутстрэп используется в непараметрической постановке. Это означает, что у нас нет никакого семейства распределений F_θ , а есть только реализация выборки x_1, \dots, x_n из некоторого неизвестного распределения F .

В этом случае бутстрэп-выборки генерируются с помощью выбора с возвращением.

Бутстрэп

Теоретически это можно обосновать с помощью понятия эмпирической функции распределения.

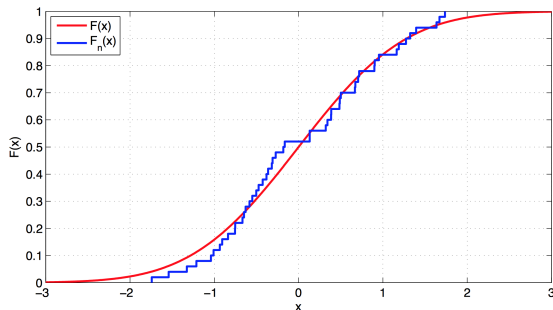
Эмпирическая функция распределения $\hat{F}_n(u)$ определяется формулой

$$\hat{F}_n(u) = \frac{1}{n} \sum_{i=1}^n \mathbf{I}_{\{x_i \leq u\}},$$

где $\mathbf{I}_{\{x_i \leq u\}}$ — индикатор события $\{x_i \leq u\}$.

Бутстрэп

График $\hat{F}_n(x)$ представляет собой ступенчатую функцию, растущую скачками высоты $1/n$. Скачки происходят в точках с координатами x_1, \dots, x_n .



Бутстрэп

Известно, что эмпирическая функция распределения является очень хорошим приближением для истинной функции распределения.

Следовательно, чтобы сгенерировать бутстрэп-выборку, можно использовать закон, соответствующий эмпирической функции распределения.

А это и будет выбором с возвращением.