

In [22]:

```
1 import numpy as np
2 import pandas as pd
3 import scipy.stats as sps
4 import seaborn as sns
5 from matplotlib import pyplot as plt
6 from collections import namedtuple
7 from tqdm import tqdm_notebook as tqdm
8
9 %matplotlib inline
```

## Задача 1

Предположим, что вы разработали лекарство от коронавируса. Перед применением оно обязательно должно пройти клинические испытания. Для начала было разрешено проверить лекарство на двух независимых группах по 10 человек. Одна группа принимает плацебо, другая -- ваш препарат. Большее количество пациентов на первом этапе брать не разрешают -- слишком велики риски отрицательного результата.

Для каждого пациента измерялось количество дней от приема препарата до выздоровления. Получились следующие результаты:

In [5]:

```
1 placebo = [6, 16, 8, 13, 9, 4, 7, 10, 3, 14]
2 medicine = [5, 10, 3, 1, 5, 3, 19, 2, 2, 5]
```

Что вы можете сказать на основе этих результатов?

- Лекарство эффективнее, подтверждается статистическими методами;
- Наверное, лекарство эффективнее, но статистическими методами это пока не подтверждено, нужно продолжить эксперимент. Подумайте, как обосновать необходимость продолжение эксперимента;
- По результатам эксперимента нельзя сделать какой-либо вывод. Стоит ли продолжать эксперименты? Если да, то четко это обоснуйте;
- Лекарство неэффективно, нужно немедленно прекращать эксперимент.

Естественно, для того, чтобы сделать выводы по этой задаче, требуется воспользоваться критериями дисперсионного анализа для сравнения средних.

Поскольку люди были разные, можем пользоваться методами для независимых выборок.

### t-test

Самый простой **ВОЗМОЖНЫЙ** способ - это критерий Стьюдента для независимых выборок

In [6]:

```
1 sps.ttest_ind(placebo, medicine, equal_var=False)
```

Out[6]:

```
Ttest_indResult(statistic=1.6077916961893066, pvalue=0.12612583305411368)
```

В данном случае нет отвержения. При этом этот критерий требует нормальности выборок, которой может здесь не быть, либо же большой выборки

### Критерий Уилкоксона-Манна-Уитни

Ещё один **ВОЗМОЖНЫЙ** способ проверки - критерий Манна-Уитни. Применим его

In [7]:

```
1 sps.mannwhitneyu(placebo, medicine, alternative='greater')
```

Out[7]:

```
MannwhitneyUResult(statistic=76.5, pvalue=0.024256557087931816)
```

Тогда гипотеза равенства средних отвергается на уровне значимости 0.05, однако надо быть внимательным.

В данном случае размер выборки всего 10, а асимптотика в критерии Манна-Уитни без использования поправок справедлива лишь тогда, когда размер обеих выборок не меньше 50

### Перестановочный критерий

Оставшийся вариант - разобранные на лекции перестановочные критерии. Если даны две выборки  $X$  и  $Y$ , то в качестве статистики можно взять  $\bar{X} - \bar{Y}$ .

В качестве группы перестановок:

$$G = \{(s_1, \dots, s_{n+m}) \mid \{s_1, \dots, s_n\} \in C_{\{1, \dots, n+m\}}^n\}$$

Посчитаем статистику для всех перестановок отсюда

In [8]:

```
1 N_SAMPLES = int(1e5) # количество бутстрепных перестановок
```

In [9]:

```
1 n = len(placebo)
2 m = len(medicine)
3 concated = np.array(placebo + medicine)
4
5 bootstrap_stats = []
6
7 for _ in tqdm(range(N_SAMPLES)):
8     # выбираем индексы для первой части
9     index_x = np.random.choice(np.arange(n + m), n, replace=False)
10    # выделяем оставшиеся
11    index_y = np.arange(n + m)[~np.in1d(np.arange(n + m), index_x)]
12
13    x_iter = concated[index_x]
14    y_iter = concated[index_y]
15
16    # добавляем статистику от перестановки
17    bootstrap_stats.append(np.mean(x_iter) - np.mean(y_iter))
```

Смотрим, в какой доле случаев статистика после перестановки оказывается больше изначальной (случай для односторонней альтернативы)

In [10]:

```
1 (np.array(bootstrap_stats) > (np.mean(placebo) - np.mean(medicine))).mean()
```

Out[10]:

0.05657

Получаем, что снова нет отвержения.

Однако это не единственная статистика, которую можно использовать. Например, можно взять статистику Манна-Уитни. Это сумма рангов срока выздоровления при приёме лекарства, то есть, если лекарство имеет эффект, то выздоровление происходит скорее, и статистика будет меньше, чем в противном случае. Так что знак в проверке стоит развернуть.

Проведём такую же процедуру для статистики Манна-Уитни.

In [11]:

```
1 n = len(placebo)
2 m = len(medicine)
3 concated = np.array(placebo + medicine)
4
5 bootstrap_stats = []
6
7 for _ in tqdm(range(N_SAMPLES)):
8     # выбираем индексы для первой части
9     index_x = np.random.choice(np.arange(n + m), n, replace=False)
10    # выделяем оставшиеся
11    index_y = np.arange(n + m)[~np.in1d(np.arange(n + m), index_x)]
12
13    x_iter = concated[index_x]
14    y_iter = concated[index_y]
15
16    # считаем статистику Манна-Уитни
17    bootstrap_stats.append(sps.mannwhitneyu(x_iter, y_iter)[0])
```

In [12]:

```
1 (np.array(bootstrap_stats) < sps.mannwhitneyu(placebo, medicine)[0]).mean()
```

Out[12]:

0.04015

### Вывод:

В критерии перестановок с использованием статистики Манна-Уитни наблюдается отвержение на уровне значимости 0.05. Ясно, что такой критерий не первым приходит в голову, но всё-таки даже на такой небольшой выборке можно найти корректную процедуру, которая подтвердит статистическую значимость. Но тем не менее можно сказать, что это отвержение "на грани", поэтому, если есть возможность, то стоит ещё поэкспериментировать.

### АА-тест

Проверим, как подобная процедура будет работать, если нулевая гипотеза верна. Запустим её на двух выборках из одинакового равномерного распределения.

Посмотрим на распределение `rvalue`. Оно должно получиться равномерным. Для этого проведём АА-тест много раз

In [ ]:

```
1  def make_aa_test(sample_size, n_samples):
2      """
3      Функция проведения АА-теста
4      Применяет перестановочный критерий со статистикой Манна-Уитни
5      к выборкам из U[0, 20] одинакового размера
6
7      Параметры
8      -----
9      sample_size : int
10         Размер каждой выборки
11      n_samples : int
12         Количество бутстрепных перестановок
13      Возвращает
14      -----
15      float : pvalue
16      """
17
18      distr = sps.randint(low=1, high=20)
19      x = distr.rvs(sample_size)
20      y = distr.rvs(sample_size)
21      concated = np.hstack([x, y])
22
23      bootstrap_stats = []
24      for _ in range(n_samples):
25         index_x = np.random.choice(np.arange(n + m), n, replace=False)
26         index_y = np.arange(n + m)[~np.in1d(np.arange(n + m), index_x)]
27
28         x_iter = concated[index_x]
29         y_iter = concated[index_y]
30
31         bootstrap_stats.append(sps.mannwhitneyu(x_iter, y_iter)[0])
32
33     return (np.array(bootstrap_stats) < sps.mannwhitneyu(x, y)[0]).mean()
```

Проведём тест один раз

In [19]:

```
1  make_aa_test(sample_size=len(medicine), n_samples=N_SAMPLES)
```

Out[19]:

0.6458

Видно, что нет отвержения на одинаковых распределениях. Но с некоторой вероятностью могло произойти и отвержение. Посмотрим на распределение pvalue.

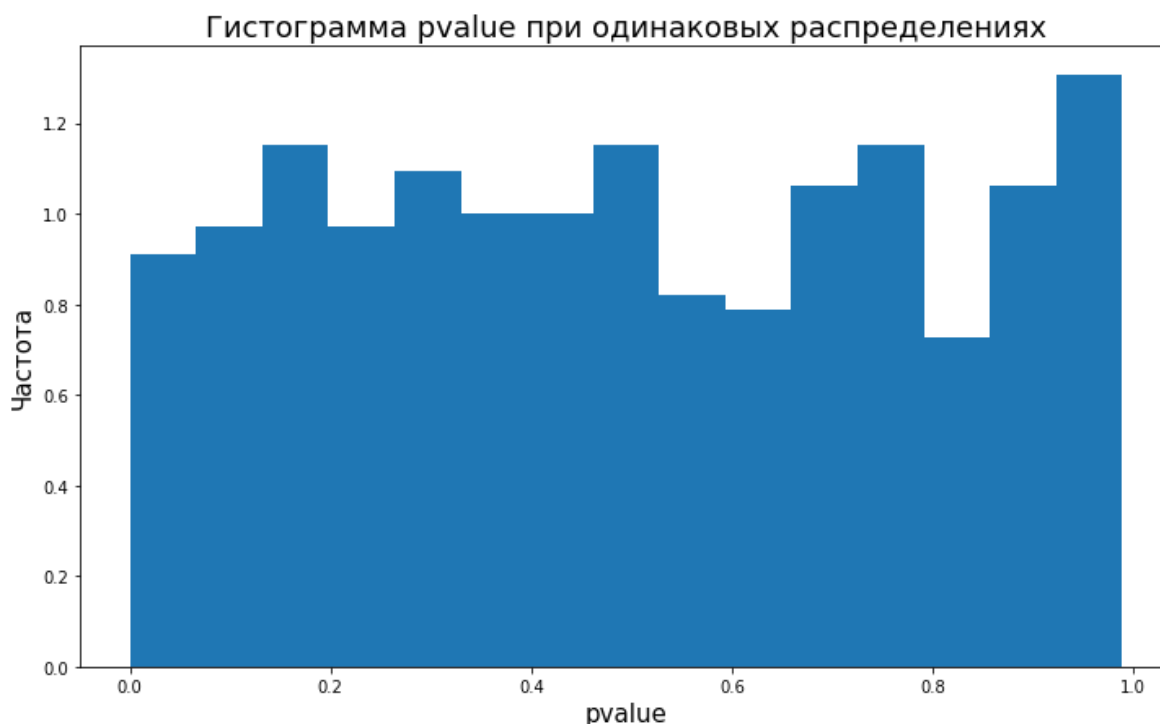
In [20]:

```
1  N_TESTS = 500 # количество АА-тестов
2
3  pvalues = []
4  for _ in tqdm(range(N_TESTS)):
5      pvalues.append(make_aa_test(sample_size=len(medicine), n_samples=10000))
```

Построим гистограмму

In [26]:

```
1 plt.figure(figsize=(12, 7))
2 plt.title('Гистограмма pvalue при одинаковых распределениях', fontsize=18)
3 plt.xlabel('pvalue', fontsize=15)
4 plt.ylabel('Частота', fontsize = 15)
5 plt.hist(pvalues, density=True, bins=15)
6 plt.show()
```



In [27]:

```
1 pvalues = np.array(pvalues)
2 n_rejects = np.sum(pvalues < 0.05)
3 print('Доля отвержений: {:.3f}'.format(n_rejects / pvalues.size))
```

Доля отвержений: 0.052

### Вывод:

По гистограмме видим, что распределение похоже на равномерное. Возможно, будет ещё больше похоже при увеличении количества запусков АА-теста. Доля отвержений близка к 0.05, но для точной оценки нужно гораздо больше запусков.

## Задача 2

В задании 6 по машинному обучению вы предсказывали цену жилья по его характеристикам, в процессе чего принимали участие на Kaggle. Пришло время внедрять разработки в продакшн. На первый взгляд может показаться, что внедрять нужно решение победителя, однако оно может повлечь множество технических сложностей при работе в продакшне, поэтому предлагается выбрать оптимум между

качеством и технической сложностью модели. В данной задаче вам предлагается исследовать модели на наличие статистически значимой разницы по их качеству. Не забывайте про практическую значимость результата.

Сравните три модели по качеству предсказания на тестовой выборке: свою и два решения, рассказанных на семинарах. В решении данного задания не нужно приводить код каждой из моделей. Достаточно прочитать 4 файла: истинные ответы на тесте и предсказания трех моделей.

---

Сравнивать модели будем по значениям функции ошибки, которые получаются у этих моделей на тестовой выборке. В данном соревновании стояла задача регрессии, поэтому можно оценивать модели по MAPE, то есть по средней относительной ошибке в процентах.

Таким образом, если есть правильные ответы  $Y_i$  и предсказания модели  $\widehat{Y}_i$ , то можно получить выборку ошибок:

$$e_i = \frac{|Y_i - \widehat{Y}_i|}{Y_i}$$

Тогда, если у двух моделей получились ошибки  $e^1$  и  $e^2$ , то нас интересует гипотеза о равенстве средних.

Модели были обучены на одной выборке, они получают один и тот же вход, значит эти выборки зависимы. Поэтому будем пользоваться критерием знаков.

Можем считать, что в данных есть шум, который не убирается моделью, и его распределение одинаково для двух моделей, потому что данные те же. А значит предположения критерия знаков выполнены.

Для начала загрузим данные: правильные ответы, а также ответы двух лучших работ и бейзлайна

In [2]:

```
1 ▼ golden_answers = pd.read_csv('solution.csv',
2                               index_col='index')
3 ▼ ivanov_answers = pd.read_csv('answer_Viacheslav_Ivanov.csv',
4                                index_col='index')
5 ▼ gracheva_answers = pd.read_csv('answer_Anastasia_Gracheva.csv',
6                                  index_col='index')
7 ▼ baseline_answers = pd.read_csv('sample_submission.csv',
8                                  index_col='index')
```

In [3]:

```
1 golden_answers.head()
```

Out[3]:

	price
index	
1	330000
2	960000
3	219900
4	430000
5	221700

Вычислим выборки ошибок. Поменяв функцию вычисления, можно перейти к другим ошибкам, например, MSE или MAE

In [4]:

```
1 def get_errors(answers_df):  
2     """Функция подсчёта ошибок"""  
3  
4     correct = golden_answers.values  
5     predicted = answers_df.values  
6     return (np.abs(correct - predicted) * 100. / correct) \  
7         .reshape(answers_df.values.size)
```

In [5]:

```
1 ivanov_errors = [] # Слава не ошибается
```

In [6]:

```
1 ivanov_errors = get_errors(ivanov_answers)  
2 gracheva_errors = get_errors(gracheva_answers)  
3 baseline_errors = get_errors(baseline_answers)
```

Теперь необходимо проверить 3 гипотезы о равенстве средних в этих ошибках. Для проверки напомним функцию, которая вычислит pvalue для критерия знаков, а также построит сравнительный график kde распределений ошибок



In [7]:

```
1 ▾ # Переиспользуем код из прошлого задания
2 SignedTestResult = namedtuple('SignedTestResult', ('statistic', 'pvalue'))
3
4 ▾ def _signedtest_n_greater_15(s, n, alternative):
5     """Внутренняя функция, считает критерий знаков для случая n > 15"""
6
7     assert n > 15
8
9     statistics = (s - n / 2 - 1 / 2) / np.sqrt(n / 4)
10
11 ▾     if alternative == 'two-sided':
12         pvalue = 2 * min(sps.norm.sf(statistics), sps.norm.cdf(statistics))
13 ▾     elif alternative == 'less':
14         pvalue = sps.norm.cdf(statistics)
15 ▾     else:
16         pvalue = sps.norm.sf(statistics)
17
18     return SignedTestResult(statistics, min(pvalue, 1))
19
20 ▾ def _signedtest_n_small_sizes(s, n, alternative):
21     """Внутренняя функция, считает критерий знаков для случая n <= 15"""
22
23     assert n <= 15
24
25     statistics = s
26     left_part = sps.binom.cdf(statistics, n, 1/2)
27     right_part = sps.binom.sf(statistics, n, 1/2)
28
29 ▾     if alternative == 'two-sided':
30         pvalue = 2 * min(left_part, right_part)
31 ▾     elif alternative == 'less':
32         pvalue = left_part
33 ▾     else:
34         pvalue = right_part
35
36     return SignedTestResult(statistics, min(pvalue, 1))
37
38 ▾ def signedtest(differences, alternative='two-sided'):
39     """
40     Критерий знаков
41     Параметры
42     -----
43 ▾     differences : array_like
44         Массив разностей между двумя связными выборками
45 ▾     alternative : {'two-sided', 'less', 'greater'}, optional
46         Определение альтернативной гипотезы.
47         'two-sided' - альтернативная гипотеза: медиана ошибок не равна 0
48         'less' - альтернативная гипотеза: медиана ошибок меньше 0
49         'greater' - альтернативная гипотеза: медиана ошибок больше 0
50     Возвращает
51     -----
52     statistic : float
53     pvalue : float
54     """
55
56     assert alternative in ['two-sided', 'less', 'greater']
57
58     z = differences
59     z = np.unique(np.array(z)) # выкидываем совпадения
```

```

60     s = np.sum(z > 0)
61     n = z.shape[0]
62
63     if n > 15:
64         return _signedtest_n_greater_15(s, n, alternative)
65     return _signedtest_n_small_sizes(s, n, alternative)

```

In [8]:

```

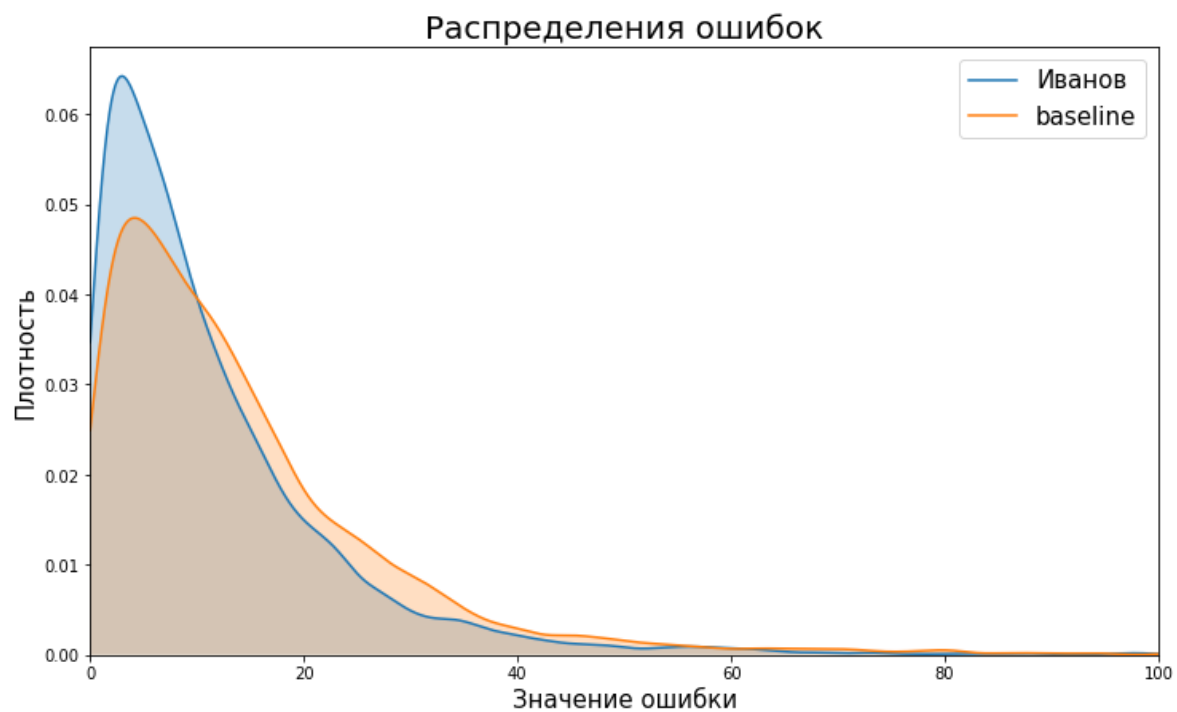
1  def make_models_comparison(x_errs, y_errs, label_x='First',
2                             label_y='Second', max_err=100):
3
4     """
5     Функция для сравнения моделей регрессии в соревновании
6     Печатает pvalue, абсолютную и относительную разницу в метриках
7     Строит график распределений ошибок
8
9     Параметры
10    -----
11    x_errs : array_like, 1d
12            Ошибки первой модели
13    y_errs : array_like, 1d
14            Ошибки второй модели
15    label_x : str
16            Название первой модели
17    label_y : str
18            Название второй модели
19    max_err : float
20            Наибольшая ошибка, правый предел графика
21    Возвращает
22    -----
23    float : pvalue
24
25    plt.figure(figsize=(12, 7))
26    plt.title('Распределения ошибок', fontsize=20)
27    plt.xlabel('Значение ошибки', fontsize=15)
28    plt.ylabel('Плотность', fontsize=15)
29    sns.kdeplot(x_errs, shade=True, label=label_x, gridsize=5000)
30    sns.kdeplot(y_errs, shade=True, label=label_y, gridsize=5000)
31    plt.xlim((0, max_err))
32    plt.legend(fontsize=15)
33    plt.show()
34
35    diff = x_errs.mean() - y_errs.mean()
36    diff_percent = diff / x_errs.mean()
37
38    _, pvalue = signedtest(x_errs - y_errs)
39    print('p-value критерия знаков: {:.5f}'.format(pvalue))
40    print('Разница средних: {:.3f}'.format(diff))
41    print('Относительная разница: {:.1%}'.format(diff_percent))
42    return pvalue

```

Теперь применим критерий 3 раза, потом применим процедуру МПГ при помощи метода Холма

In [9]:

```
1  ivanov_base = make_models_comparison(  
2      ivanov_errors,  
3      baseline_errors,  
4      label_x='Иванов',  
5      label_y='baseline'  
6  )
```



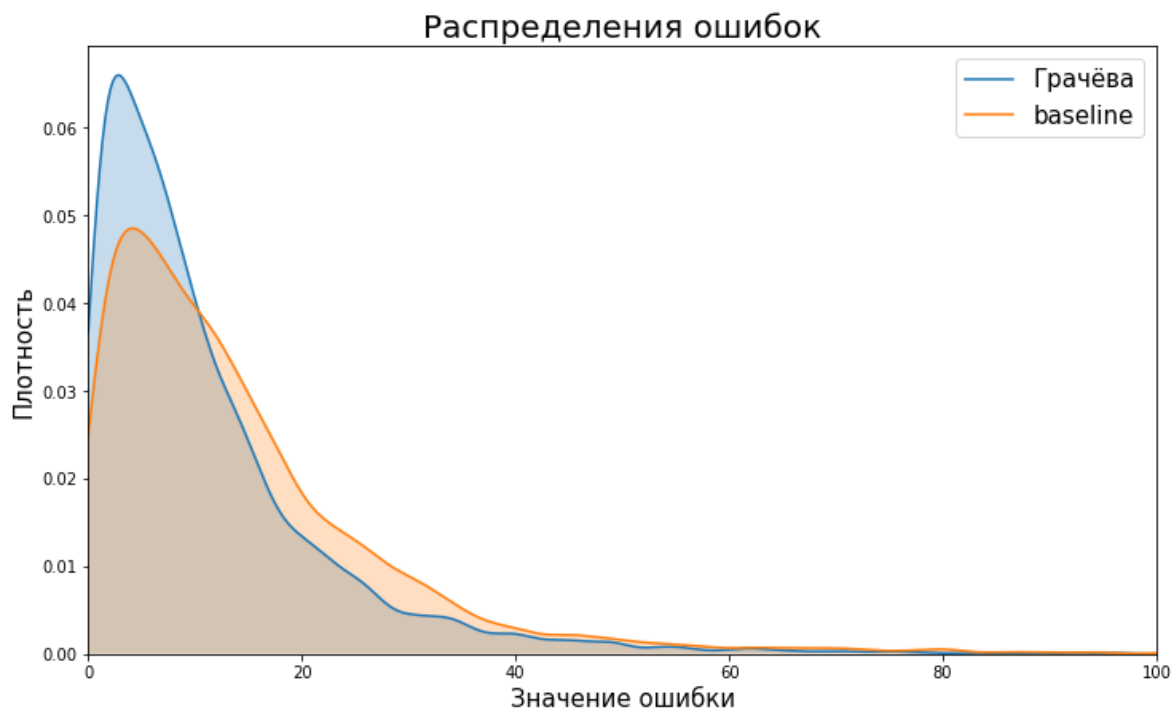
p-value критерия знаков: 0.00000

Разница средних: -3.019

Относительная разница: -26.0%

In [10]:

```
1 ▾ gracheva_base = make_models_comparison(  
2     gracheva_errors,  
3     baseline_errors,  
4     label_x='Грачёва',  
5     label_y='baseline'  
6 )
```



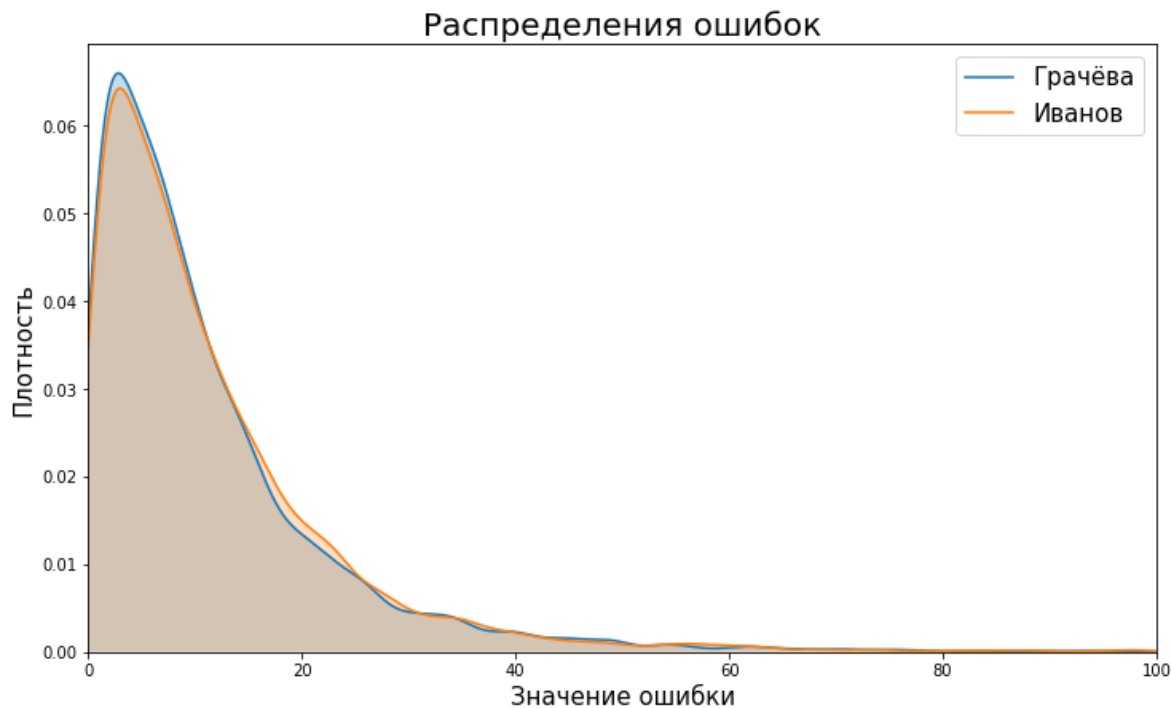
p-value критерия знаков: 0.00000

Разница средних: -3.353

Относительная разница: -29.7%

In [11]:

```
1 ▾ gracheva_ivanov = make_models_comparison(  
2     gracheva_errors,  
3     ivanov_errors,  
4     label_x='Грачёва',  
5     label_y='Иванов'  
6 )
```



p-value критерия знаков: 0.00021  
Разница средних: -0.334  
Относительная разница: -3.0%

### Вывод:

По графикам видно, что модели участников оказались сильно лучше бейзлайна. Это подтверждают и критерии, pvalue крайне мал, что говорит о высокой уверенности.

При этом оказалось, что и результат двух участников различается значимо. Лучше результат у Насти Грачёвой. Однако здесь pvalue уже не такой мизерный.

Но тем не менее и этого достаточно, чтобы отвергнуть все гипотезы даже после коррекции методом Бонферрони. Оказывается, что критерий знаков позволяет отвергнуть незначимость разницы, когда в метрике MAPE абсолютная разница составила 0.3, что было на 3 процента лучше, чем во второй по рейтингу работе.

## Задача 3

Сессией в интернете называется промежуток времени, охватывающий работу пользователя с момента открытия первой страницы и до закрытия последней. В каждой сессии пользователь может кликнуть на целевой объект. Предположим, базовая конверсия сессии в клик равна 0.08 (т.е. вероятность клика за сессию). С целью проведения АВ-тестирования для 5 процентов аудитории выкатывается новый дизайн, от которого ожидается улучшение конверсии на 0.01, т.е. вероятность клика станет равна 0.09. Оцените двумя способами срок такого АВ-теста, считая, что все сессии независимы, а на сайте происходит 1000 сессий в день.

In [2]:

```
1 control_conversion = 0.08
2 test_conversion = 0.09
3
4 control_fraction = 0.95
5 test_fraction = 0.05
6
7 events_per_day = 1000
```

Для того, чтобы определить в данном случае срок АВ-теста, воспользуемся двумя предложенными на лекциях методами. А именно посчитаем срок:

- Из соображений мощности
- Из соображений эффекта

Пусть оценки вероятности конверсий оказались в контрольной и тестовой выборке соответственно  $p_a$  и  $p_b$ , количество сессий аналогично  $n_a$  и  $n_b$

Искомое общее количество сессий обозначим за  $n = n_a + n_b$

Долю пользователей, на которых работает тестовая версия, обозначим за  $f$

### Из соображений эффекта

Мы хотим, чтобы произошло отвержение при заданных конверсиях. Для этого надо, чтобы величина, которая, согласно тесту, распределена нормально, была больше 2.

$$\frac{p_b - p_a}{\hat{\sigma}} \geq 2$$

$$\hat{\sigma} = \sqrt{\frac{p_b(1-p_b)}{n_b} + \frac{p_a(1-p_a)}{n_a}} = \sqrt{\frac{p_b(1-p_b)}{fn} + \frac{p_a(1-p_a)}{(1-f)n}} = \frac{1}{\sqrt{n}} \sqrt{\frac{p_b(1-p_b)}{f} + \frac{p_a(1-p_a)}{1-f}}$$

Подставляя выражение для  $\hat{\sigma}$ , получаем следующее неравенство:

$$\sqrt{n} \frac{p_b - p_a}{\sqrt{\frac{p_b(1-p_b)}{f} + \frac{p_a(1-p_a)}{1-f}}} \geq 2$$

Возведя в квадрат, можем получить неравенство на  $n$

$$n \geq 4 \frac{\frac{p_b(1-p_b)}{f} + \frac{p_a(1-p_a)}{1-f}}{(p_b - p_a)^2}$$

Остаётся только не забыть про то, что число дней должно быть в 1000 раз меньше числа сессий, потому что в день на сайте наблюдается 1000 сессий

In [24]:

```
1  def get_effect_ab_length(control_conversion, test_conversion, test_fraction):
2      """
3      Функция для определения срока АБ-теста из соображений эффекта
4      Параметры
5      -----
6      control_conversion : float
7          Текущая конверсия
8      test_conversion : float
9          Конверсия, на которой должно быть отвержение
10     test_fraction : float
11         Доля сессий в тестовой выборке
12     Возвращает
13     -----
14     int : количество дней проведения теста
15     """
16
17     control_fraction = 1 - test_fraction
18
19     test_var = test_conversion * (1 - test_conversion) / test_fraction
20     control_var = control_conversion * (1 - control_conversion) / control_fraction
21
22     n_sessions = 4 * (test_var + control_var) \
23                 / (control_conversion - test_conversion) ** 2
24     n_days = int(np.ceil(n_sessions / events_per_day))
25
26     return n_days
```

In [26]:

```
1  get_effect_ab_length(control_conversion, test_conversion, test_fraction)
```

Out[26]:

69

Получаем, что необходимо проводить АБ-тест в течение 69 дней

### Из соображений мощности

Здесь будем действовать проще. При помощи бинарного поиска будем фиксировать количество сессий, семплировать бернуллиевские выборки из контроля и теста, затем на каждой считать критерий Вальда, чтобы оценить мощность.

Мощность нужно считать с точностью хотя бы до 2 знака, поэтому нужно 10000 семплирований

In [29]:

```
1 ▼ def wald_test_greater(x, y):
2     """
3     Критерий Вальда для бернуллиевских выборок
4     Альтернатива: mean(y) > mean(x)
5     Параметры
6     -----
7 ▼    x, y : array-like
8         Выборки из распределения Бернулли
9     Возвращает
10    -----
11    float : pvalue критерия
12    """
13
14    p_x, p_y = x.mean(), y.mean()
15
16    var_x = p_x * (1 - p_x)
17    var_y = p_y * (1 - p_y)
18
19    sigma = np.sqrt(var_x / x.size + var_y / y.size)
20    stat = (p_y - p_x) / sigma
21
22    return sps.norm.sf(stat)
```



In [30]:

```
1  ▾ def get_power_ab_length(  
2      control_conversion, test_conversion, test_fraction,  
3      max_sessions=1000000, n_samples=10000, demanded_power=0.8  
4  ▾ ):  
5      """  
6      Функция для определения срока АБ-теста из соображений мощности  
7      с помощью бинарного поиска  
8      Параметры  
9      -----  
10     ▾ control_conversion : float  
11         Текущая конверсия  
12     ▾ test_conversion : float  
13         Конверсия, на которой должно быть отвержение  
14     ▾ test_fraction : float  
15         Доля сессий в тестовой выборке  
16     ▾ max_sessions : int  
17         Максимально возможное число сессий в эксперименте  
18     ▾ n_samples : int  
19         Количество семплирования для определения мощности  
20     ▾ demanded_power : float  
21         Требуемая мощность  
22     Возвращает  
23     -----  
24     int : количество дней проведения теста  
25     """  
26  
27     left, right = 1, max_sessions  
28     ▾ while (right - left) > events_per_day:  
29         # делим пополам, пока различия больше одного дня  
30         middle = (left + right) / 2  
31  
32         n_rejects = 0  
33     ▾ for _ in range(n_samples):  
34         n_control = int(middle * control_fraction)  
35         n_test = int(middle * test_fraction)  
36         # семплируем из распределения бернулли нужное количество примеров  
37         control = sps.bernoulli(p=control_conversion).rvs(n_control)  
38         test = sps.bernoulli(p=test_conversion).rvs(n_test)  
39         pvalue = wald_test_greater(control, test) # применяем критерий Ва.  
40         n_rejects += pvalue < 0.05  
41         power_est = n_rejects / n_samples  
42  
43     ▾ print('Число сессий: {:.0f}, оценка мощности: {:.2f}'\  
44         .format(middle, power_est))  
45  
46     ▾ if power_est > demanded_power:  
47         right = middle  
48     ▾ else:  
49         left = middle  
50  
51     return int(np.ceil(middle / events_per_day))
```

In [33]:

```
1 ▾ get_power_ab_length(control_conversion, test_conversion,  
2                               test_fraction, n_samples=10000)
```

Число сессий: 500000, оценка мощности: 1.00  
Число сессий: 250001, оценка мощности: 0.99  
Число сессий: 125001, оценка мощности: 0.86  
Число сессий: 62501, оценка мощности: 0.61  
Число сессий: 93751, оценка мощности: 0.77  
Число сессий: 109376, оценка мощности: 0.82  
Число сессий: 101563, оценка мощности: 0.79  
Число сессий: 105470, оценка мощности: 0.81  
Число сессий: 103517, оценка мощности: 0.80  
Число сессий: 102540, оценка мощности: 0.79

Out[33]:

103

Получаем, что данный метод даёт большее количество дней, а именно 103

---

## Задача 4

Имеются две задачи:

1. Пусть  $X_i = (X_{i1}, \dots, X_{id}), i = 1, \dots, n$  -- выборка из  $d$ -мерных объектов. Для всех пар признаков требуется проверить гипотезу о независимости этих признаков, т.е.  $H_{jk}: X_{ij}, X_{ik}$  независимы. Для проверки используется критерий на основе коэффициента корреляции Спирмена.
2. Пусть решается задача регрессии и  $X_i = (X_{i1}, \dots, X_{id}), i = 1, \dots, n$  -- выборка из  $d$ -мерных регрессоров, а  $Y_1, \dots, Y_n$  -- соответствующие значения отклика. Требуется отобрать значимые признаки, для чего проверяются гипотезы  $H_j: X_{ij}, Y_i$  независимы. Для проверки используется критерий на основе коэффициента корреляции Спирмена.

В обеих задачах предполагается использовать перестановочный критерий, в котором статистикой является коэффициент корреляции Спирмена.

Задание:

- Предложите группу перестановок для реализации данного критерия, а так же метод генерации бутстрепных выборок.
- Можно ли в данных задачах использовать метод множественной проверки гипотезы на основе перестановок (maxT-статистика)? Если нет, то нужно привести пример, для которого нарушается какое-либо требование. Если да, то нужно привести пару примеров, для которых свойство выполняется. Примеры можно привести с помощью семплирования.

---

### Группа перестановок

Нужно предложить такую группу перестановок, которая не меняет распределение коэффициента корреляции Спирмена, поскольку именно на нём основан критерий.

Если гипотеза верна, то есть выборки независимы, при перемешивании элементы внутри каждой выборки, распределение коэффициента корреляции Спирмена не меняется. Аналогичная ситуация для мультигипотез, если все выборки независимы.

Таким образом в качестве группы перестановок для критерия мы берём совокупности перестановок элементов каждой из выборок.

## **Применимость $\max T$**

Рассмотрим на примерах применимость МПГ с помощью перестановок. Чтобы подтвердить её или опровергнуть, проверим распределения корреляции Спирмена.

В качестве модели, из которой генерируются данные, возьмём просто многомерное нормальное распределение. Будем задавать матрицу ковариаций и вектор средних и смотреть на распределения коэффициента корреляции.

Нужно проверить, например, отличается ли совместное распределение коэффициентов корреляции для некоторых гипотез в зависимости от верности некоторого набора других гипотез.

Напишем функцию, которая генерирует многомерные нормальные данные для матрицы ковариаций в случае верных и неверных гипотез, а потом возвращает набор коэффициентов корреляции для заданных индексов (регрессоров или отклика)

In [35]:

```
1  ▼ def compare_correlation_distributions(cov_true, cov_false, inds1, inds2,
2                                     sample_size=100, n_samples=20000):
3      """
4      Функция для сравнения коэффициентов корреляции
5      в проверке гипотез о независимости выборок
6      Параметры
7      -----
8  ▼   cov_true : 2d array
9       Матрица ковариаций распределения при справедливости гипотез
10  ▼   cov_false : 2d array
11       Матрица ковариаций распределения при несправедливости гипотез
12  ▼   inds1 : Tuple(int, int)
13       Пара индексов выборок первой гипотезы из набора
14  ▼   inds2 : Tuple(int, int)
15       Пара индексов выборок второй гипотезы из набора
16  ▼   sample_size : int
17       Размер многомерной выборки
18  ▼   n_samples : int
19       Количество семплирования = размер выборки коэффициентов
20      """
21
22     cov_true = np.array(cov_true)
23     cov_false = np.array(cov_false)
24     mean = np.zeros(cov_true.shape[0])
25
26     T_true = []
27     T_false = []
28
29  ▼   for _ in tqdm(range(n_samples)):
30       # генерируем многомерную выборку
31  ▼       sample = sps.multivariate_normal(mean=mean,
32                                       cov=cov_true).rvs(size=sample_size)
33
34       # считаем коэффициенты корреляции для нужных регрессоров/отклика
35  ▼       T_true.append([
36           sps.spearmanr(sample[:, inds1[0]], sample[:, inds1[1]])[0],
37           sps.spearmanr(sample[:, inds2[0]], sample[:, inds2[1]])[0]
38       ])
39
40  ▼       sample = sps.multivariate_normal(mean=mean,
41                                       cov=cov_false).rvs(size=sample_size)
42
43  ▼       T_false.append([
44           sps.spearmanr(sample[:, inds1[0]], sample[:, inds1[1]])[0],
45           sps.spearmanr(sample[:, inds2[0]], sample[:, inds2[1]])[0]
46       ])
47
48     return T_true, T_false
```

Также требуется функция, которая изобразит распределения полученных статистик

In [37]:

```
1 ▾ def show_distributions(T_true, T_false, x_label='Корреляция 1', y_label='Корр'
2     """
3     Функция для построения совместных распределений коэффициентов корреляции
4     Параметры
5     -----
6 ▾   T_true : list
7       Список попарных значений коэффициентов при справедливости набора гипотез
8 ▾   T_false : list
9       Список попарных значений коэффициентов при несправедливости набора гипотез
10 ▾  x_label, y_label : str
11     Метки для осей
12     """
13
14     T_true = np.array(T_true)
15     T_false = np.array(T_false)
16     plt.figure(figsize=(10, 10))
17     plt.title('Сравнение совместных распределений', fontsize=18)
18
19 ▾     sns.kdeplot(T_true[:, 0], T_true[:, 1], gridsize=100, shade=True,
20                 shade_lowest=False, alpha=0.5)
21 ▾     sns.kdeplot(T_false[:, 0], T_false[:, 1], gridsize=100,
22                 color='red', shade=True,
23                 shade_lowest=False, alpha=0.5)
24
25     plt.xlabel(x_label, fontsize=15)
26     plt.ylabel(y_label, fontsize=15)
27     plt.xlim((-0.3, 0.3))
28     plt.ylim((-0.3, 0.3))
29     plt.show()
```

## Пункт 1

Рассматриваем 3 нормальные выборки. Пусть  $T_{jk}$  - статистика для проверки  $H_{jk}$

Покажем, что совместное распределение статистик  $T_{12}$  и  $T_{13}$  при справедливости  $H_{12}$  и  $H_{13}$  зависит от справедливости  $H_{23}$

Соответственно, рассматриваем матрицы ковариаций, в одной из которых будет зависимость между 2 и 3 признаком, а в другой - не будет

In [47]:

```
1 ▾ cov_true = [[1, 0, 0],
2              [0, 1, 0],
3              [0, 0, 1]]
4
5 ▾ cov_false = [[1, 0, 0],
6               [0, 1, 0.5],
7               [0, 0.5, 1]]
```

In [48]:

```
1 ▼ T_true, T_false = compare_correlation_distributions(  
2     cov_true, cov_false, (0, 1), (0, 2)  
3 )
```

In [49]:

```
1 ▼ show_distributions(T_true, T_false, x_label="$Corr(X_1, X_2)$",  
2     y_label="$Corr(X_1, X_3)$")
```



Видим, что распределения отличаются, поэтому сразу видим невозможность использования  $\max T$

## Пункт 2

Рассматриваем три нормальных выборки регрессоров и одну откликов. Пусть  $T_j$  - статистика для проверки  $H_j$ . Нужно, например, показать, что совместное распределение статистик  $T_1$  и  $T_2$  при справедливости  $H_1$  и  $H_2$  не зависит от справедливости  $H_3$ .

В качестве первого примера рассмотрим случай, когда все признаки независимы. А зависимость у третьего признака с откликом либо есть, либо нет.

In [41]:

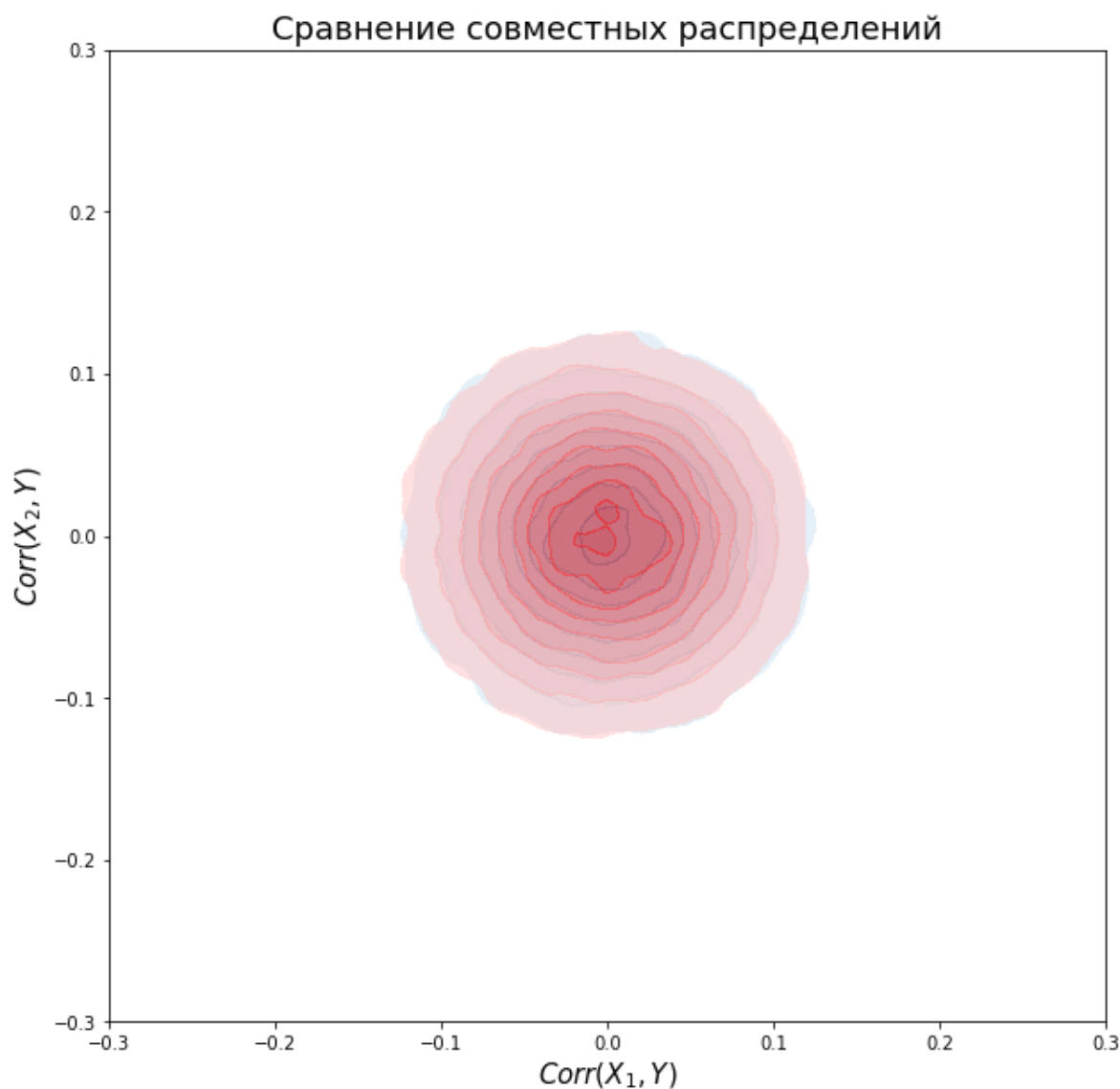
```
1 ▼ cov_true = [[1, 0, 0, 0],
2               [0, 1, 0, 0],
3               [0, 0, 1, 0],
4               [0, 0, 0, 1]]
5
6 ▼ cov_false = [[1, 0, 0, 0],
7                [0, 1, 0, 0],
8                [0, 0, 1, 0.9],
9                [0, 0, 0.9, 1]]
```

In [42]:

```
1 ▼ T_true, T_false = compare_correlation_distributions(
2     cov_true, cov_false, (0, 3), (1, 3),
3     n_samples=100000, sample_size=300
4 )
```

In [44]:

```
1 ▼ show_distributions(T_true, T_false, x_label="$Corr(X_1, Y)$",
2                     y_label="$Corr(X_2, Y)$")
```



По графику можно видеть, что распределения совпадают. Добавим также второй пример, в котором между признаками есть зависимость

In [50]:

```
1 ▼ cov_true = [[1, 0.3, 0.4, 0],
2               [0.3, 1, 0.1, 0],
3               [0.4, 0.1, 1, 0],
4               [0, 0, 0, 1]]
5
6 ▼ cov_false = [[1, 0.3, 0.4, 0],
7                [0.3, 1, 0.1, 0],
8                [0.4, 0.1, 1, 0.9],
9                [0, 0, 0.9, 1]]
```

In [51]:

```
1 ▼ T_true, T_false = compare_correlation_distributions(
2     cov_true, cov_false, (0, 3), (1, 3),
3     n_samples=100000, sample_size=300
4 )
```



In [52]:

```
1 show_distributions(T_true, T_false, x_label="$Corr(X_1, Y)$",  
2                   y_label="$Corr(X_2, Y)$")
```



Снова видим совпадение распределений

### Вывод:

Таким образом, в случае с проверкой независимости признаков точно нельзя использовать МПГ на основе перестановок. В случае со значимостью признаков и проверкой независимости с откликом мы получаем, что на рассмотренных примерах совместные распределения совпадают, а значит наверняка применить такую схему возможно.