

Занятие №7

# Непараметрические тесты и методы

Сергей Москвин



- 
1. Некоторые факты из прошлой лекции
  2. Проверка на нормальность
  3. Последствия отклонения от нормальности
  4. Непараметрические тесты
  5. Resampling

# Некоторые важные факты из прошлой лекции

---



- Доверительные интервалы
- Тесты: t-test
- Ошибки первого и второго рода
- P-value

# Доверительный интервал



- $\alpha$ -процентный двусторонний доверительный интервал для выборочного среднего

$$\bar{X} - T_{\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + T_{1-\frac{\alpha}{2}} * \frac{\sigma}{\sqrt{n}}$$

- Подразумевает нормальное распределение  $X$

# Тесты на соответствие распределению



## Нормальное распределение

- Шапиро-Уилка

## Произвольное распределение

- Колмогорова
- Критерий Пирсона (для дискретных величин)
- QQ-plot



Notebook Goodness of Fit

# Отклонения от нормальности: выбросы

---



Игнорирование выбросов

- Trimming
- Windzorizing

Использование робастных статистик



**Робастность** – устойчивость статистики к наличию выбросов в данных

**Breakdown point** – минимальное количество выбросов, необходимое для того, чтобы тестовая статистика приняла сколь угодно большое/малое значение

BP для среднего – 0%

BP для медианы – 50%

# Трансформация данных



Если ненормальность данных состоит не в наличии выбросов, а в ненормальной форме распределения, можно попробовать произвести трансформацию данных

Например:

- Логарифм
- Box-Cox transform





# Непараметрические тесты



- Sign test
- Wilcoxon signed rank test
- Mann-Whitney rank sum test



Notebook non\_parametric\_tests

# Предпосылки тестов



- То, что тест непараметрический, еще не означает, что у него нет набора предпосылок

Wilcoxon sign-rank

Симметричные  
распределения

Mann-Whitney

Одинаковые  
распределения

- Непараметрические тесты чувствительны к выборкам с разными дисперсиями



Notebook Parametric\_vs\_nonparametric

# Аргументы за использование параметрических тестов

---



- В силу ЦПТ статистики (например, среднее), построенные на выборках достаточно большого размера будут асимптотически иметь нормальное распределение
- Параметрические тесты можно применять для выборок с разной дисперсией
- Мощность тестов – при выполнении предпосылок, параметрические тесты, как правило, обладают большей мощностью



# Аргументы за использование непараметрических тестов

---



- Более адекватны для задачи, в которых распределение лучше характеризует медиана
- Маленький объем выборки, для которого не выполняется ЦПТ
- Наличие выбросов в исходных данных
- Ранговые или ординальные данные

# Resampling



Методы для получения качественных оценок распределения статистик без предпосылок о типе распределения исходной случайной величины

- Jackknife
- Bootstrap



Notebook Bootstrap

# Восстановление плотности распределения



Позволяет:

- Визуализировать данные
- Дифференцировать функцию плотности



Notebook KDE

# Множественная проверка гипотез



	Product A	Product B
Criterion 1	50	55
Criterion 2	45	40
Criterion 3	90	90
..		
Criterion n	30	60

# Множественная проверка гипотез



FWER (Family-wise error rate)

FDR (False Discovery Rate)



Notebook Multiple HT





Факт возникновения проблем с множественной проверкой гипотез можно предвидеть заранее. Его нужно учитывать при разработке дизайна исследования

- Создание предварительного плана исследования
- Создание индексов
- Дизайн исследования устойчивый к FWER

# Дополнительные соображения

---



- Data dredging
- P-hacking

# Домашнее задание



Файл с данными salaries.csv

1. Сформулируйте нулевые и альтернативные гипотезы и проведите следующие тесты

1.1 Одновыборочный параметрический тест

1.2 Двухвыборочный параметрический и непараметрический тесты

1.3 Параметрический и непараметрический тесты для парных наблюдений

2. Сформулируйте точную альтернативную гипотезу для одного из проведенных тестов и оцените мощность теста

3. Рассматривая проведенные тесты как единую группу, сделайте поправку на множественное тестирование (любую). Изменились ли результаты?

4. Постройте параметрические доверительные интервалы для среднего и непараметрические интервалы на основе бутстрепа



**Спасибо за  
внимание!**

**Сергей Москвин**

`smos@list.ru`