

# Дисперсионный анализ

ПМИ ФКН ВШЭ, 2 ноября 2019 г.

Денис Деркач, Алексей Артёмов

ФКН ВШЭ

Денис Деркач

# Оглавление

Мотивация

Тестируемые гипотезы

Дисперсионный анализ (ANalysis Of VAriances)

Анализ Краскела — Уоллиса

Апостериорные тесты

Анализ контрастов

Мотивация

# Мотивация

- › Ранее, мы рассматривали одно- и двухвыборочные тесты (например,  $t$ -тест);
- › что делать, если мы хотим сравнить сразу несколько выборок?

# Мотивирующий пример

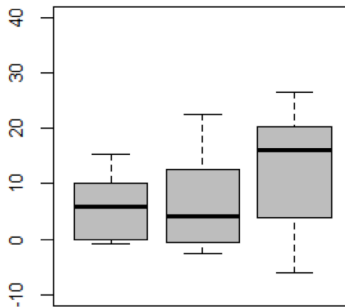
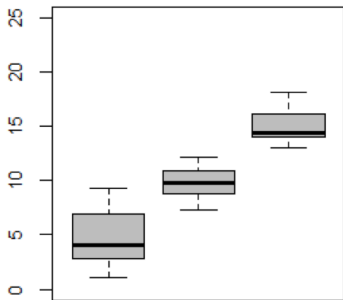
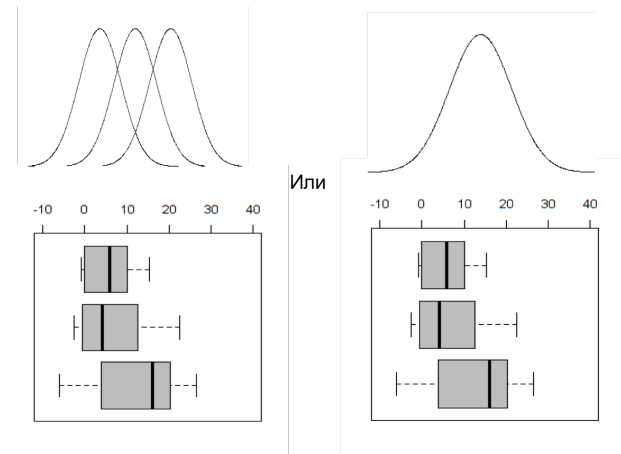


График слева: мы можем сказать, что средние, видимо, отличаются.  
Что мы можем сказать про график справа?

# Мотивирующий пример

Иными словами, пришли ли выборки из одного распределения или из разных?



# Обсуждение мотивирующего примера

- › Если события во всех выборках происходят из одного распределения, то все события должны быть распределены одинаково, как внутри группы, так между группами.
- › Это, в частности, означает асимптотическое равенство дисперсий и средних.
- › Необходимо сравнить дисперсию внутри группы с дисперсией между группами, для того, чтобы получить вывод о равенстве средних.

# Тестируемые гипотезы



# Тестирование гипотез: нулевая гипотеза

Сформулируем гипотезы:

›  $H_0: \mu_1 = \mu_2 = \mu_3.$

›  $H_1: \mu_1 \neq \mu_2 \text{ или } \mu_2 \neq \mu_3 \text{ или } \mu_1 \neq \mu_3.$

# Какой тест предпочесть?

Мы можем взять несколько попарных  $t$ -тестов, проверяя:

›  $H_0: \mu_1 = \mu_2;$

$H_1: \mu_1 \neq \mu_2.$

›  $H_0: \mu_2 = \mu_3;$

$H_1: \mu_2 \neq \mu_3.$

› и т.д.

Проблема: вероятность ошибки первого рода резко увеличивается.

**TABLE 1:** Probability of Committing at Least One Type I Error by Using Two-Sample  $t$  Tests for All  $C$  Pairwise Comparisons of  $k$  Means\*

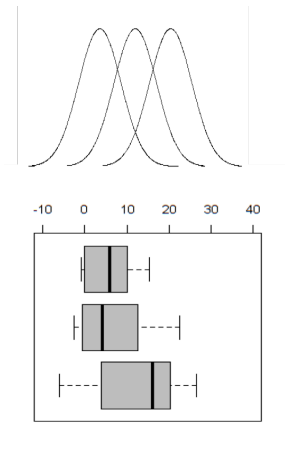
$k$	$C$	Level of Significance, $\alpha$ , Used in the $t$ Tests				
		0.10	0.05	0.01	0.005	0.001
2	1	0.10	0.05	0.01	0.005	0.001
3	3	0.27	0.14	0.03	0.015	0.003
4	6	0.47	0.26	0.06	0.030	0.006
5	10	0.65	0.40	0.10	0.049	0.010
6	15	0.79	0.54	0.14	0.072	0.015
10	45	0.99	0.90	0.36	0.202	0.044
	$\infty$	1.00	1.00	1.00	1.000	1.000

\*There are  $C = k(k - 1)/2$  pairwise comparisons of  $k$  means. This is the number of combinations of  $k$  items taken two at a time.

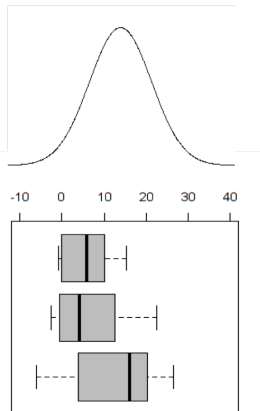
# Идея ANOVA

- › Заметим, что при верной  $H_0$  все группы получены из популяций с одинаковыми средним  $\mu$  и дисперсией  $\sigma^2$ .

# Вспомним пример



Или



Как мы можем отличить два случая?

# Идея ANOVA

- › Заметим, что при верной  $H_0$  все группы получены из популяций с одинаковыми средним  $\mu$  и дисперсией  $\sigma^2$ .
- › Давайте оценим дисперсию разными независимыми способами и сравним!

# Двухвыборочный F-тест

Предположим, что у нас есть две iid выборки

$X_1, \dots, X_n \sim \mathcal{N}(\mu_X, \sigma_X^2)$  и  $Y_1, \dots, Y_m \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ . Необходимо проверить, что  $X$  и  $Y$  имеют одинаковую дисперсию.

$$H_0 : \sigma_X^2 = \sigma_Y^2;$$

$$H_a : \sigma_X^2 \neq \sigma_Y^2;$$

Мы можем оценить дисперсии:

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

$$S_Y^2 = \frac{1}{m-1} \sum_{i=1}^m (Y_i - \bar{Y})^2$$

# Распределение Фишера

Тогда статистика:

$$\frac{S_X^2}{S_Y^2} \sim F(m-1, n-1).$$

Где  $F$  - распределение Фишера, характерной функцией плотности вероятности для которого является:

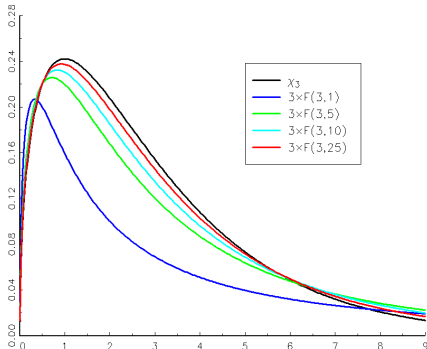
$$f(x; d_1, d_2) = \frac{1}{B\left(\frac{d_1}{2}, \frac{d_2}{2}\right)} \left(\frac{d_1}{d_2}\right)^{\frac{d_1}{2}} x^{\frac{d_1}{2}-1} \left(1 + \frac{d_1}{d_2} x\right)^{-\frac{d_1+d_2}{2}},$$

где  $B$  - бета-функция.

# Распределение Фишера

$$\frac{S_X^2}{S_Y^2} \sim F(m-1, n-1).$$

Распределение часто используется для характеристики отношений случайных величин  $\chi_k/\chi_m$ ,  $\Gamma_1/\Gamma_2$  и пр.

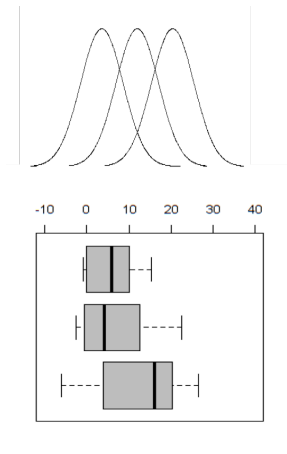




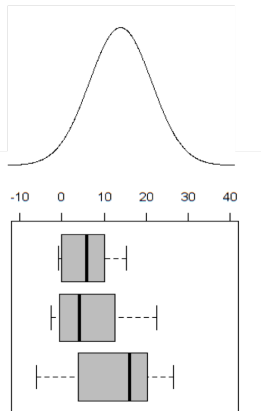
# Свойства F-теста

- › односторонний;
- › относительно неустойчив к ненормальности распределений;
- › в случае больших  $\alpha$  или близких размеров выборок, может проявлять устойчивость, при этом всё равно теряет мощность.

# Вспомним пример



Или



Как мы можем отличить два случая?

# Идея ANOVA

- › Заметим, что при верной  $H_0$  все группы получены из популяций с одинаковыми средним  $\mu$  и дисперсией  $\sigma^2$ .
- › Давайте оценим дисперсию разными независимыми способами и сравним!
- › Можем оценить исходя из вариативности внутри группы и вариативности между групп.

# Дисперсионный анализ (ANalysis Of VAriances)

# Данные

Данные состоят из  $N = \sum_{j=1}^k n_j$  наблюдений  $x_{ij}$ , по  $n_j$  наблюдений в  $j$ -й выборке,  $j = 1, \dots, k$ . Фактически, мы анализируем такую таблицу:

Измерения			
1	2	...	$k$
$x_{11}$	$x_{12}$	...	$x_{1k}$
$x_{21}$	$x_{22}$	...	$x_{2k}$
...	...	...	...
$x_{n_1 1}$	$x_{n_2 2}$	...	$x_{n_k k}$

Для каждого столбца характерно среднее значение  $\mu_j$ , оцениваемое  $\bar{X}_j$

# Тестирование гипотез: нулевая гипотеза

Сформулируем нулевую гипотезу:

$$> H_0: \mu_1 = \mu_2 = \mu_3.$$

Это сложная гипотеза, то есть она содержит в себе много парных и непарных простых (например  $\mu_1 = \mu_2$  или  $\mu_2 = \frac{\mu_1 + \mu_3}{2}$ ).

NB: есть другие способы сформулировать нулевую гипотезу ANOVA.

NB2: на этой лекции рассматриваем только однофакторную (one-way) ANOVA, то есть анализ с одной независимой переменной.

# Тестирование гипотез: альтернативная гипотеза

Сформулируем альтернативную гипотезу:

$$> H_1: \mu_1 \neq \mu_2 \text{ или } \mu_2 \neq \mu_3 \text{ или } \mu_1 \neq \mu_3.$$

Заметим, что  $H_0$  отвергается, если верна хотя бы одна из маленьких частных альтернативных гипотез (парных или комплексных).

NB: ANOVA не говорит какая.

# Вариативность между группами

Оценим  $\sigma^2$  на основе дисперсии средних между группами (посчитаем ошибку среднего, как будто это выборочные средние, и из неё вычислим дисперсию):

$$s_{\bar{x}} = \frac{\sum_j (\bar{X}_j - \bar{\bar{X}})}{k - 1}$$

Тогда mean square between groups ( $MS_B$ ):

$$MS_B = \frac{\sum_j (\bar{X}_j - \bar{\bar{X}})n_j}{k - 1}$$

Количество степеней свободы при этом:

$$DF_B = k - 1,$$

где  $k$  — число групп.



# Вариативность внутри группы

Mean square within groups = error MS

$$MS_W = \frac{s_1^2 + \dots + s_k^2}{k}$$

Количество степеней свободы при этом:

$$DF_B = N - k,$$

где  $k$  — число групп,  $N$  - полное число семплов.

# F-статистика

$$F = \frac{\text{оценка дисперсии между группами}}{\text{оценка дисперсии внутри групп}} = \frac{MS_B}{MS_W}.$$

Тестирование  $H_0$

- › для заданных df рассчитывается критическое значение  $F$ ;
- › на основе групп считается  $F$  и сравнивается с критическим значением;
- › если  $F$  больше критического —  $H_0$  о равенстве средних в группах отвергается;
- ›  $F$  - это отношение дисперсий, оно имеет особое распределение, оно всегда положительно; ANOVA — принципиально односторонний тест.

# Sum of squares

SS - это суммы квадратов отклонений (sum of squared deviations):

- ›  $SS_{\text{Between}}$  - сумма квадратов отклонений каждого среднего в группе от общего среднего = Effect;
- ›  $SS_{\text{Within}}$  — сумма квадратов отклонений каждого измерения от среднего в соответствующей группе = Error;
- ›  $SS_{\text{Total}}$  — сумма квадратов отклонений каждого измерения от общего среднего = Total.

При этом:

$$SS_T = SS_W + SS_B.$$

# ANOVA достигнутый эффект

Для того, чтобы понять насколько значим полученный результат в тесте строят два типа переменных:

- ›  $R^2 = \eta^2 = \frac{SS_B}{SS_T}$ , чем выше  $R^2$ , тем больше полученный эффект.
- ›  $f = \frac{s_{\bar{x}}}{\sqrt{MS_W}}$ , чем выше, тем больше полученный эффект.

# Глоссарий ANOVA

Типы переменных:

- › Группирующая переменная, фактор (factor, predictor).
- › Зависимая переменная (dependent variable, response).

Мы пока разбираем случай с одним фактором (one-way). В ANOVA одна зависимая переменная, а факторов может быть несколько, и они могут составлять довольно сложные конструкции.

# Дизайн эксперимента

Факторы могут быть двух видов:

- › fixed. Рассматриваются именно эти значения фактора. Другие значения не существуют или не интересуют. Пример: пол, время суток и тд.
- › random. Рассматриваются случайно выбранные значения фактора из многих возможных. За пределами исследования существуют другие значения фактора. Пример: происхождение семпла, процент лекарства.

Для этих типов факторов по-разному оценивается межгрупповая изменчивость. Когда фактор один, это не важно, но в сложных моделях с несколькими факторами эти различия очень важны!

# Допущения ANOVA

- › Выборки должны быть случайными, измерения — независимыми.
- › Размеры групп должны различаться как можно меньше.
- › Нормальность в каждой группе по отдельности.
- › Равенство дисперсий в группах.

# Требование нормальности

Возможные проверки:

- › Сделать тест по методу моментов (Обычно достаточно проверить эксцесс и асимметрию).
- › Построить гистограмму распределения остатков ( $x_i - \bar{x}$ ) внутри каждого семпла (и проверить goodness-of-fit тесты).
- › Тест Шапиро-Уилка.
- › Тест Д'Агостино-Пирсона.
- › Тест Андресона-Дарлинга.

Подробный разбор тестов в этой ссылке.



# Тест Шапиро-Уилка

Критерий Шапиро-Уилка основан на оптимальной линейной несмещённой оценке дисперсии к её обычной оценке методом максимального правдоподобия. Статистика критерия имеет вид:

$$W = \frac{1}{s^2} \left[ \sum_{i=1}^n a_{n-i+1} (x_{n-i+1} - x_i) \right]^2 ,$$

где  $s^2 = \sum_{i=1}^n (x_i - \bar{x})^2$ ,  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ , а коэффициенты  $a_{n-i+1}$  берутся из таблиц.

На практике, следует проверять применимость таблиц в том или ином софте. Можно также использовать тест Шапиро-Франча.

# Ненормальные данные

- › Если семплы достаточно большие, можно оставить как есть.
- › Провести преобразование к нормальным:
  - › стандартизовать распределение;
  - › метод Бокса-Кокса.
- › Использовать непараметрический ANOVA:
  - › односторонний дисперсионный анализ Краскела—Уоллиса.

# Метод Бокса-Кокса

Для последовательности:  $\{y_1, \dots, y_n\}$ ,  $y_i > 0$  однопараметрическое преобразование Бокса-Кокса с параметром  $\lambda$  определяется следующим образом:

$$y_i^\lambda = \begin{cases} \frac{y_i^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0, \\ \log(y_i), & \text{if } \lambda = 0. \end{cases}$$

Где  $\lambda$  - свободный параметр.

# Требование равенства дисперсий

Для проверки можно:

- › Использовать тест Левена (Ливиня).
- › Использовать F-тест.
- › Построить зависимость остатков (residuals) от средних.

# Тест Ливиня

Проверяет равенство дисперсий всех семплов.

$$W = \frac{(N - k)}{(k - 1)} \cdot \frac{\sum_{i=1}^k N_i (Z_{i\cdot} - Z_{..})^2}{\sum_{i=1}^k \sum_{j=1}^{N_i} (Z_{ij} - Z_{i\cdot})^2},$$

Здесь  $Z$  может центрироваться на:

- › среднее выборки (для симметричных распределений);
- › медиану выборки (для асимметричных распределений);
- › усечённое среднее выборки (для распределений с тяжёлыми хвостами).

В общем случае рекомендуют использовать медиану.

Значение  $W$  затем сравнивается с соответствующим  $F$  распределением.

# Свойства ANOVA

- › Возможно провести one-way ANOVA в случае, если у нас в руках есть только средние значения, показатели разброса (SD, SE,  $s^2$ ) и размер выборок (например, из какой-нибудь статьи).
- › В случае двух выборок ANOVA эквивалентна  $t$ -тесту.

# Анализ Краскела — Уоллиса

# Мотивация

В некоторых случаях невозможно добиться данных с нормальным распределением. Потому необходимо использовать непараметрические тесты или линейные модели.



# Анализ Краскела — Уоллиса

- › непараметрический тест, который не требует нормальности данных (но чувствителен к разным дисперсиям);
- › менее мощный, чем параметрические тесты, в случае однофакторного анализа мощность около 95%, в остальных случаях ниже 80%.
- › для двух групп эквивалентен тесту Манна-Уитни.

NB: если данные гетероскедастичны, следует использовать дисперсионный анализ Уелча.

# Алгоритм

1. выставить ранг  $r_{ij}$  согласно значению всех  $x_{ij}$ ;
2. положив:

$$S_j = \sum_{i=1}^{n_j} r_{ij}, r_{.j} = S_j/n_j,$$

тогда

$$r_{..} = 1/N \sum r_{ij} = (N + 1)/2$$

подсчитать статистику:

$$W = \left( \frac{12}{N(N + 1)} \sum_{j=1}^k S_j^2/n_j \right) - 3(N + 1),$$

3. найти критическую область для  $W$ -статистики.

# Пост-Мотивация

$W$  мотивирована сравнением рангов между группами и внутри группы (то есть, анализ дисперсий на рангах):

$$W = (N - 1) \frac{\sum_{j=1}^k n_j (r_{j\cdot} - r_{\cdot\cdot})^2}{\sum_{j=1}^k \sum_{i=1}^{n_j} (r_{ij} - r_{\cdot\cdot})^2},$$

# Совпадения

Мы использовали, что  $\sum_i r_i = \frac{N(N+1)}{2}$ . В случае, если есть совпадающие значения  $x_{ij}$ , это не так. Для вычисления  $W$  в этом случае следует пользоваться коррекцией:

- › при подсчёте  $W$  использовать среднее значение рангов для совпадающих элементов;
- › для тестирования использовать  $W' = W/\gamma$ , где  $\gamma = 1 - \frac{1}{N(N^2-1)} \sum_{m=1}^g l_m(l_m^2 - 1)$ , где  $g$  - число групп совпадений,  $l_m$  - количество элементов в  $m$ -й группе

# Наблюдения

- › для  $k > 5$  можно сравнивать с  $\chi^2_{k-1}$ ;
- › для нормальной работы теста рекомендуют  $n_i > 5$ ;
- ›  $H_0$ : все группы происходят из одного распределения (Отличается от обычного дисперсионного анализа!).

# Многофакторные непараметрические задачи

В случае бОльшего количества факторов можно использовать улучшения теста Краскела-Уоллеса:

- › Scheirer-Ray-Hare Test
- › Jonckheere-Terpstra Test

NB: относительная мощность непараметрических тестов быстро падает с ростом количества факторов.

# Апостериорные тесты

# Апостериорные (post-hoc) тесты

ANOVA не называет причину, по которой была отвергнута гипотеза  $H_0$ . Потому используют апостериорные тесты:

- › Сначала сравнить все группы между собой с помощью ANOVA.
- › Если различия есть, использовать методы множественного сравнения (сравнивают группы попарно, сохраняя общую  $\alpha = 0.05$ ).
- › Если различий нет, анализ следует считать завершённым (и не проводить post-hoc тесты).

NB: проведение апостериорных тестов может испортить весь анализ.



# Тест Тьюки

Он же honestly significant difference test (HSD test) or wholly significant difference test (WSD test).

- › Выстраиваем средние по выборке по возрастанию.
- › Строим статистику  $q = \frac{Y_A - Y_B}{SE}$ , где  $Y$  - среднее (причём  $Y_A > Y_B$ , SE - стандартное отклонение.
- › Ищем значимость  $q_{\alpha, N-k, k}$  для нужного  $\alpha$ .

# Тест Тьюки

- › Наиболее распространённый и рекомендуемый в литературе тест;
- › строго контролирует  $\alpha$  (0.05);
- › проверяет все парные гипотезы сразу;
- › плохо работает, если размер групп сильно различается;
- › чувствителен к неравенству дисперсий;
- › считает статистику ( $q$ ) на основе  $MS_{within}$  и  $df$ .

## Другие post-hoc тесты

- › Тьюки-Крамера, решает проблему теста Тьюки для неравных выборок.
- › Критерий Ньюмена-Кейлса. Все средние упорядочивают по возрастанию и пошагово вычисляют статистики; начинают от сравнения наибольшего с наименьшим. Сравнивают с  $q_{\alpha, N-k, p}$ , где  $p$  — диапазон средних. Мощнее теста Тьюки, но плохо контролирует ошибку 1-го рода.
- › Критерий Шеффе (Scheffe test) — очень консервативный, мощность меньше, чем у теста Тьюки (но см. ниже).
- › Критерий Даннетта (Dunnett test) — используется для сравнения нескольких групп с контрольной группой, мощнее, чем тест Тьюки. Размер контрольной группы рекомендуется делать больше, чем размеры остальных групп в  $\sqrt{k-1}$  раз.

# Failed post-hoc

Бывает так, что в ANOVA нулевая гипотеза отвергается, а пост-хок тесты не обнаруживают различий, так как их мощность ниже. В этом случае необходимо увеличивать размер выборки.

# Анализ контрастов

# Анализ контрастов (planned comparisons)

- › Проводится вместо ANOVA.
- › Важно: то, какие группы сравнивать, выбирают заранее, до проведения какого-либо анализа. В идеале — ещё при постановке исследования.
- › В тесте проверяется только одна гипотеза;
- › Можно провести 2-3 таких теста в пределах одного «набора» групп, только надо следить, чтобы сравнения не сильно перекрывались, не были избыточными.
- › Мощнее post-hoc тестов.

# Пример

У нас 4 группы тигров, их кормят: овощами; фруктами; рыбой; мясом.

Вопрос: отличается ли масса тигров, питающихся животной и растительной едой?

# Построение контрастов

Контраст — линейная комбинация средних значений.

Коэффициенты сравнения — константы, на которые умножены средние. Таким образом гипотезы формулируются:

$$H_0 : \sum_i C_i \mu_i = 0;$$

$$H_1 : \sum_i C_i \mu_i \neq 0.$$

При этом  $\sum_i C_i = 0$ .

Если тестируется несколько гипотез:  $\sum_i C_{1,i} C_{2,i} = 0$ . В этом случае статистика строится:  $t = \frac{\sum_i C_i \mu_i}{SE}$  и имеет  $t$  распределение.



# Пример

У нас 4 группы тигров, их кормят: овощами; фруктами;рыбой; мясом.

Вопрос: отличается ли масса тигров, питающихся животной и растительной едой?

Мы строим контраст:  $\frac{1}{2}\mu_1 + \frac{1}{2}\mu_2 - \frac{1}{2}\mu_3 - \frac{1}{2}\mu_4$ .