

Теория и практика применения бутстрэпа

Москва
2023

Содержание

1. Введение	3
1.1. История метода бутстрэпа	3
1.2. Классическая теория построения доверительных интервалов	3
1.3. Построение псевдовыборок	4
1.4. Суть метода бутстрэпа	5
2. Теория основных версий бутстрэпа	5
2.1. Перцентильный бутстрэп	5
2.2. Обратный перцентильный бутстрэп	6
2.3. Стьюдентизированный бутстрэп (бутстрэп t-статистики)	6
2.4. Бутстрэп в бутстрэпе	7
2.5. Графическая иллюстрация точности методов	8
2.6. Бутстрэп с коррекцией смещения и ускорением (Bias-corrected and accelerated bootstrap, BCa)	13
3. Бутстрэп в линейной регрессии	14
3.1. Парный бутстрэп	14
3.2. Параметрический бутстрэп	14
4. Бутстрэп в решающих деревьях и ансамблях	14
4.1. Бэггинг	14

1. Введение

1.1. История метода бутстрэпа

Словом "бутстрэп" обозначают некоторое семейство или множество статистических методов или алгоритмов, предназначенных для нахождения оценок параметров распределения и доверительных интервалов. Эти методы основаны на генерации большого числа выборок из исходной, имеющейся в распоряжении исследователя выборки с возвращением с использованием датчика (псевдо-) случайных чисел. Таким образом, из-за объема вычислений реализация бутстрэпа предполагается на компьютере.

Одной из первых работ, в которой была изложена суть метода бутстрэпа, стала статья Брэдли Эфрона «Bootstrap methods: another look at the jackknife» (1979 г.), вдохновленная более ранней работой о методе Jackknife. Позже разрабатывалась теория для нахождения более точных оценок дисперсий параметров, повышения устойчивости метода на выборках небольших размеров, с включением байесовского подхода и т.д.

1.2. Классическая теория построения доверительных интервалов

Пусть имеется выборка $\mathbf{X} = (X_1, \dots, X_n)$ и известно, что она была сгенерирована из нормального распределения $\mathcal{N}(\mu, \sigma^2)$ с неизвестным математическим ожиданием μ и известной дисперсией σ^2 . В таком случае известным результатом является построение следующего доверительного интервала:

$$\mathbb{P} \left(\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right) = 1 - \alpha$$

где z_α — квантиль стандартной нормальной случайной величины порядка α .

Кроме того, похожий результат может быть получен в случае неизвестного параметра σ^2 . В таком случае используется несмещенная оценка дисперсии

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$$

и получается, что

$$n \cdot \frac{\bar{X} - \mu}{\hat{\sigma}} \sim t_{n-1}$$

Тогда точный доверительный интервал приобретает вид

$$\mathbb{P} \left(\bar{X} - t_{1-\frac{\alpha}{2}; n-1} \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{1-\frac{\alpha}{2}; n-1} \frac{\hat{\sigma}}{\sqrt{n}} \right) = 1 - \alpha$$

При большом числе наблюдений стандартизированная случайная величина $n \cdot \frac{\bar{X} - \mu}{\hat{\sigma}}$, согласно центральной предельной теореме, по распределению сходится к стандартному

нормальному распределению вне зависимости от вида закона распределения X (при условии существования конечного математического ожидания μ и конечной дисперсии σ^2), т. е.

$$n \cdot \frac{\bar{X} - \mu}{\hat{\sigma}} \xrightarrow[n \rightarrow \infty]{dist} \mathcal{N}(0, 1)$$

Поэтому может быть получен следующий асимптотический доверительный интервал:

$$\mathbb{P} \left(\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}}{\sqrt{n}} \right) \approx 1 - \alpha.$$

Итак, было показано, что в случае построения доверительного интервала для параметра среднего в случае известного закона распределения имеется хорошо разработанная теория и могут быть получены точные или асимптотические формулы. Кроме того, такие результаты широко известны в случае построения доверительных интервалов для разницы математических ожиданий, дисперсий и отношения дисперсий.

Однако возникает вопрос, какую схему действий предпринять в случае, если закон распределения, порождающий данные, неизвестен, а также если стоит задача построить доверительный интервал для такой не менее важной характеристики, как, например, медиана распределения, асимптотический закон распределения которой не является широко известным.

По-прежнему можно построить $\hat{\theta}$ как точечную оценку неизвестного параметра θ . Но теоретических знаний о распределении $\hat{\theta}$ у нас может не быть, либо же эта теория недоступна исследователю. В таком случае можно использовать бутстрэп для построения доверительного интервала. Однако прежде стоит разобраться с тем, каким образом можно создавать новые выборки на основе исходной.

1.3. Построение псевдовыборок

Существует два распространенных метода построения случайной выборки (пример взят из [Конспекта](#)).

- Выборка без возвращения (without replacement, simple random sampling):

Предположим, мы поочередно берем 10 карт наугад из колоды из 52 карт, не возвращая ни одну из карт обратно в колоду между взятиями. Это называется выборкой без возвращению или простой случайной выборкой. При таком методе в нашей выборке из 10 карт не будет дубликатов карт. То есть, в данном подходе, если мы захотим создать выборку, размер которой совпадает с исходной, будет получена сама исходная выборка, что ограничивает возможности в исследовании распределения данных.

- Выборка с возвращением (with replacement):

Теперь предположим, что мы берем 10 карт наугад из колоды, но после каждого взятия мы кладем карту обратно в колоду и перемешиваем карты. Это называется выборкой с возвращением. При использовании этого метода выборка из 10 карт мо-

жет иметь дубликаты. Возможно и такое, что была вытянута шестерка червей все 10 раз. Такая процедура позволяет создавать отличные от исходной новые выборки, размер которых при этом совпадает с размером исходной выборки.

1.4. Суть метода бутстрэпа

На основе имеющейся выборки может быть рассчитана выборочная или эмпирическая функция распределения $\hat{F}(x)$, которая для каждого x показывает долю наблюдений в выборке, не превосходящих x :

$$\hat{F}(x) = \frac{\sum_{i=1}^n \mathbb{I}\{X_i \leq x\}}{n}$$

Известно, что эта функция обладает рядом "хороших" свойств: состоятельность, эффективность, асимптотическая нормальность, поэтому является "хорошей" оценкой для истинной функции распределения $F(x)$, порождающей данные.

Итак, основной идеей бутстрэпа в данном случае будет генерация n_{boot} бутстрэпированных выборок $\mathbf{X}_1^*, \dots, \mathbf{X}_{n_{boot}}^*$ с повторениями из исходной выборки \mathbf{X} и подсчет оценки неизвестного параметра распределения $\hat{\theta}_j^*$ для j -ой бутстрэпированной выборки, $j \in \{1, \dots, n_{boot}\}$, для получения распределения $\hat{\theta}$.

2. Теория основных версий бутстрэпа

2.1. Перцентильный бутстрэп

Одной из основных и наиболее простых техник бутстрэпа является перцентильный бутстрэп (Percentile Bootstrap). Его алгоритм следующий:

1. Из исходной выборки \mathbf{X} генерируется n_{boot} бутстрэпированных выборок $\mathbf{X}_1^*, \mathbf{X}_2^*, \dots, \mathbf{X}_{n_{boot}}^*$;
2. Для каждой бутстрэпированной выборки оценивается параметр $\hat{\theta}_j^*$, $1 \leq j \leq n_{boot}$, затем оцененные параметры сортируются по возрастанию: $\hat{\theta}_{(1)}^* \leq \hat{\theta}_{(2)}^* \leq \dots \leq \hat{\theta}_{(n_{boot})}^*$, и собирается вектор $\hat{\theta}^* = (\hat{\theta}_{(1)}^*, \hat{\theta}_{(2)}^*, \dots, \hat{\theta}_{(n_{boot})}^*)$;
3. Оцененным $(1 - \alpha)$ -процентным доверительным интервалом истинного параметра θ будет отрезок от $\frac{\alpha}{2}$ квантиля до $1 - \frac{\alpha}{2}$ квантиля вектора $\hat{\theta}^*$: $\widehat{CI}_\theta = [\hat{\theta}_{\frac{\alpha}{2}}^*; \hat{\theta}_{1-\frac{\alpha}{2}}^*]$.

Насколько же точен данный метод построения доверительных интервалов? Введем обозначение $CL_\alpha = 1 - \alpha$ — теоретический уровень значимости доверительного интервала, \widehat{CL}_n — фактический (выборочный) уровень значимости доверительного интервала. Разность между теоретическим и фактическим уровнем значимости может быть оценена следующим образом: $|CL_\alpha - \widehat{CL}_n| = O\left(\frac{1}{\sqrt{n}}\right)$.

Оцененный таким образом доверительный интервал в среднем будет занижать вероятность накрытия.

Кроме того, если выборочное распределение было асимметричным, это также приведет к асимметрии доверительного интервала, полученного с помощью перцентильного бутстрэпа.

2.2. Обратный перцентильный бутстрэп

Модификацией первого метода можно считать обратный перцентильный бутстрэп (Reverse Percentile Bootstrap). Его алгоритм:

1. Из исходной выборки \mathbf{X} генерируется n_{boot} бутстрэпированных выборок \mathbf{X}_1^* , \mathbf{X}_2^* , ..., $\mathbf{X}_{n_{boot}}^*$;
2. Для каждой бутстрэпированной выборки оценивается параметр $\hat{\theta}_j^*$, $1 \leq j \leq n_{boot}$, далее считается отклонение полученной оценки от выборочной оценки $d_j^* = \hat{\theta}_j^* - \hat{\theta}$, полученные разницы сортируются по возрастанию: $d_{(1)}^* \leq d_{(2)}^* \leq \dots \leq d_{(n_{boot})}^*$, и собирается вектор $d^* = (d_{(1)}^*, d_{(2)}^*, \dots, d_{(n_{boot})}^*)$;
3. Оцененным $(1 - \alpha)$ -процентным доверительным интервалом разницы $d = \hat{\theta} - \theta$ будет отрезок от квантиля $\frac{\alpha}{2}$ до квантиля $1 - \frac{\alpha}{2}$ вектора d^* : $\widehat{CI}_d = [d_{\frac{\alpha}{2}}^*; d_{1-\frac{\alpha}{2}}^*]$. Тогда для получения оценки доверительного интервала для параметра θ необходимо выполнить ряд преобразований:

$$d_{\frac{\alpha}{2}}^* \leq d \leq d_{1-\frac{\alpha}{2}}^* \Leftrightarrow d_{\frac{\alpha}{2}}^* \leq \hat{\theta} - \theta \leq d_{1-\frac{\alpha}{2}}^* \Leftrightarrow \hat{\theta} - d_{1-\frac{\alpha}{2}}^* \leq \theta \leq \hat{\theta} - d_{\frac{\alpha}{2}}^*$$

$$\text{Таким образом, } \widehat{CI}_\theta = [\hat{\theta} - d_{1-\frac{\alpha}{2}}^*; \hat{\theta} - d_{\frac{\alpha}{2}}^*].$$

Скорость сходимости обратного перцентильного бутстрэпа совпадает со скоростью сходимости обычного перцентильного бутстрэпа: $|CL_\alpha - \widehat{CI}_n| = O\left(\frac{1}{\sqrt{n}}\right)$.

2.3. Стьюдентизированный бутстрэп (бутстрэп t-статистики)

Далее рассмотрим метод, который позволит рассчитывать стандартную ошибку оценки, и учитывать ее влияние на распределение. Идея этого подхода заключается в том, чтобы приблизить величину $\frac{\hat{\theta} - \theta}{se(\hat{\theta})}$ величиной $\frac{\hat{\theta}^* - \hat{\theta}}{se(\hat{\theta}^*)}$. Алгоритм стьюдентизированного бутстрэпа:

1. Из исходной выборки \mathbf{X} генерируется n_{boot} бутстрэпированных выборок \mathbf{X}_1^* , \mathbf{X}_2^* , ..., $\mathbf{X}_{n_{boot}}^*$;
2. Для каждой бутстрэпированной выборки оценивается параметр $\hat{\theta}_j^*$, $1 \leq j \leq n_{boot}$, далее считается стандартная ошибка полученной оценки. В случае построения доверительного интервала для математического ожидания можно воспользоваться формулой

$$\widehat{se}(\hat{\theta}_j^*) = \sqrt{\frac{\sum_{i=1}^n ((X_i^*)_j - \bar{X}_j^*)^2}{n \cdot (n - 1)}} = \frac{1}{\sqrt{n}} \sigma_{X_j^*}$$

где \bar{X}_j^* — среднее по бутстрэпированной выборке j , $\sigma_{X_j^*}$ — среднеквадратичное

отклонение бутстрэпированной выборки j . Далее вычисляются n_{boot} бутстрэпированных t -статистик:

$$t_j^* = \frac{\hat{\theta}_j^* - \hat{\theta}}{\widehat{se}(\hat{\theta}_j^*)}$$

полученные t -статистики сортируются по возрастанию: $t_{(1)}^* \leq t_{(2)}^* \leq \dots \leq t_{(n_{boot})}^*$, и собирается вектор $t^* = (t_{(1)}^*, t_{(2)}^*, \dots, t_{(n_{boot})}^*)$;

3. Оценным $(1 - \alpha)$ -процентным доверительным интервалом стандартизированного параметра $\theta^{st} = \frac{\hat{\theta} - \theta}{se(\hat{\theta})}$ будет отрезок от квантиля $\frac{\alpha}{2}$ до квантиля $1 - \frac{\alpha}{2}$ вектора t^* : $\widehat{CI}_{\theta^{st}} = [t_{\frac{\alpha}{2}}^*; t_{1-\frac{\alpha}{2}}^*]$. Тогда для получения оценки доверительного интервала для параметра θ необходимо выполнить ряд преобразований:

$$t_{\frac{\alpha}{2}}^* \leq \theta^{st} \leq t_{1-\frac{\alpha}{2}}^* \Leftrightarrow t_{\frac{\alpha}{2}}^* \leq \frac{\hat{\theta} - \theta}{se(\hat{\theta})} \leq t_{1-\frac{\alpha}{2}}^* \Leftrightarrow \hat{\theta} - t_{1-\frac{\alpha}{2}}^* \cdot se(\hat{\theta}) \leq \theta \leq \hat{\theta} - t_{\frac{\alpha}{2}}^* \cdot se(\hat{\theta})$$

Таким образом, $\widehat{CI}_{\theta} = [\hat{\theta} - t_{1-\frac{\alpha}{2}}^* \cdot se(\hat{\theta}); \hat{\theta} - t_{\frac{\alpha}{2}}^* \cdot se(\hat{\theta})]$.

Эта версия бутстрэпа имеет ряд преимуществ над предыдущими версиями. Прежде всего, построенный с помощью студентизированного бутстрэпа доверительный интервал быстрее сходится к теоретическому: $|CL_{\alpha} - \widehat{CI}_n| = O(\frac{1}{n})$. Также стоит добавить, что такой доверительный интервал будет более корректно учитывать асимметрию выборочного распределения.

2.4. Бутстрэп в бутстрэпе

Однако точная формула стандартной ошибки, например, медианы либо других выборочных характеристик может быть неизвестна. В такой ситуации для подсчета стандартной ошибки может быть применен бутстрэп в бутстрэпе:

1. Из каждой бутстрэпированной выборки \mathbf{X}_j^* , $1 \leq j \leq n_{boot}$, генерируется n_{bb} псевдо-выборок $\mathbf{X}_{j1}^{**}, \mathbf{X}_{j2}^{**}, \dots, \mathbf{X}_{jn_{bb}}^{**}$;
2. Для каждой бутстрэпированной выборки второго уровня оценивается параметр $\hat{\theta}_{jk}^{**}$, $1 \leq k \leq n_{bb}$, далее вычисляется стандартная ошибка по формуле

$$\widehat{se}(\hat{\theta}_j) = \sqrt{\frac{\sum_{k=1}^{n_{bb}} (\hat{\theta}_{jk}^{**} - \bar{\hat{\theta}}_j^{**})^2}{n_{bb} - 1}}$$

где $\bar{\hat{\theta}}_j^{**} = \frac{1}{n_{bb}} \sum_{k=1}^{n_{bb}} \hat{\theta}_{jk}^{**}$ — среднее значение параметра по бутстрэпированным выборкам второго уровня.

2.5. Графическая иллюстрация точности методов

В данном разделе приводится графическая иллюстрация точности построения доверительных интервалов для математического ожидания (среднего) с помощью ЦПТ, перцентильного бутстрэпа и стьюдентизированного бутстрэпа. При использовании обеих версий бутстрэпа количество генерируемых выборок равняется 10^5 . В качестве распределений использовались нормальное с параметрами $\mu = 10$, $\sigma^2 = 25$ и равномерное на отрезке $[0; 10]$.

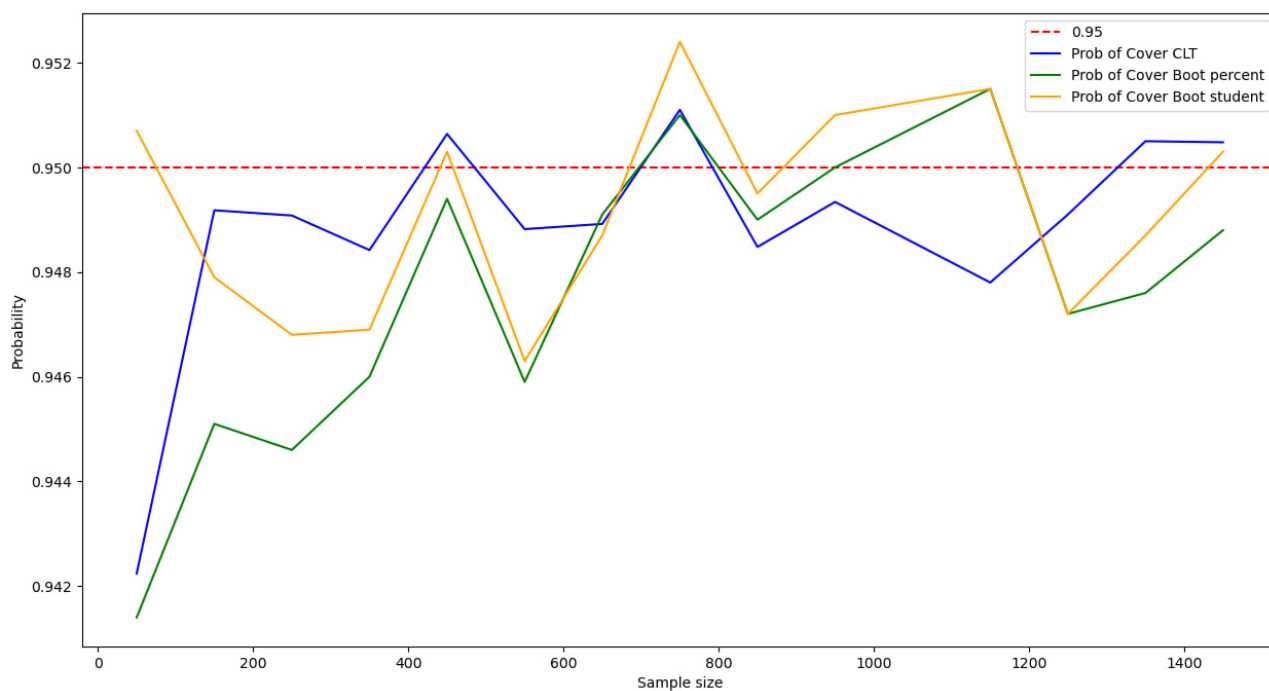


Рисунок 2.1. Зависимость оцениваемой вероятности накрытия от размера выборки, нормальное распределение. Число экспериментов для ЦПТ = $5 \cdot 10^5$, для бутстрэпа = 10^5

Рис. 2.1 демонстрирует, что доверительный интервал, построенный при помощи бутстрэпа t-статистики, при небольших размерах выборки оказывается точнее доверительного интервала на основе ЦПТ и перцентильного бутстрэпа. Однако при объемах выборки в несколько сотен наблюдений в среднем более точным оказывается результат применения ЦПТ.

Для равномерного распределения (Рис. 2.2) вероятности накрытия на основе бутстрэпа также оказываются ближе к 0.95 только при размерах выборки до 200 — 300 наблюдений. Для больших размеров выборки д. и., оцененные с помощью ЦПТ, оказываются стабильно ближе к теоретическим.

Далее более внимательно сравним точность работы бутстрэпа по сравнению с ЦПТ для выборок маленького разброса (30, 50 и 100).

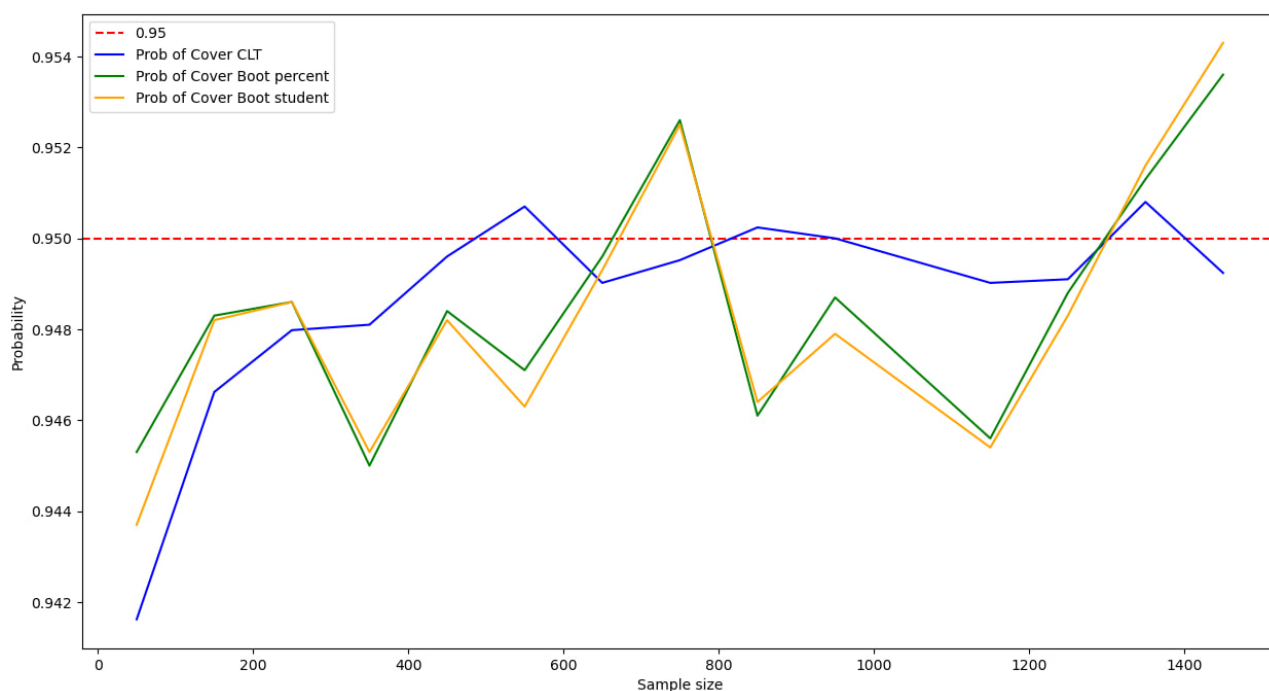


Рисунок 2.2. Зависимость оцениваемой вероятности накрытия от размера выборки, равномерное распределение. Число экспериментов для ЦПТ = $5 \cdot 10^5$, для бутстрэпа = 10^5

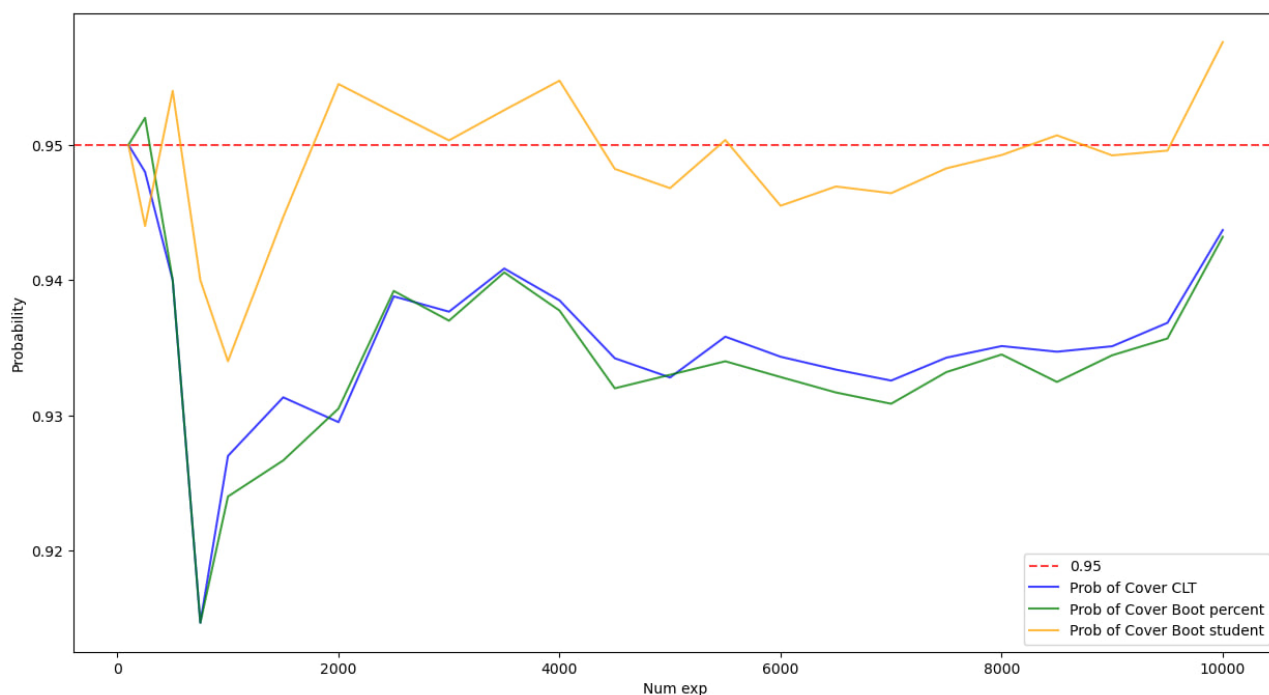


Рисунок 2.3. Зависимость оцениваемой вероятности накрытия от числа экспериментов, нормальное распределение, 30 наблюдений

Все три графика (Рис. 2.3, 2.4, 2.5) демонстрируют, что поведение вероятности накрытия истинного параметра (математического ожидания) очень схоже для ЦПТ и перцентильного бутстрэпа. При этом они оба систематически уступают студентизированно-

му бутстрэпу в точности оценивания, особенно в ситуации работы с выборкой размером 30 наблюдений.

В случае оценки доверительного интервала для математического ожидания равномерного распределения результаты менее однозначны. Точность построения д. и. схожа для всех трех методов, причем все они склонны занижать размер оцененного доверительного интервала по сравнению с теоретическим (Рис. 2.6, 2.7, 2.8). Это происходит из-за низкого эксцесса равномерного распределения (случайные величины не концентрируются вблизи математического ожидания).

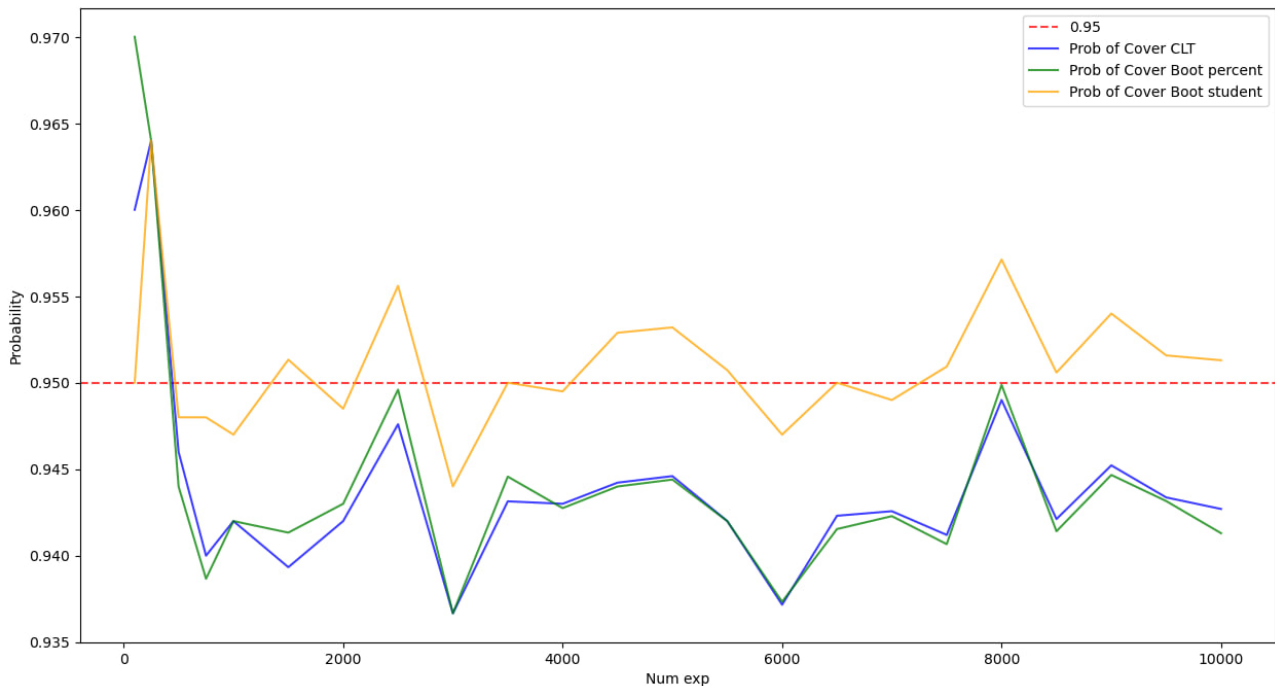


Рисунок 2.4. Зависимость оцениваемой вероятности накрытия от числа экспериментов, нормальное распределение, 50 наблюдений

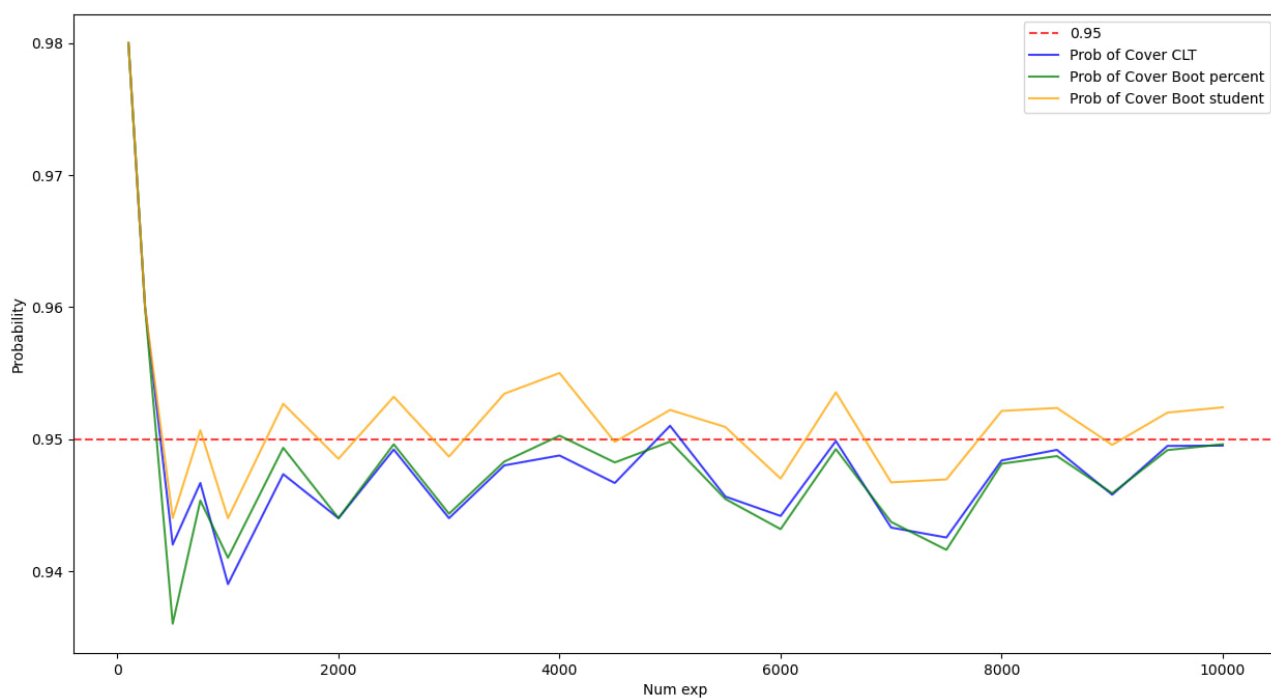


Рисунок 2.5. Зависимость оцениваемой вероятности накрытия от числа экспериментов, нормальное распределение, 100 наблюдений

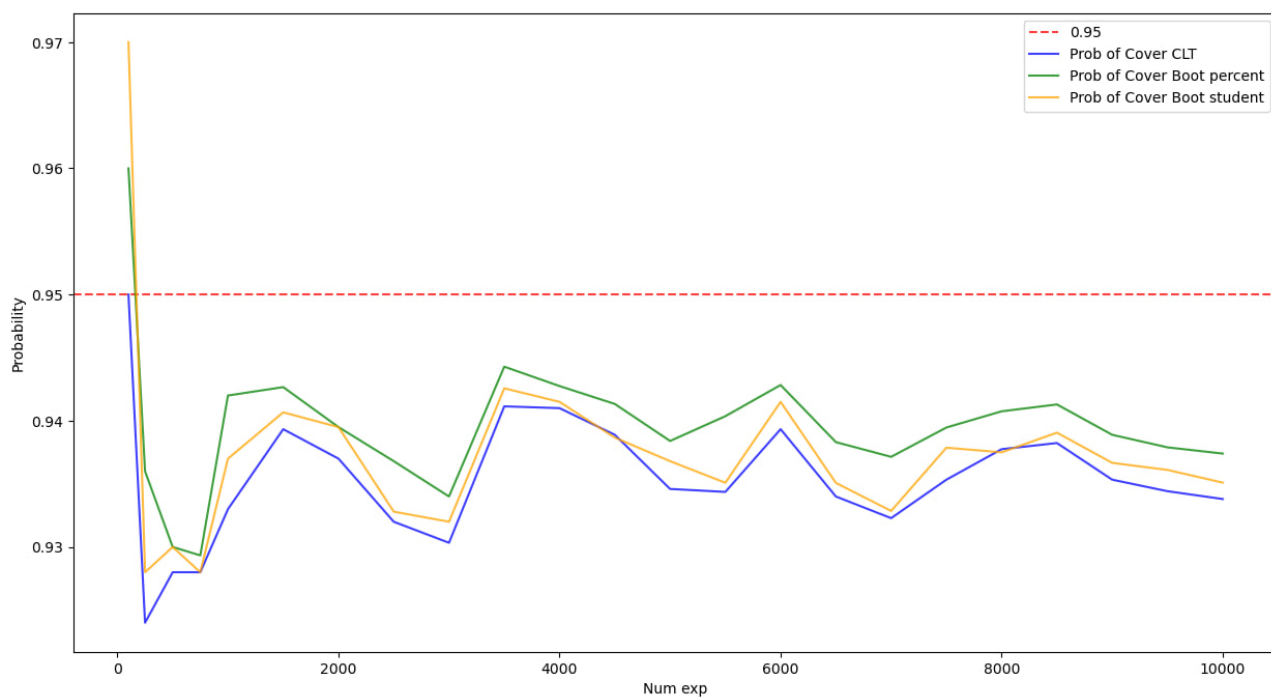


Рисунок 2.6. Зависимость оцениваемой вероятности накрытия от числа экспериментов, равномерное распределение, 30 наблюдений

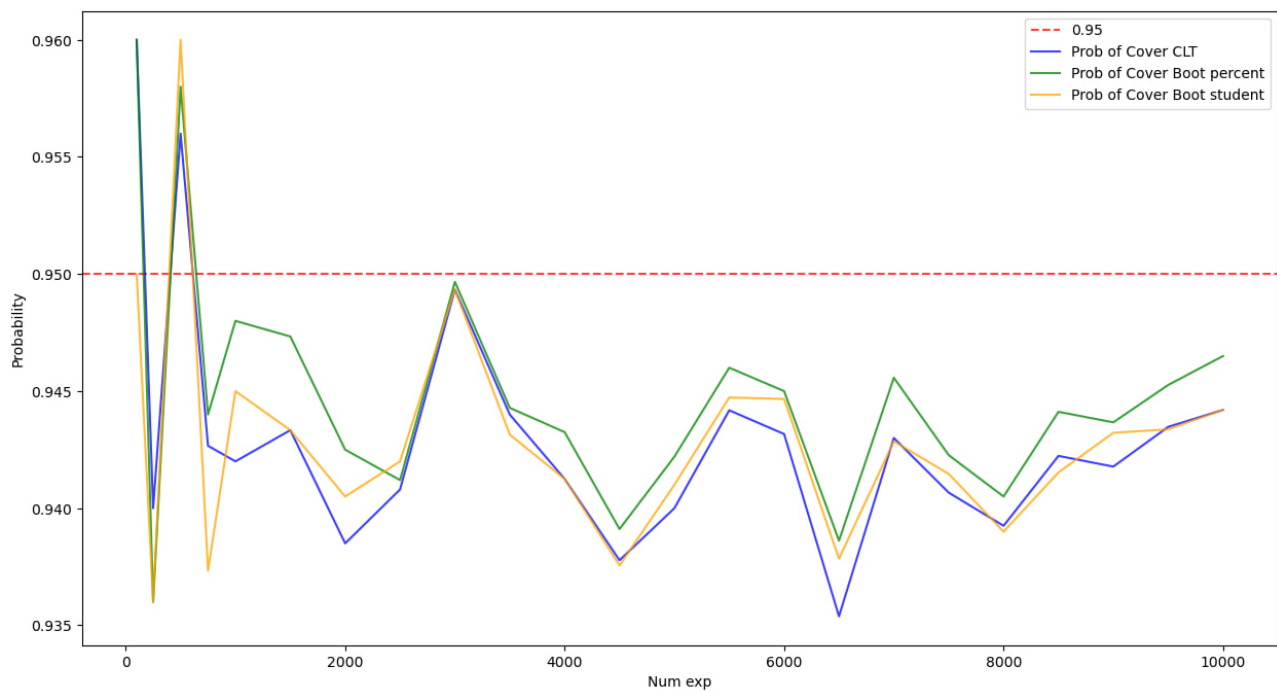


Рисунок 2.7. Зависимость оцениваемой вероятности накрытия от числа экспериментов, равномерное распределение, 50 наблюдений

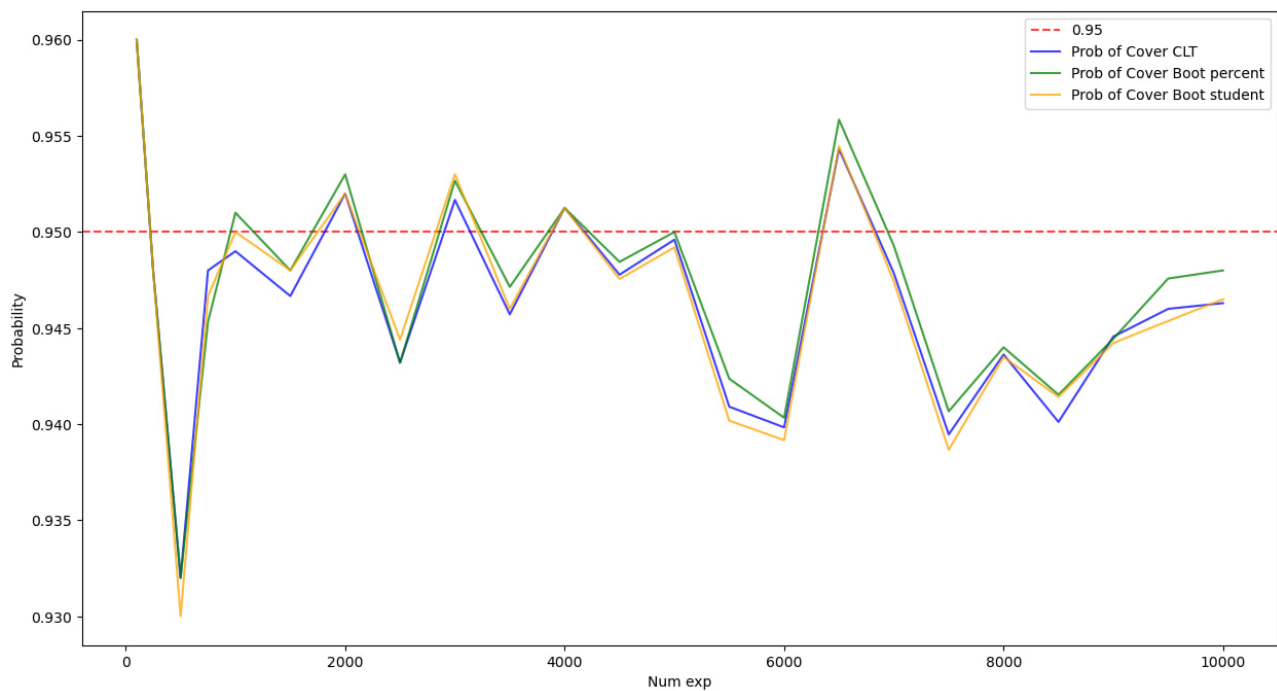


Рисунок 2.8. Зависимость оцениваемой вероятности накрытия от числа экспериментов, равномерное распределение, 100 наблюдений

2.6. Бутстрэп с коррекцией смещения и ускорением (Bias-corrected and accelerated bootstrap, BCa)

Методом, модифицирующим перцентильный бутстрэп и позволяющий корректировать асимметрию выборочного распределения, является ВСа-бутстрэп. В нем дополнительно оцениваются два параметра: \hat{a} — параметр "ускорения", и \hat{z}_0 — фактор коррекции смещения. Алгоритм ВСа-бутстрэпа выглядит следующим образом:

1. Повторяются шаги 1, 2 из метода обычного перцентильного бутстрэпа;
2. Оценивается фактор коррекции смещения:

$$\hat{z}_0 = \Phi^{-1} \left(\frac{1}{n_{boot}} \sum_{j=1}^{n_{boot}} \mathbb{I}\{\hat{\theta}_j^* < \hat{\theta}\} \right)$$

где Φ^{-1} — обратная функция от функции распределения стандартной нормальной случайной величины, а ее аргументом является доля бутстрэпированных оценок, меньших выборочной оценки.

3. Далее необходимо оценить параметр "ускорения". Введем обозначения: $\hat{\theta}_{(i)}$ — оценка параметра θ без учета наблюдения X_i , $\hat{\theta}_{(.)} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{(i)}$ — среднее таких оценок (jackknife оценка). Тогда параметр "ускорения" оценивается следующим образом:

$$\hat{a} = \frac{\sum_{i=1}^n \left(\hat{\theta}_{(.)} - \hat{\theta}_{(i)} \right)^3}{6 \left[\sum_{i=1}^n \left(\hat{\theta}_{(.)} - \hat{\theta}_{(i)} \right)^2 \right]^{3/2}}$$

4. Следующим шагом вычисляются квантили бутстрэпированного распределения $\hat{\theta}$, которые станут границами доверительного интервала:

$$\alpha_1 = \Phi \left(\hat{z}_0 + \frac{\hat{z}_0 + z_\alpha}{1 - \hat{a}(\hat{z}_0 + z_\alpha)} \right)$$

$$\alpha_2 = \Phi \left(\hat{z}_0 + \frac{\hat{z}_0 + z_{1-\alpha}}{1 - \hat{a}(\hat{z}_0 + z_{1-\alpha})} \right)$$

где Φ — функция распределения стандартного нормального распределения, z_α — квантиль стандартного нормального распределения порядка α .

5. Оценным $(1 - \alpha)$ -процентным доверительным интервалом истинного параметра θ будет отрезок: $\widehat{CI}_\theta = [\hat{\theta}_{\alpha_1}^*; \hat{\theta}_{\alpha_2}^*]$.

3. Бутстрэп в линейной регрессии

3.1. Парный бутстрэп

3.2. Параметрический бутстрэп

4. Бутстрэп в решающих деревьях и ансамблях

4.1. Бэггинг