

Математическая статистика

Практическое задание 1

В данном задании рассматриваются различные способы генерации выборки из некоторых стандартных распределений, а так же рассматриваются некоторые свойства эмпирической функции распределения и ядерной оценки плотности.

Правила:

- Баллы за каждую задачу указаны далее. Если сумма баллов за задание меньше 25%, то все задание оценивается в 0 баллов.
- Выполненную работу нужно отправить на почту `probability.diht@yandex.ru`, указав тему письма "[номер группы] Фамилия Имя - Задание 1". Квадратные скобки обязательны.
- Прислать нужно ноутбук и его pdf-версию. Названия файлов должны быть такими: `1.N.ipynb` и `1.N.pdf`, где N - ваш номер из таблицы с оценками.
- Никакой код из данного задания при проверке запускаться не будет.

Баллы за задание:

- Задача 1 - 13 баллов
- Задача 2 - 3 балла
- Задача 3 - 5 баллов
- Задача 4 - 3 балла
- Задача 5 - 2 балла
- Задача 6 - 1 балл
- Задача 7 - 3 балла

```
In [81]: import numpy as np
import scipy.stats as sps
import matplotlib.pyplot as plt

%matplotlib inline
```

Задача 1. Имеется симметричная монета. С ее помощью напишите функцию генерации выборки из многомерного нормального распределения с заданными параметрами.

Часть 1. Напишите сначала функцию генерации равномерного распределения на отрезке $[0, 1]$ с заданной точностью. Это можно сделать, записав случайную величину $\xi \sim U[0, 1]$ в двоичной системе счисления $\xi = 0, \xi_1 \xi_2 \xi_3 \dots$. Тогда $\xi_i \sim \text{Bern}(1/2)$ и независимы в совокупности. Приближение заключается в том, что вместо генерации бесконечного количества ξ_i мы полагаем $\xi = 0, \xi_1 \xi_2 \xi_3 \dots \xi_n$.

Для получения максимального балла реализовать функцию нужно так, чтобы она могла принимать на вход в качестве параметра `size` объект `tuple` любой размерности, и возвращать объект `numpy.array` соответствующей размерности. Например, если `size=(10, 1, 5)`, то функция должна вернуть объект размера $10 \times 1 \times 5$. Кроме того, функцию `coin` можно вызвать только один раз, и, конечно же, не использовать какие-либо циклы.

Решение: Подкинем монетку (произведение всех координат кортежа `size`) * `precision`) раз. (Именно столько, так как произведение координат - это "размер" нашей будущей выборки $\xi_i \sim U[0, 1]$, и для каждого необходима точность `precision`).

После применения функции `reshape` у нас получится список из /количество наблюдений/ (это = `sample_size`) строк, и в каждой строке - выборка из распределения Бернулли размера `precision`.

Далее создаем массив из отрицательных степеней двойки. И с его помощью переведем каждое $\xi = 0, \xi_1 \xi_2 \xi_3 \dots \xi_n$ в десятичную систему счисления. Согласно условию задачи, это и будет необходимое значение. Реализуем это через матричное умножение - оно как раз умножит i -ый разряд после запятой на двойку в $-i$ степени и просуммирует результаты.

Полученному одномерному массиву с выборкой $\xi_i \sim U[0, 1]$ зададим нужную размерность.

```
In [82]: coin = sps.bernoulli(0.5).rvs # симметричная монета
# coin(size=10) --- реализация 10 бросков монеты

def uniform(size=1, precision=30):

    # определяем размер выборки и подкидываем монету необходимое число раз
    sample_size = np.prod(size)
    bern_sample = coin(sample_size * precision)
    bern_sample = bern_sample.reshape(sample_size, precision)

    # создаем массив отрицательных степеней двоек и осуществляем перевод в другую систему счисления
    deg_two_arr = np.power((np.ones(precision) * 2), np.arange(1,precision + 1) * -1)
    uni_sample = (bern_sample @ deg_two_arr)

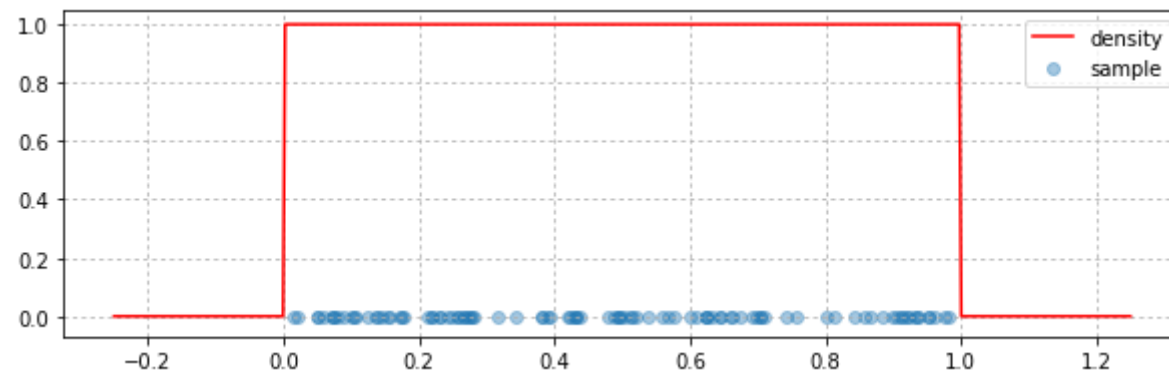
    return uni_sample.reshape(size)

# Решение в одну строчку:
# return (coin(np.prod(size) * precision).reshape(np.prod(size), precision) @ np.power((np.ones(precision)
```

Для $U[0, 1]$ сгенерируйте выборку и постройте график плотности.

```
In [83]: size = 100
grid = np.linspace(-0.25, 1.25, 500)

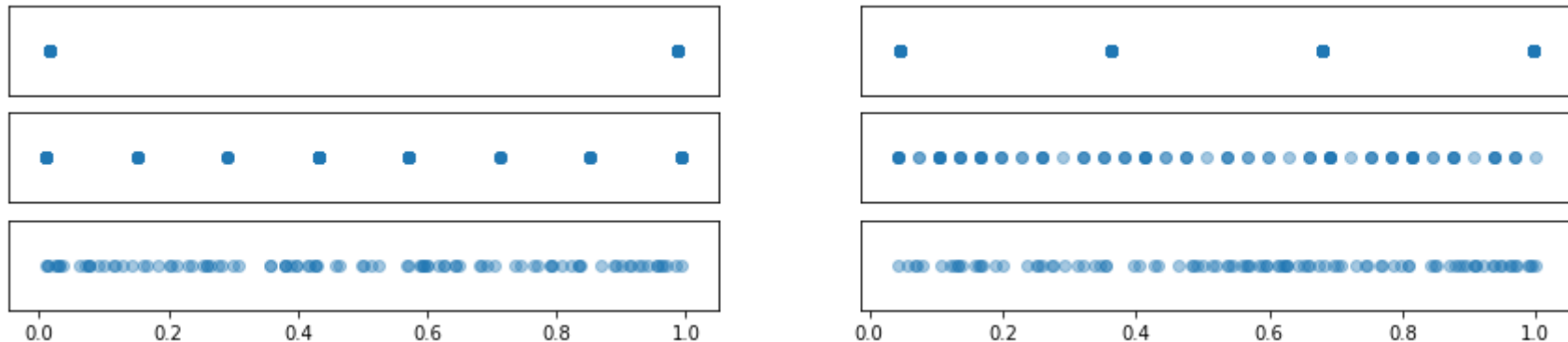
plt.figure(figsize=(10, 3))
plt.scatter(uniform(size, 50), np.zeros(size), alpha=0.4, label='sample')
plt.plot(grid, sps.uniform.pdf(grid), color='red', label='density')
plt.legend()
plt.grid(ls=':')
plt.show()
```



Исследуйте, как меняется выборка в зависимости от precision.

```
In [84]: size = 100

plt.figure(figsize=(15, 3))
for i, precision in enumerate([1, 2, 3, 5, 10, 30]):
    plt.subplot(3, 2, i + 1)
    plt.scatter(uniform(size, precision), np.zeros(size), alpha=0.4)
    plt.yticks([])
    if i < 4: plt.xticks([])
plt.show()
```



Вывод:

Чем большее число precision мы указываем в качестве параметра, тем большее количество значений принимает наша выборка.

Часть 2. Напишите функцию генерации выборки размера size (как и раньше, тут может быть tuple) из распределения $\mathcal{N}(loc, scale^2)$ с помощью преобразования Бокса-Мюллера (задача 7.12 из книги по теории вероятностей).

Для получения полного балла реализация должна быть без циклов.

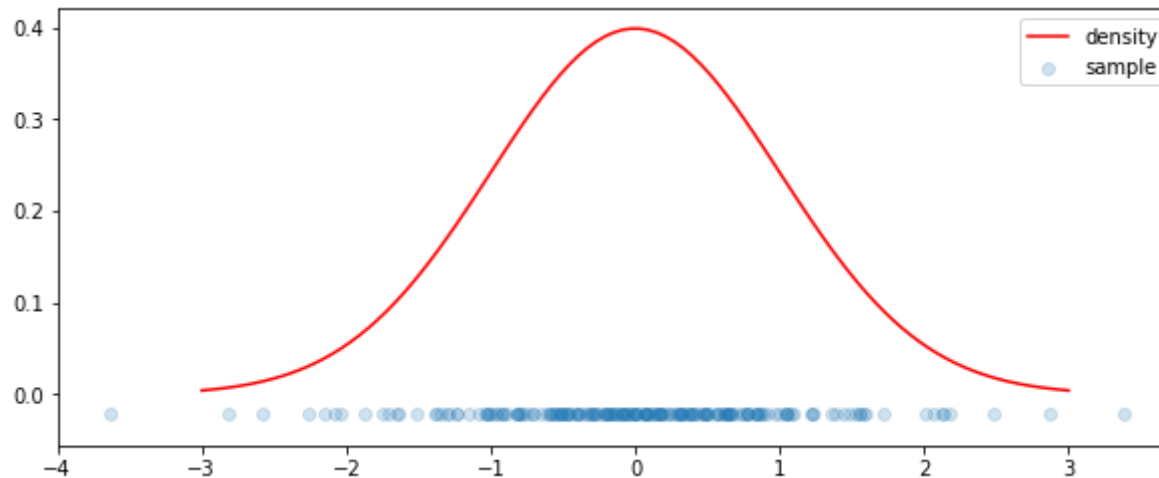
```
In [85]: def normal(size=1, loc=0, scale=1, precision=30):
    uni_sample_1 = uniform(size, precision)
    uni_sample_2 = uniform(size, precision)

    return np.cos(2 * np.pi * uni_sample_1) * np.sqrt(-2 * np.log(uni_sample_2)) * scale + loc
```

Для $\mathcal{N}(0, 1)$ сгенерируйте выборку и постройте график плотности этого распределения на отрезке $[-3, 3]$.

```
In [86]: grid = np.linspace(-3, 3, 100)
sample = normal(200, 0, 1)

plt.figure(figsize=(10, 4))
plt.scatter(sample, np.zeros(200) - 0.02, alpha=0.2, label='sample')
plt.plot(grid, sps.norm.pdf(grid), color='red', label='density')
plt.legend()
plt.show()
```



Пусть P --- некоторое распределение на $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. Числа a и σ называются параметрами сдвига и масштаба соответственно для семейства распределений $\{P_{a,\sigma} \mid a \in \mathbb{R}, \sigma \in \mathbb{R}_+\}$, где $P_{a,\sigma}(B) = P\left(\frac{B-a}{\sigma}\right)$ и $\frac{B-a}{\sigma} = \left\{\frac{x-a}{\sigma} \mid x \in B\right\}$

Вопрос: Найдите плотность $P_{a,\sigma}$, если P имеет плотность $p(x)$.

Ответ: $p_{a,\sigma}(x) = F'_{a,\sigma}(x) = F'\left(\frac{x-a}{\sigma}\right) = p\left(\frac{x-a}{\sigma}\right) * \left(\frac{x-a}{\sigma}\right)' = \frac{1}{\sigma} p\left(\frac{x-a}{\sigma}\right)$

Вопрос: Пусть P --- стандартное нормальное распределение. Выпишите параметрическое семейство распределений, параметризованное параметрами сдвига и масштаба по отношению к распределению P . Какая связь между параметрами и характеристиками распределения (например, математическое ожидание)?

Ответ: \mathcal{F}_θ --- параметрическое семейство распределений, параметризованное вектором параметров $\theta = (a, \sigma^2)$, где $\mathcal{F}_\theta = \mathcal{N}(a, \sigma^2) \mid a \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+$.

Тогда: математическое ожидание = параметр сдвига, дисперсия = параметр масштаба в квадрате. То есть между параметром сдвига и мат. ожиданием наблюдаем линейную зависимость, а между параметром масштаба и дисперсией --- квадратичную.

Постройте на одном графике разными цветами плотности стандартного нормального распределения, а так же для параметров $a = 3, \sigma = 1$ и $a = 0, \sigma = 2$. Интервал по оси икс $[-7, 7]$.

Ниже графика теми же цветами изобразите также точку a и 3σ -интервал, используя шаблон, приведенный ниже.

```
In [87]: grid = np.linspace(-7, 7, 300)

plt.figure(figsize=(10, 5))

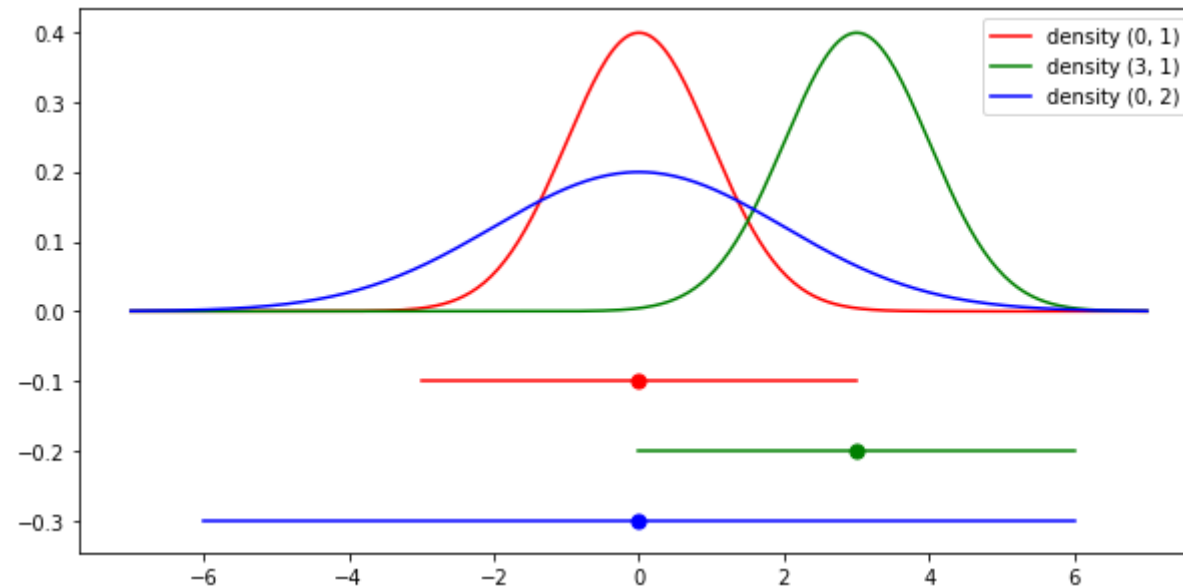
params = [ [0,1,'red'], [3,1,'green'], [0,2,'blue'] ]
for i, par in enumerate(params):
    plt.plot(grid, sps.norm.pdf(grid, loc=par[0], scale=par[1]), color=par[2],
             label='density (' + str(par[0]) + ', ' + str(par[1]) + ')')
    plt.plot([par[0] - 3 * par[1], par[0] + 3 * par[1]], [-0.1 * (i + 1), -0.1 * (i + 1)], color=par[2])
    plt.scatter(par[0], -0.1 * (i + 1), color=par[2], s=50)

# аналогичные действия, но без цикла for
# plt.plot(grid, sps.norm.pdf(grid, loc=0, scale=1), color='red', label='density (0,1)')
# plt.plot(grid, sps.norm.pdf(grid, loc=3, scale=1), color='green', label='density (3,1)')
# plt.plot(grid, sps.norm.pdf(grid, loc=0, scale=2), color='blue', label='density (0,2)')

# plt.plot([0 - 3 * 1, 0 + 3 * 1], [-0.1, -0.1], color='red')
# plt.plot([3 - 3 * 1, 3 + 3 * 1], [-0.2, -0.2], color='green')
# plt.plot([0 - 3 * 2, 0 + 3 * 2], [-0.3, -0.3], color='blue')

# plt.scatter(0, -0.1, color='red', s=50)
# plt.scatter(3, -0.2, color='green', s=50)
# plt.scatter(0, -0.3, color='blue', s=50)

plt.legend()
plt.show()
```

Вывод:

Выполненные исследования наглядно демонстрируют нам на графике Правило Трех Сигм. Оно гласит, что приблизительно с вероятностью 0,9973 значение нормально распределённой случайной величины лежит в интервале $(a - 3\sigma; a + 3\sigma)$

Часть 3. Теперь напишите функцию генерации выборки из многомерного нормального распределения с заданным вектором средних `mean` и матрицей ковариаций `cov_matrix`. Помочь в этом может теорема об эквивалентных определениях гауссовского вектора. Для извлечения квадратного корня из матрицы может пригодиться следующая функция, которая вычисляет собственные значения и векторы матрицы.

```
In [88]: from scipy.linalg import eigh
```

На этот раз достаточно, чтобы функция корректно работала в случае, когда `size` является числом.

```
In [89]: mean = np.array([1,2,3,4])
         mean.shape[0]
```

```
Out[89]: 4
```

```
In [90]: def gauss(mean, cov_matrix, size=1, precision=30):
# Преобразование типов
mean = np.array(mean)
cov_matrix = np.array(cov_matrix)

# Проверка на корректность входа
assert mean.ndim == 1 and cov_matrix.ndim == 2
assert mean.shape[0] == cov_matrix.shape[0]
assert cov_matrix.shape[0] == cov_matrix.shape[1]

eigenvals, eigenvecs = eigh(cov_matrix)

# Получим выборку из "нашего" одномерного нормального распределения
norm_sample = normal((mean.shape[0], size), precision=precision)

# Вычислим, используя формулу из курса Теории Вероятностей
matrix = eigenvecs.T @ np.diag(np.sqrt(eigenvals)) @ norm_sample

return (matrix + mean.reshape((mean.shape[0], 1))).T
```

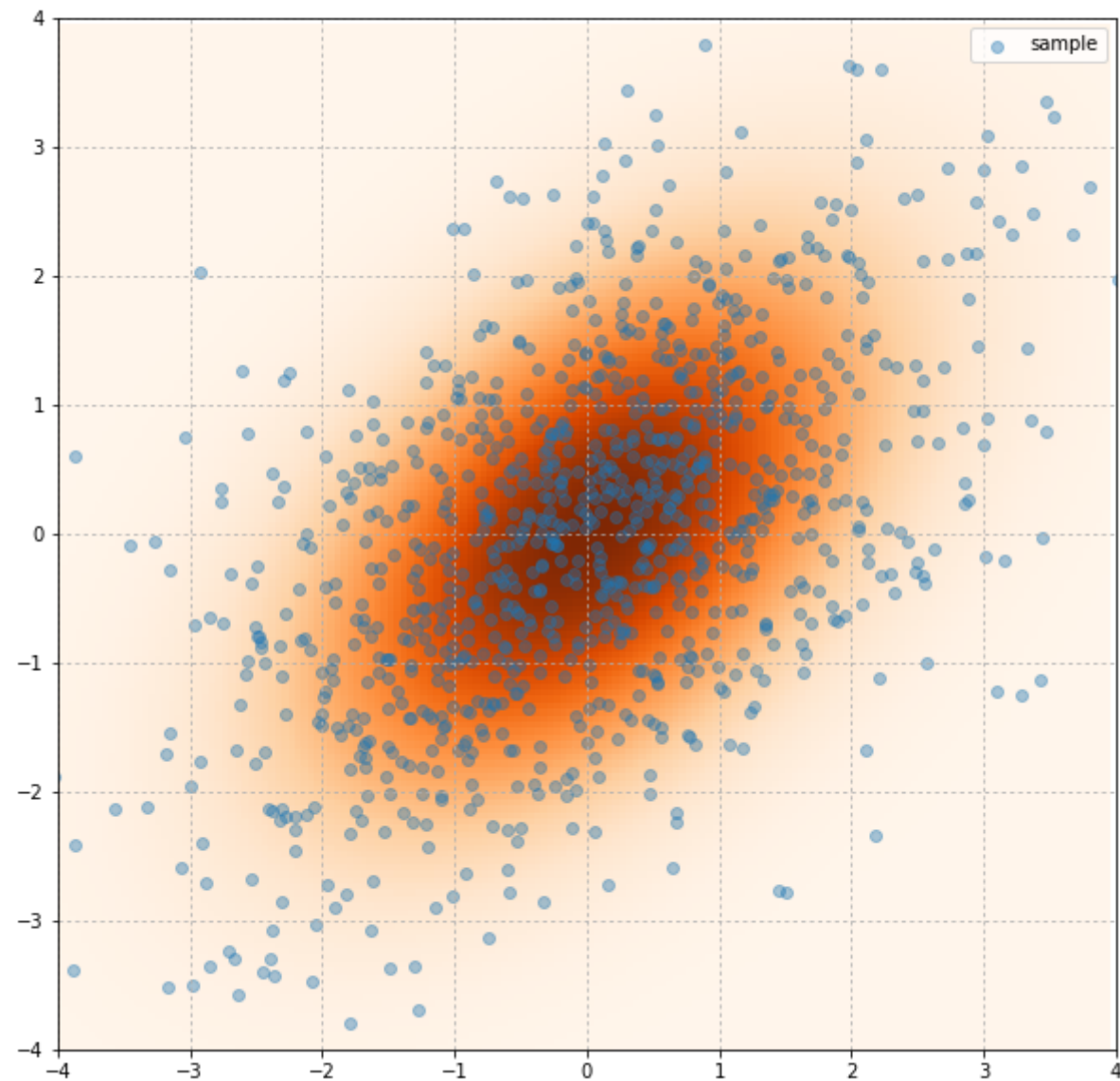
Сгенерируйте выборку размера `size` из двумерного нормального распределения с нулевым вектором средних и матрицей ковариаций $\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$. Нанесите точки выборки на график и отметьте цветом значение плотности.

В инструкциях по Питону плотность вычислялась с помощью неэффективного кода. Подумайте, как можно написать эффективный короткий код, не использующий циклы.

```
In [91]: size = 1000
mean = [0,0]
cov_matrix = [[2, 1], [1, 2]]
sample = gauss(mean, cov_matrix, size)

grid = np.mgrid[-4:4:0.05, -4:4:0.05]
density = sps.multivariate_normal.pdf(grid.T, mean, cov_matrix) # Вычисление плотности

plt.figure(figsize=(10, 10))
plt.pcolormesh(grid[0], grid[1], density, cmap='Oranges')
plt.scatter(sample[:, 0], sample[:, 1], alpha=0.4, label='sample')
plt.legend()
plt.grid(ls=':')
plt.xlim((-4, 4))
plt.ylim((-4, 4))
plt.show()
```



Вывод: на графике мы можем заметить, что значительная часть точек выборки сконцентрирована в окрестности точки с координатами из вектора mean, и что "правило трёх сигм" работает и в этом случае.

Задача 2. Вы уже научились генерировать выборку из равномерного распределения. Напишите функцию генерации выборки из экспоненциального распределения, используя результат задачи 6.9 из книги по теории вероятностей.

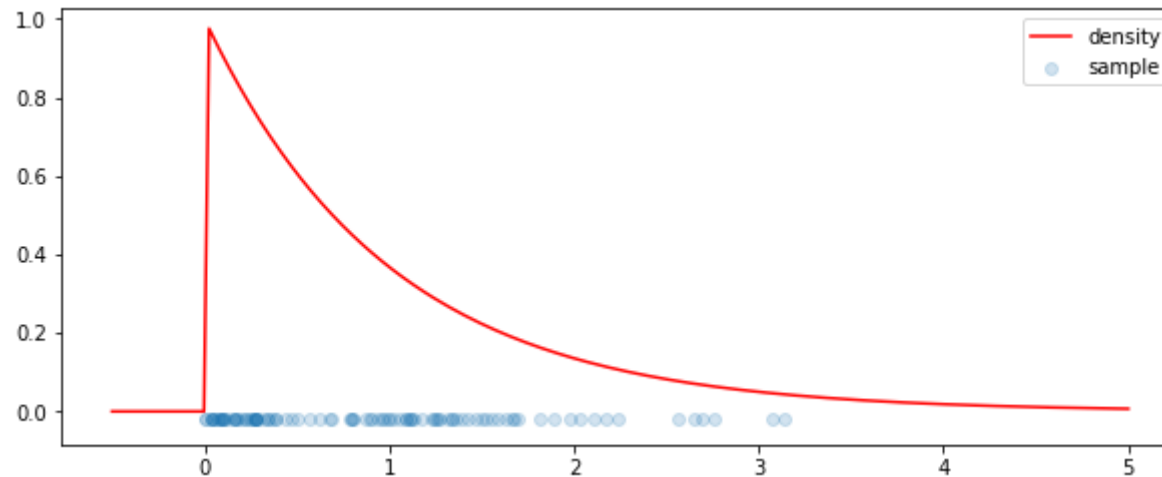
Для получения полного балла реализация должна быть без циклов, а параметр size может быть типа tuple.

Решение: Из задачи 6.9: Пусть ξ --- случайная величина с функцией распределения F . Тогда случайная величина $F(\xi)$ распределена равномерно на отрезке $[0, 1]$. Значит, зная функцию распределения экспоненциального распределения: $\eta = (1 - e^{-\lambda\xi}) \sim U[0, 1]$,
 $\xi \sim -\frac{\ln(1-\eta)}{\lambda}$.

```
In [92]: def expon(size=1, lambd=1, precision=30):  
         return -np.log(1 - uniform(size, precision)) / lambd
```

Для $Exp(1)$ сгенерируйте выборку размера 100 и постройте график плотности этого распределения на отрезке $[-0.5, 5]$.

```
In [93]: grid = np.linspace(-0.5, 5, 200)  
exp_sample = expon(100, 1)  
  
plt.figure(figsize=(10, 4))  
plt.scatter(exp_sample, np.zeros(100) - 0.02, alpha=0.2, label='sample')  
plt.plot(grid, sps.expon.pdf(grid), color='red', label='density')  
plt.legend()  
plt.show()
```



Вывод:

Используя задачу из книги по Теории Вероятностей, мы научились получать функцию распределения через равномерное распределение на отрезке $[0, 1]$.

Задача 3. Для каждого распределения постройте эмпирическую функцию распределения (ЭФР), гистограмму и ядерную оценку плотности. Сделать это помогут следующие функции.

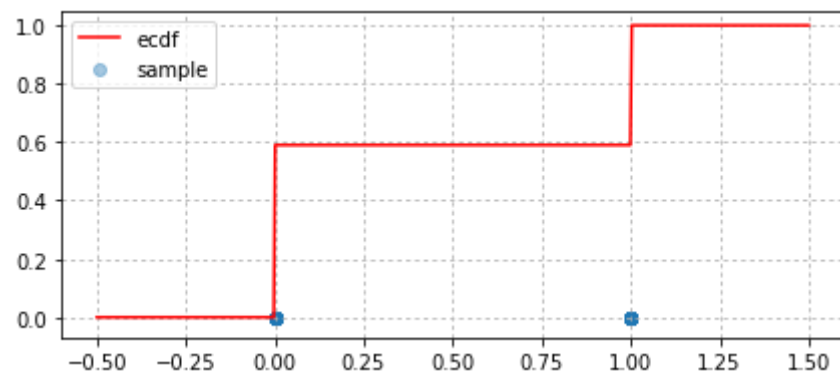
```
In [94]: from statsmodels.distributions.empirical_distribution import ECDF
        from statsmodels.nonparametric.kde import KDEUnivariate
```

1. Бернуллиевское.

Тут приведен пример построения ЭФР, просто запустите эту ячейку.

```
In [95]: sample = coin(size=100)
        ecdf = ECDF(sample)
        grid = np.linspace(-0.5, 1.5, 500)

        plt.figure(figsize=(7, 3))
        plt.scatter(sample, np.zeros(100), alpha=0.4, label='sample')
        plt.plot(grid, ecdf(grid), color='red', label='ecdf')
        plt.legend()
        plt.grid(ls=':')
        plt.show()
```



Далее, чтобы не копировать несколько раз один и тот же код, напомним некоторую функцию.

В третьей функции нужно построить ядерную оценку плотности, о которой будет рассказано на лекциях. В частности, формула была на презентации на первой лекции. Пример построения можно посмотреть тут

<http://statsmodels.sourceforge.net/0.6.0/generated/statsmodels.nonparametric.kde.KDEUnivariate.html>

(<http://statsmodels.sourceforge.net/0.6.0/generated/statsmodels.nonparametric.kde.KDEUnivariate.html>)

```

In [96]: def draw_ecdf(sample, grid, cdf=None):
    ''' По сетке grid строит графики эмпирической функции распределения
    и истинной (если она задана) для всей выборки и для 1/10 ее части.
    '''

    plt.figure(figsize=(16, 3))
    for i, size in enumerate([len(sample) // 10, len(sample)]):
        plt.subplot(1, 2, i + 1)

        plt.scatter(sample[:size], np.zeros(size), alpha=0.4, label='sample')

        if cdf is not None:
            plt.plot(grid, cdf(grid), color='green', alpha=0.3, lw=2, label='true cdf')

        plt.plot(grid, ECDF(sample[:size])(grid), color='red', label='ecdf')

        plt.legend()
        plt.grid(ls=':')
        plt.title('sample size = {}'.format(size))
    plt.show()

def draw_hist(sample, grid, pdf=None):
    ''' Строит гистограмму и по сетке grid график истинной плотности
    (если она задана) для всей выборки и для 1/10 ее части.
    '''

    plt.figure(figsize=(16, 3))
    for i, size in enumerate([len(sample) // 10, len(sample)]):
        plt.subplot(1, 2, i + 1)
        plt.hist(sample[:size], bins=20, range=(grid.min(), grid.max()), normed=True)

        if pdf is not None:
            plt.plot(grid, pdf(grid), color='green', alpha=0.3, lw=2)
    plt.show()

def draw_pdf(sample, grid, pdf=None):
    ''' По сетке grid строит графики ядерной оценки плотности
    и истинной плотности (если она задана) для всей выборки и для 1/10 ее части.
    '''

```



```

plt.figure(figsize=(16, 3))
for i, size in enumerate([len(sample) // 10, len(sample)]):
    plt.subplot(1, 2, i + 1)
    kernel_density = KDEUnivariate(sample[:size])
    kernel_density.fit()

    plt.scatter(sample[:size], np.zeros(size), alpha=0.4, label='sample')

    if pdf is not None:
        plt.plot(grid, pdf(grid), color='green', alpha=0.3, lw=2, label='true pdf')

    plt.plot(grid, kernel_density.evaluate(grid), color='red', label='kde')

plt.legend()
plt.grid(ls=':')
plt.show()

```

При использовании KDEUnivariate могут возникать разные проблемы. Можно попробовать их решить следующими способами:

1. В режиме суперюзера в файле `/usr/local/lib/python3.5/dist-packages/statsmodels/nonparametric/kdetools.py` замените строку 20 на

```
y = X[:int(m/2+1)] + np.r_[0,X[int(m/2+1):],0]*1j
```

В файле `/usr/local/lib/python3.5/dist-packages/statsmodels/nonparametric/kde.py` замените строку 327 на

```
nobs = len(X) # after trim
```

2. Попробуйте скачать с гитхаба <https://github.com/statsmodels/statsmodels/> (<https://github.com/statsmodels/statsmodels/>), установить руками. При этом должен быть установлен cython.

Можно также воспользоваться другой реализацией <http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KernelDensity.html#sklearn.neighbors.KernelDensity> (<http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KernelDensity.html#sklearn.neighbors.KernelDensity>)

Теперь примените реализованные выше функции к выборкам размера 500 из распределений $U[0, 1]$, $\mathcal{N}(0, 1)$ и $Exp(1)$. Графики (ф.р., плотностей) стройте на интервалах $(-0.2, 1.2)$, $(-3, 3)$ и $(-0.5, 5)$ соответственно.

In [97]: size = 500

```
# Равномерное на [0, 1]
uni_sample = uniform(size)
uni_grid = np.linspace(-0.2, 1.2, 500)

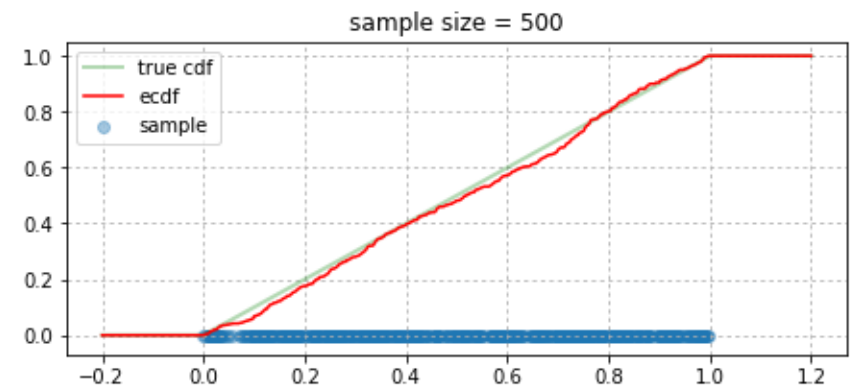
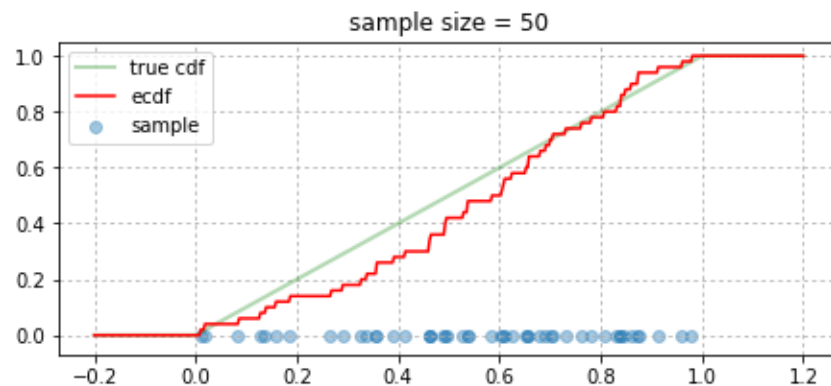
draw_ecdf(uni_sample, uni_grid, cdf=sps.uniform.cdf)
draw_hist(uni_sample, uni_grid, pdf=sps.uniform.pdf)
draw_pdf(uni_sample, uni_grid, pdf=sps.uniform.pdf)

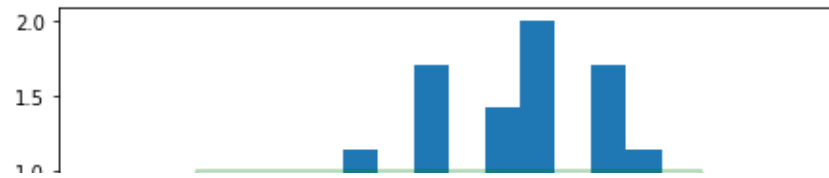
# Нормальное (0, 1)
norm_sample = normal(size)
norm_grid = np.linspace(-3, 3, 500)

draw_ecdf(norm_sample, norm_grid, cdf=sps.norm.cdf)
draw_hist(norm_sample, norm_grid, pdf=sps.norm.pdf)
draw_pdf(norm_sample, norm_grid, pdf=sps.norm.pdf)

# Экспоненциальное (1)
exp_sample = expon(size, 1)
exp_grid = np.linspace(-0.5, 5, 500)

draw_ecdf(exp_sample, exp_grid, cdf=sps.expon.cdf)
draw_hist(exp_sample, exp_grid, pdf=sps.expon.pdf)
draw_pdf(exp_sample, exp_grid, pdf=sps.expon.pdf)
```





Вывод:

На построенных нами графиках видно, что эмпирическая функция распределения и ядерная оценка плотности довольно хорошо приближают настоящие функцию распределения и функцию плотности соответственно. Причем отметим, что, чем больше размер выборки, тем точнее.

Задача 4. Сгенерируйте выборку X_1, \dots, X_{10000} из стандартного нормального распределения. Для каждого $n \leq 10000$ постройте эмпирическую функцию распределения F_n^* и посчитайте **точное** значение статистики

$$D_n = \sup_{x \in \mathbb{R}} |F_n^*(x) - F(x)|.$$

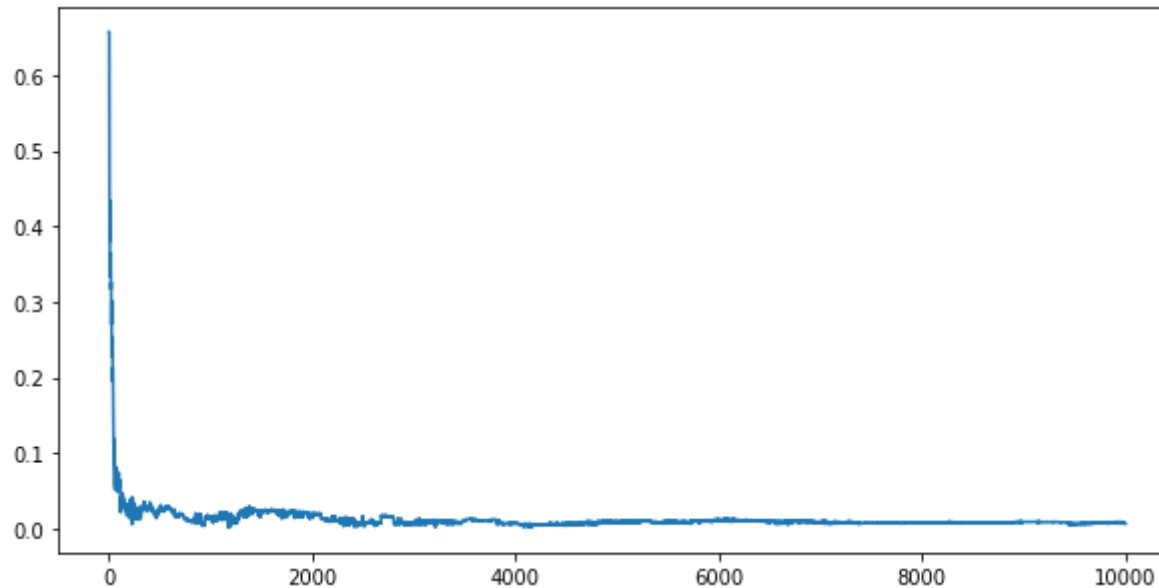
Постройте график зависимости статистики D_n от n . Верно ли, что $D_n \rightarrow 0$ и в каком смысле? Не забудьте сделать вывод.

```
In [98]: from scipy.optimize import minimize_scalar

size = 10000
norm_sample = normal(size)
Dn = []
F = sps.norm(0, 1).cdf

for n in range(1, size + 1):
    Fn = ECDF(norm_sample[:n])
    Dn.append(-minimize_scalar(lambda x: -abs(Fn(x) - F(x))).fun)
```

```
In [99]: plt.figure(figsize=(10,5))
plt.plot(np.linspace(1, size, size), Dn);
#plt.ylim([0, 0.5]);
#plt.grid(True)
plt.show()
```



Вывод:

Из построенного графика хорошо видно, что при увеличении n , действительно, $D_n \rightarrow 0$, то есть при увеличении размера выборки эмпирическая функция распределения принимает всё более и более близкие к истинной функции распределения значения.

Задача 5. 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 y -координатой должны быть нанесены точки выборки. Для экономии места стройте графики в два столбца.

Не забудьте сделать вывод.

Задача 6. В файле `countries.csv` дан список стран и территорий с указанием их площади. Нанести значения площади на график и постройте эмпирическую функцию распределения и ядерную оценку плотности. Поскольку некоторые страны слишком большие, ограничьте график по оси *икс*. Не забудьте сделать вывод.

Задача 7. Проведите небольшое исследование. Выберите случайных n человек в социальной сети. Вы можете выбирать их случайно из всех зарегистрированных в этой социальной сети, либо по какому-то *одному* критерию (укажите его). Составьте выборку X_1, \dots, X_n , где X_i --- количество друзей у i -го человека. Постройте по этой выборке эмпирическую функцию распределения. Можете ли вы сказать, какому закону подчиняется распределение количества друзей?

Выборка должна быть из не менее 30 человек, ограничений сверху нет. Вы можете также написать программу, которая будет автоматически собирать данные. Не забудьте сделать вывод.

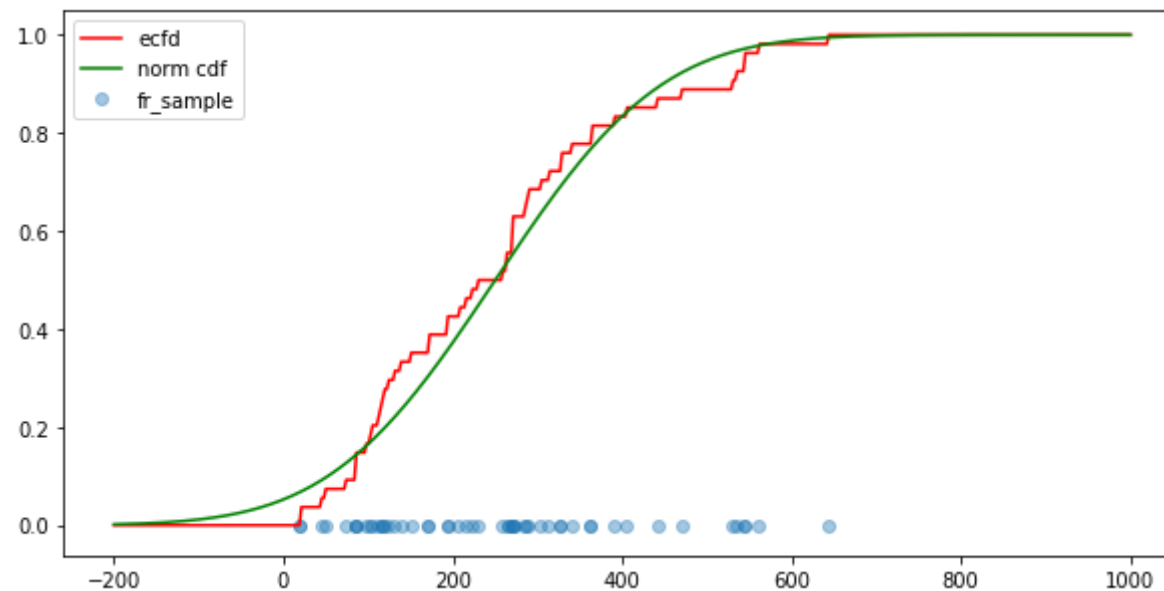
```
In [100]: fr_sample = np.array([97,74,222,20,230,207,45,105,289,391,151,284,363,340,287,442,404,530,470,215,264,171,545])
size = fr_sample.size
ecdf = ECDF(fr_sample)
grid = np.linspace(-200, 1000, 500)
plt.figure(figsize=(10,5))

plt.scatter(fr_sample, np.zeros(size), alpha=0.4, label='fr_sample')
plt.plot(grid, ecdf(grid), color='red', label='ecdf')

loc = np.mean(fr_sample)
scale = np.sqrt(np.var(fr_sample))
plt.plot(grid, sps.norm(loc=loc, scale=scale).cdf(grid), color='green', label='norm cdf')

plt.legend()
plt.show
```

Out[100]: <function matplotlib.pyplot.show>



Вывод: Взяли выборку из друзей в Вконтакте по критерию Мужчина, Москвич, Студент МФТИ (54 человека). На графике был дополнительно проведен график функции нормального распределения. Благодаря этому хорошо видно, что функция распределения практически совпадает с функцией распределения нормального распределения. Значит, сделаем вывод, что количество друзей

подчиняется нормальному закону распределения.