

Введение в анализ данных

Домашнее задание 1. Numpy, matplotlib, scipy.stats

Правила:

- Дедлайн **25 марта 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 "[номер группы] Фамилия Имя - Задание 1". Квадратные скобки обязательны.
- Прислать нужно ноутбук и его pdf-версию (без архивов). Названия файлов должны быть такими: `1.N.ipynb` и `1.N.pdf`, где `N` -- ваш номер из таблицы с оценками. *pdf-версию можно сделать с помощью Ctrl+P. Пожалуйста, посмотрите ее полностью перед отправкой. Если что-то существенное не напечатается в pdf, то баллы могут быть снижены.*
- 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.
- Пропущенные описания принимаемых аргументов дописать на русском.
- Если код будет не понятен проверяющему, оценка может быть снижена.

Баллы за задание:

Легкая часть (достаточно на "хор"):

- Задача 1.1 -- 3 балла
- Задача 1.2 -- 3 балла
- Задача 2 -- 3 балла

Сложная часть (необходимо на "отл"):

- Задача 1.3 -- 3 балла
- Задача 3.1 -- 3 балла
- Задача 3.2 -- 3 балла
- Задача 3.3 -- 3 балла
- Задача 4 -- 4 балла

Баллы за разные части суммируются отдельно, нормируются впоследствии также отдельно. Иначе говоря, 1 балл за легкую часть может быть не равен 1 баллу за сложную часть.

In [19]:

```
import numpy as np
import scipy.stats as sps

import matplotlib.pyplot as plt
import matplotlib.cm as cm
from mpl_toolkits.mplot3d import Axes3D
import ipywidgets as widgets
import pandas

import typing

%matplotlib inline
```

Легкая часть: генерация

В этой части другие библиотеки использовать запрещено. Шаблоны кода ниже менять нельзя.

Задача 1

Имеется симметричная монета. Напишите функцию генерации независимых случайных величин из нормального и экспоненциального распределений с заданными параметрами.

In [20]:

```
# Эта ячейка -- единственная в задаче 1, в которой нужно использовать
# библиотечную функцию для генерации случайных чисел.
# В других ячейках данной задачи используйте функцию coin.

# симметричная монета
#coin = sps.norm(loc=0, scale=1).rvs
#coin = sps.expon(scale=1).rvs
coin = sps.bernoulli(p = 1/2).rvs
```

Проверьте работоспособность функции, сгенерировав 10 бросков симметричной монеты.

In [21]:

```
coin(size=10)
```

Out[21]:

```
array([0, 1, 0, 0, 1, 0, 1, 1, 1, 1])
```

Часть 1. Напишите сначала функцию генерации случайных величин из 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 $\xi = 0, \xi_1 \xi_2 \xi_3 \dots \xi_n$.

Нужно реализовать функцию так, чтобы она могла принимать на вход в качестве параметра `size` как число, так и 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. Аргумент `precision` отвечает за число n .

In [22]:

```
def uniform(size=1, precision=30):  
    count_elements = np.array([size]).prod()  
    arr_of_dec = np.array(coin(count_elements*precision)).reshape(count_elements, p  
    matr_of_dec = 0.5 ** np.arange(1, precision + 1, 1)  
    return np.array(arr_of_dec[:]) @ matr_of_dec
```

In [23]:

```
print(uniform((2, 2), 5))
```

```
[0.875    0.71875 0.09375 0.34375]
```

Для $U[0, 1]$ сгенерируйте 200 независимых случайных величин, постройте график плотности на отрезке $[-0.25, 1.25]$, а также гистограмму по сгенерированным случайным величинам.

In [24]:

```

size = 200
grid = np.linspace(-0.25, 1.25, 500)
sample = uniform(size, 50)

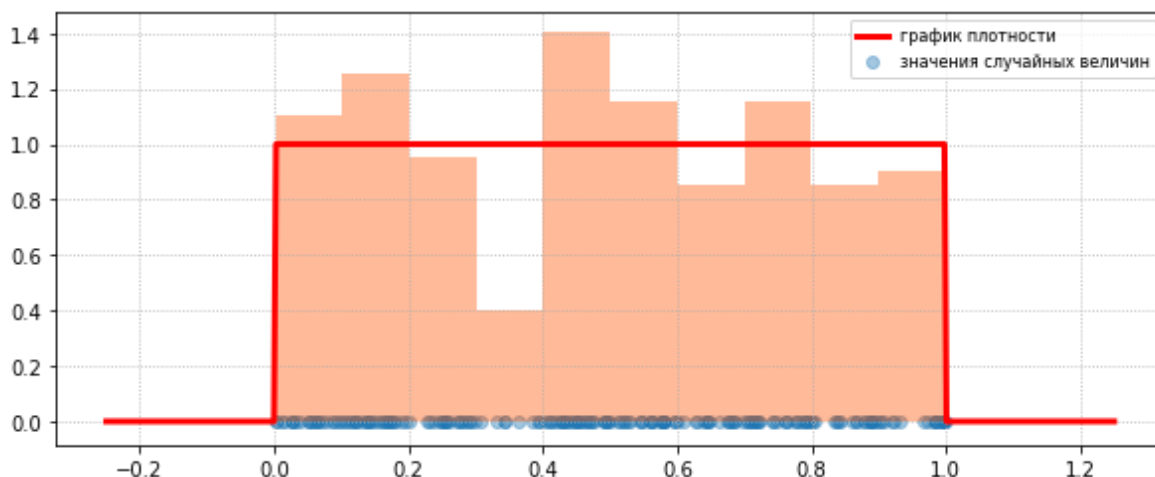
# Отрисовка графика
plt.subplots(figsize=(10, 4))#<определите график размера 10 на 4>

# отображаем значения случайных величин полупрозрачными точками
points = plt.scatter(
    sample,
    np.zeros(size),
    alpha=0.4,
    label='значения случайных величин'
)

# по точкам строим нормированную полупрозрачную гистограмму
histogram = plt.hist(
    sample,
    bins=10,
    density=True,
    alpha=0.4,
    color='#FF5300'
)

# рисуем график плотности
line = plt.plot(
    grid,
    sps.uniform.pdf(grid),
    color='#FF0000',
    linewidth=3,
    label='график плотности'
)
plt.legend(fontsize=8, loc=1)
plt.grid(ls=':')
plt.show()

```



Исследуйте, как меняются значения случайных величин в зависимости от `precision` .

In [25]:

```

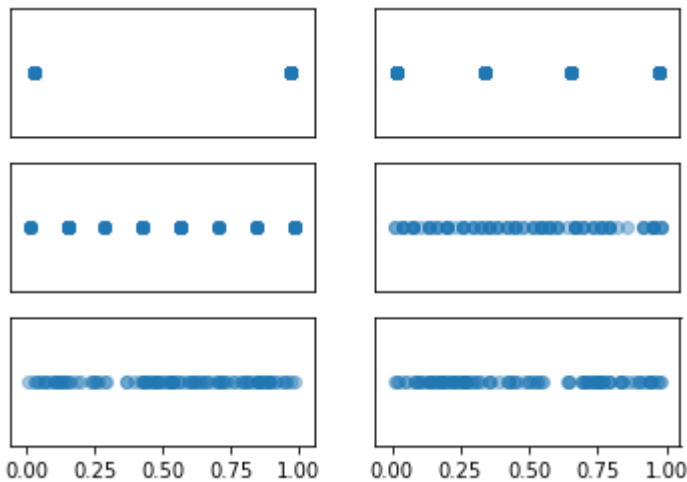
size = 100

plt.plot(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:**

Используя генератор случайной выборки из непрерывного равномерного распределения, мы можем построить генератор выборки любого непрерывного равномерного распределения.

Часть 2. Напишите функцию генерации случайных величин в количестве `size` штук (как и раньше, тут может быть `tuple`) из распределения $\mathcal{N}(loc, scale^2)$ с помощью преобразования Бокса-Мюллера, которое заключается в следующем. Пусть ξ и η -- независимые случайные величины, равномерно распределенные на $(0, 1]$. Тогда случайные величины $X = \cos(2\pi\xi)\sqrt{-2\ln\eta}$, $Y = \sin(2\pi\xi)\sqrt{-2\ln\eta}$ являются независимыми нормальными $\mathcal{N}(0, 1)$.

Реализация должна быть без циклов. Желательно использовать как можно меньше бросков монеты.

In [26]:

```

def normal(size=1, loc=0, scale=1, precision=30):
    return ((np.sin(2*np.pi*uniform(size, precision))*np.sqrt(-2)*np.log(uniform(s

```

Для $\mathcal{N}(0, 1)$ сгенерируйте 200 независимых случайных величин, постройте график плотности на отрезке $[-3, 3]$, а также гистограмму по сгенерированным случайным величинам.

In [30]:

```
size = 200
grid = np.linspace(-3, 3, 500)

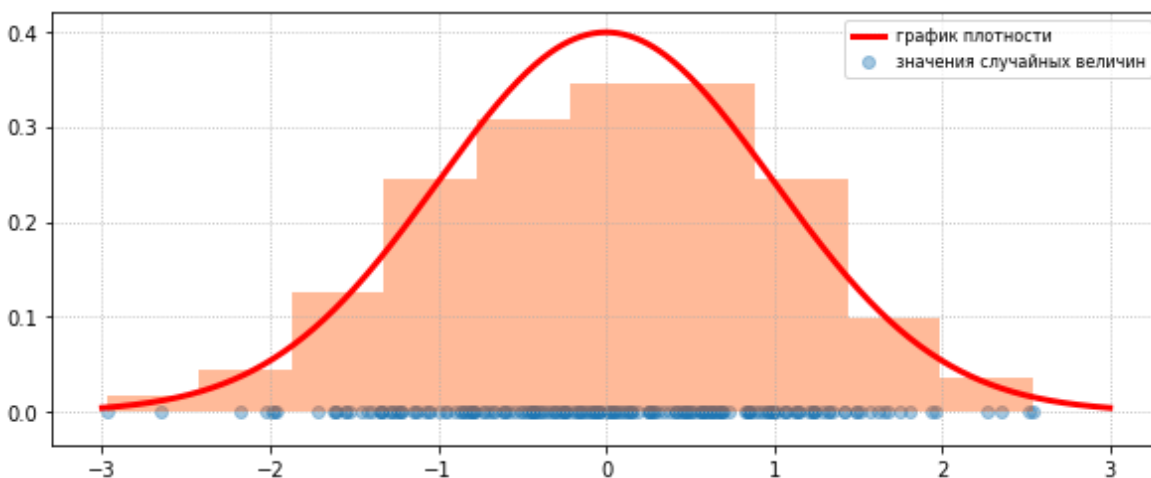
sample = normal(size, 0, 1, 30)

plt.figure(figsize=(10,4))

# отображаем значения случайных величин
points = plt.scatter(
    sample,
    np.zeros(size),
    alpha=0.4,
    label='значения случайных величин'
)

# по точкам строим нормированную полупрозрачную гистограмму
histogram = plt.hist(
    sample,
    bins=10,
    density=True,
    alpha=0.4,
    color='#FF5300'
)

# рисуем график плотности
line = plt.plot(
    grid,
    sps.norm.pdf(grid),
    color='#FF0000',
    linewidth=3,
    label='график плотности'
)
plt.legend(fontsize=8, loc=1)
plt.grid(ls=':')
plt.show()
```



Сложная часть: генерация

Часть 3. Вы уже научились генерировать выборку из равномерного распределения. Напишите функцию генерации выборки из экспоненциального распределения, используя из теории вероятностей:

Если ξ --- случайная величина, имеющая абсолютно непрерывное распределение, и F --- ее функция распределения, то случайная величина $F(\xi)$ имеет равномерное распределение на $[0, 1]$.

Какое преобразование над равномерной случайной величиной необходимо совершить?

$F(x) = -\frac{1}{\lambda} \ln(1 - x)$ --- обратная функция экспоненциального распределения, значит необходимы такое преобразование: $-\frac{1}{\lambda} \ln(1 - \xi)$

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

In [33]:

```
def expon(size=1, lambd=1, precision=30):  
    return ((-1) / lambd * np.log(1 - uniform(size, precision))).reshape(size)
```

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

In [42]:

```

size = 100
grid = np.linspace(-0.5, 5, 500)

sample = expon(size, lambd=1, precision=50)

plt.figure(figsize=(10,4))

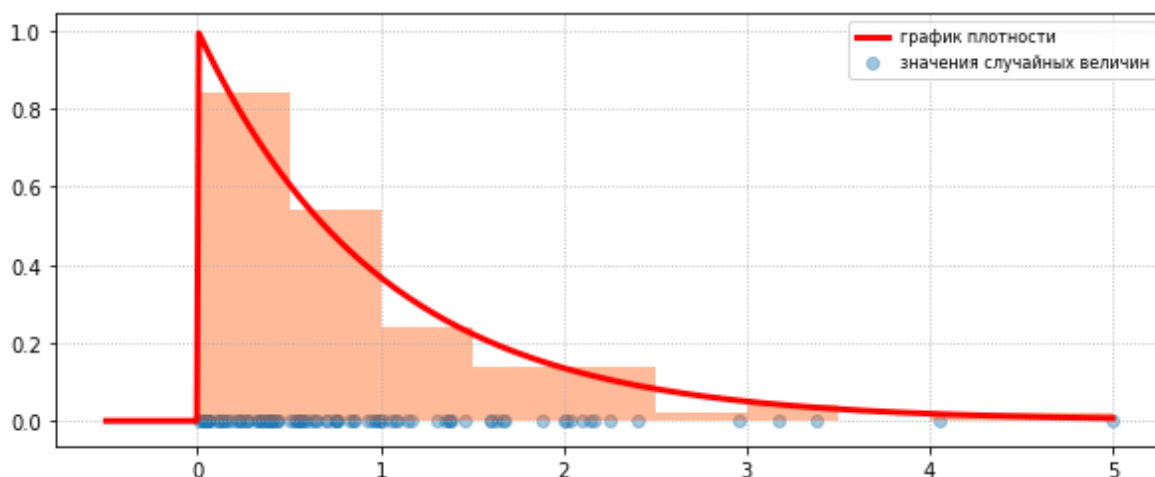
# отображаем значения случайных величин
points = plt.scatter(
    sample,
    np.zeros(size),
    alpha=0.4,
    label='значения случайных величин'
)

# по точкам строим нормированную полупрозрачную гистограмму
histogram = plt.hist(
    sample,
    bins=10,
    density=True,
    alpha=0.4,
    color='#FF5300'
)

# рисуем график плотности
line = plt.plot(
    grid,
    sps.expon.pdf(grid),
    color='#FF0000',
    linewidth=3,
    label='график плотности'
)

plt.legend(fontsize=8, loc=1)
plt.grid(ls=':')
plt.show()

```



Вывод по задаче:

Используя генератор случайной выборки из непрерывного равномерного распределения, мы можем построить генератор выборки любого непрерывного равномерного распределения, используя функцию обратную к функции распределения случайной величины, методом обратного преобразования.

Легкая часть: матричное умножение

Задача 2

Напишите функцию, реализующую матричное умножение. При вычислении разрешается создавать объекты размерности три. Запрещается пользоваться функциями, реализующими матричное умножение (`numpy.dot`, операция `@`, операция умножения в классе `numpy.matrix`). Разрешено пользоваться только простыми векторно-арифметическими операциями над `numpy.array`, а также преобразованиями осей. Авторское решение занимает одну строчку.

In [11]:

```
def matrix_multiplication(A, B):  
    return np.apply_along_axis(lambda my_str: my_str.sum(), 1, np.apply_along_axis(
```

Проверьте правильность реализации на случайных матрицах. Должен получиться ноль.

In [12]:

```
A = sps.uniform.rvs(size=(10, 20))  
B = sps.uniform.rvs(size=(20, 30))  
np.abs(matrix_multiplication(A, B) - A @ B).sum()
```

Out[12]:

1.3589129821411916e-13

На основе опыта: вот в таком стиле многие из вас присылали бы нам свои работы, если не стали бы делать это задание :)

In [13]:

```
def stupid_matrix_multiplication(A, B):  
    C = [[0 for j in range(len(B[0]))] for i in range(len(A))]  
    for i in range(len(A)):  
        for j in range(len(B[0])):  
            for k in range(len(B)):  
                C[i][j] += A[i][k] * B[k][j]  
    return C
```

Проверьте, насколько быстрее работает ваш код по сравнению с неэффективной реализацией `stupid_matrix_multiplication`. Эффективный код должен работать почти в 200 раз быстрее. Для примера посмотрите также, насколько быстрее работают встроенные `numpy`-функции.

In [14]:

```
A = sps.uniform.rvs(size=(400, 200))
B = sps.uniform.rvs(size=(200, 300))

%time C1 = matrix_multiplication(A, B)
%time C2 = A @ B # python 3.5
%time C3 = np.matrix(A) * np.matrix(B)
%time C4 = stupid_matrix_multiplication(A, B)
%time C5 = np.einsum('ij,jk->ik', A, B)
```

```
CPU times: user 640 ms, sys: 76 ms, total: 716 ms
Wall time: 716 ms
CPU times: user 4.24 ms, sys: 0 ns, total: 4.24 ms
Wall time: 2.33 ms
CPU times: user 4.52 ms, sys: 0 ns, total: 4.52 ms
Wall time: 1.36 ms
CPU times: user 18.3 s, sys: 67.3 ms, total: 18.4 s
Wall time: 18.3 s
CPU times: user 9.51 ms, sys: 0 ns, total: 9.51 ms
Wall time: 9.21 ms
```

Ниже для примера приведена полная реализация функции. Вас мы, конечно, не будем требовать проверять входные данные на корректность, но документации к функциям нужно писать.

In [17]:

```
def matrix_multiplication(A, B):
    '''Возвращает матрицу, которая является результатом
    матричного умножения матриц A и B.

    ...

    # Если A или B имеют другой тип, нужно выполнить преобразование типов
    A = np.array(A)
    B = np.array(B)

    # Проверка данных входных данных на корректность
    assert A.ndim == 2 and B.ndim == 2, 'Размер матриц не равен 2'
    assert A.shape[1] == B.shape[0], \
        ('Матрицы размерностей {} и {} неперемножаемы'.format(A.shape, B.shape))

    C = np.apply_along_axis(lambda my_str: my_str.sum(), 1, np.apply_along_axis(lambda
    return C
```

Сложная часть: броуновское движение

Задача 3

Познавательная часть задачи (не пригодится для решения задачи)

Абсолютное значение скорости движения частиц идеального газа, находящегося в состоянии ТД-равновесия, есть случайная величина, имеющая распределение Максвелла и зависящая только от одного термодинамического параметра — температуры T .

В общем случае плотность вероятности распределения Максвелла для n -мерного пространства имеет вид:

$$p(v) = C e^{-\frac{mv^2}{2kT}} v^{n-1},$$

где $v \in [0, +\infty)$, а константа C находится из условия нормировки $\int_0^{+\infty} p(v) dv = 1$.

Физический смысл этой функции таков: вероятность того, что скорость частицы входит в промежуток $[v_0, v_0 + dv]$, приближённо равна $p(v_0)dv$ при достаточно малом dv . Тут надо оговориться, что математически корректное утверждение таково:

$$\lim_{dv \rightarrow 0} \frac{P\{v \mid v \in [v_0, v_0 + dv]\}}{dv} = p(v_0).$$

Поскольку это распределение не ограничено справа, определённая доля частиц среды приобретает настолько высокие скорости, что при столкновении с макрообъектом может происходить заметное отклонение как траектории, так и скорости его движения.

Мы предполагаем идеальность газа, поэтому компоненты вектора скорости частиц среды v_i можно считать независимыми нормально распределёнными случайными величинами, т.е.

$$v_i \sim \mathcal{N}(0, s^2),$$

где s зависит от температуры и массы частиц и одинаково для всех направлений движения.

При столкновении макрообъекта с частицами среды происходит перераспределение импульса в соответствии с законами сохранения энергии и импульса, но в силу большого числа подобных событий за единицу времени, моделировать их напрямую достаточно затруднительно. Поэтому для выполнения этого ноутбука сделаем следующие предположения:

- Приращение компоненты координаты броуновской частицы за фиксированный промежуток времени (или за шаг) Δt имеет вид $\Delta x_i \sim \mathcal{N}(0, \sigma^2)$.
- σ является конкретным числом, зависящим как от Δt , так и от параметров броуновской частицы и среды.
- При этом σ не зависит ни от координат, ни от текущего вектора скорости броуновской частицы.

Если говорить формальным языком, в этом ноутбуке мы будем моделировать [Винеровский случайный процесс](https://ru.wikipedia.org/wiki/%D0%92%D0%B8%D0%BD%D0%B5%D1%80%D0%BE%D0%B2%D1%81%D0%F) с фиксированным шагом.

Задание

1. Разработать функцию симуляции броуновского движения

Функция должна вычислять приращение координаты частицы на каждом шаге как $\Delta x_{ijk} \sim \mathcal{N}(0, \sigma^2) \forall i, j, k$, где i — номер частицы, j — номер координаты, а k — номер шага. Функция принимает в качестве аргументов:

- Параметр σ ;
- Количество последовательных изменений координат (шагов), приходящихся на один процесс;
- Число процессов для генерации (количество различных частиц);

- Количество пространственных измерений для генерации процесса.

Возвращаемое значение:

- 3-х мерный массив `result`, где `result[i,j,k]` — значение j -й координаты i -й частицы на k -м шаге.

Общее требование

- Считать, что все частицы в начальный момент времени находятся в начале координат.

Что нужно сделать

- Реализовать функцию для произвольной размерности, не используя циклы.
- Дописать проверки типов для остальных аргументов.

Обратите внимание на использование аннотаций для типов аргументов и возвращаемого значения функции. В новых версиях Питона подобные возможности синтаксиса используются в качестве подсказок для программистов и статических анализаторов кода, и никакой дополнительной функциональности не добавляют.

Например, `typing.Union[int, float]` означает "или `int`, или `float`".

Что может оказаться полезным

- Генерация нормальной выборки: `scipy.stats.norm`. [Ссылка](https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.norm.html) (<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.norm.html>)
- Кумулятивная сумма: метод `cumsum` у `np.ndarray`. [Ссылка](https://docs.scipy.org/doc/numpy/reference/generated/numpy.ndarray.cumsum.html) (<https://docs.scipy.org/doc/numpy/reference/generated/numpy.ndarray.cumsum.html>)

In []:

```
def generate_brownian(sigma: typing.Union[int, float] = 1,
                      *,
                      n_proc: int = 10,
                      n_dims: int = 2,
                      n_steps: int = 100) -> np.ndarray:
    """
    :param sigma: стандартное отклонение нормального распределения,
                  генерирующего пошаговые смещения координат
    :param n_proc: <ДОПИСАТЬ>
    :param n_dims: <ДОПИСАТЬ>
    :param n_steps: <ДОПИСАТЬ>

    :return: np.ndarray размера (n_proc, n_dims, n_steps), содержащий
             на позиции [i,j,k] значение j-й координаты i-й частицы
             на k-м шаге.
    """
    if not np.issubdtype(type(sigma), np.number):
        raise TypeError("Параметр 'sigma' должен быть числом")
    # <ДОПИСАТЬ ПРОВЕРКИ ТИПОВ>

    return <...>
```

Символ `*` в заголовке означает, что все аргументы, объявленные после него, необходимо определять только по имени.

Например,

```
generate_brownian(323, 3)           # Ошибка
generate_brownian(323, n_steps=3)   # ОК
```

При проверке типов остальных аргументов, по аналогии с `np.number`, можно использовать `np.integer`. Конструкция `np.issubdtype(type(param), np.number)` используется по причине того, что стандартная питоновская проверка `isinstance(sigma, (int, float))` не будет работать для `numpy`-чисел `int64`, `int32`, `float64` и т.д.

In []:

```
brownian_2d = generate_brownian(2, n_steps=12000, n_proc=500, n_dims=2)
assert brownian_2d.shape == (500, 2, 12000)
```

2. Визуализируйте траектории для 9-ти первых броуновских частиц

Что нужно сделать

- Нарисовать 2D-графики для `brownian_2d`.
- Нарисовать 3D-графики для `brownian_3d = generate_brownian(2, n_steps=12000, n_proc=500, n_dims=3)`.

Общие требования

- Установить соотношение масштабов осей, равное 1, для каждого из подграфиков.

Что может оказаться полезным

- [Тьюториал \(https://matplotlib.org/devdocs/gallery/subplots_axes_and_figures/subplots_demo.html\)](https://matplotlib.org/devdocs/gallery/subplots_axes_and_figures/subplots_demo.html) по построению нескольких графиков на одной странице.
- Метод `plot` у `AxesSubplot` (переменная `ax` в цикле ниже).
- Метод `set_aspect` у `AxesSubplot`.

In []:

```
fig, axes = plt.subplots(3, 3, figsize=(18, 10))
fig.suptitle('Траектории броуновского движения', fontsize=20)

for ax, (xs, ys) in zip(axes.flat, brownian_2d):
    <...>
    pass
```

3. Постройте график среднего расстояния частицы от начала координат в зависимости от времени (шага)

- Постройте для `n_dims` от 1 до 5 включительно.
- Кривые должны быть отрисованы на одном графике. Каждая кривая должна иметь легенду.
- Для графиков подписи к осям обязательны.

Вопросы

- Как вы думаете, какой функцией может описываться данная зависимость?
- Сильно ли её вид зависит от размерности пространства?

- Можно ли её линеаризовать? Если да, нарисуйте график с такими же требованиями.

In []:

```
plt.figure(figsize=(12, 6))

for n_dims in range(1, 6):
    plt.plot(
        <...>
        label=f'Размерность: {n_dims}'
    )

plt.ylabel('Ср. раст. частицы от нач. координат')
plt.xlabel('Шаг')
plt.legend(loc='best')
plt.tight_layout()
plt.show()
```

Сложная часть: визуализация распределений

Задача 4

В этой задаче вам нужно исследовать свойства дискретных распределений и абсолютно непрерывных распределений.

Для перечисленных ниже распределений нужно

- 1) На основе графиков дискретной плотности (функции массы) для различных параметров пояснить, за что отвечает каждый параметр.
- 2) Сгенерировать набор независимых случайных величин из этого распределения и построить по ним гистограмму.
- 3) Сделать выводы о свойствах каждого из распределений.

Распределения:

- Бернулли
- Биномиальное
- Равномерное
- Геометрическое

Для выполнения данного задания можно использовать код с лекции.

In []:

```
<...>
```