



# Вероятностное моделирование

## лекция 2

**Артур Игнатъев**

CS Space, МХН СПбГУ, ИТМО

**Владимир Евменов**

CS Space, Huawei



- А если мы хотим оценивать характеристику не только одним числом?
- Среднее число клиентов в день находится в каком-то интервале с большой вероятностью
- Бизнес вопрос. Сколько проект заработает за месяц?
  - Можно дать оценку суммарного дохода за месяц
  - А можно сказать, что с большой вероятностью доход будет от  $a$  до  $b$



- Уже умеем оценивать характеристики
  - Если можем сэмплировать из  $X$
1. Генерируем  $N$  выборок  $x_{[n],1}, \dots, x_{[n],N}$  размера  $n$ .
  2. Для каждой из них считаем нужную нам характеристику  $\phi_i^* = \phi^*(x_{[n],i})$ .
  3. У полученной выборки  $\phi_{[N]}^*$  считаем дисперсию, доверительный интервал.
- Выборки  $x_n^*$  и  $\phi_N^b$  называются **бутстраповскими**.



- Всегда ли мы можем сэмплировать из  $X$ ?
- Если можем проводить эксперимент много раз, то круто
- А если с.в. это данные, например, из бизнеса
- **Бутстрап (bootstrap)** — это метод Монте-Карло, примененный к какой-либо аппроксимации  $\mathcal{P}_X$  (например, к  $\mathcal{P}_{[n]}^*$ ).
- Усредненный рецепт такой:
  1. Генерируем  $N$  выборок  $x_{[n],1}^*, \dots, x_{[n],N}^*$  размера  $n$  из  $X^*$ . Для генерации одной выборки  $x_{[n],i}^*$  нужно взять  $n$  случайных элементов исходной выборки  $x_{[n]}$ .
  2. Для каждой из них считаем нужную нам характеристику  $\phi_i^b = \phi^*(x_{[n],i}^*)$ .
  3. У полученной выборки  $\phi_{[N]}^b$  считаем дисперсию, доверительный интервал.
- Выборки  $x_n^*$  и  $\phi_N^b$  называются **бутстраповскими**.



Ошибка бутстрапа складывается из двух слагаемых:

- **устраняемая ошибка** — из-за того, что мы взяли  $N$  выборок, а не  $\infty$ ,
- **неустраняемая ошибка** — из-за того, что мы взяли  $\mathcal{P}_{[n]}^*$  вместо  $\mathcal{P}_X$ .

# Интервальные оценки



Вместо того, чтобы приблизить значение  $\phi$  точно, мы хотим оценить область, в которой лежит истинное значение.

# Доверительный интервал



- Пара статистик  $(\phi_L^*, \phi_R^*)$  называется **доверительным интервалом** для  $\phi(\mathcal{P}_X)$  с уровнем доверия  $\gamma$  если

$$P(\phi_L^*(X_{[n]}) \leq \phi(\mathcal{P}_X) \leq \phi_R^*(X_{[n]})) = \gamma$$

- Если сгенерировать  $N$  выборок из распределения  $\mathcal{P}_X$  и для каждой из них построить доверительный интервал, то примерно  $\gamma N$  из них накроют истинное значение  $\phi(\mathcal{P}_X)$ .
- Конкретный доверительный интервал  $(\phi_L^*(x_{[n]}), \phi_R^*(x_{[n]}))$  либо покрывает истинное значение, либо нет, никакой вероятности  $\gamma$  этого события нет: нельзя говорить «Этот интервал с вероятностью 90% содержит истину».



- Пусть  $\phi_L^*$  —  $\alpha_1$ -квантиль бутстраповской выборки  $\phi_{[N]}^*$ , а  $\phi_R^* = 1 - \alpha_2$ -квантиль. Тогда  $(\phi_L^*, \phi_R^*)$  это асимптотический доверительный интервал с уровнем доверия  $1 - \alpha_1 - \alpha_2$ .



# Пример



Контрольная группа	А	Б
Размер группы	893	923
Число купивших	34	28
Конверсия	3.81%	3.03%

- Ого, выигрыш в 26%, время идти просить премию, или нет?
- Аналитический можно показать, что разница мат. ожиданий лежит в интервале  $(-0.93\%, 2.48\%)$  с вероятностью 0.95
- Давайте предположим, что у нас нету возможности дать ответ аналитический. Воспользуемся бутстрапом.
- 10 000 раз сэмплируем из обеих групп, считаем разницу средних, далее строим доверительный интервал по распределению разницы.
- Получаем  $(-0.89\%, 2.4\%)$