

# Статистика, прикладной поток

## Практическое задание 6

В данном задании вы исследуете некоторые свойства непараметрических методов, а также проверки статистических гипотез.

### Правила:

- Дедлайн **2 декабря 23:59**. После дедлайна работы не принимаются кроме случаев наличия уважительной причины.
- Выполненную работу нужно отправить на почту `mipt.stats@yandex.ru`, указав тему письма "[applied] Фамилия Имя - задание 6". Квадратные скобки обязательны. Если письмо дошло, придет ответ от автоответчика.
- Прислать нужно ноутбук и его pdf-версию (без архивов). Названия файлов должны быть такими: `6.N.ipynb` и `6.N.pdf`, где `N` - ваш номер из таблицы с оценками.
- Решения, размещенные на каких-либо интернет-ресурсах не принимаются. Кроме того, публикация решения в открытом доступе может быть приравнена к предоставлению возможности списать.
- Для выполнения задания используйте этот ноутбук в качестве основы, ничего не удаляя из него.
- Никакой код из данного задания при проверке запускаться не будет.

### Баллы за задание:

- Задача 1 - 2 балла **О3**
- Задача 2 - 10 баллов **О2**
- Задача 3 - 5 баллов **О3**
- Задача 4 - 5 баллов **О3**
- Задача 5 - 5 баллов **О3**
- Задача 6 - 12 баллов **О3**
- Задача 7 - 20 баллов **О3**
- Задача 8 - 15 баллов **О3**

In [2]:

```
1 import numpy as np
2 import scipy.stats as sps
3 import pandas as pd
4
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 import warnings
8 from statsmodels.nonparametric.kde import kernel_switch, KDEUnivariate
9 from statsmodels.distributions.empirical_distribution import ECDF
10
11 from datetime import datetime
12
13 %matplotlib inline
```

## Непараметрический подход

## Задача 1.

Напишите определение ядра, используемого для построения ядерных оценок плотности.

- Ядро  $q$ , используемое для построения ядерных оценок плотности - это некоторая симметричная плотность.

Прежде чем начать работу с ядерными оценками плотности, изучим виды ядер. В библиотеке `statsmodels` реализованы следующие ядра:

In [2]:

```
1 list(kernel_switch.keys())
```

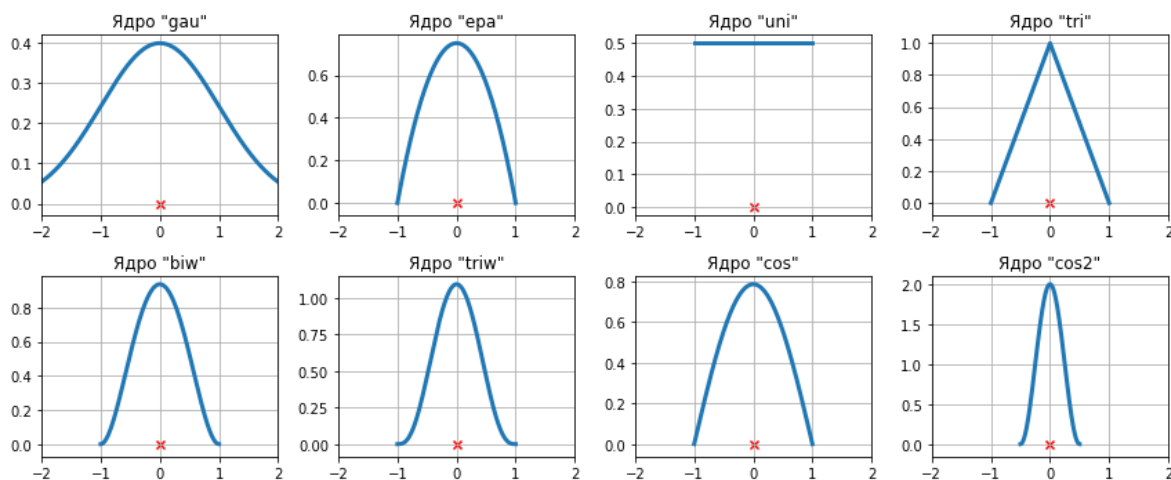
Out[2]:

```
['gau', 'epa', 'uni', 'tri', 'biw', 'triw', 'cos', 'cos2']
```

Нарисуем эти ядра. Запустите код в ячейке.

In [3]:

```
1 fig = plt.figure(figsize=(12, 5))
2
3 for i, (ker_name, ker_class) in enumerate(kernel_switch.items()):
4
5     kernel = ker_class() # ядро
6     domain = kernel.domain or [-2, 2] # носитель
7     x_vals = np.linspace(*domain, 2**10)
8     y_vals = kernel(x_vals)
9
10    ax = fig.add_subplot(2, 4, i + 1)
11    ax.set_title('Ядро "{}"'.format(ker_name))
12    ax.plot(x_vals, y_vals, lw=3, label='{}'.format(ker_name))
13    ax.scatter([0], [0], marker='x', color='red')
14    plt.grid(True, zorder=-5)
15    ax.set_xlim((-2, 2))
16
17 plt.tight_layout()
```



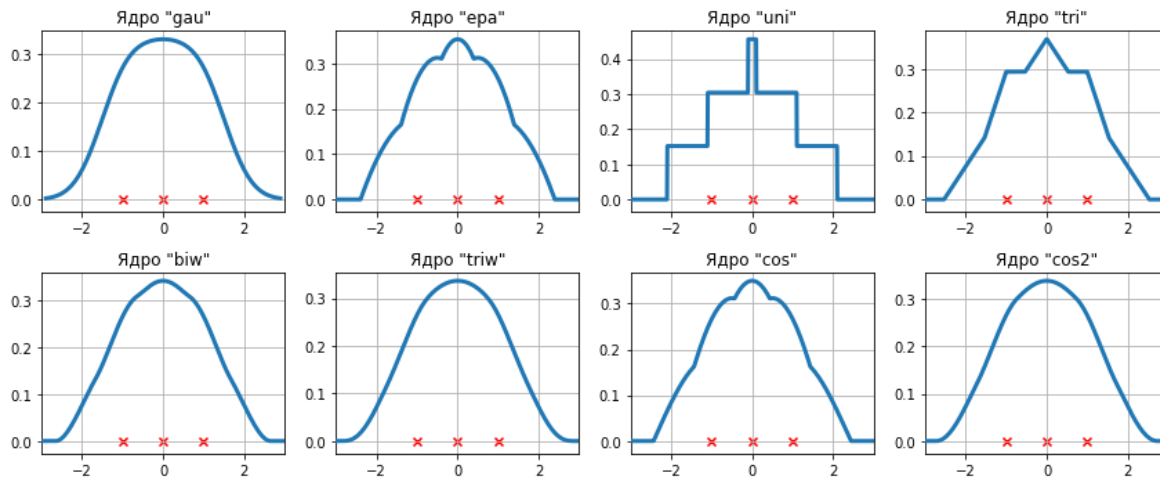
Посмотрим, как будет выглядеть ядерная оценка плотности в зависимости от типа ядра, взяв в качестве реализации выборки набор точек  $[-1, 0, 1]$ . Запустите код в ячейке.

In [4]:

```

1 data = np.linspace(-1, 1, 3) # выборка
2 kde = KDEUnivariate(data) # объект, выполняющий построение оценки
3
4 fig = plt.figure(figsize=(12, 5))
5
6 for i, kernel in enumerate(kernel_switch.keys()):
7
8     ax = fig.add_subplot(2, 4, i + 1)
9     ax.set_title('Ядро "{}"'.format(kernel))
10
11     # построение ядерной оценки плотности с заданным ядром
12     kde.fit(kernel=kernel, fft=False, gridsize=2**10)
13     # отрисовка полученной оценки
14     ax.plot(kde.support, kde.density, lw=3, zorder=10)
15     # отрисовка выборки
16     ax.scatter(data, np.zeros_like(data), marker='x', color='red')
17     plt.grid(True, zorder=-5)
18     ax.set_xlim([-3, 3])
19
20

```



Какие вы можете назвать преимущества и недостатки для первых трех ядер как с точки зрения самой оценки плотности, так и с вычислительной точки зрения?

Мы видим, что первое ядро - то есть гауссовское - хорошо приближает наши три точки - отличается на графике гладкостью, видно что кривая не претерпевает резкие скачки при переходе от одной точки выборки до другой - действительно по виду похоже на некоторую плотность распределения. А, например, ядро Епаненчикова уже гораздо существенней реагирует на отдельные точки выборки. У прямоугольного ядра kde вообще разрывная, хорошо видно, что претерпевает резкие скачки при переходе от одного элемента выборки до другого.

Что касается вычисления этих функций, то понятно, что вычисление прямоугольной kde - легче вычислить чем остальные два. Затем по степени сложности вычисления стоит ядро Епаненчикова (квадратичная функция), а сложнее всего вычислять Гауссовское, так там экспонента.

Другие примеры работы с ядерными оценками плотности можно посмотреть по ссылке

[http://www.statsmodels.org/dev/examples/notebooks/generated/kernel\\_density.html](http://www.statsmodels.org/dev/examples/notebooks/generated/kernel_density.html)

([http://www.statsmodels.org/dev/examples/notebooks/generated/kernel\\_density.html](http://www.statsmodels.org/dev/examples/notebooks/generated/kernel_density.html))

## Задача 2.

1. Сгенерируйте выборку  $X_1, \dots, X_{10000}$  из стандартного нормального распределения. Для каждого  $n \leq 10000$  постройте эмпирическую функцию распределения  $\widehat{F}_n$  и посчитайте **точное** значение статистики

$$D_n = \sup_{x \in \mathbb{R}} \left| \widehat{F}_n(x) - F(x) \right|.$$

Постройте график зависимости статистики  $D_n$  от  $n$ . Верно ли, что  $D_n \rightarrow 0$  и в каком смысле?

Для выполнения задания можно использовать следующую функцию:

In [5]:

```
1 from statsmodels.distributions.empirical_distribution import ECDF
2 help(ECDF) # В случае затруднений раскомментировать и выполнить
```

Help on class ECDF in module statsmodels.distributions.empirical\_distribution:

```
class ECDF(StepFunction)
|   Return the Empirical CDF of an array as a step function.
|
|   Parameters
|   -----
|   x : array-like
|       Observations
|   side : {'left', 'right'}, optional
|       Default is 'right'. Defines the shape of the intervals constit
|   using the
|       steps. 'right' correspond to [a, b) intervals and 'left' to
|   (a, b].
|
|   Returns
|   -----
|   Empirical CDF as a step function.
|
|   Examples
|   -----
|   >>> import numpy as np
|   >>> from statsmodels.distributions.empirical_distribution import E
CDF
|   >>>
|   >>> ecdf = ECDF([3, 3, 1, 4])
|   >>>
|   >>> ecdf([3, 55, 0.5, 1.5])
|   array([ 0.75,  1.   ,  0.   ,  0.25])
|
|   Method resolution order:
|       ECDF
|       StepFunction
|       builtins.object
|
|   Methods defined here:
|
|   __init__(self, x, side='right')
|       Initialize self. See help(type(self)) for accurate signature.
|
|   -----
|
|   Methods inherited from StepFunction:
|
|   __call__(self, time)
|       Call self as a function.
|
|   -----
|
|   Data descriptors inherited from StepFunction:
|
|   __dict__
|       dictionary for instance variables (if defined)
|
|   __weakref__
```

| list of weak references to the object (if defined)

In [2]:

```
1 # получаю выборку из нормального
2 sample = sps.norm().rvs(10000)
3 cdfs = sps.norm().cdf(sample)
4 ones = np.ones(len(sample))
5
6 # значения cdf на выборке sample, продублированные
7 # на квадратичную матрицу
8 cdfs = np.einsum("i, j->ij", ones, cdfs)
```

In [3]:

```
1 cdfs = np.tril(cdfs) + \
2 np.triu(np.ones((len(sample), len(sample))), k=1)
```

In [4]:

```
1 cdfs = np.sort(cdfs)
```

In [5]:

```
1 grid = np.arange(1, len(sample) + 1)
2
3 # в нижнем треугольнике значения эмпирической cdf -
4 # в точках sample[:i], 1<=i<=n
5 # i-ая строка - значения на i-ом префиксе
6 ecdfs = np.einsum("i, j->ij", 1/grid, grid)
```

In [6]:

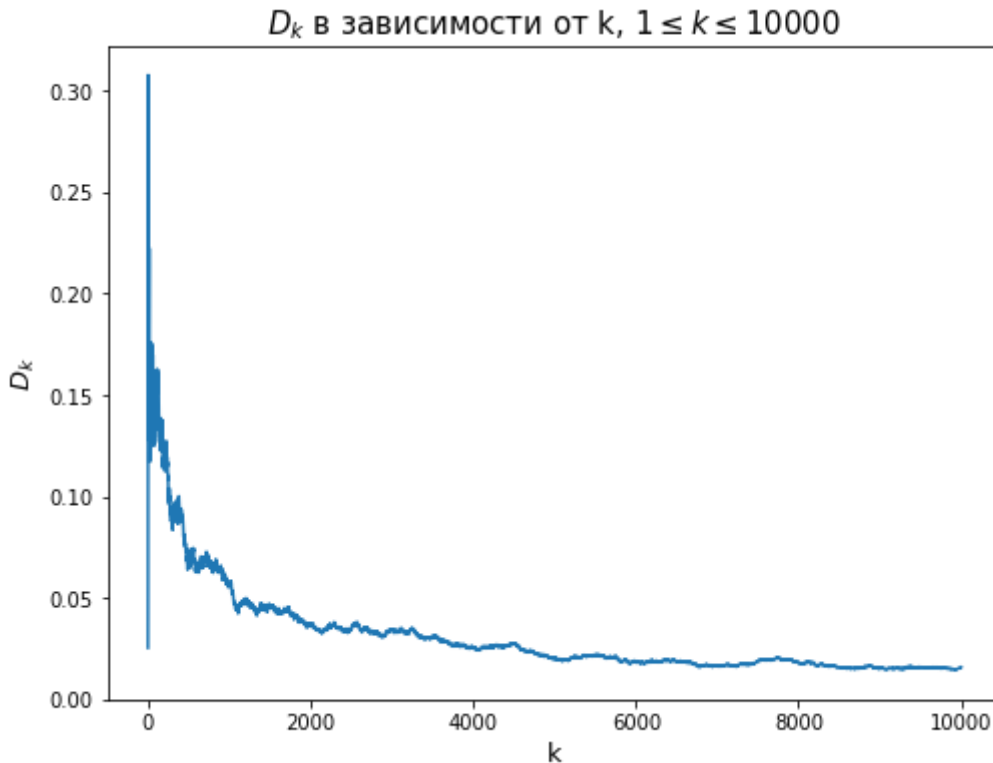
```
1 # сама D_k
2 Ds = np.max(np.abs(np.tril(cdfs - ecdfs)), axis=1)
```

In [7]:

```

1 plt.figure(figsize=(8,6))
2
3 plt.title(r"$D_k$ в зависимости от $k$, $1 \leq k \leq 10000$",
4           fontsize=15)
5
6 plt.plot(grid, Ds)
7
8 plt.xlabel("$k$", fontsize=13)
9 plt.ylabel(r"$D_k$", fontsize=13)
10
11 plt.show()

```

**Вывод:**

Действительно мы можем видеть, что последовательность  $D_n \rightarrow 0$  при  $n \rightarrow +\infty$ , причем сходимость к нулю довольно равномерная. Это подтверждает теорему Гливенко-Кантелли о сходимости  $D_n$  к нулю. То есть мы можем сказать, что эмперическая функция распределения действительно хорошо описывает приближает истинную.

2. Для  $n = 10000$  и  $k = 10000$  (значение  $k$  можно взять больше) выполните следующее:

Вычислите  $D_n^1, \dots, D_n^k$  для независимых выборок  $(X_1^1, \dots, X_n^1), \dots, (X_1^k, \dots, X_n^k)$  из стандартного нормального распределения. Постройте график гистограммы значений  $\sqrt{n}D_n^1, \dots, \sqrt{n}D_n^k$  и ядерной оценки плотности распределения этих величин.

Для выполнения задания можно воспользоваться как функцией `seaborn.distplot` (желательно), так и реализацией в `statsmodels`.

In [3]:

```
1 def DX2(sample):
2     ones = np.ones(len(sample))
3     grid = np.arange(1, len(sample) + 1)
4     ecdfs = np.einsum("i,j->ij", ones, 1/grid[::-1])
5     cdfs = np.sort(sps.norm().cdf(sample), axis=1)
6
7     return np.max(np.abs(ecdfs - cdfs), axis=1)
8
9
10 sample = sps.norm().rvs(size = (10000, 10000))
11
12
13 # DD = DX2(sample)
```

In [6]:

```
1 ones = np.ones(len(sample))
2 grid = np.arange(1, len(sample) + 1)
```

In [16]:

```
1 ecdfs = np.einsum("i,j->ij", ones, grid/10000)
```

In [4]:

```
1 cdfs = sps.norm().cdf(sample)
```

In [5]:

```
1 cdfs = np.sort(cdfs, axis=1)
```

In [17]:

```
1 DD = np.max(np.abs(ecdfs - cdfs), axis=1)
```



In [30]:

```

1 ax = sns.distplot(100*DD)
2
3
4 ax.set_title(r"Гистограммы значений  $\sqrt{n}D_n^k$  и KDE",
5             fontsize=14)
6
7 plt.show()

```

**Вывод:**

Как показывает график, мы можем наблюдать что  $\sqrt{n}D_n$  - похоже на конкретное распределение - а именно на распределение Колмогорова. Это подтверждает теорему Колмогорова-Смирнова о сходимости  $\sqrt{n}D_n$  по распределению к случайной величине из распределения Колмогорова.

**Задача 3.**

Для каждого распределения постройте график эмпирической функции распределения, гистограмму и график ядерной оценки плотности.

Вам выдется почти готовый код для выполнения задания с некоторыми пропусками. Код предполагает использование реализации ядерных оценок плотности из `statsmodels`. При желании вы можете написать аналогичный код, используя реализацию в `seaborn`.

In [7]:

```
1 def draw_ecdf(sample, grid, cdf=None):
2     """По сетке grid строит графики эмпирической функции распределения
3     и истинной (если она задана) для всей выборки и для 1/10 ее части.
4     """
5     plt.figure(figsize=(16, 3))
6     for i, size in enumerate([len(sample) // 10, len(sample)]):
7         plt.subplot(1, 2, i + 1)
8
9         plt.scatter(sample[:size], np.zeros(size),
10                    alpha=0.4, label='sample')
11
12         if cdf is not None:
13             plt.plot(grid,
14                    cdf(grid),
15                    color='green', alpha=0.3, lw=2, label='true cdf')
16
17         plt.plot(grid,
18                    ECDF(sample[:size])(grid),
19                    color='red', label='ecdf')
20
21         plt.legend()
22         plt.grid(ls=':')
23         plt.title('sample size = {}'.format(size))
24     plt.show()
```

In [12]:

```
1 def draw_hist(sample, grid, pdf=None):
2     """Строит гистограмму и, по сетке grid, график истинной плотности
3     (если она задана) для всей выборки и для 1/10 ее части.
4     """
5     plt.figure(figsize=(16, 3))
6     for i, size in enumerate([len(sample) // 10, len(sample)]):
7         plt.subplot(1, 2, i + 1)
8
9         plt.hist(sample[:size],
10                bins=20,
11                range=(grid.min(), grid.max()),
12                density=True, label='sample hist')
13
14         if pdf is not None :
15             plt.plot(grid,
16                    pdf(grid),
17                    color='green', alpha=0.3, lw=2, label='true pdf')
18
19         plt.legend()
20     plt.show()
```

In [21]:

```

1  def draw_pdf(sample, grid, pdf=None):
2      """По сетке grid строит графики ядерной оценки плотности
3      и истинной плотности (если она задана) для всей выборки
4      и для 1/10 ее части.
5      """
6      plt.figure(figsize=(16, 3))
7      for i, size in enumerate([len(sample) // 10, len(sample)]):
8          plt.subplot(1, 2, i + 1)
9
10         plt.scatter(sample[:size], np.zeros(size),
11                     alpha=0.4, label='sample')
12
13         if pdf is not None:
14             plt.plot(grid,
15                     pdf(grid), color='green',
16                     alpha=0.3, lw=2, label='true pdf')
17
18         kde = KDEUnivariate(sample[:size])
19         kde.fit(kernel='gau')
20         plt.plot(grid,
21                 kde.evaluate(grid),
22                 color='red', label='kde')
23
24         plt.legend()
25         plt.grid(ls=':')
26     plt.show()

```

Теперь примените реализованные выше функции к выборкам размера 500 для следующих распределений:

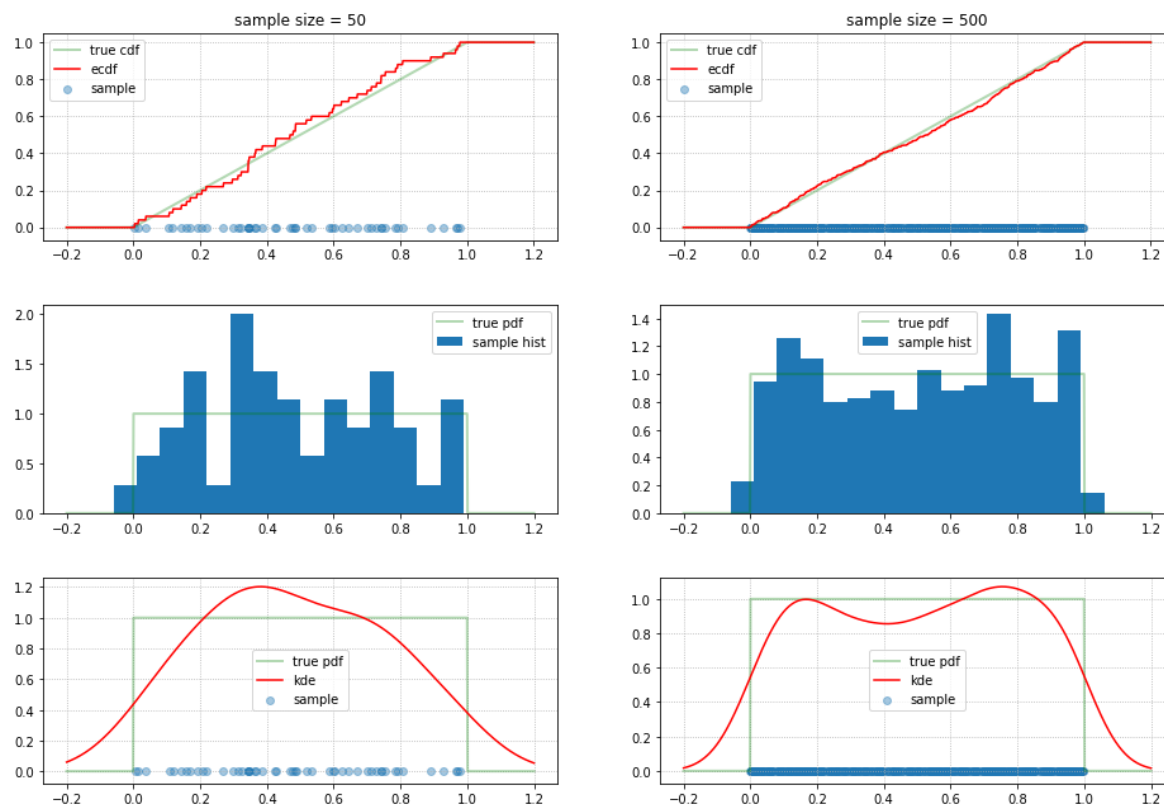
1. Равномерное распределение  $U[0, 1]$ . Графики (ф.р., плотностей) строить на интервале  $(-0.2, 1.2)$ .

In [22]:

```

1 n = 500
2 sample = sps.uniform().rvs(500)
3 grid = np.linspace(-0.2, 1.2, 1000)
4
5 draw_ecdf(sample, grid, sps.uniform().cdf)
6
7 draw_hist(sample, grid, sps.uniform().pdf)
8
9 draw_pdf(sample, grid, sps.uniform().pdf)

```



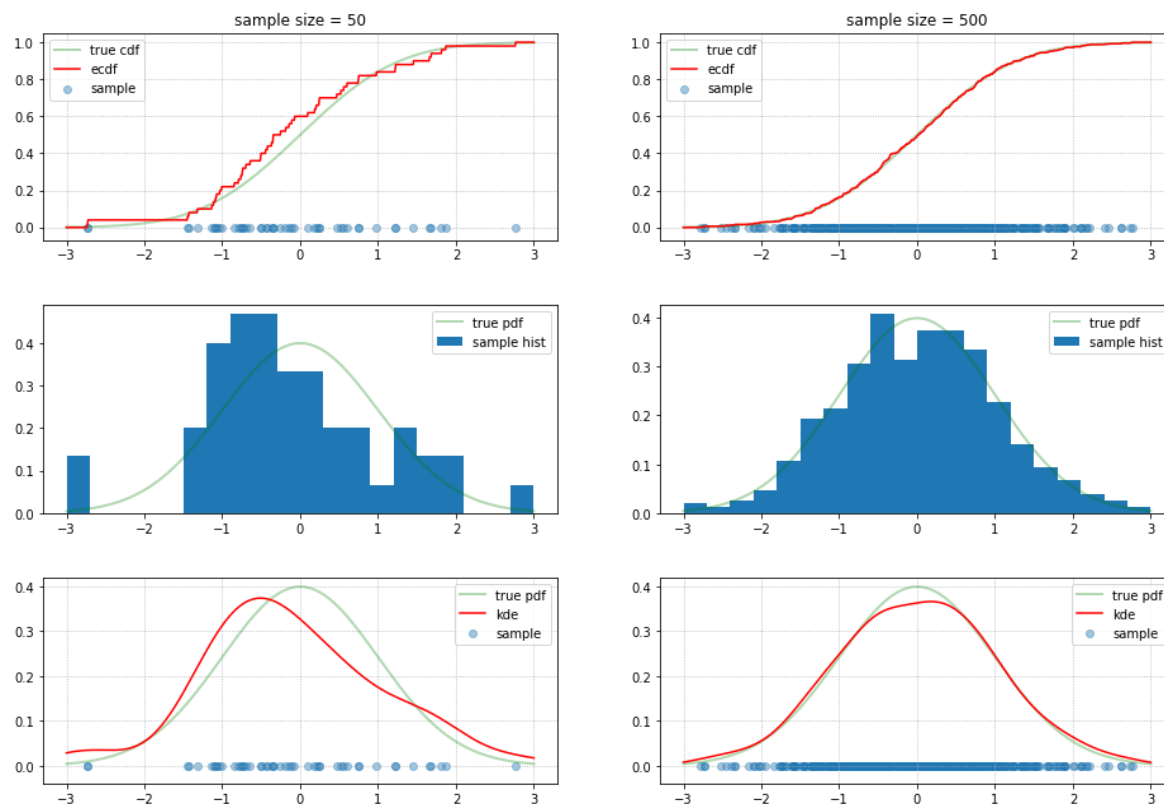
2. Нормальное распределение  $\mathcal{N}(0, 1)$ . Графики строить на интервале  $(-3, 3)$ .

In [14]:

```

1 sample = sps.norm().rvs(500)
2 grid = np.linspace(-3, 3, 1000)
3
4 draw_ecdf(sample, grid, sps.norm().cdf)
5
6 draw_hist(sample, grid, sps.norm().pdf)
7
8 draw_pdf(sample, grid, sps.norm().pdf)

```



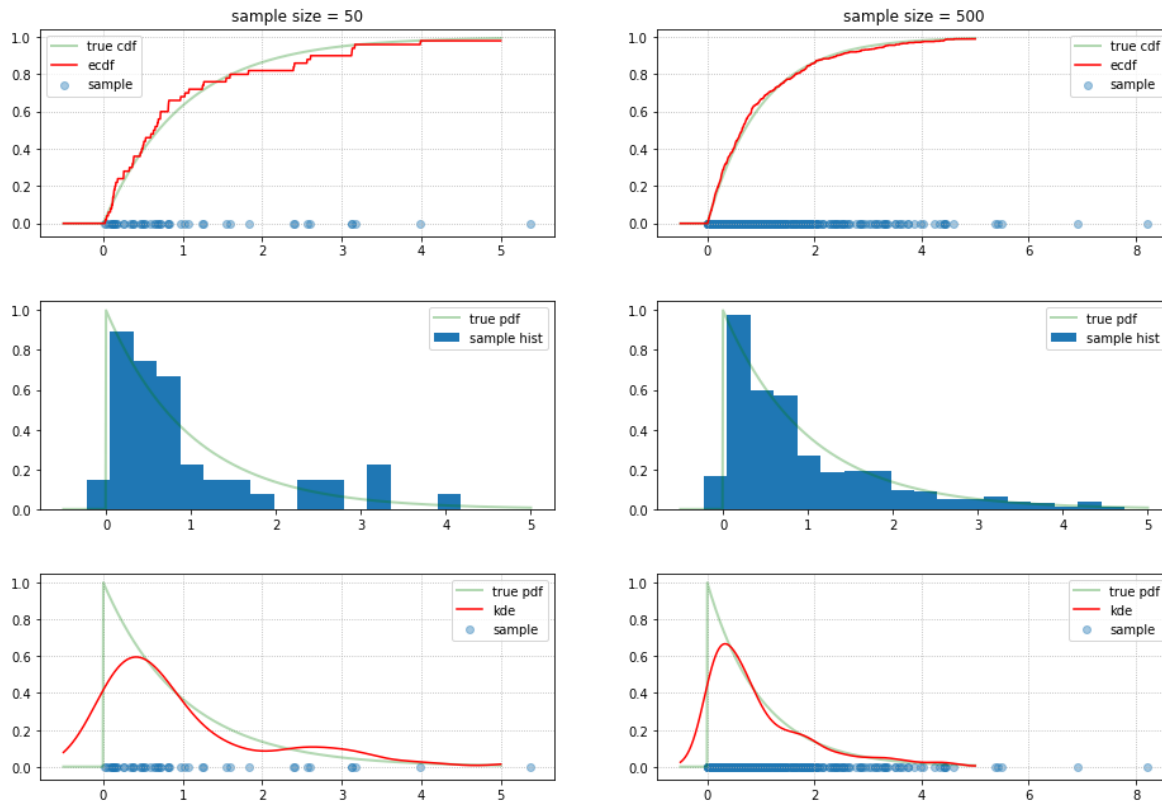
3. Экспоненциальное распределение  $Exp(1)$ . Графики строить на интервале  $(-0.5, 5)$ .

In [16]:

```

1 sample = sps.expon().rvs(500)
2 grid = np.linspace(-0.5, 5, 1000)
3
4 draw_ecdf(sample, grid, sps.expon().cdf)
5
6 draw_hist(sample, grid, sps.expon().pdf)
7
8 draw_pdf(sample, grid, sps.expon().pdf)

```

**Вывод:**

Мы видим, что уже на  $\frac{1}{10}$  части выборки чаще всего угадывается распределение, но можно заметить, что ядерная оценка плотности на небольшой выборке часто имеет лишние "горбы", так как на этом этапе значение одного элемента выборки имеет большой вес. Но на большой выборке все лишние скачки выравниваются, и оценка плотности все больше становится похожа на истинную.

Если брать гауссовское ядро, то любая оценка плотности получается довольно гладкой. Это наиболее хорошо подходит для плотности нормального распределения, также довольно неплохо оценивается экспоненциальная плотность, хотя она не успевает достать до "гребня" в силу того, что оценка старается все сгладить. Наиболее плохо она оценивает равномерную плотность, опять же потому что старается сгладить края.

**Задача 4.**

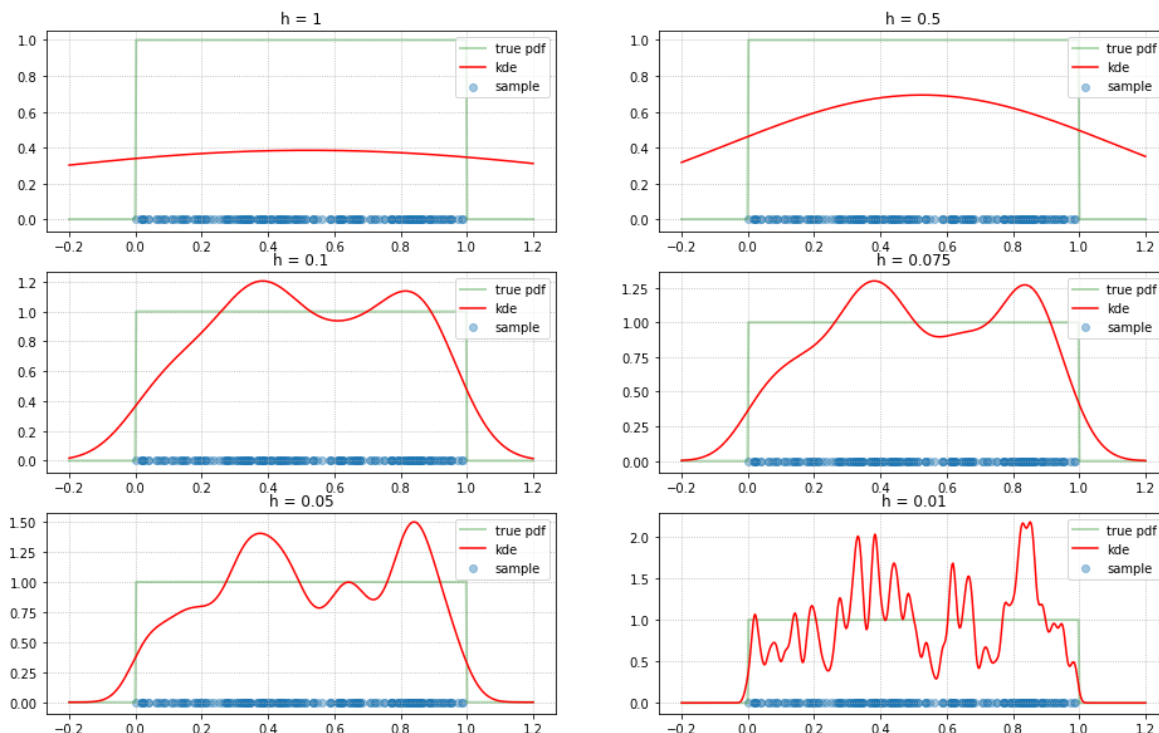
Исследуйте вид ядерной оценки плотности в зависимости от вида ядра и его ширины.

Для этого сгенерируйте выборку  $X_1, \dots, X_{200}$  из распределения  $U[0, 1]$  и постройте серию графиков для различной ширины гауссовского ядра, а затем другую серию графиков для различных типов ядер при фиксированной ширине. На каждом графике на отрезке  $[-0.2, 1.2]$  должны быть изображены истинная плотность (полупрозрачным цветом) и ее ядерная оценка, а так же с нулевой у-координатой должны быть нанесены точки выборки. Для экономии места стройте графики в два столбца.

Вам выдется почти готовый код для выполнения задания с некоторыми пропусками. Код предполагает использование реализации ядерных оценок плотности из `statsmodels`. При желании вы можете написать аналогичный код, используя реализацию в `seaborn`.

In [23]:

```
1 size = 200
2 sample = sps.uniform().rvs(size)
3 grid = np.linspace(-0.2, 1.2, 500)
4
5 plt.figure(figsize=(16, 10))
6 for i, bw in enumerate([1, 0.5, 0.1, 0.075, 0.05, 0.01]):
7     plt.subplot(3, 2, i + 1)
8     kernel_density = KDEUnivariate(sample)
9     kernel_density.fit(bw=bw)
10
11     plt.scatter(sample, np.zeros(size), alpha=0.4, label='sample')
12     plt.plot(grid, sps.uniform.pdf(grid), color='green',
13              alpha=0.3, lw=2, label='true pdf')
14     plt.plot(grid, kernel_density.evaluate(grid),
15              color='red', label='kde')
16     plt.legend(loc=1)
17     plt.grid(ls=':')
18     plt.title('h = {}'.format(bw))
19 plt.show()
```

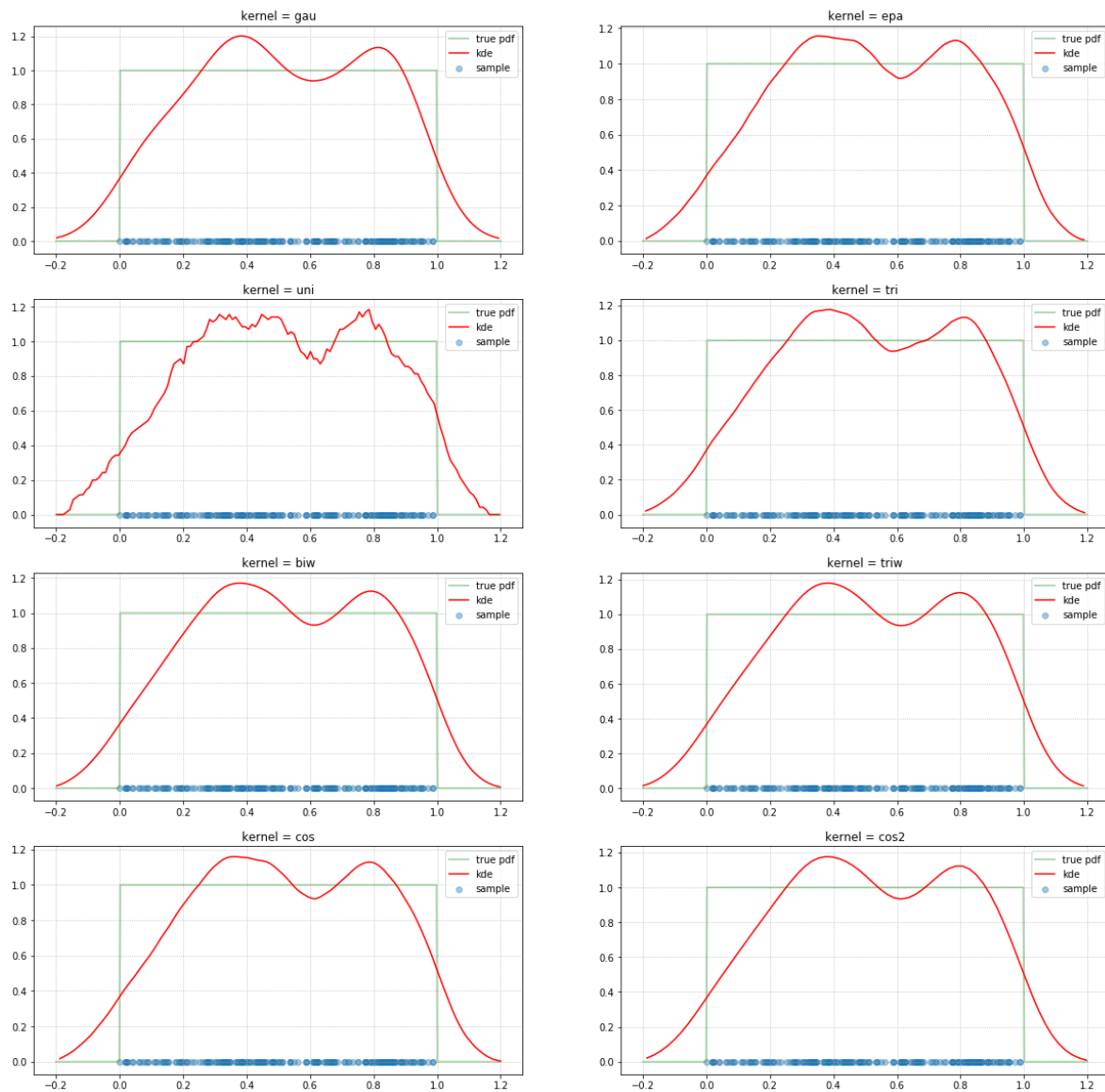


Во втором случае графики постройте аналогичным образом, проведя итерации по типу ядра.

In [52]:

```
1 plt.figure(figsize=(20, 20))
2 for i, kernel in enumerate(kernel_switch.keys()):
3     plt.subplot(4, 2, i + 1)
4     kernel_density = KDEUnivariate(sample)
5     kernel_density.fit(kernel=kernel, fft=False)
6
7     plt.scatter(sample, np.zeros(size), alpha=0.4, label='sample')
8     plt.plot(grid, sps.uniform.pdf(grid), color='green',
9             alpha=0.3, lw=2, label='true pdf')
10
11     support = kernel_density.support
12     k_density = kernel_density.density[np.logical_and(support >= -0.2, support
13 support = support[np.logical_and(support >= -0.2, support <= 1.2)]
14
15 #     print(len(kernel_density.support))
16 #     print(len(kernel_density.density[np.logical_and(support >= -0.2, support
17
18
19     plt.plot(support, k_density,
20             color='red', label='kde')
21     plt.legend(loc=1)
22     plt.grid(ls=':')
23     plt.title('kernel = {}'.format(kernel))
24
25 plt.show()
```





### Вывод:

Все ядра на выборке размером 200 дают очень похожие результаты, разве что равномерное немного отличается, так как оно наименее гладкое среди всех, хотя во основном форма почти одна и та же.

А вот вариация с  $h$  дает различные результаты. При слишком большом значении  $h$  - плотность получается слишком гладкой - можно ошибочно подумать, что данные выглядят как единая группа, на всем рассматриваемом отрезке, плотность никак не отражает отдельные скопления в данных. Наоборот, если взять слишком маленькую  $h$ -ку, то в плотности будут присутствовать слишком много скачков. Поэтому стоит вдумчиво подходить к выбору  $h$

## Задача 5.

Скачайте данные [wine dataset](http://archive.ics.uci.edu/ml/datasets/wine) (<http://archive.ics.uci.edu/ml/datasets/wine>) и выберите произвольные 7 столбцов с действительными числами. С помощью `seaborn.PairGrid` постройте таблицу графиков, состоящую из

- одномерных ядерных оценок плотности по диагонали;
- двумерных ядерных оценок плотности ниже диагонали;
- scatter-plot выше диагонали (`plt.scatter`).

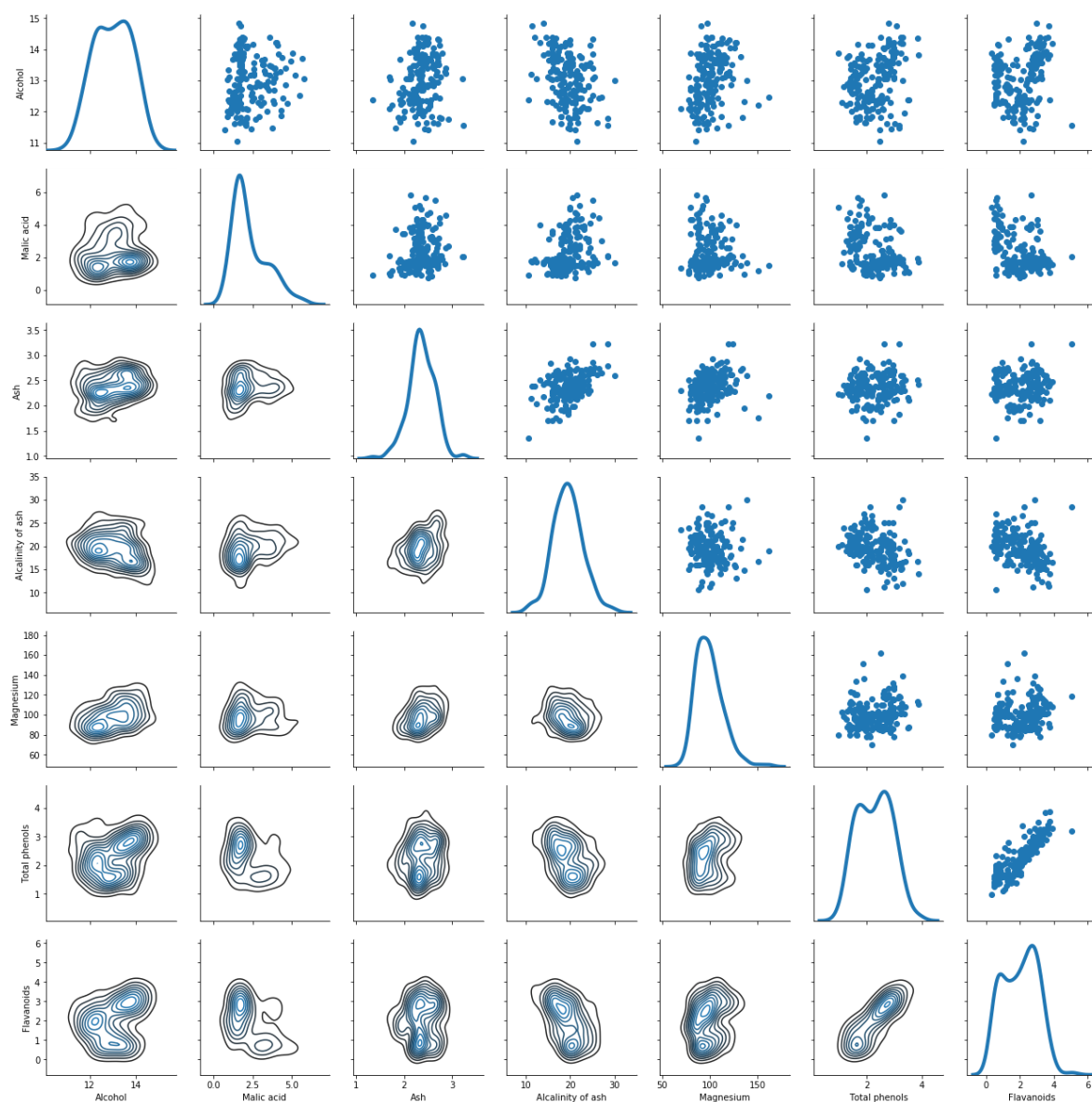
При возникновении затруднений посмотрите обучающий ноутбук по `seaborn`.

In [5]:

```
1 data = pd.read_csv("wine.data")
2
3 columns = ["Alcohol", "Malic acid", "Ash",
4           "Alcalinity of ash", "Magnesium",
5           "Total phenols", "Flavanoids"]
6
7 data = data.loc[:, columns]
8
9 # sns.set(style="darkgrid")
10
11 g = sns.PairGrid(data, diag_sharey=False)
12
13 g.map_diag(sns.kdeplot, lw=4)
14 g.map_lower(sns.kdeplot)
15 g.map_upper(plt.scatter)
16
```

Out[5]:

<seaborn.axisgrid.PairGrid at 0x7f84e8e124e0>



**Вывод:**

Судя по двумерным ядерным оценкам плотностей в данных - мы можем делать какие-то выводы о зависимости признаков в данных, а также можем судить о выбросах в данных.

Судя по одномерным оценкам плотностей мы можем сказать кое-что о распределении данных. Например, то, что

- Alcohol, Total phenols и Flavanoids - довольно похоже на равномерное распределение
- Ash, Alcanity of Ash - больше смахивает на нормальное
- Malic Acid - похоже на экспотенциальное или даже гамма
- Magnesium - возможно тоже экспотенциальное

---

**Задача 6.**

Около месяца назад (24-25 октября) в Краснодарском крае сильные ливни привели наводнению, из-за которого сильно пострадали города Сочи и Туапсе. Вам выданы данные об уровне воды за 2014-2018 год по следующим рекам Краснодарского края:

- Мзымта (Сочи, Адлерский район, Роза Хутор)
- Сочи (Сочи, Центральный район)
- Туапсе (Туапсе)
- Херота (Сочи, Адлерский район)
- Хоста (Сочи, Хостинский район)

В файлах используйте столбец `Уровень воды (по БСВ)`. Это уровень воды по [Балтийской системе высот](https://ru.wikipedia.org/wiki/Балтийская_система_высот) ([https://ru.wikipedia.org/wiki/Балтийская\\_система\\_высот](https://ru.wikipedia.org/wiki/Балтийская_система_высот)) — принятой в СССР системе нормальных высот, отсчёт которых ведётся от нуля Кронштадтского футштока. Для каждой реки нарисуйте график уровня воды.

Данные собраны за каждые 10 минут, что достаточно тяжело обрабатывать. Преобразуйте данные, рассмотрев максимальное значение уровня воды за сутки. Вам нужно по данным до октября 2018 года не включительно построить верхнюю границу предсказательного интервала уровня воды и сравнить ее с максимальным значением, достигавшимся в октябре 2018 года.

Предсказательный интервал постройте в три этапа:

1. Бутстрепный доверительный интервал для среднего значения максимального уровня воды за сутки;
2. Бутстрепный доверительный интервал для стандартного отклонения максимального уровня воды за сутки;
3. Сложите границу доверительного интервала для среднего с границей доверительного интервала для стандартного отклонения, домноженной на 2.

Рассмотрите три способа построения бутстрепных доверительных интервалов, рассказанные на лекции.

Сделайте выводы.

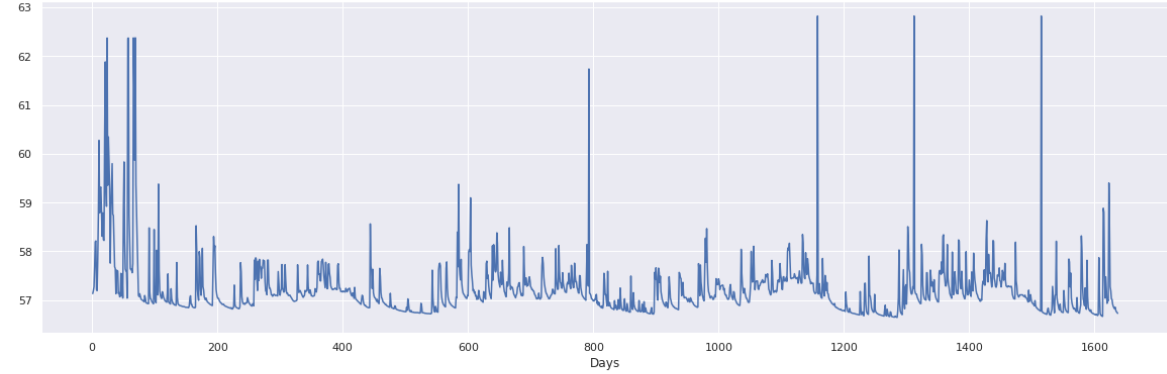
In [197]:

```
1  #преобразование данных в посуточные
2  def get_mod_data(name):
3      data = pd.read_csv("data/" + name + ".csv", '\t')
4      data = data.iloc[:, [0,2]]
5
6      data['Время'] = pd.to_datetime(data['Время'])
7
8      data['День'] = data['Время'].dt.date
9      data['Время'] = data['Время'].dt.time
10
11     october_begin = datetime.strptime('2018-10-01',
12                                       "%Y-%m-%d").date()
13
14     october_end = datetime.strptime('2018-10-31',
15                                    "%Y-%m-%d").date()
16
17
18     past_data = data.where(data['День'] < \
19 october_begin).dropna()
20
21     october_data = data.where((data['День'] >= \
22 october_begin) & (data['День'] <= october_end)).dropna()
23
24     return (past_data.loc[:, ['Уровень воды (по БСВ)',
25                               'День']].groupby(by='День').max(),
26            october_data.loc[:, ['Уровень воды (по БСВ)',
27                                 'День']].groupby(by='День').max())
28
29
30 Sochi, october_sochi = get_mod_data("Сочи")
31 Tuapse, october_tuapse = get_mod_data("Туапсе")
32 Herota, october_herota = get_mod_data("Херота")
33 Hosta, october_hosta = get_mod_data("Хоста")
34 Msiymta, october_msiymta = get_mod_data("Мзымта")
35
```

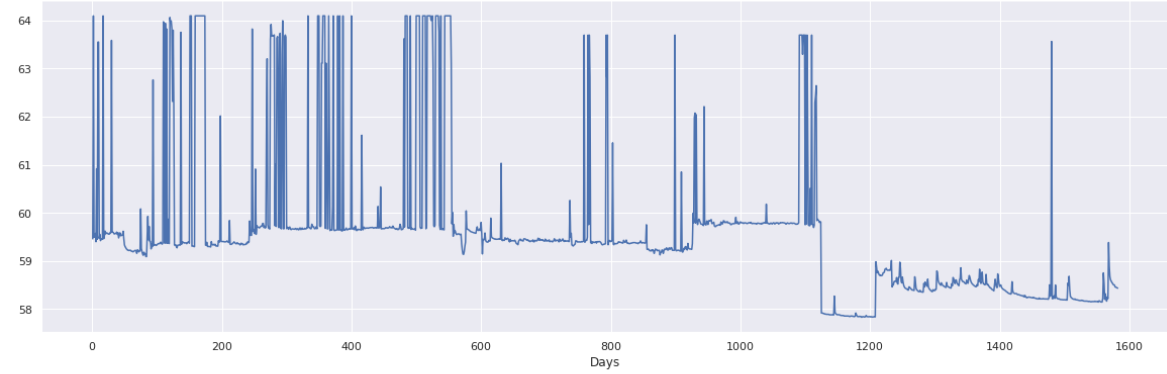
In [198]:

```
1 def get_plot(data, name, ind):
2     plt.subplot(5,1,ind)
3     plt.title(name, fontsize=14)
4     grid = np.arange(1, len(data) + 1)
5
6     plt.xlabel('Days', fontsize=12)
7
8     plt.plot(grid, data['Уровень воды (по БСВ)'],
9              label='Уровень воды (по БСВ)')
10
11
12 plt.figure(figsize=(20,35))
13
14 get_plot(Sochi, "Сочи", 1)
15 get_plot(Tuapse, "Туапсе", 2)
16 get_plot(Herota, "Херота", 3)
17 get_plot(Hosta, "Хоста", 4)
18 get_plot(Msiymta, "Мзымта", 5)
19
```

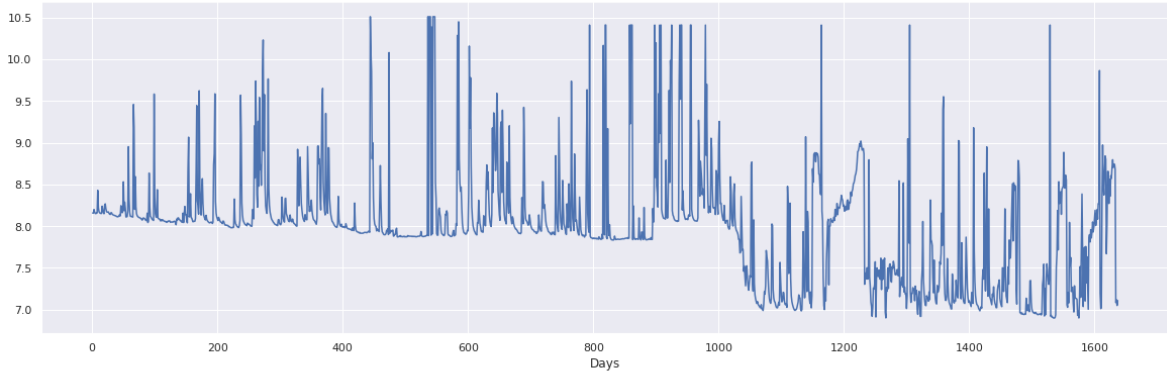
Сочи



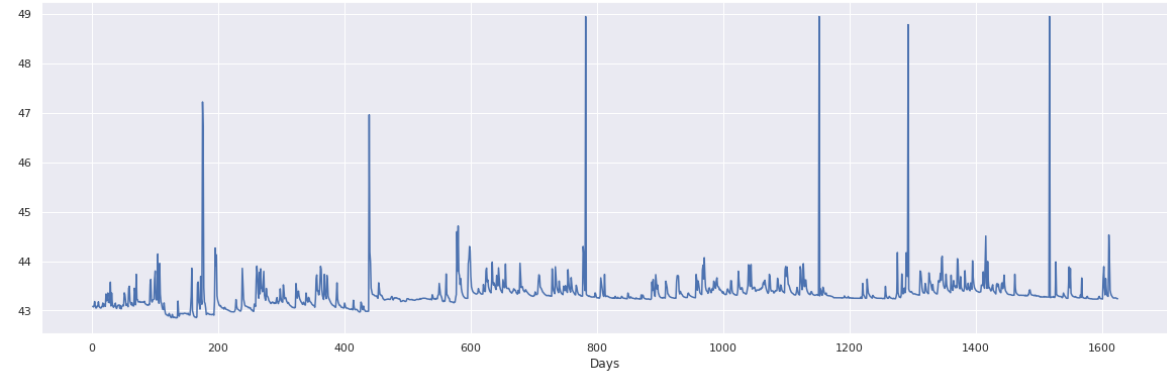
Туанце



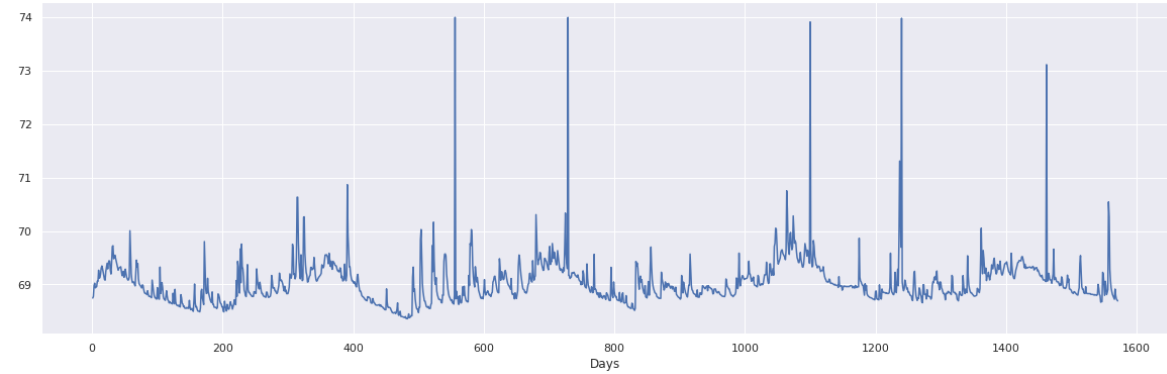
Херота



Хоста



Мзымта



In [199]:

```
1  # реализация бутстрепных интервалов
2
3  def bootstrap(T, sample, count):
4      """
5      возвращает множество бутстрепных выборок
6      из sample
7      в количестве count
8      """
9      indices = np.random.choice(len(sample), size=(count, len(sample)))
10     bootstraps = np.take(sample, indices)
11     return T(bootstraps, axis=-1)
12
13
14 def norm_interval(T, data, alpha=0.05, count=10000):
15     """
16     реализует бутстрепный нормальный интервал параметра
17     по выборке sample, где статистика T - а.н.о. параметра,
18     alpha - уровень доверия
19     count - количество бутстрепных выборок
20     возвращает верхнюю границу интервала
21     """
22     quantile = sps.norm().ppf((1 + alpha)/2)
23
24     sample = data['Уровень воды (по БСВ)'].get_values()
25     bstraps = bootstrap(T, sample, count)
26     boot_var = (bstraps**2).mean() - (bstraps.mean())**2
27     return T(sample, axis=0) + quantile*np.sqrt(boot_var)
28
29
30 def ref_interval(T, data, alpha=0.05, count=10000):
31     """
32     реализует опорный интервал
33     """
34     sample = data['Уровень воды (по БСВ)'].get_values()
35     bstraps = bootstrap(T, sample, count)
36
37     k = int(np.floor(count*(1 - alpha)/2))
38     q = np.partition(bstraps, k)[k]
39
40     return 2*T(sample, axis=0) - q
41
42
43 def q_interval(T, data, alpha=0.05, count=10000):
44     """
45     реализует квантильный интервал
46     """
47     sample = data['Уровень воды (по БСВ)'].get_values()
48     bstraps = bootstrap(T, sample, count)
49
50     k = int(np.ceil(count*(1 + alpha)/2) )
51     return np.partition(bstraps, k)[k]
52
53
```

In [203]:

```

1 datas = [Sochi, Tuapse, Herota, Hosta, Msiymta]
2 october_datas = [october_sochi, october_tuapse, october_herota,
3                  october_hosta, october_msiymta]
4 names = ["Сочи", "Туапсе", "Херота", "Хоста", "Мзымта"]
5
6 table = list()
7
8 for data, october_data, name in zip(datas, october_datas, names):
9     norm = norm_interval(np.mean, data) + 2*norm_interval(np.std, data)
10    ref = ref_interval(np.mean, data) + 2*ref_interval(np.std, data)
11    quant = q_interval(np.mean, data) + 2*q_interval(np.std, data)
12
13
14    raw = (norm, ref, quant,
15           october_data['Уровень воды (по БСВ)'].max())
16    table.append(raw)
17
18
19 result = pd.DataFrame(table, columns=["Нормальный инт-л",
20                                     "Опорный инт-л",
21                                     "Квантильный инт-л",
22                                     "Октябрьское макс. значение"
23                                     ],
24                       index=names)
25
26 result

```

Out[203]:

	Нормальный инт-л	Опорный инт-л	Квантильный инт-л	Октябрьское макс. значение
Сочи	58.476398	58.479912	58.471512	59.603
Туапсе	62.547360	62.549826	62.545613	62.346
Херота	9.327940	9.328774	9.326740	10.063
Хоста	44.119223	44.126584	44.111835	47.513
Мзымта	69.899866	69.905802	69.892356	71.990

**Вывод:** Как мы видим, в реке Сочи уровень воды выше правой границы доверительного интервала на 1 метр, как и река Мзымта. У реки Херота, выше примерно на 0.68 м. Макс. значение в октябре у реки Туапсе попадает в дов. интервал, то есть в ее уровне воды нет ничего необычного. А вот уровень реки Хоста превышает "допустимый" больше чем на 3 метра, что довольно существенно. Значит, вполне вероятно, что разлив именно реки Хосты из этих 5 рек больше всего повлиял на наводнение.

Что касается дов. интервалов, то мы рассмотрели 3 их вида, и видим, что значения их границы практически одинаковы с точностью до тысячных. То есть доверительные интервалы, полученные с помощью различных методов, при  $n \rightarrow +\infty$ , где  $n$  — размер выборки, сходятся к похожим значениям.

## Проверка статистических гипотез

### Задача 7.



Существует примета, что если перед вами дорогу перебегают черные коты, то скоро случится неудача. Вы же уже достаточно хорошо знаете статистику и хотите проверить данную примету. Сформулируем задачу на математическом языке.

Пусть  $X_1, \dots, X_n \sim \text{Bern}(p)$  --- проведенные наблюдения, где  $X_i = 1$ , если в  $i$ -м испытании случилась неудача после того, как черный кот перебежал дорогу, а  $p$  --- неизвестная вероятность такого события. Вы хотите проверить гипотезу  $H_0: p = 1/2$  (отсутствие связи между черным котом и неудачей) против альтернативы  $H_1: p > 1/2$  (неудача происходит чаще если черный кот перебегает дорогу).

Известно, что  $S = \{T(X) \geq c_\alpha\}$ , где  $T(X) = \sum X_i$ , является равномерно наиболее мощным критерием для проверки этих гипотез. Чему при этом равно  $c_\alpha$  и как определяется p-value?

- В силу того, что  $T(X) = \sum X_i \sim \text{Bin}(n, p)$ , то для критерия уровня значимости  $\alpha$ ,  $c_\alpha$  будет определяться как  $(1 - \alpha)$  квантиль распределения  $\text{Bin}(n, \frac{1}{2})$
- Пусть  $t$  - реализация статистики  $T(x)$ , где  $x$  - реализация выборки. Тогда p-value мы назовем вероятностью при справедливости гипотезы  $H_0$  - получить такое значение статистики или еще более экстремальное, то есть в данном случае  $P_0(T(X) \geq t)$  - при условии, что  $p = 1/2$

Для начала проверьте, что критерий работает. Возьмите несколько значений  $n$  и реализаций статистики  $T(X)$ . В каждом случае найдите значение  $c_\alpha$  и p-value. Оформите это в виде таблицы (можно через `pandas.DataFrame`).

Пользуйтесь функциями из `scipy.stats`. Внимательно проверьте правильность строгих и нестрогих знаков.

In [223]:

```

1 ns = [5, 10, 15, 20, 30, 50, 100]
2 alpha = 0.05
3
4 data = list()
5
6 for n in ns:
7     sample = sps.bernoulli(p=0.5).rvs(n)
8     t = sample.sum()
9     c = sps.binom(n=n, p=0.5).ppf(1-alpha)
10    p_value = 1 - sps.binom(n=n, p=0.5).cdf(t)
11    data.append((t, c, p_value))
12
13 pd.DataFrame(data, columns = ["Реализация статистики",
14                               r"$c_\alpha$",
15                               "p-value"],
16               index=ns)
17

```

Out[223]:

	Реализация статистики	$c_\alpha$	p-value
5	2	4.0	0.500000
10	6	8.0	0.171875
15	7	11.0	0.500000
20	12	14.0	0.131588
30	13	19.0	0.707668
50	26	31.0	0.335906
100	47	58.0	0.691350

Для каких истинных значений  $p$  с точки зрения практики можно считать, что связь между черным котом и неудачей есть?

- Я думаю с точки практики стоит считать, что связь есть, если  $p$  - существенно отличается от 0.5; Если же  $p$  - находится в некоторой небольшой окрестности  $\frac{1}{2}$  (к примеру, если  $p=0.51$ ) - то с точки зрения практики связь неудачи с котом практически неотличима от отсутствия связи.

Теперь сгенерируйте 10 выборок для двух случаев: 1).  $n = 5, p = 0.75$ ; 2).  $n = 10^5, p = 0.51$ . В каждом случае в виде таблицы выведите реализацию статистики  $T(X)$ , соответствующее p-value и 0/1 -- отвергается ли  $H_0$  (выводите 1, если отвергается).

In [240]:

```

1 ns = [5, 10**5]
2 ps = [0.75, 0.51]
3
4 table = list()
5
6 for n,p in zip(ns, ps):
7     sample = sps.bernoulli(p=p).rvs((10, n))
8     t = sample.sum(axis=1)
9     p_value = 1 - sps.binom(n=n, p=0.5).cdf(t)
10    reject = p_value <= alpha
11    table.append((t, p_value,
12                reject))
13    table = list(zip(t, p_value, reject))
14    print("n = %d"% (n))
15    print(pd.DataFrame(table, columns=["Реализация статистики",
16                                      "p-value",
17                                      "Отвергается?"]))
18    print("\n-----\n")

```

```

n = 5
Реализация статистики  p-value  Отвергается?
0                      4  0.03125      True
1                      4  0.03125      True
2                      3  0.18750     False
3                      1  0.81250     False
4                      4  0.03125      True
5                      3  0.18750     False
6                      5  0.00000      True
7                      5  0.00000      True
8                      5  0.00000      True
9                      5  0.00000      True

```

-----

```

n = 100000
Реализация статистики      p-value  Отвергается?
0          50870  1.838940e-08      True
1          50679  8.632950e-06      True
2          50912  3.932831e-09      True
3          51059  1.034195e-11      True
4          51105  1.354250e-12      True
5          51014  6.973233e-11      True
6          50914  3.648066e-09      True
7          50559  2.011102e-04      True
8          51172  6.039613e-14      True
9          51087  3.030243e-12      True

```

-----

**Вывод:**

На этом примере мы можем видеть, что при маленькой выборке ( $n = 5$ ), чаще всего гипотеза не отвергается, несмотря на то, что истинное значение существенно отличается от значения гипотезы  $H_0$ .

А при очень большой выборке ( $n = 1e5$ ), наша гипотеза всегда отвергается, хотя истинное значение параметра ( $p = 0.51$ ) практически неотличимо от значения в  $H_0$ .

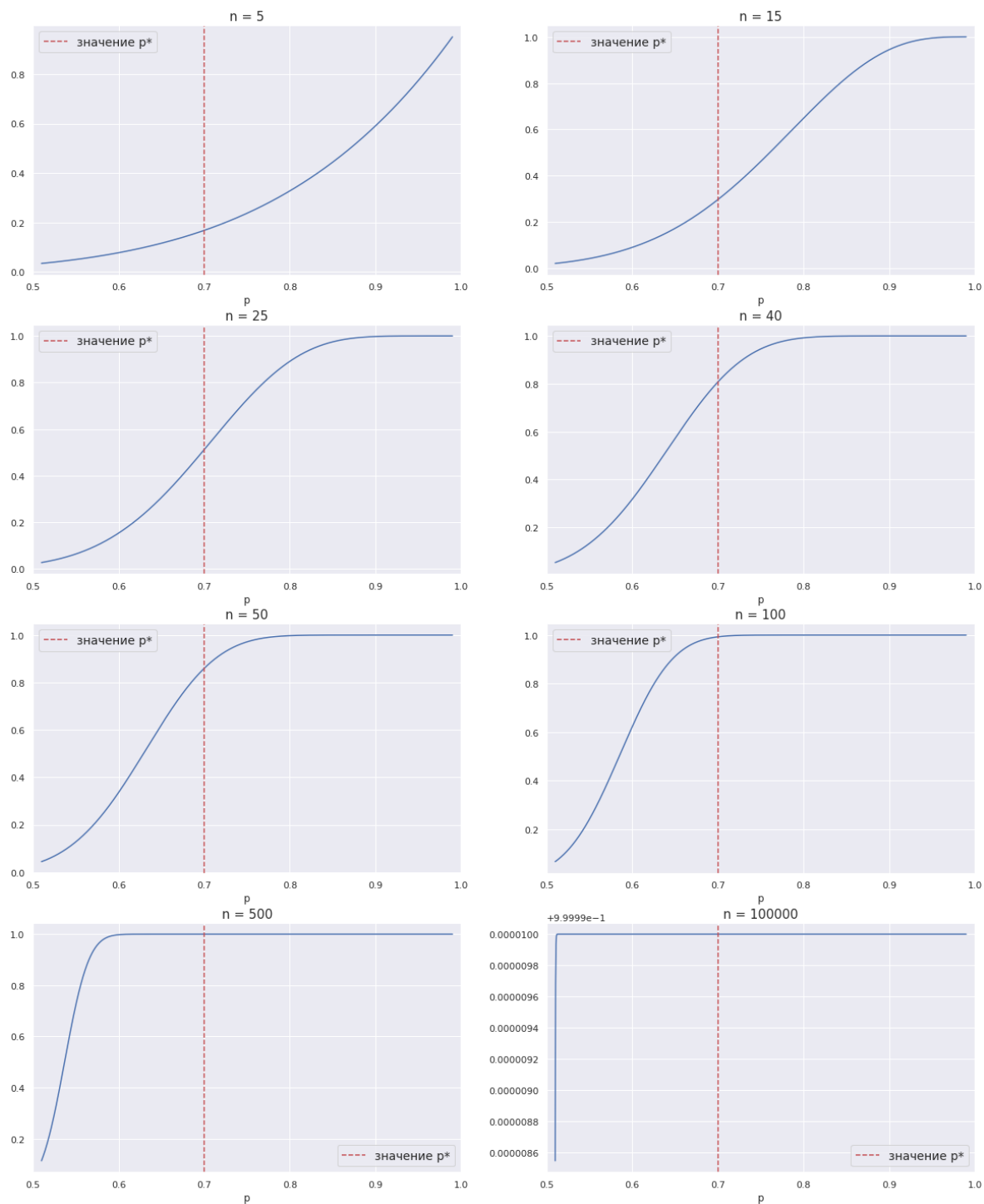
Возникает задача подбора оптимального размера выборки.

Для этого сначала зафиксируйте значение  $p^* > 1/2$ , которое будет обладать следующим свойством. Если истинное  $p > p^*$ , то такое отклонение от  $1/2$  с практической точки зрения признается существенным, то есть действительно чаще случается неудача после того, как черный кот перебегает дорогу. В противном случае отклонение с практической точки зрения признается несущественным.

Теперь для некоторых  $n$  постройте графики функции мощности критерия при  $1/2 < p < 1$  и уровне значимости 0.05. Выберите такое  $n^*$ , для которого функция мощности дает значение 0.8 при  $p^*$ .

In [261]:

```
1 p_thrd = 0.7
2 alpha = 0.05
3 ns = [5, 15, 25, 40, 50, 100, 500, 10**5]
4
5 plt.figure(figsize = (20, 25))
6
7 for i, n in enumerate(ns):
8     plt.subplot(4, 2, i+1)
9     grid = np.linspace(0.51, 0.99, 1000)
10
11     c = sps.binom(n=n, p=0.5).ppf(1-alpha)
12
13     power = 1 - sps.binom(n=n, p=grid).cdf(c)
14
15     plt.title("n = %d" % (n), fontsize=15)
16     plt.xlabel("p", fontsize=12)
17
18     plt.xlim((1/2, None))
19     plt.plot(grid, power)
20     plt.axvline(p_thrd, linestyle='--', c='r',
21                 label='значение p*')
22
23     plt.legend( fontsize=14)
24
25
26
27 plt.show()
```



In [260]:

```
1 c = sps.binom(n=np.arange(30, 50), p=0.5).ppf(1-alpha)
2 1 - sps.binom(n=np.arange(30, 50), p=0.7).cdf(c)
```

Out[260]:

```
array([0.73037039, 0.68790188, 0.64401775, 0.73338179, 0.69316797,
       0.77292538, 0.73651439, 0.80709569, 0.77445184, 0.73971878,
       0.80744825, 0.77618849, 0.83599882, 0.80808733, 0.86059582,
       0.83584182, 0.80895385, 0.85988282, 0.83594621, 0.88062789])
```

Наиболее подходит  $n^* = 40$  - выберем его

Для выбранного  $n^*$  проведите эксперимент, аналогичный проведенным ранее экспериментам, сгенерировав выборки для следующих истинных значений  $p$ : 1).  $1/2 < p < p^*$ ; 2).  $p > p^*$ .

In [269]:

```

1  n_thrd = 40
2  ps = [0.85, 0.55]
3
4  for p in ps:
5      sample = sps.bernoulli(p=p).rvs((10, n_thrd))
6      t = sample.sum(axis=1)
7
8      p_value = 1 - sps.binom(n=n_thrd, p=0.5).cdf(t)
9      reject = p_value <= alpha
10     table.append((t, p_value,
11                  reject))
12
13     table = list(zip(t, p_value, reject))
14     print("n = %d, p=%.2f" % (n_thrd, p))
15     print(pd.DataFrame(table, columns=["Реализация статистики",
16                                       "p-value",
17                                       "Отвергается?"]))
18
19     print("\n-----\n")

```

n = 40, p=0.85

	Реализация статистики	p-value	Отвергается?
0	36	9.732503e-09	True
1	35	9.285122e-08	True
2	36	9.732503e-09	True
3	38	3.728928e-11	True
4	31	9.108291e-05	True
5	31	9.108291e-05	True
6	33	4.182292e-06	True
7	33	4.182292e-06	True
8	32	2.113851e-05	True
9	35	9.285122e-08	True

-----

n = 40, p=0.55

	Реализация статистики	p-value	Отвергается?
0	18	0.682086	False
1	23	0.134094	False
2	22	0.214795	False
3	26	0.019239	True
4	19	0.562685	False
5	24	0.076930	False
6	22	0.214795	False
7	26	0.019239	True
8	18	0.682086	False
9	25	0.040345	True

-----

**Вывод:**

Для выбранного нами  $n^*$  с помощью предварительно выбранного  $p^*$  - ситуация выглядит гораздо лучше. Так - если у нас истинное значение  $p$  - больше порогового, то гипотеза почти наверняка отвергается. А если значение  $p$  - меньше, (то находится относительно близко к значению  $p = 1/2$ ), то гипотеза чаще

всего не отвергается. То есть теперь благодаря хорошему выбору  $n^*$  - ситуация выглядит гораздо правдоподобней и логичней с точки зрения практики.

Мы наблюдали, что во время проверки гипотезы - зачастую нужно быть осторожным с выбором размера выборки, а также то, что наши конечные выводы должны опираться не только на математическую основу, но и на логику, на реальную практическую значимость результата.

## Задача 8.

В Долгопрудном крупная торговая сеть Y10 имеет 100 магазинов и планирует открыть еще 5 магазинов. Сеть Y10 считает магазин успешным, если его дневная выручка в 27 днях из 30 превышает некоторый установленный порог. Благодаря модели машинного обучения, обученной на предыдущих 100 магазинах, было выбрано 5 потенциальных точек, в которых и решено было открыть новые магазины.

Вам предоставлена чистая выручка по каждому из 5 магазинов за день в течении трех месяцев (31 + 30 + 31 день) работы этих пяти магазинов. Считается, что магазин успешен в течении дня, если его выручка за этот день превышает 50000 рублей. По этому правилу сопоставьте каждому магазину набор бернуллиевских случайных величин, которые принимают значение 1, если магазин успешен в течении дня. Для простоты будем считать, что выручка за день не зависит от аналогичных показателей за предыдущие дни, и ее распределение не меняется во времени, то есть данные образуют выборку.

Модель машинного обучения выдает также параметры для априорного бета-распределения распределения по каждому из магазинов:

1.  $Beta(1, 1)$ ;
2.  $Beta(2900, 100)$ ;
3.  $Beta(29, 1)$ ;
4.  $Beta(29, 1)$ ;
5.  $Beta(1, 1)$ .

Проанализируйте данные и ответьте на следующие вопросы.

1. Для каких из магазинов можно утверждать, что с вероятностью 0.95 их можно считать успешными?  
Для ответа на этот вопрос нужно выполнить байесовскую проверку гипотез  $H_0: p = 27/30$  vs.  $H_1: p > 27/30$ , посчитав апостериорную вероятность события  $\{p > 27/30\}$ . Если эта вероятность не меньше 0.95, магазин можно считать успешным.
2. Про какие из магазинов можно сказать, что их лучше закрыть? Закрывать магазин затратно, для этого он должен быть убыточным в 27 днях из 30. Для ответа на этот вопрос нужно выполнить байесовскую проверку гипотез  $H_0: p = 3/30$  vs.  $H_1: p < 3/30$ , посчитав апостериорную вероятность события  $\{p < 3/30\}$ . Если эта вероятность не меньше 0.95, магазин можно закрыть.
3. Что можно сказать об остальных магазинах?
4. Выполнены ли предположения модели на практике? Имеет ли смысл собирать данные сразу после открытия магазина?



In [6]:

```

1 data = pd.read_csv("data/y10.csv", sep='\t')
2 data = data > 50000
3
4 data.head()

```

Out[6]:

	Shop 1	Shop 2	Shop 3	Shop 4	Shop 5
0	True	False	False	False	False
1	True	False	False	False	False
2	True	True	False	True	False
3	True	False	False	False	False
4	True	False	False	True	False

In [7]:

```

1 #проверим магазины на успешность
2 As = [1, 2900, 29, 29, 1]
3 Bs = [1, 100, 1, 1, 1]
4 n = 92
5
6 for i, (a, b) in enumerate(zip(As, Bs)):
7     sample = data["Shop %d" % (i+1)].get_values().astype(int)
8
9     S = 1 - sps.beta(a=a+sample.sum(),
10                    b=b + n - sample.sum() ).cdf(27/30)
11
12     if (S >= 0.95):
13         print("Shop %d - успешный, вероятность (p > 27/30) : %f" % (i + 1, S))
14     else:
15         print("Shop %d - не успешный, вероятность (p > 27/30) : %f" % (i + 1, S))
16     print("-----")
17

```

Shop 1 - успешный, вероятность (p > 27/30) : 0.999371

-----

Shop 2 - успешный, вероятность (p > 27/30) : 1.000000

-----

Shop 3 - не успешный, вероятность (p > 27/30) : 0.000000

-----

Shop 4 - не успешный, вероятность (p > 27/30) : 0.000000

-----

Shop 5 - не успешный, вероятность (p > 27/30) : 0.000000

-----

In [10]:

```

1  #проверим магазины на затратность
2
3  for i, (a, b) in enumerate(zip(As, Bs)):
4      sample = data["Shop %d" % (i+1)].get_values().astype(int)
5
6      S = sps.beta(a=a+sample.sum(),
7                  b=b + n - sample.sum() ).cdf(3/30)
8
9      if (S >= 0.95):
10         print("Shop %d - затратный, вероятность (p < 3/30) - %f" % (i + 1, S))
11     else:
12         print("Shop %d - не затратный, вероятность (p < 3/30) - %f" % (i + 1, S))
13     print("-----")

```

```

Shop 1 - не затратный, вероятность (p < 3/30) - 0.000000
-----
Shop 2 - не затратный, вероятность (p < 3/30) - 0.000000
-----
Shop 3 - не затратный, вероятность (p < 3/30) - 0.000001
-----
Shop 4 - не затратный, вероятность (p < 3/30) - 0.000000
-----
Shop 5 - затратный, вероятность (p < 3/30) - 0.961839
-----

```

Мы ничего не сказали о магазине 3 и магазине 4. У них у обоих было априорное распределение  $Beta(29, 1)$ . Давайте посмотрим на это распределение:

In [308]:

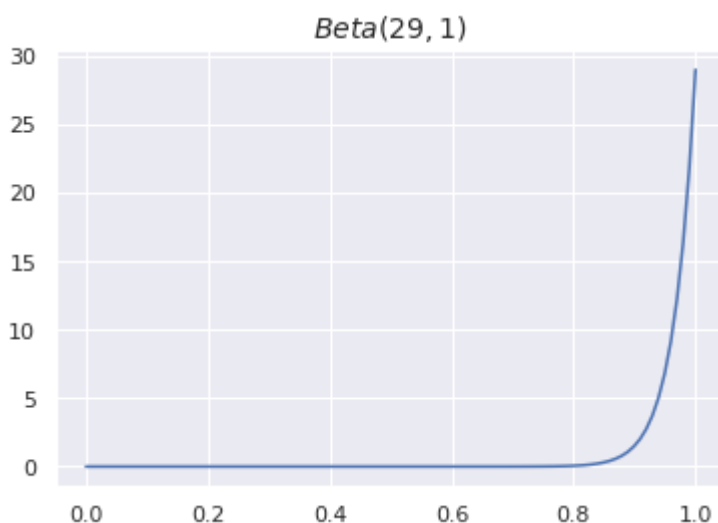
```

1  grid=np.linspace(0,1,100)
2  plt.title(r"$Beta(29,1)$", fontsize=14)
3  plt.plot(grid, sps.beta(29, 1).pdf(grid))

```

Out[308]:

[&lt;matplotlib.lines.Line2D at 0x7fb1ccb27f60&gt;]



То есть модель машинного обучения - при подсчете этих параметров выразила сильную уверенность в успешных продажах у этих магазинов.

Теперь давайте посмотрим на апостериорные распределения.

In [306]:

```

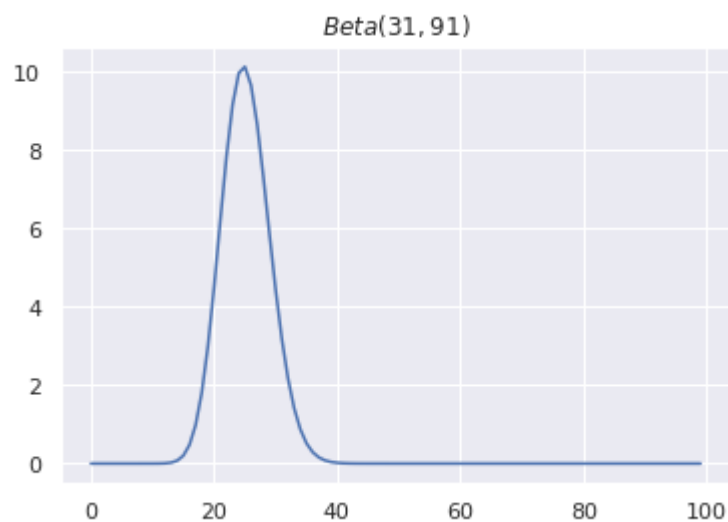
1 shops = ["Shop 3", "Shop 4"]
2
3 for shop in shops:
4     print("----- " + shop + " -----")
5     sample = data[shop].get_values().astype(int)
6     print("кол-во успешных дней : ", sample.sum())
7     print("кол-во неуспешных : ", 92-sample.sum())
8
9     new_a = 29+sample.sum()
10    new_b = 1 + 92 - sample.sum()
11    plt.title(r"$Beta(%d, %d)$" % (new_a, new_b))
12    plt.plot(sps.beta(a=new_a, b=new_b).pdf(grid))
13    plt.show()
14    print("-----")

```

```

----- Shop 3 -----
кол-во успешных дней : 2
кол-во неуспешных : 90

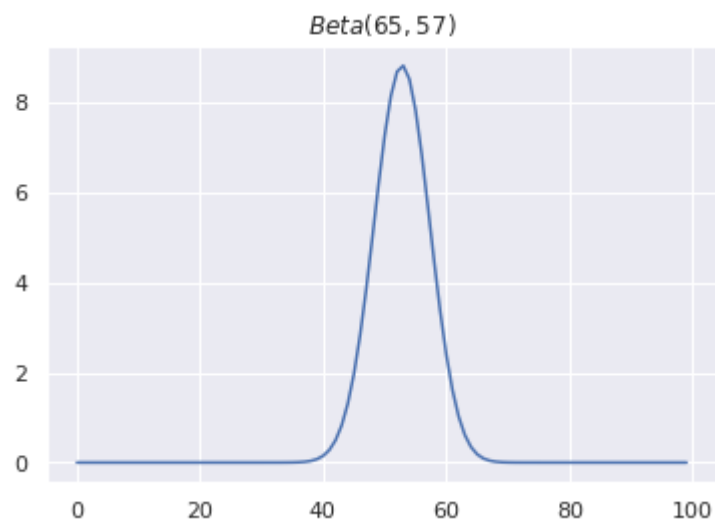
```



```

----- Shop 4 -----
кол-во успешных дней : 36
кол-во неуспешных : 56

```



У обоих магазинов апостериорное распределение получилось больше смещенное к середине, чем к какому-либо краю. А так как пик у этих распределений довольно высок и разброс относительно пика очень маленький, поэтому у краев вероятность практически нулевая. Поэтому ничего не получается сказать - о прибыльности или убыточности этих мест.

In [325]:

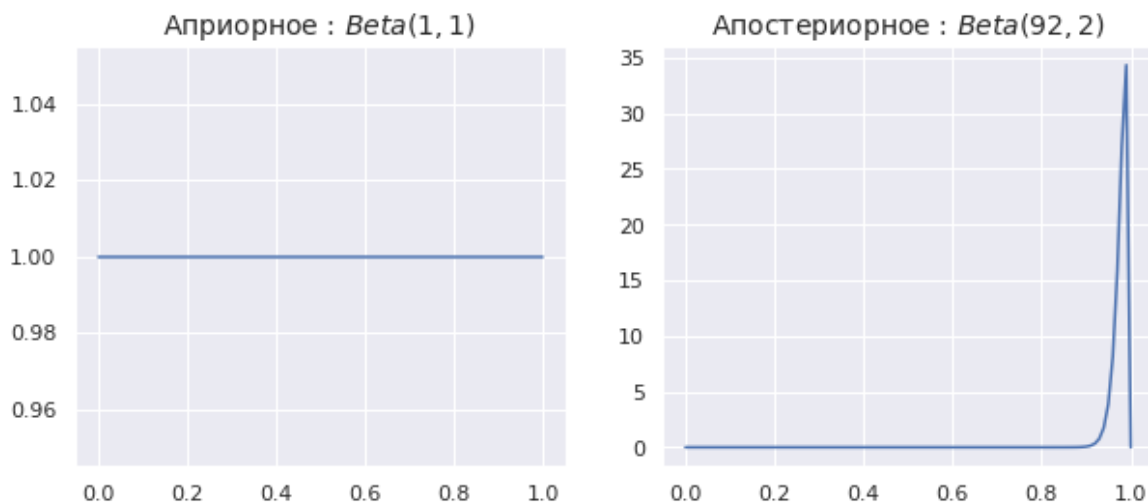
```
1 def experiments(shop, a, b):
2     print("----- " + shop + " -----")
3     sample = data[shop].get_values().astype(int)
4     print("кол-во успешных дней : ", sample.sum())
5     print("кол-во неуспешных : ", 92-sample.sum())
6
7     new_a = a+sample.sum()
8     new_b = b + 92 - sample.sum()
9     plt.figure(figsize=(10,4))
10
11     plt.subplot(1,2,1)
12     plt.title(r"Априорное : $Beta(%d, %d)$" % (a, b),
13             fontsize=14)
14     plt.plot(grid, sps.beta(a=a, b=b ).pdf(grid))
15
16
17     plt.subplot(1,2,2)
18     plt.title(r"Апостериорное : $Beta(%d, %d)$" % (new_a, new_b),
19             fontsize=14)
20     plt.plot(grid, sps.beta(a=new_a, b=new_b ).pdf(grid))
21     plt.show()
22     print("-----")
```

Рассмотрим, действительно ли предположения модели соответствуют наблюдениям.

In [326]:

```
1 experiments("Shop 1", 1, 1)
```

```
----- Shop 1 -----
кол-во успешных дней : 91
кол-во неуспешных : 1
```



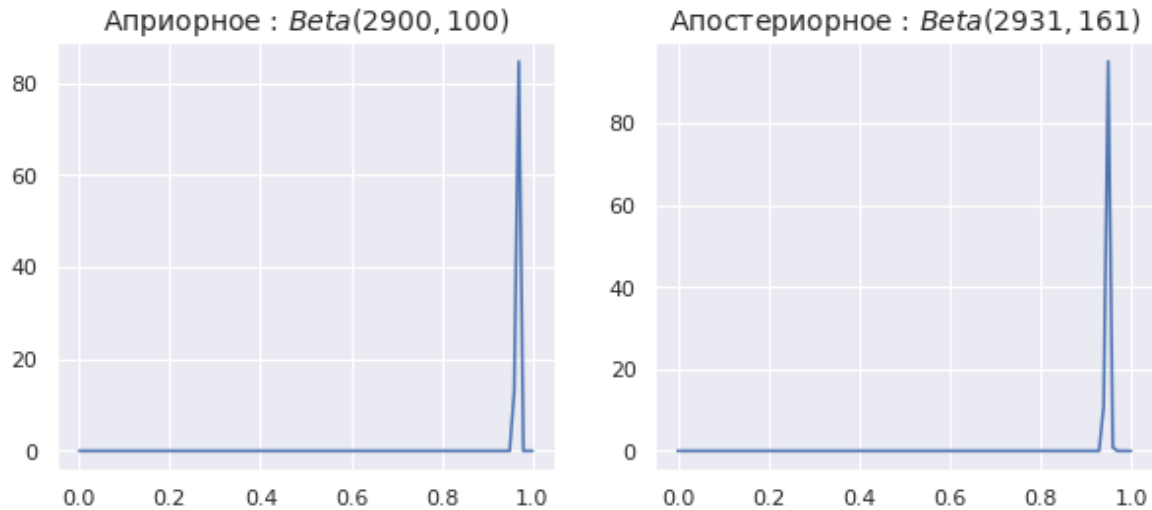
В 1ом магазине модели видимо ничего не было известно об априорном распределении прибыльности за

день, поэтому оставила априорное как равномерное на  $[0, 1]$ . В итоге, судя по выборке, оказалось, что магазин прибыльный, поэтому апостериорное и получилось таким.

In [327]:

```
1 experiments("Shop 2", 2900, 100)
```

```
----- Shop 2 -----  
кол-во успешных дней : 31  
кол-во неуспешных : 61
```



-----

Для второго магазина модель выразила очень большую уверенность в том, что прибыльность будет очень высокой - примерно один неуспешный день к 30 успешным. Но на практике оказалось, что успешных дней в половину меньше, чем плохих дней. Но благодаря "излишней уверенности в успехе" размера выборки не хватило, чтобы не только придти к соответствующему парктике результату, но даже отдалить пик распределения от 1. Поэтому мы совершенно неправильно (исходя из данных) посчитали магазин успешным.

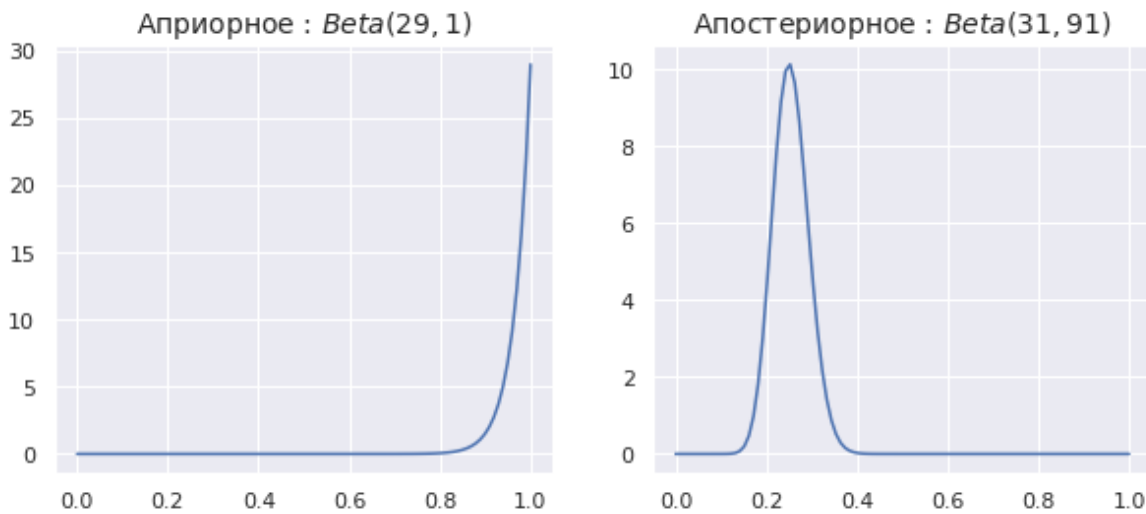
In [328]:

```
1 experiments("Shop 3", 29, 1)
```

----- Shop 3 -----

кол-во успешных дней : 2

кол-во неуспешных : 90



-----

В 3ем магазине ситуация похожа на 2ой магазин, только немного лучше. Модель также выразила уверенность, что прибыльность будет 1 плохой день к 30ти прибыльным, только сказала об этом не так уверенно как в предыдущем пункте. Мы видим, что в третьем магазине данные совершенно не соответствуют априорному распределению - всего лишь 2 успешных дня. Судя по данным вероятность успешного дня этого магазина должна быть очень мала, и апостериорному распределению видимо не хватило размера выборки, чтобы "уйти" вплотную к левому краю, но зато пик распределения все-таки теперь смещен ближе к 0 чем к 1, в отличие от 2ого магазина.

Поэтому, хоть и начальное предположение модели не соответствует выборке, но апостериорное распределение, хоть все еще плохо описывает данные, но показывает перевес в сторону неуспешных дней, что уже неплохо.

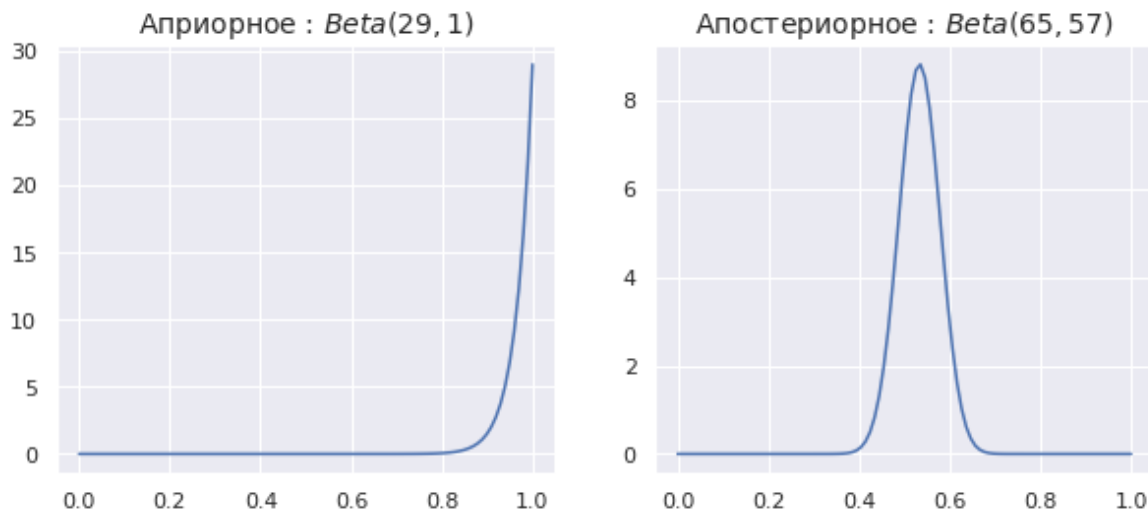
In [329]:

```
1 experiments("Shop 4", 29, 1)
```

----- Shop 4 -----

кол-во успешных дней : 36

кол-во неуспешных : 56



-----

Здесь ситуация похожа на предыдущие два, только здесь априорное распределение уже больше соответствует данным. Опять же, параметры априорного распределения плохо соответствует практике, но апостериорное распределение получилось довольно разумное - оно не относит магазин ни к успешным, ни к затратным.

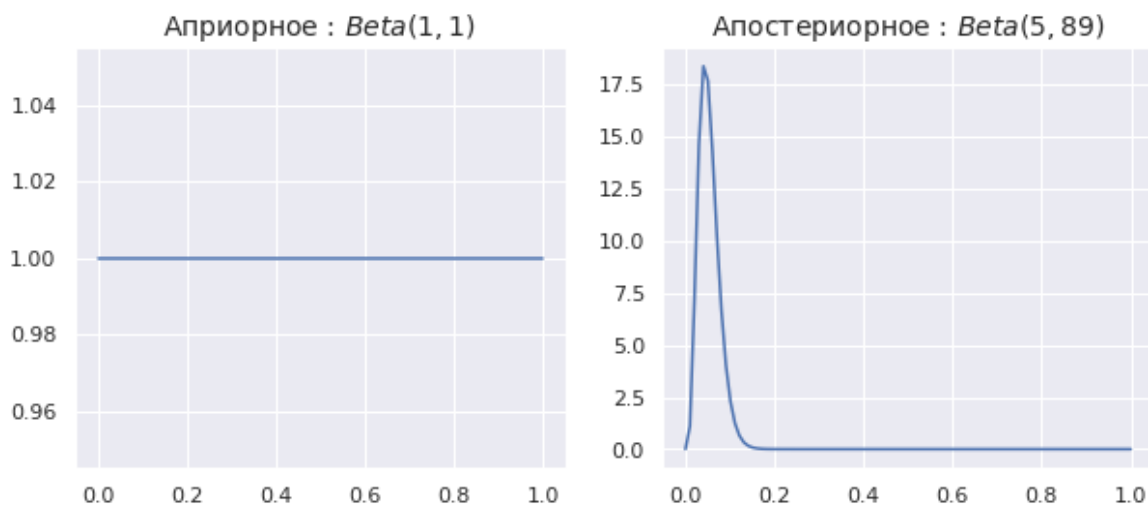
In [330]:

```
1 experiments("Shop 5", 1, 1)
```

----- Shop 5 -----

кол-во успешных дней : 4

кол-во неуспешных : 88



-----

Здесь модель, как и в 1ом магазине, не делала никаких предположений относительно доли успешных дней. В результате кол-во неуспешных получилось гораздо больше чем успешных, поэтому



апостериорное распределение получилось сильно близким к нулю, и магазин закономерно был признан затратным.

Вообще говоря, данные не стоит собирать сразу после открытия магазина, потому что они, вполне вероятно, не будут соответствовать дальнейшим данным. Например, о магазине сразу после открытия еще может не знать большая часть потенциальных покупателей, которая затем о нем узнает. Или, наоборот, была большая рекламная кампания по открытию этого магазина, и в первые дни туда ходило очень много людей, но затем магазин не оправдывает ожидания части населения и его прибыльность снизится.

### **Вывод:**

Нужно быть осторожней с выбором параметров априорного распределения и не выражать слишком большой уверенности в результате, так как если выбрать неправильные априорные параметры, которые к тому же выражают довольно экстремальное распределение (например близкое к какому-либо краю с высоким пиком), то размера реальных данных может не хватить, чтобы "исправить" апостериорное распределение на более правдивое. В итоге получаем результаты, не соответствующие реальности.

Также хотелось бы отметить, что в байесовском подходе при больших данных, зачастую получается очень маленький разброс апостериорного распределения, что не всегда хорошо, так как получается, что вероятность значения параметра, по логике не сильно отличающегося от например Матюжа апостериорного распределения, будет почти нулевой, что может не соответствовать реальной ситуации.

In [ ]:

1	
---	--