

In [1]:

```
1 import numpy as np
2 import scipy.stats as sps
3 import matplotlib.pyplot as plt
4
5 %matplotlib inline
```

In [2]:

```
1 import warnings
2 warnings.filterwarnings("ignore", category=RuntimeWarning)
```

In [3]:

```
1 red = '#FF3300'
2 blue = '#0099CC'
3 green = '#00CC66'
4 orange = 'orange'
```

Корреляционный анализ

Коэффициенты корреляции

[pearsonr](#)

(<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.pearsonr.html#scipy.stats.pearsonr>),

[spearmanr](#)

(<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.spearmanr.html#scipy.stats.spearmanr>),

[kendalltau](#)

(<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.kendalltau.html#scipy.stats.kendalltau>) -

коэффициенты корреляции Пирсона, Спирмана, Кендалла.

Общий принцип: $f(x_1, x_2) = (\text{correlation}, \text{pvalue})$

Упорядоченные по возрастанию наборы данных

In [4]:

```
1 x1, x2 = np.arange(5), np.arange(5) + 6
2 print('Выборки:', x1, x2)
3 sps.pearsonr(x1, x2), sps.spearmanr(x1, x2), sps.kendalltau(x1, x2)
```

Выборки: [0 1 2 3 4] [6 7 8 9 10]

Out[4]:

```
((1.0, 0.0),
 SpearmanrResult(correlation=0.9999999999999999, pvalue=1.404265422054
 3672e-24),
 KendalltauResult(correlation=0.9999999999999999, pvalue=0.016666666666
 6666666))
```

Одна по возрастанию, другая по убыванию.

In [5]:

```
1 print('Выборки:', x1, -x2)
2 sps.pearsonr(x1, -x2), sps.spearmanr(x1, -x2), sps.kendalltau(x1, -x2)
```

Выборки: [0 1 2 3 4] [-6 -7 -8 -9 -10]

Out[5]:

```
((-1.0, 0.0),
 SpearmanrResult(correlation=-0.9999999999999999, pvalue=1.40426542205
 43672e-24),
 KendalltauResult(correlation=-0.9999999999999999, pvalue=0.0166666666
 66666666))
```

Корреляция с каким-то другим набором

In [6]:

```
1 x2 = [4, 8, 2, 5, 1]
2 print('Выборки:', x1, x2)
3 sps.pearsonr(x1, x2), sps.spearmanr(x1, x2), sps.kendalltau(x1, x2)
```

Выборки: [0 1 2 3 4] [4, 8, 2, 5, 1]

Out[6]:

```
((-0.5196152422706632, 0.36951722839383205),
 SpearmanrResult(correlation=-0.49999999999999994, pvalue=0.3910022189
 5577053),
 KendalltauResult(correlation=-0.39999999999999997, pvalue=0.483333333
 33333334))
```

Если у одного набора поменять знак, то коэффициенты корреляции сменят знак

In [7]:

```
1 print('Выборки:', -x1, x2)
2 sps.pearsonr(-x1, x2), sps.spearmanr(-x1, x2), sps.kendalltau(-x1, x2)
```

Выборки: [0 -1 -2 -3 -4] [4, 8, 2, 5, 1]

Out[7]:

```
((0.5196152422706632, 0.36951722839383205),
 SpearmanrResult(correlation=0.49999999999999994, pvalue=0.39100221895
 577053),
 KendalltauResult(correlation=0.39999999999999997, pvalue=0.483333333
 33333334))
```

Некоторые вспомогательные функции для отрисовки графиков.

In [8]:

```
1 ▾ def autolabel(rects):
2     """
3     Attach a text label above each bar displaying its height
4     """
5     for rect in rects:
6         height = rect.get_height()
7         y = rect.get_y()
8         plt.text(rect.get_x() + rect.get_width()/2.,
9                 -0.2 if y >= 0 else 0.07,
10                '%.2f' % (height * (-1 if y < 0 else 1)),
11                ha='center', va='bottom', fontsize=16)
```

In [9]:

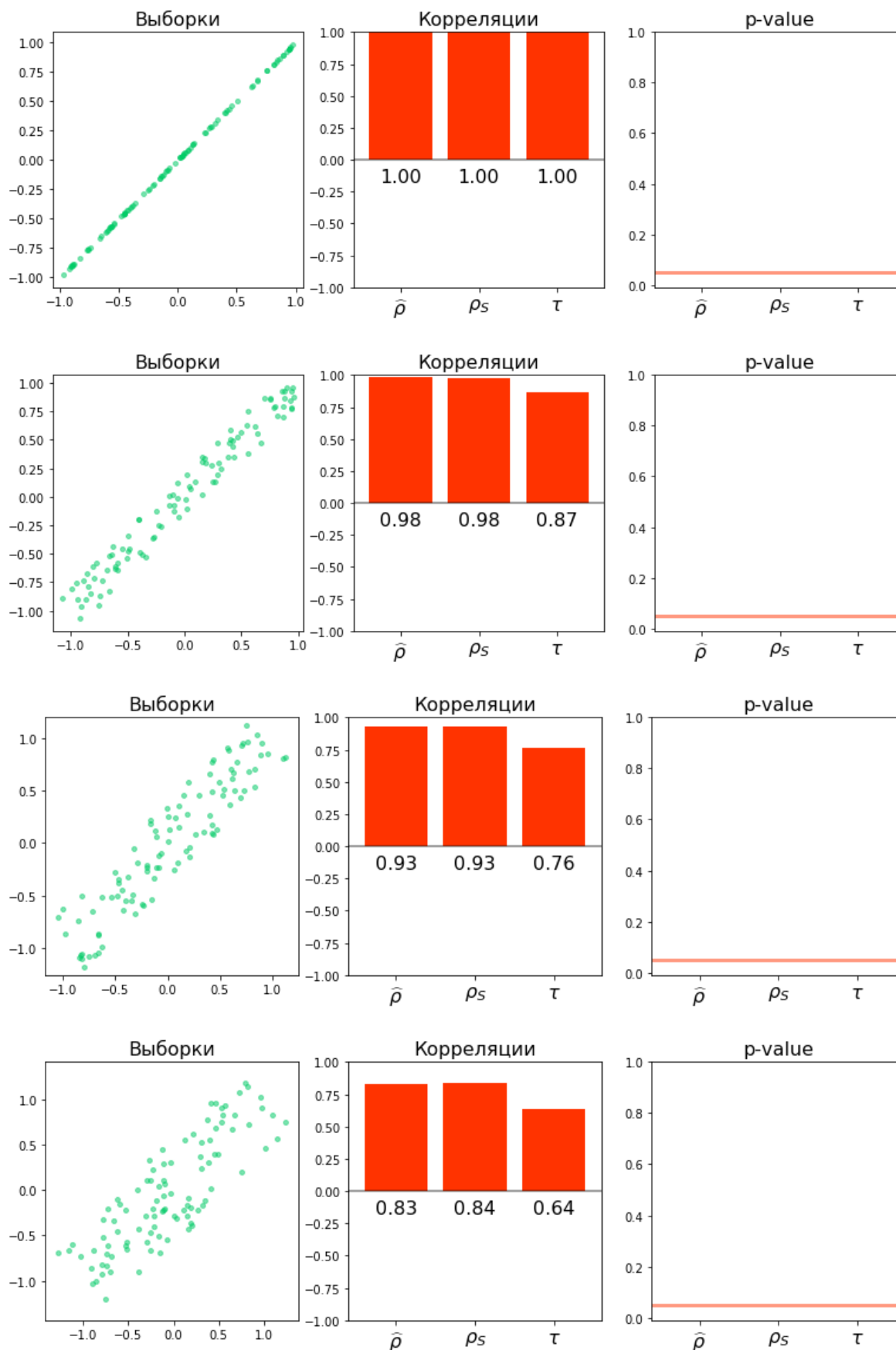
```
1 ▾ def draw_graphics(x1, x2):
2     r, pr = sps.pearsonr(x1, x2)
3     rho, prho = sps.spearmanr(x1, x2)
4     tau, ptau = sps.kendalltau(x1, x2)
5     colors = [(red if pr < 0.05 else blue),
6              (red if prho < 0.05 else blue),
7              (red if ptau < 0.05 else blue)]
8     titles = ['$\\widehat{\\rho}$', '$\\rho_S$', '$\\tau$']
9
10    plt.figure(figsize=(13, 4))
11
12    plt.subplot(1, 3, 1)
13    plt.scatter(x1, x2, alpha=0.5, s=15, color=green)
14    plt.axis('equal')
15    plt.title('Выборки', fontsize=16)
16
17    plt.subplot(1, 3, 2)
18    rects = plt.bar([1, 2, 3], [r, rho, tau], color=colors)
19    plt.hlines(0, 0.4, 3.6, color='black', alpha=0.5)
20    autolabel(rects)
21    plt.xticks([1, 2, 3], titles, fontsize=16)
22    plt.title('Корреляции', fontsize=16)
23    plt.xlim((0.4, 3.6)), plt.ylim((-1, 1))
24
25    plt.subplot(1, 3, 3)
26    plt.bar([1, 2, 3], [pr, prho, ptau], color=colors)
27    plt.hlines(0.05, 0.4, 3.6, color=red, alpha=0.5, lw=3)
28    plt.xticks([1, 2, 3], titles, fontsize=16)
29    plt.title('p-value', fontsize=16)
30    plt.xlim((0.4, 3.6)), plt.ylim((-0.01, 1))
31
32    plt.show()
```

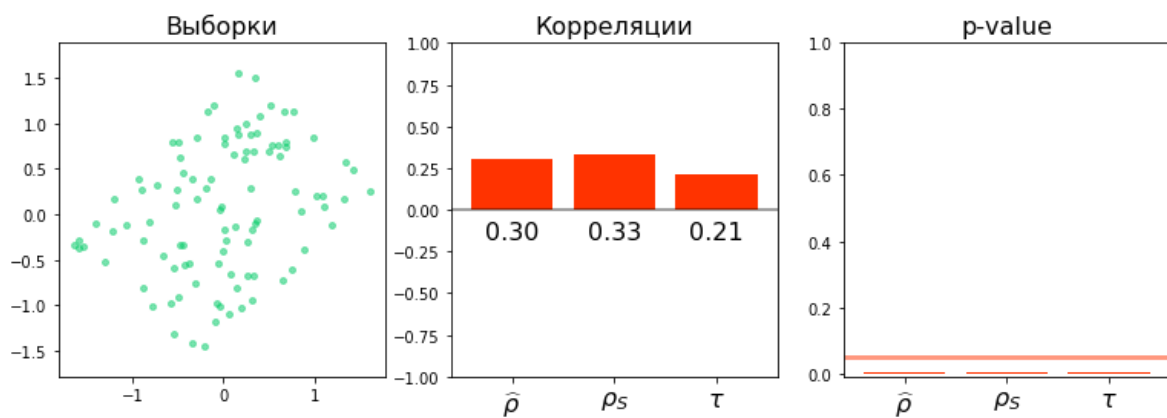
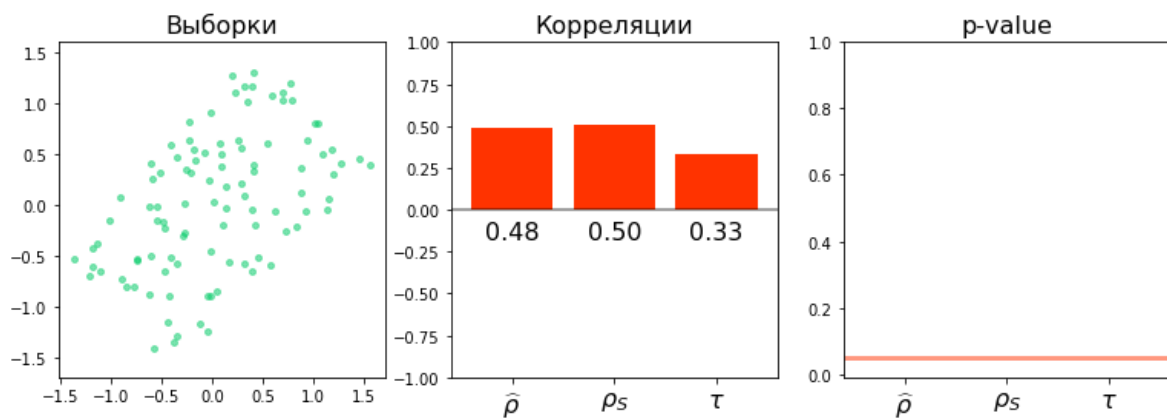
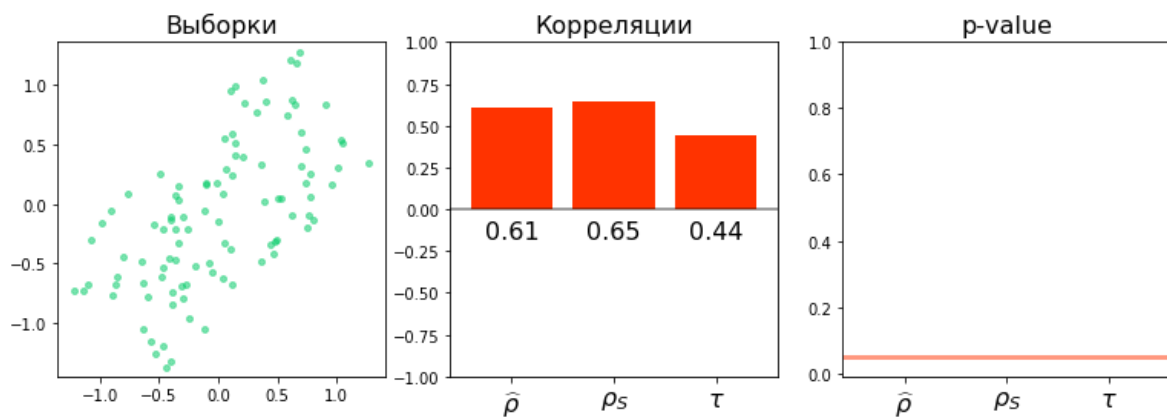
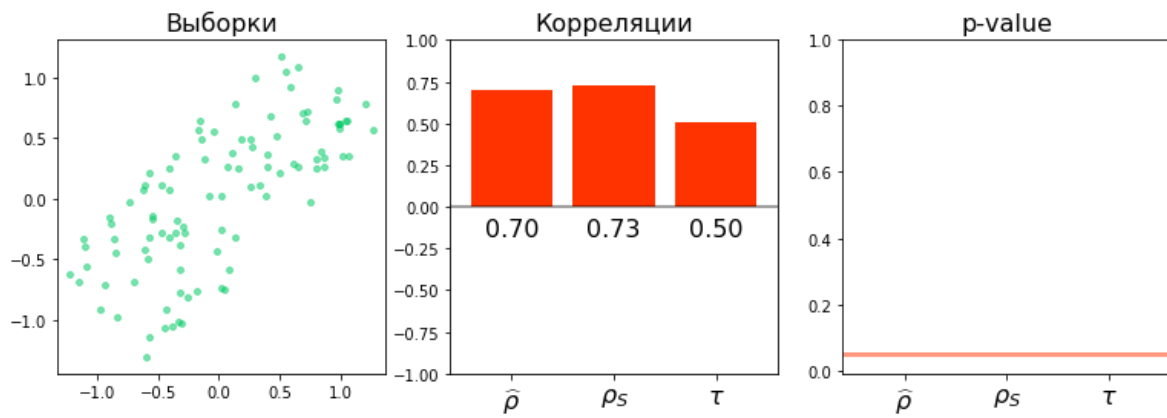
Везде ниже выборка изображается на графике без искажений, то есть масштаб по обеим осям совпадает. Столбцы отвечают за коэффициенты корреляции Пирсона, Спирмена и Кенделла соответственно. Если столбец красный, то гипотеза о независимости (точнее, некоррелированности) отвергается.

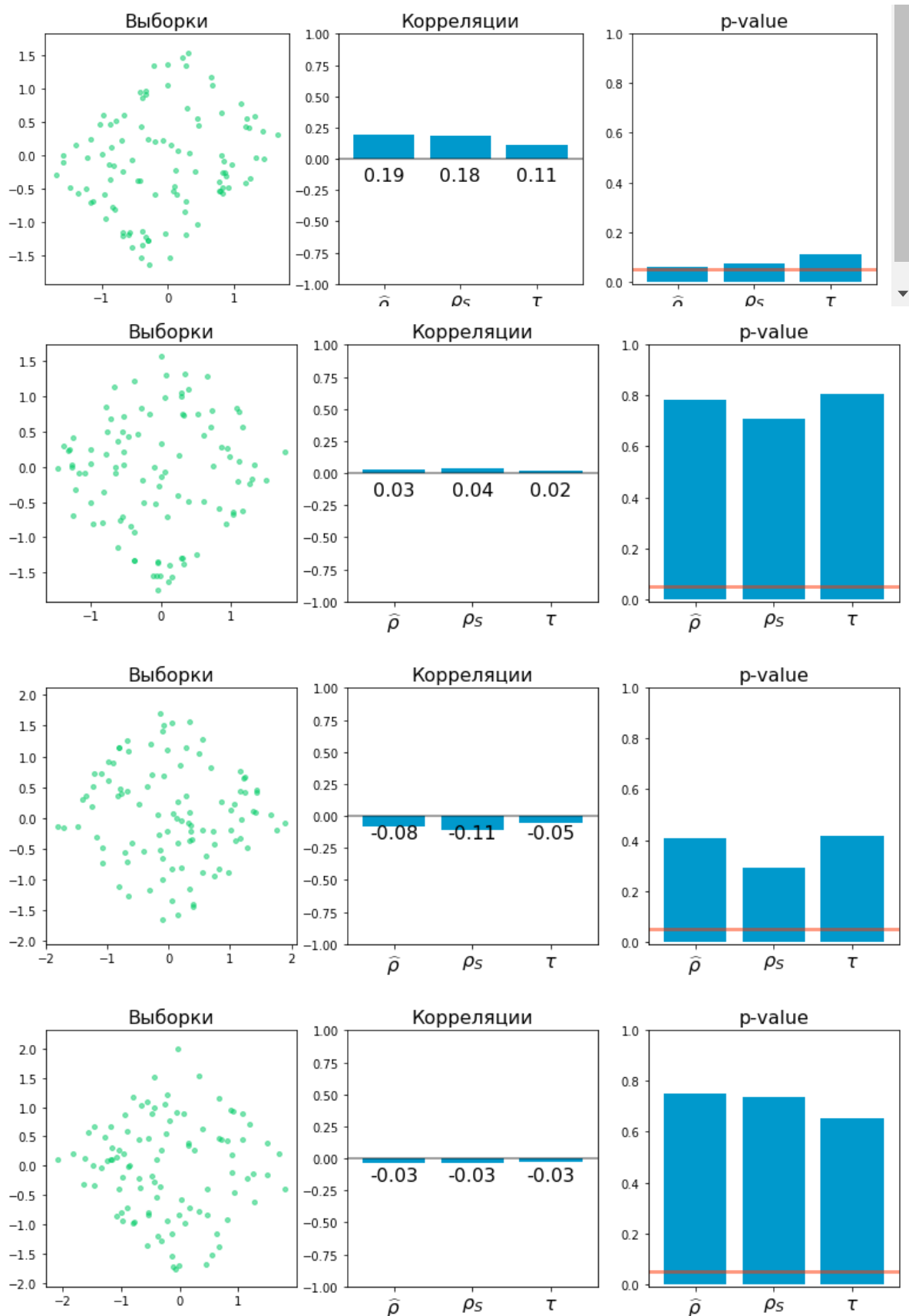
Зависимость коэффициентов корреляции от выборок. Выборки из равномерного распределения по прямой $y = x$ при $x \in [-1, 1]$ размазывается вдоль прямой $y = -x$.

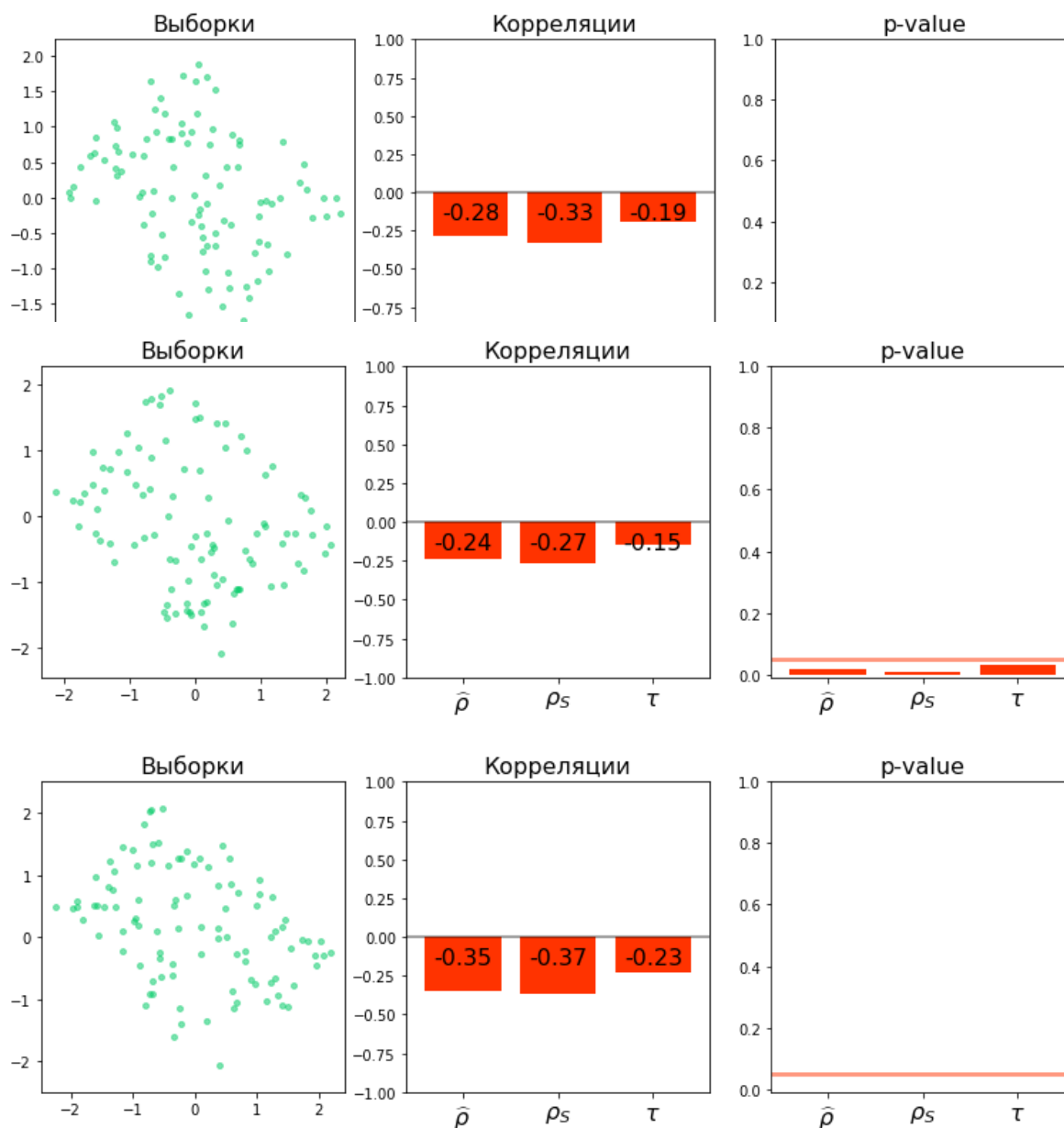
In [10]:

```
1 for i in range(15):
2     x = sps.uniform(loc=-1, scale=2).rvs(size=100)
3     y = sps.uniform(loc=-0.1*i, scale=0.2*i).rvs(size=100)
4     draw_graphics(x + y, x - y)
```





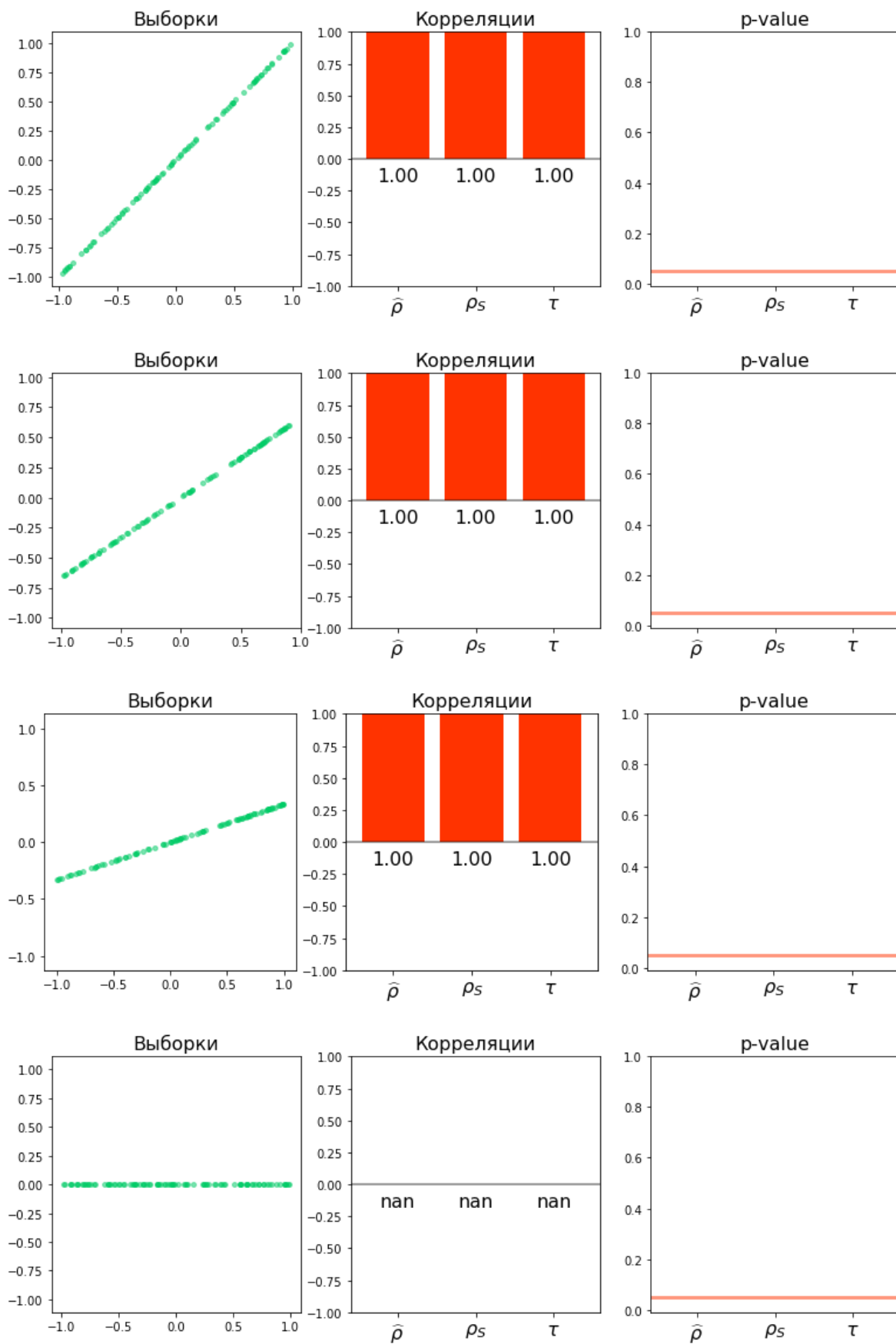


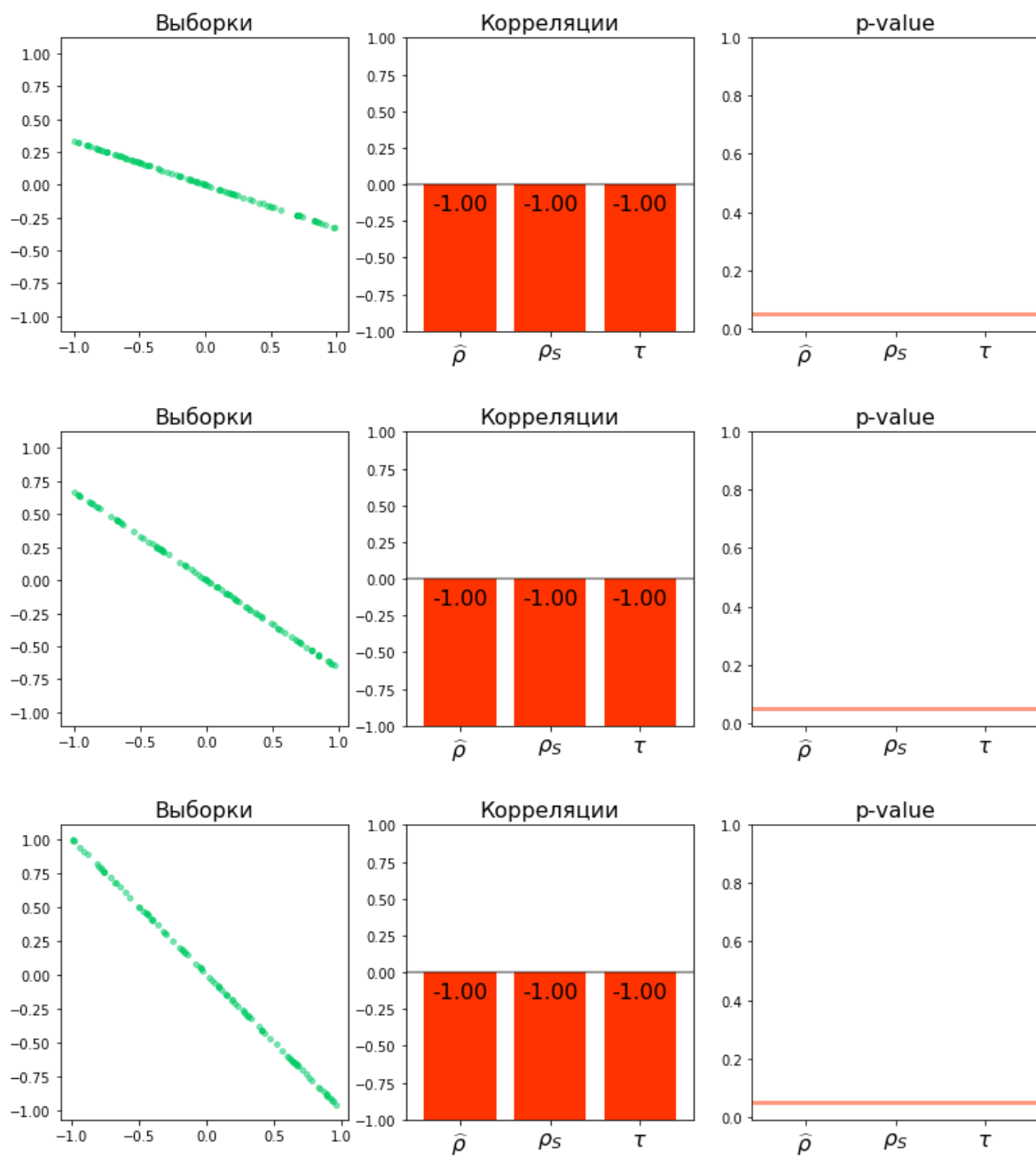


Выборки расположены вдоль поворачивающейся прямой. Если одна из выборок принимает только одно значение, то коэффициент корреляции неопределен.

In [11]:

```
1 for i in range(7):  
2     x = sps.uniform(loc=-1, scale=2).rvs(size=100)  
3     draw_graphics(x, x * (1 - i/3))
```



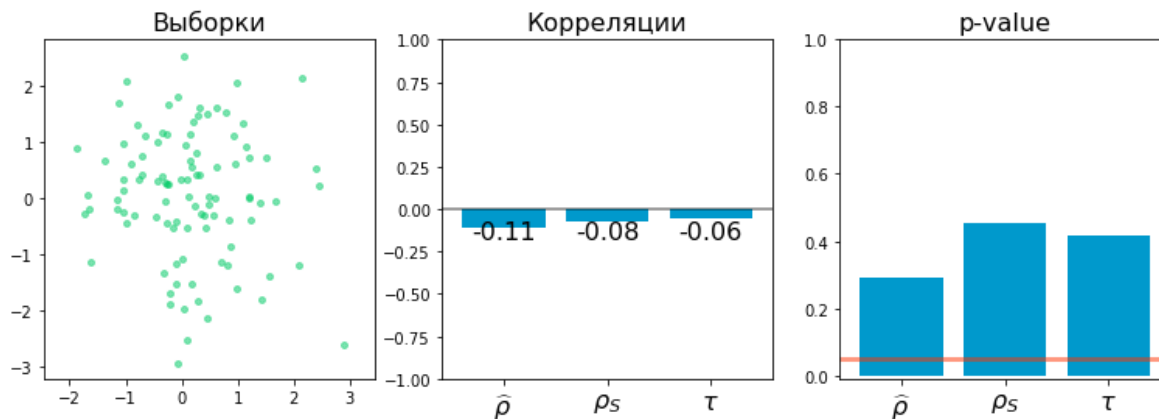


Выборка размера 100 из двумерного нормального распределения.

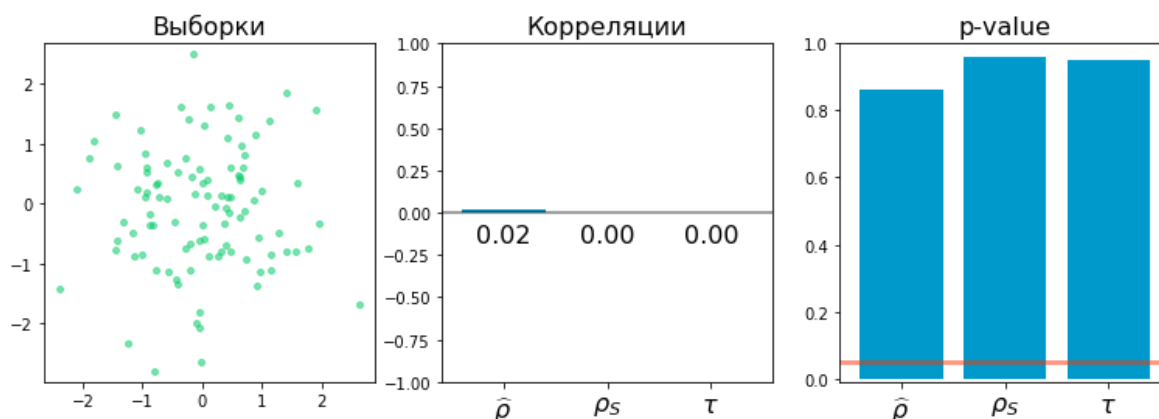
In [12]:

```
1 for i in range(11):
2     cov = 0.1 * i if i < 10 else 0.9999
3     print('Истинная корреляция: %.1f' % cov)
4     x1, x2 = sps.multivariate_normal(cov=[[1, cov], [cov, 1]]).rvs(size=100).
5     draw_graphics(x1, x2)
```

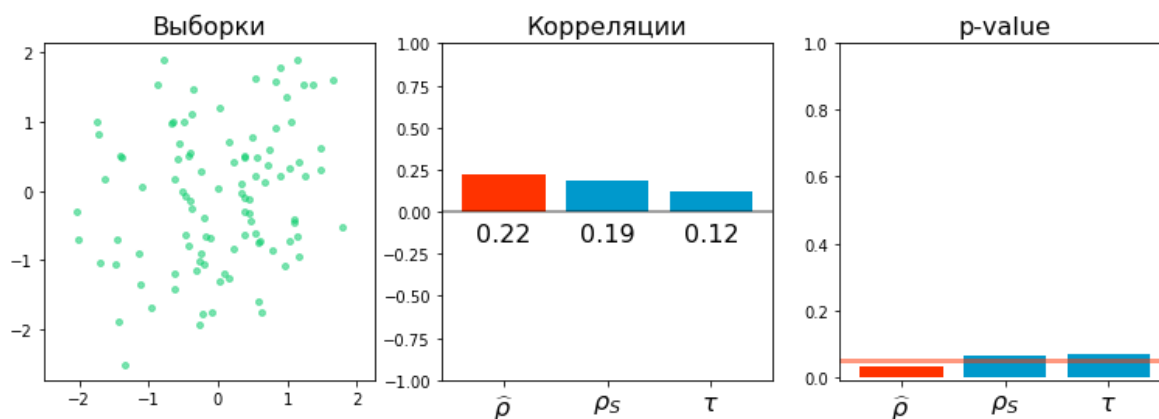
Истинная корреляция: 0.0



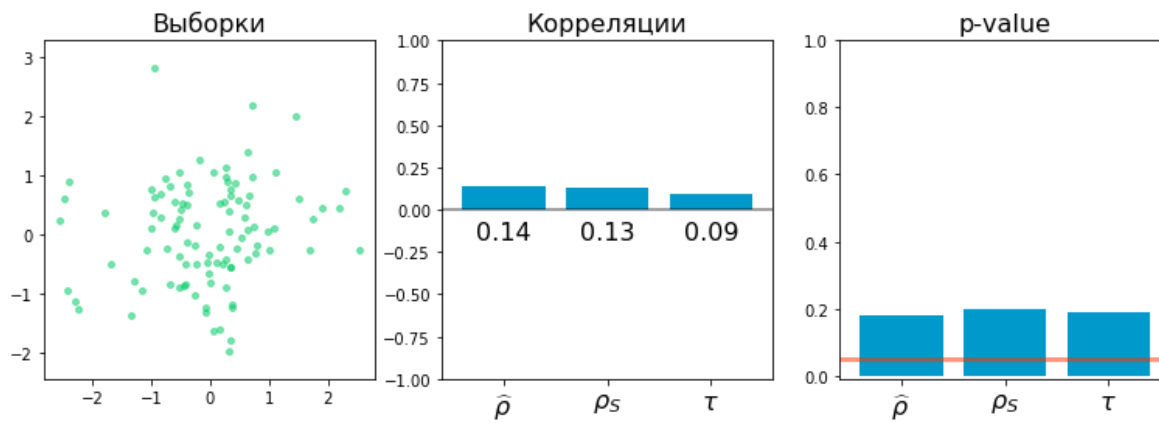
Истинная корреляция: 0.1



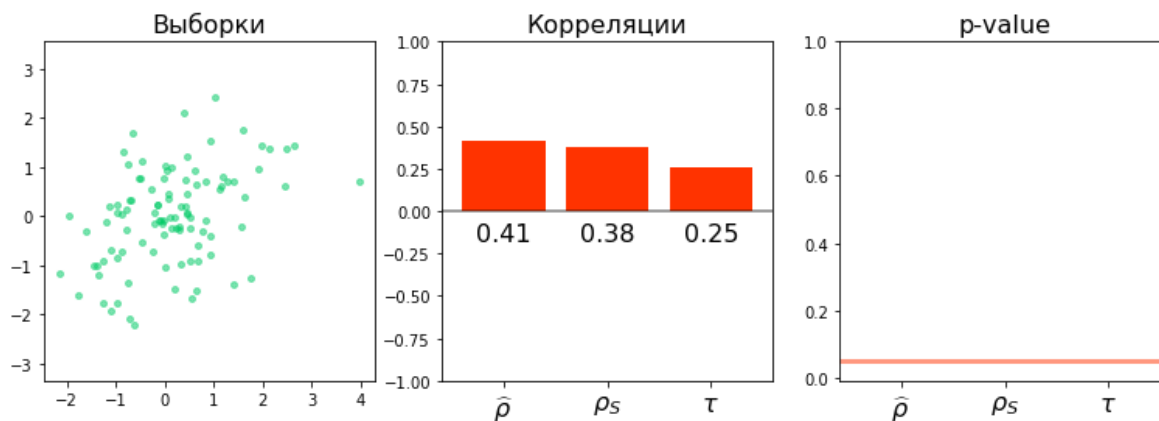
Истинная корреляция: 0.2



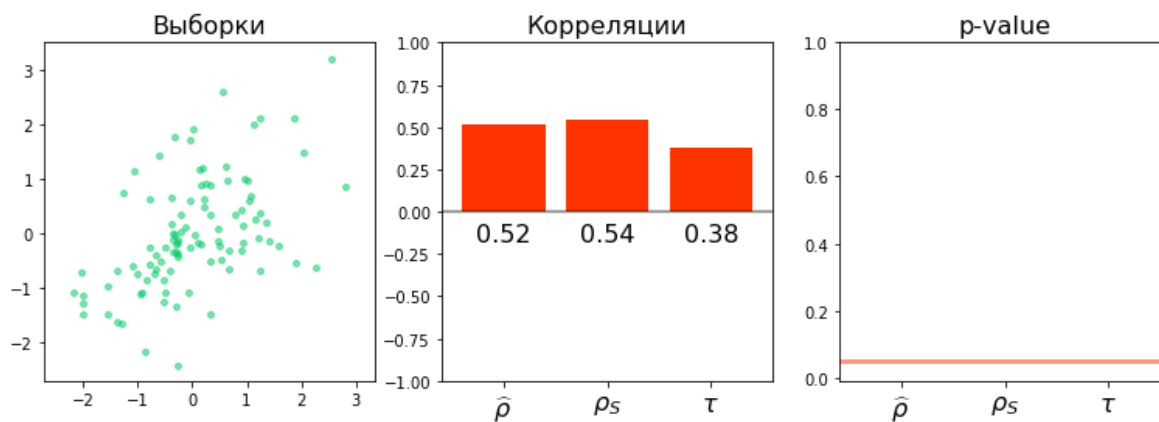
Истинная корреляция: 0.3



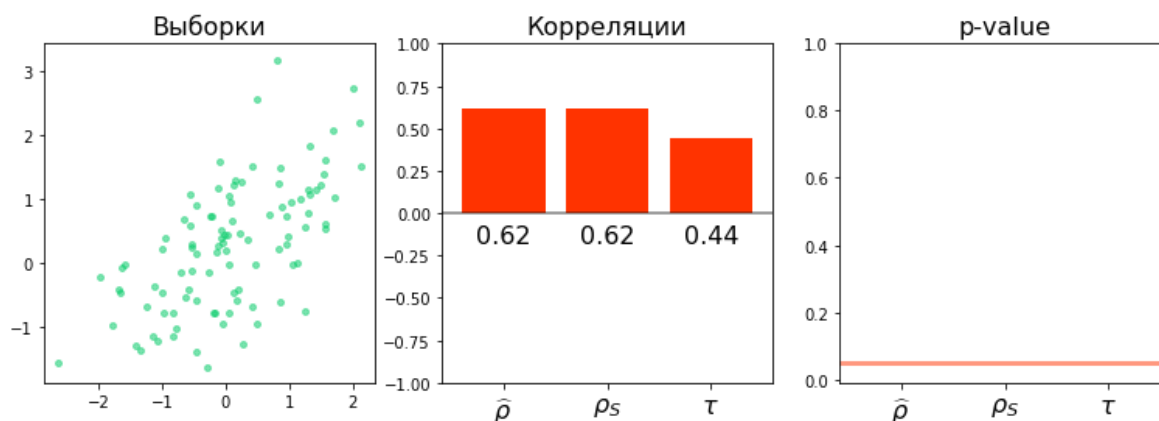
Истинная корреляция: 0.4



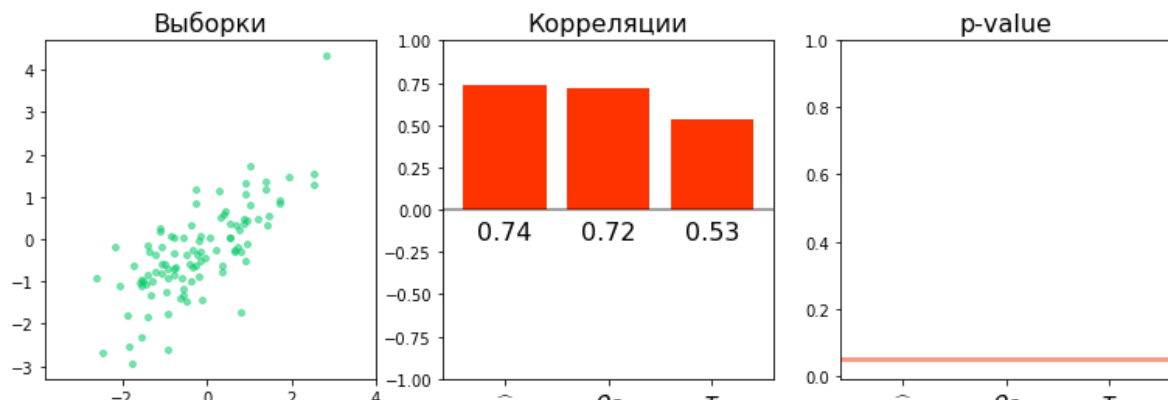
Истинная корреляция: 0.5



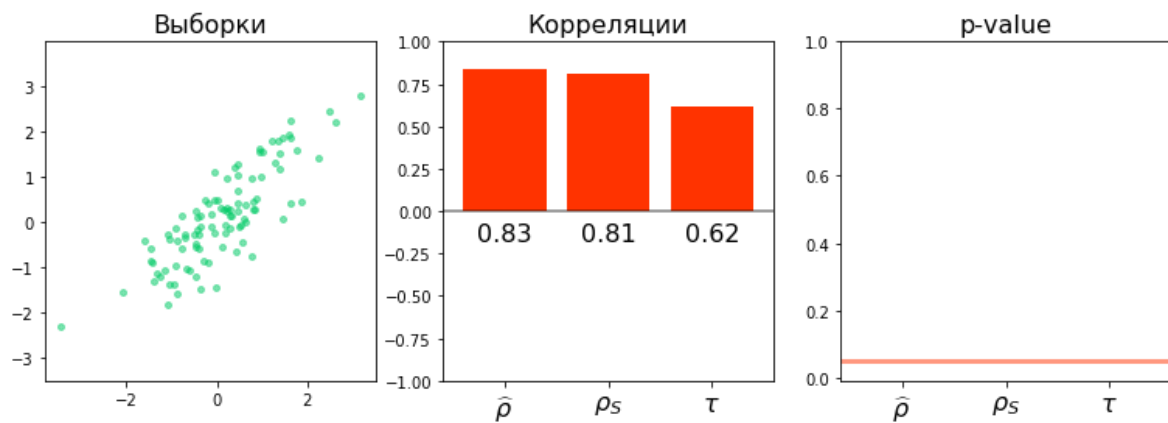
Истинная корреляция: 0.6



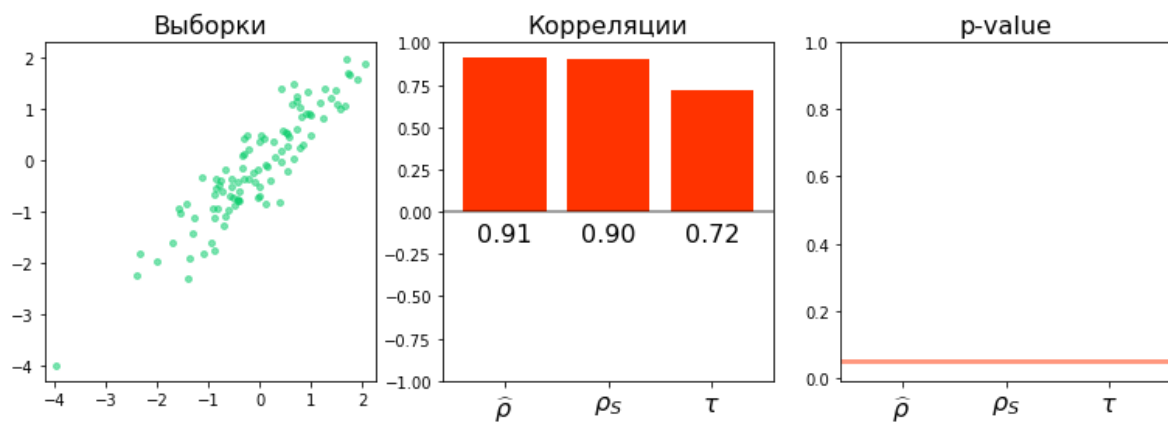
Истинная корреляция: 0.7



Истинная корреляция: 0.8

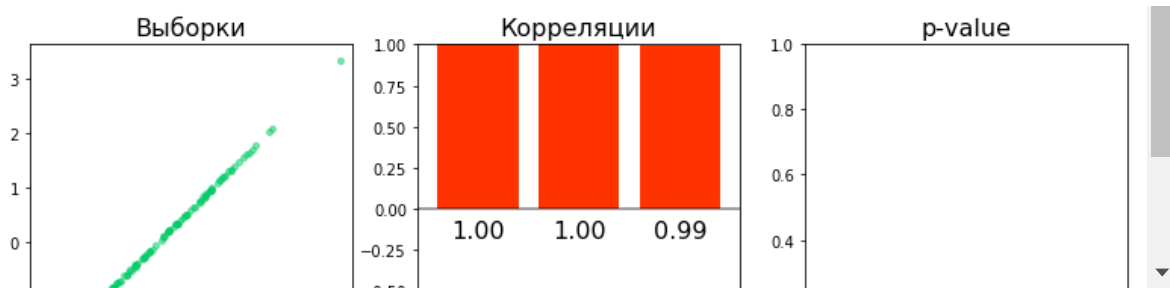


Истинная корреляция: 0.9



Истинная корреляция: 1.0



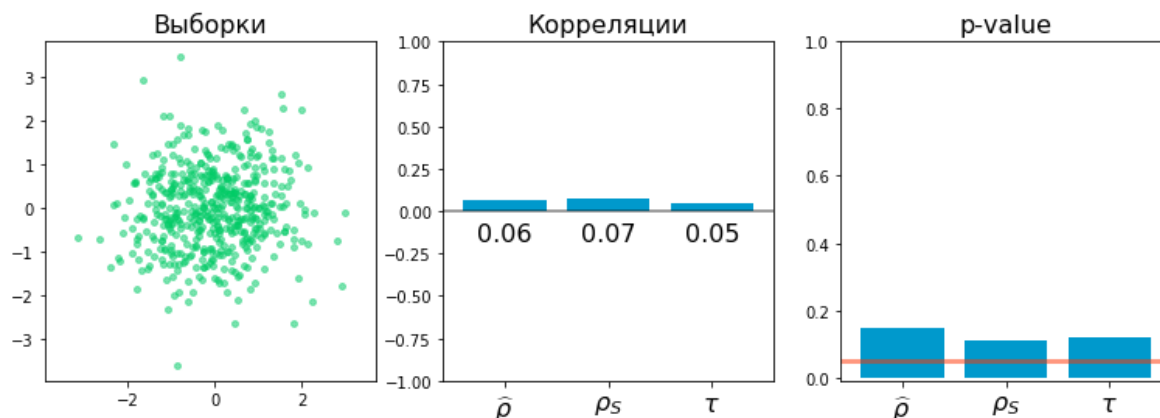


Выборка размера 500 из двумерного нормального распределения.

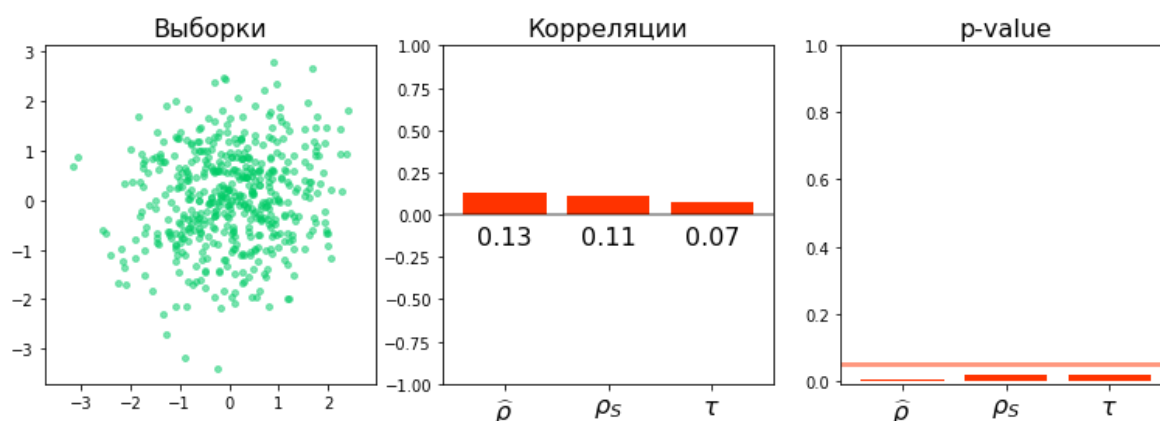
In [13]:

```
1 ▼ for i in range(11):
2     cov = 0.1 * i if i < 10 else 0.9999
3     print('Истинная корреляция: %.1f' % cov)
4     x1, x2 = sps.multivariate_normal(cov=[[1, cov], [cov, 1]]).rvs(size=500).
5     draw_graphics(x1, x2)
```

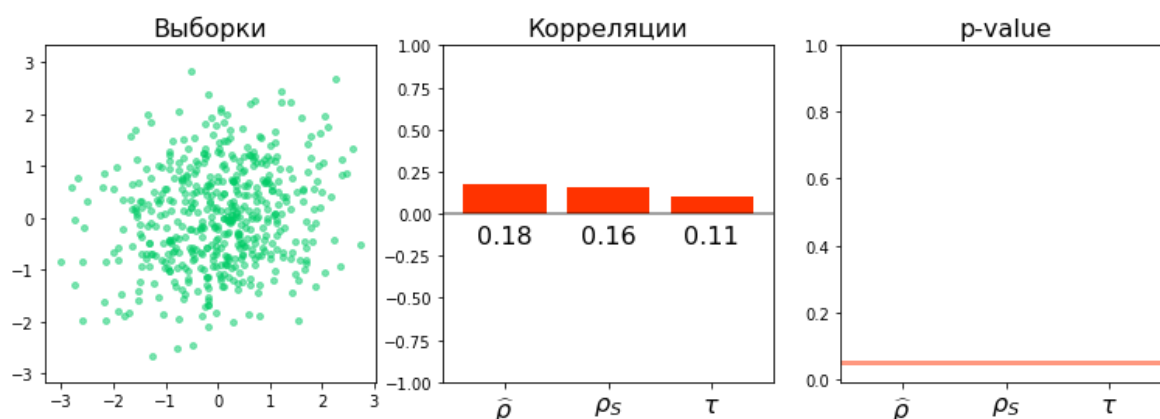
Истинная корреляция: 0.0



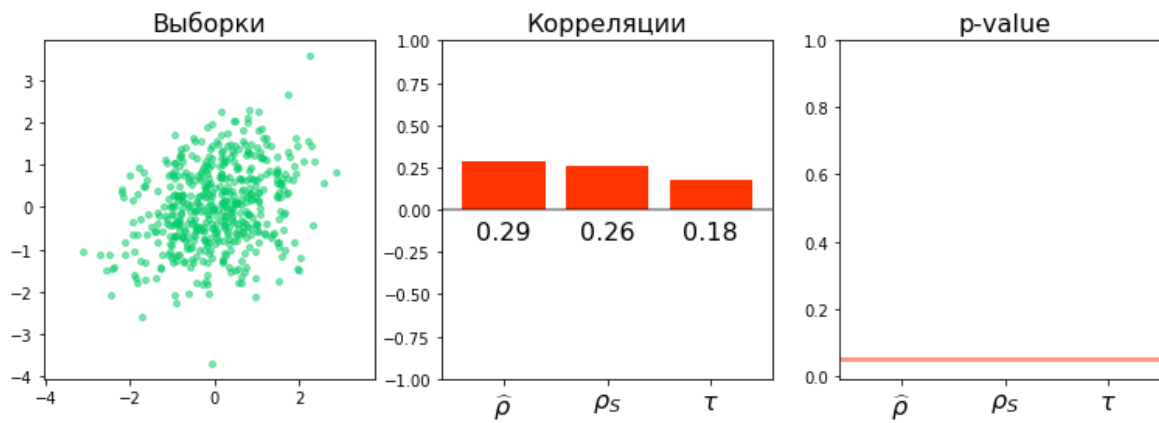
Истинная корреляция: 0.1



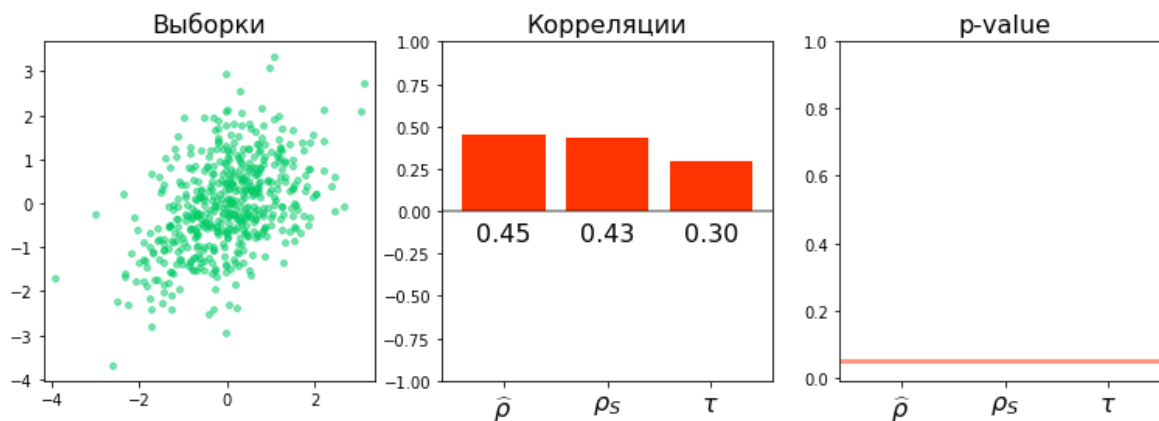
Истинная корреляция: 0.2



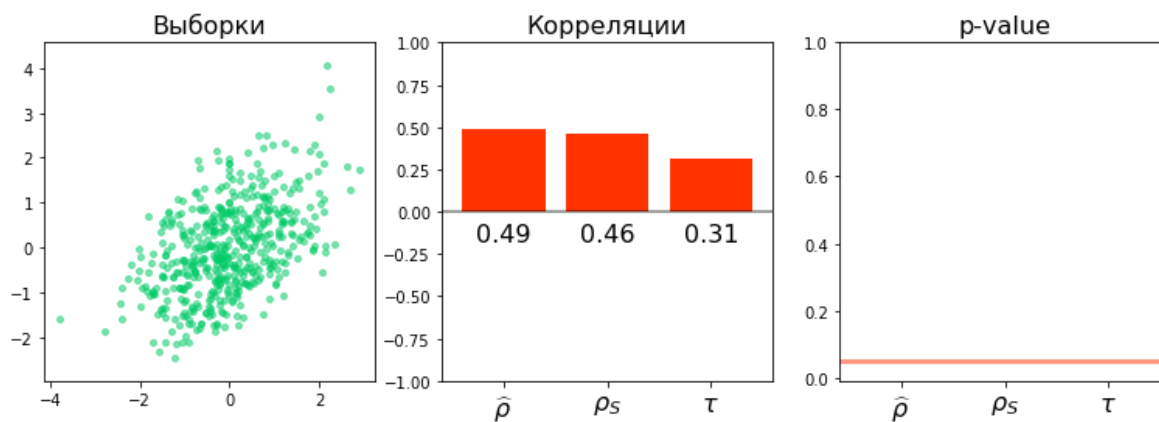
Истинная корреляция: 0.3



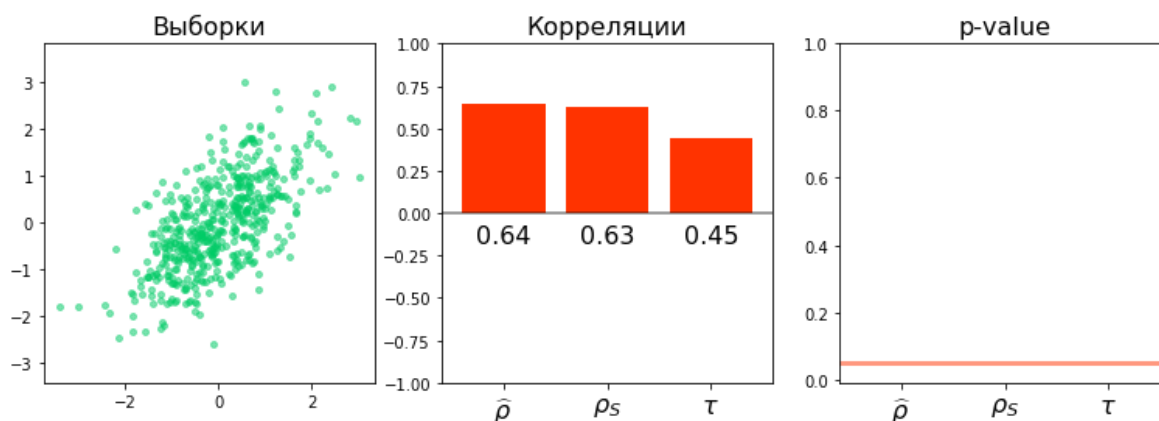
Истинная корреляция: 0.4



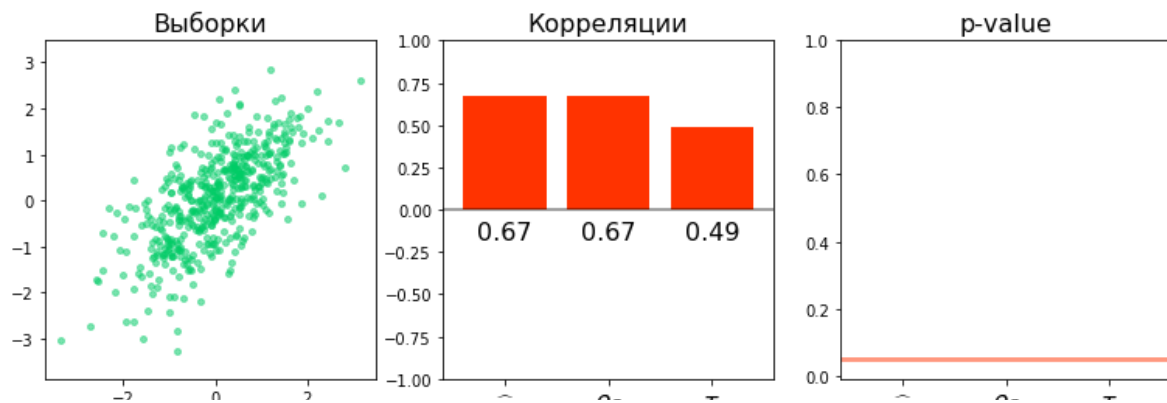
Истинная корреляция: 0.5



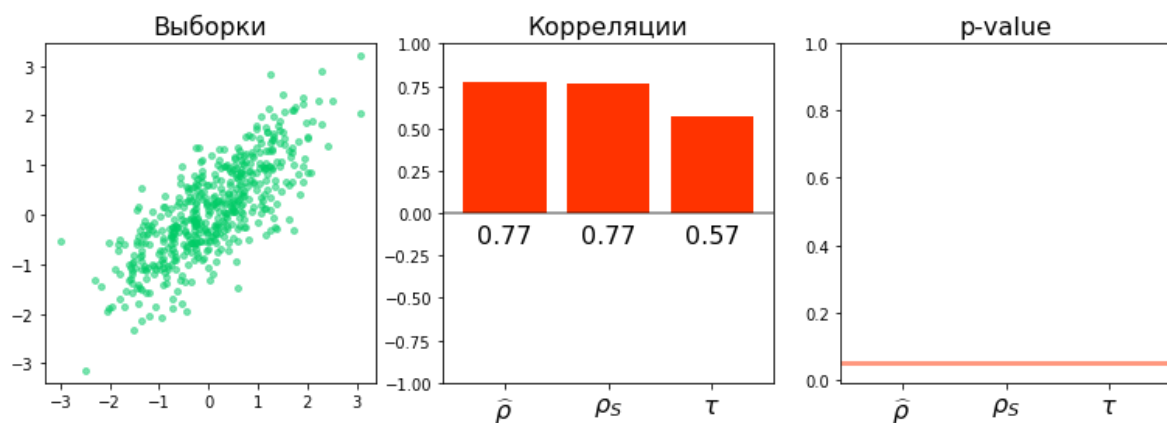
Истинная корреляция: 0.6



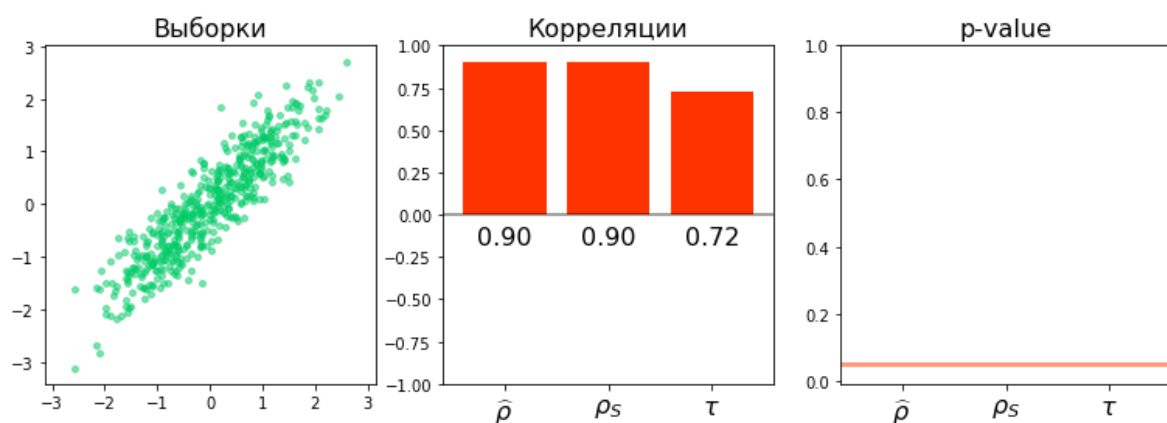
Истинная корреляция: 0.7



Истинная корреляция: 0.8

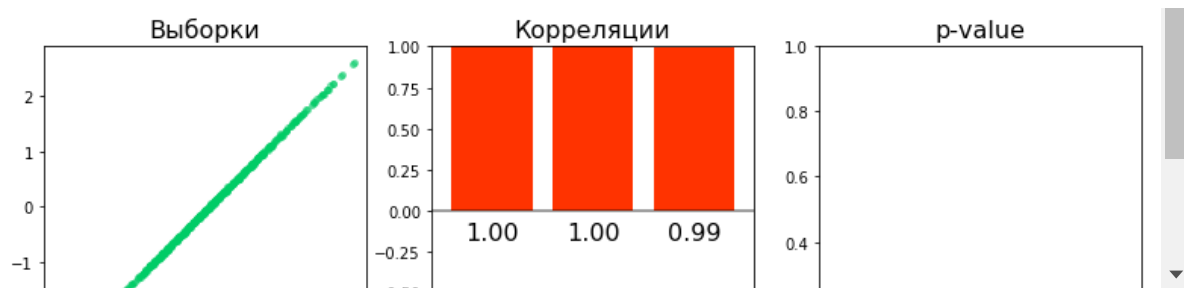


Истинная корреляция: 0.9



Истинная корреляция: 1.0



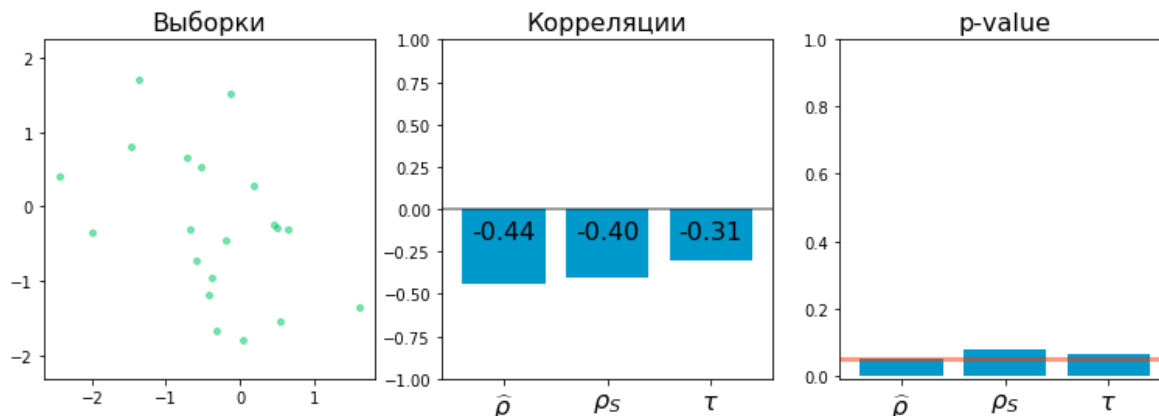


Выборка размера 20 из двумерного нормального распределения. При малых значениях корреляции гипотеза о независимости не отвергается в отличие от выборки большего размера.

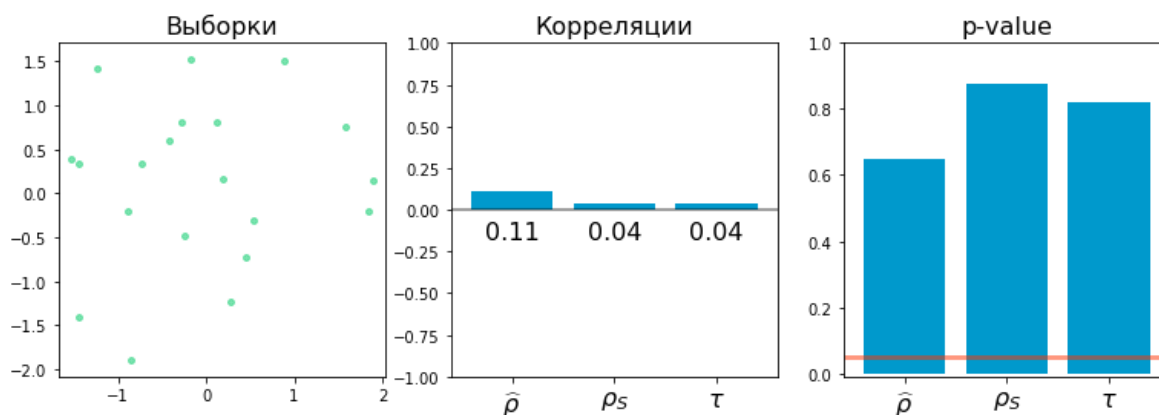
In [14]:

```
1 for i in range(11):
2     cov = 0.1 * i if i < 10 else 0.9999
3     print('Истинная корреляция: %.1f' % cov)
4     x1, x2 = sps.multivariate_normal(cov=[[1, cov], [cov, 1]]).rvs(size=20).T
5     draw_graphics(x1, x2)
```

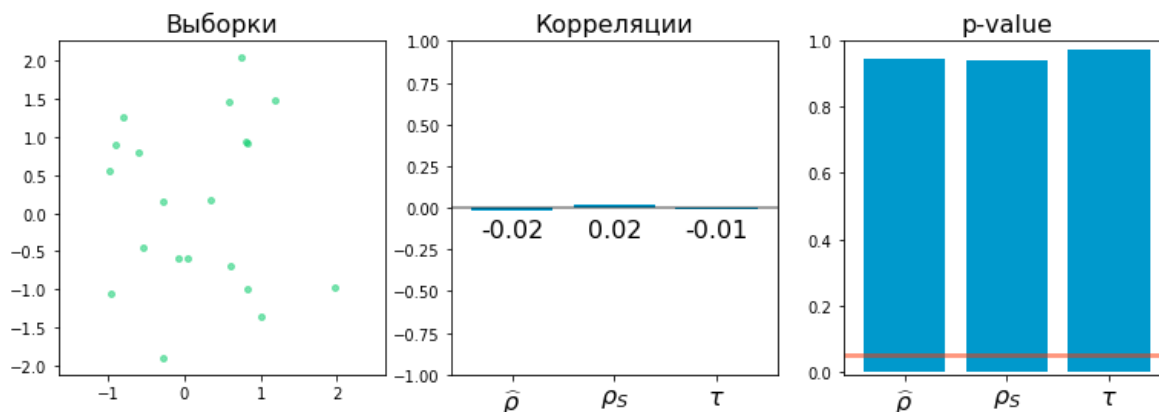
Истинная корреляция: 0.0



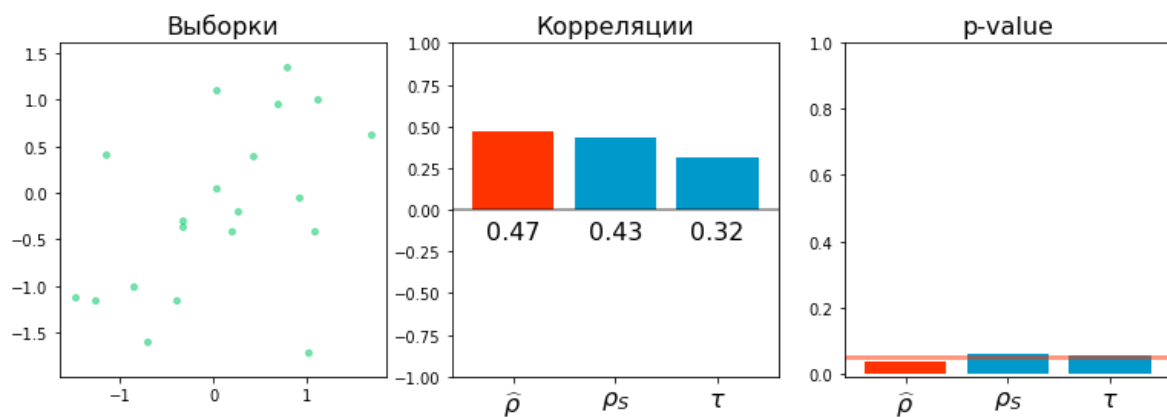
Истинная корреляция: 0.1



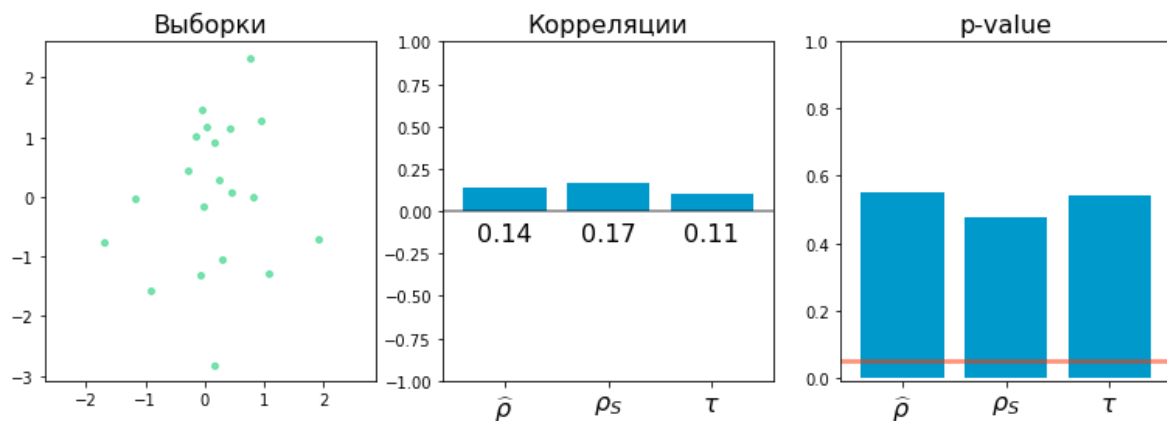
Истинная корреляция: 0.2



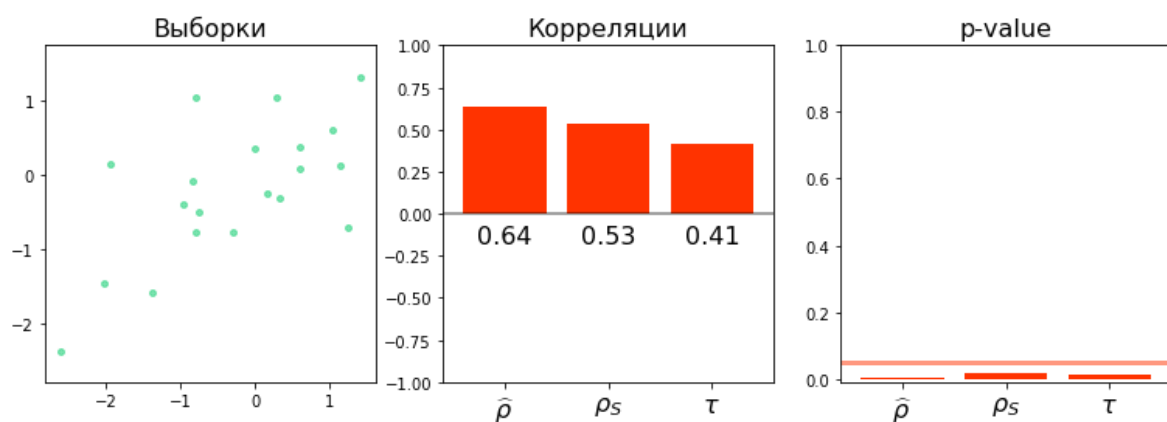
Истинная корреляция: 0.3



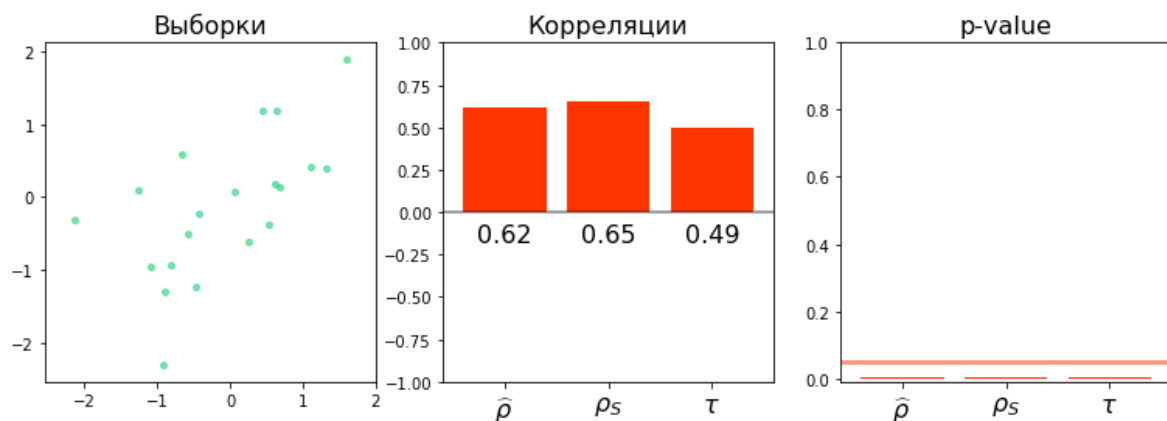
Истинная корреляция: 0.4



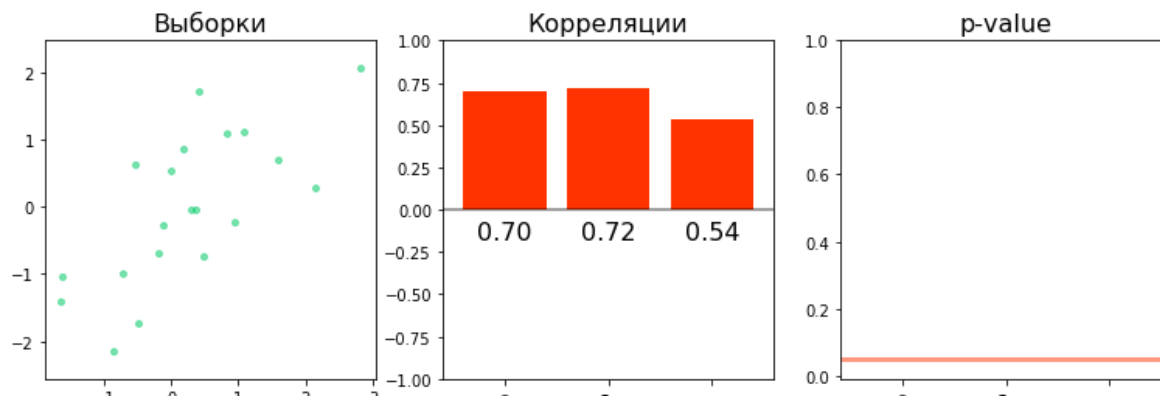
Истинная корреляция: 0.5



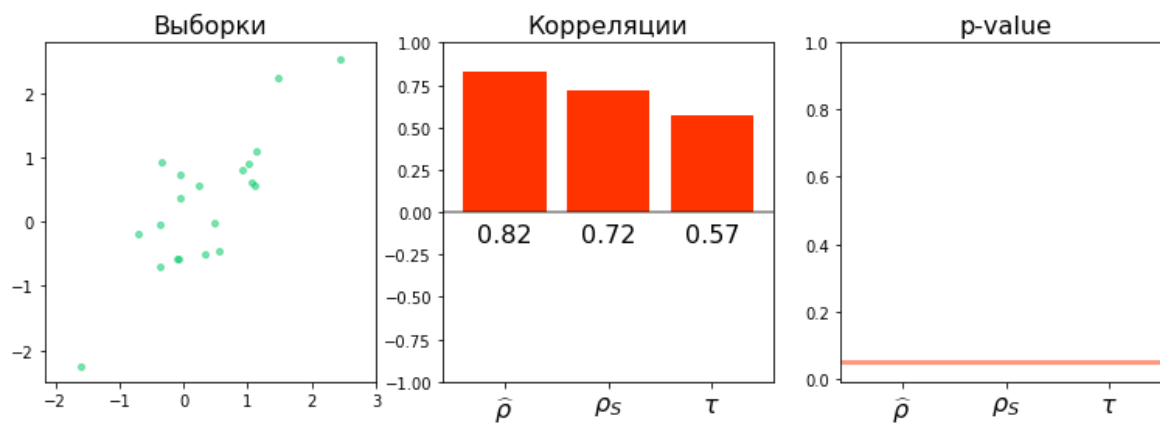
Истинная корреляция: 0.6



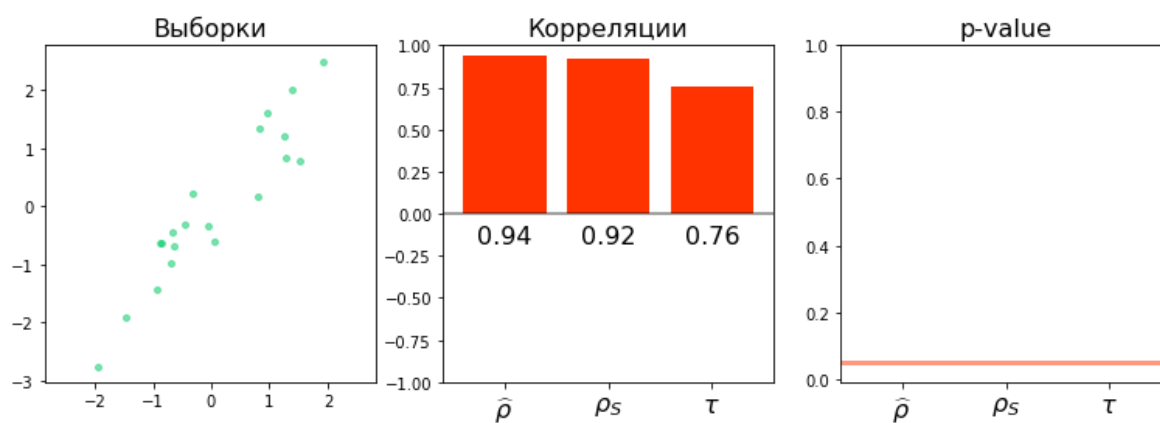
Истинная корреляция: 0.7



Истинная корреляция: 0.8

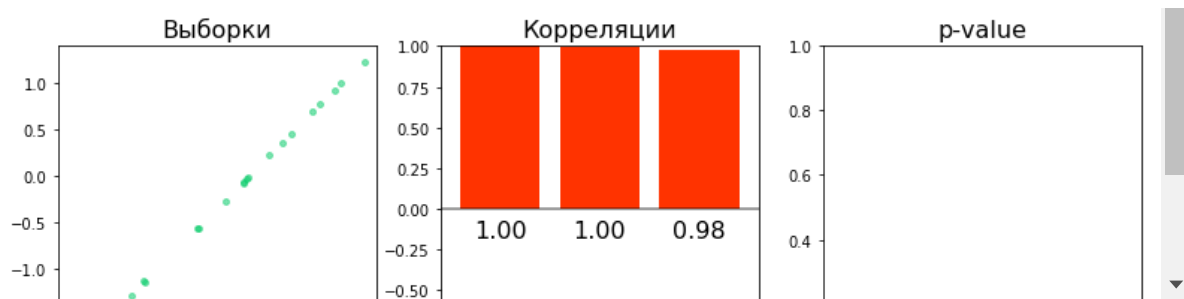


Истинная корреляция: 0.9



Истинная корреляция: 1.0

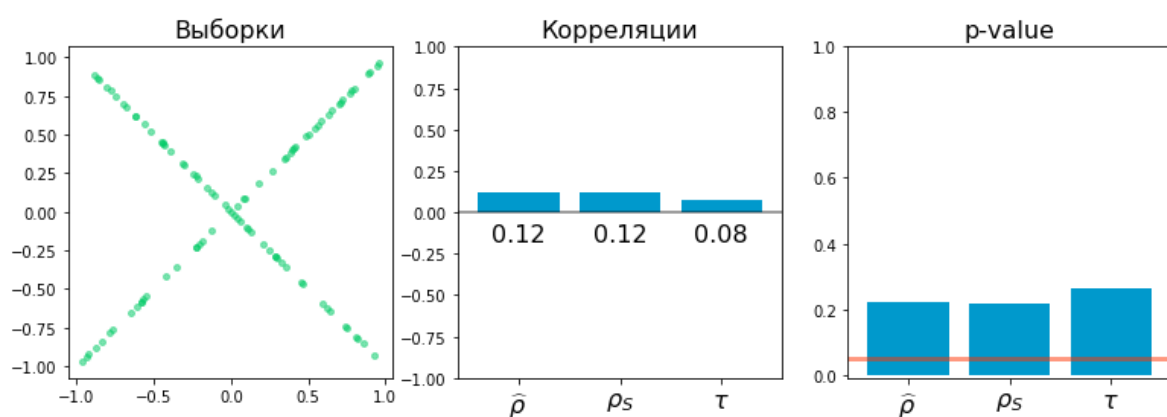




Выборки "X". Очевидно, они зависимы, но коэффициенты корреляции близки к нулю.

In [15]:

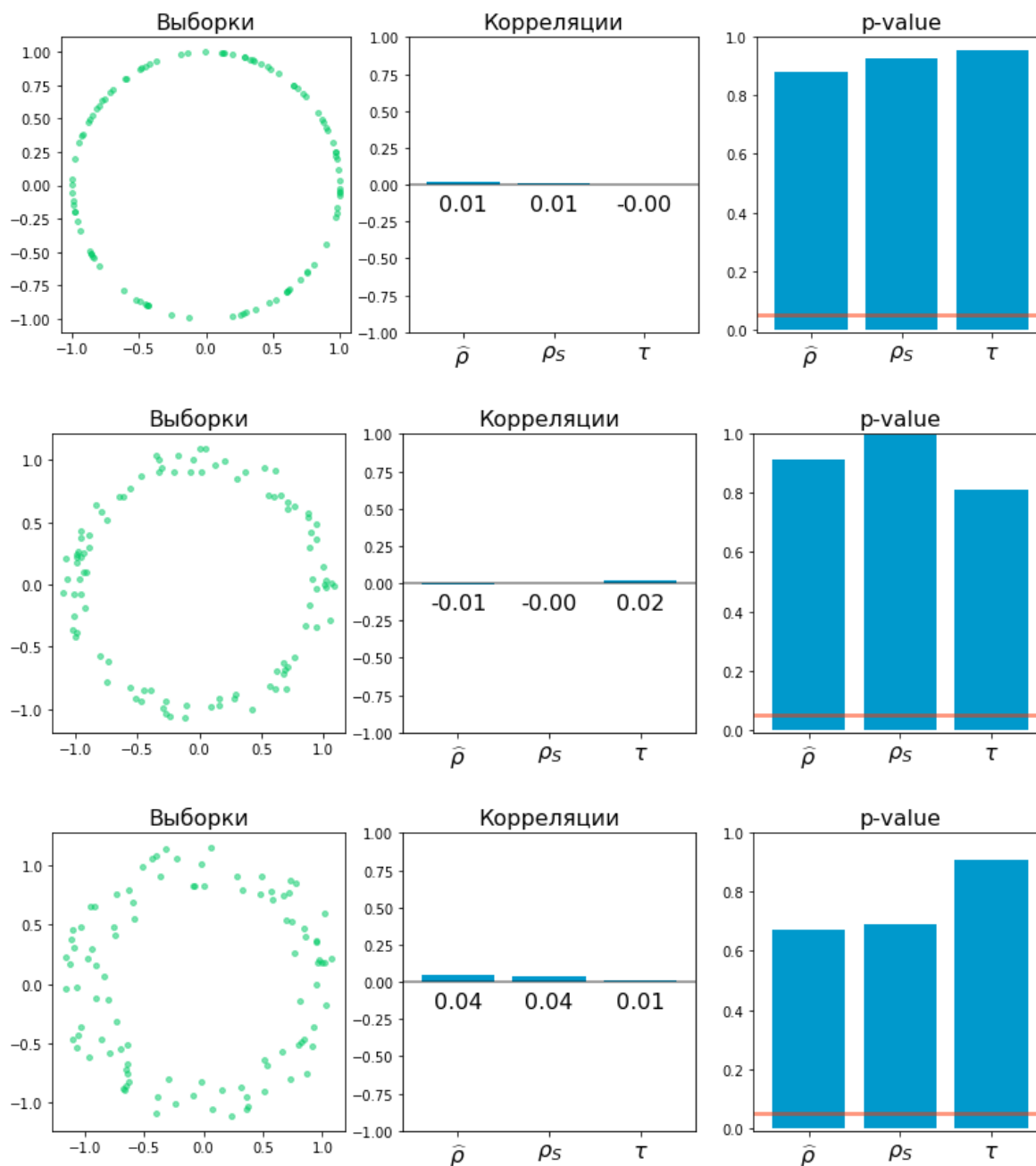
```
1 x1 = sps.uniform(loc=-1, scale=2).rvs(size=100)
2 x2 = x1 * (1 - 2 * sps.bernoulli(0.5).rvs(size=100))
3 draw_graphics(x1, x2)
```



При круговой зависимости выполняется аналогичное свойство.

In [16]:

```
1 for i in range(3):
2     phi = sps.uniform(loc=-1, scale=2).rvs(size=100)
3     r = 1 + sps.uniform(loc=-0.1*i, scale=0.2*i).rvs(size=100)
4     draw_graphics(r * np.cos(np.pi * phi), r * np.sin(np.pi * phi))
```



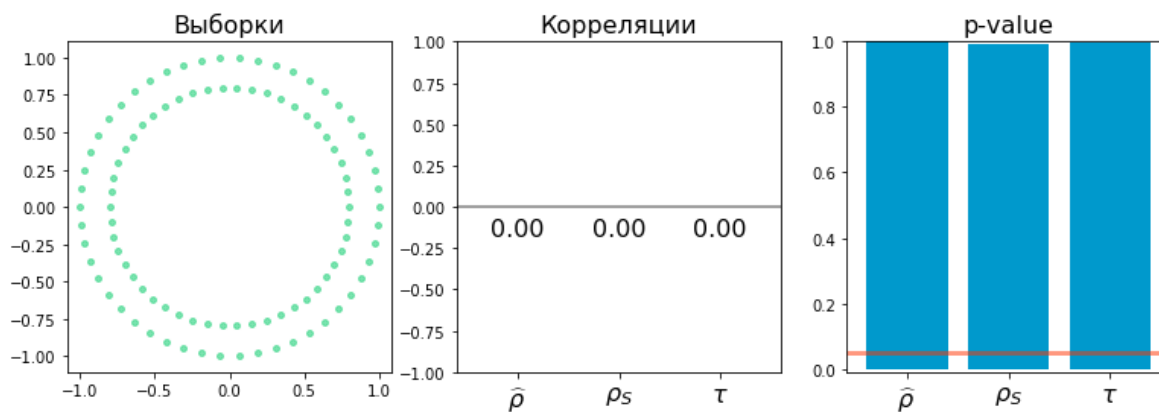
Несколько других примеров

In [17]:

```
1 from sklearn.datasets import make_circles, make_blobs
```

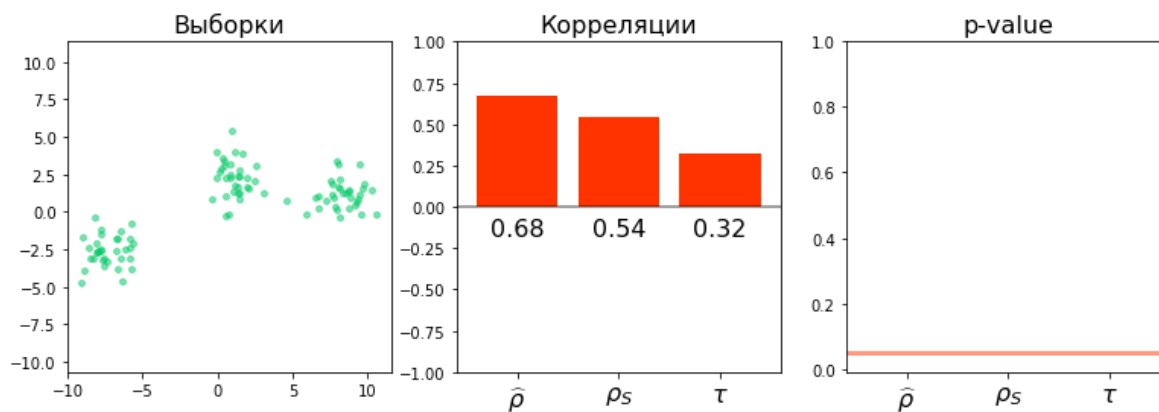
In [18]:

```
1 x, y = make_circles(n_samples=100)[0].T
2 draw_graphics(x, y)
```



In [19]:

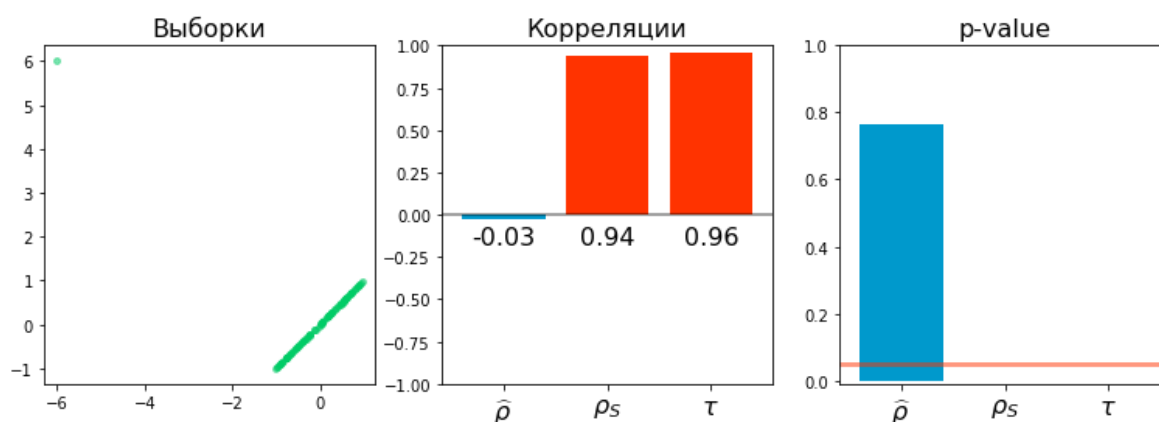
```
1 x, y = make_blobs(n_samples=100)[0].T
2 draw_graphics(x, y)
```



Пусть выборки линейно зависимы, но при этом случился один выброс (в левом верхнем углу). Коэффициент корреляции Пирсона близок к нулю, несмотря на очевидную зависимость данных. Остальные коэффициенты корреляции не сильно влияют на выброс.

In [20]:

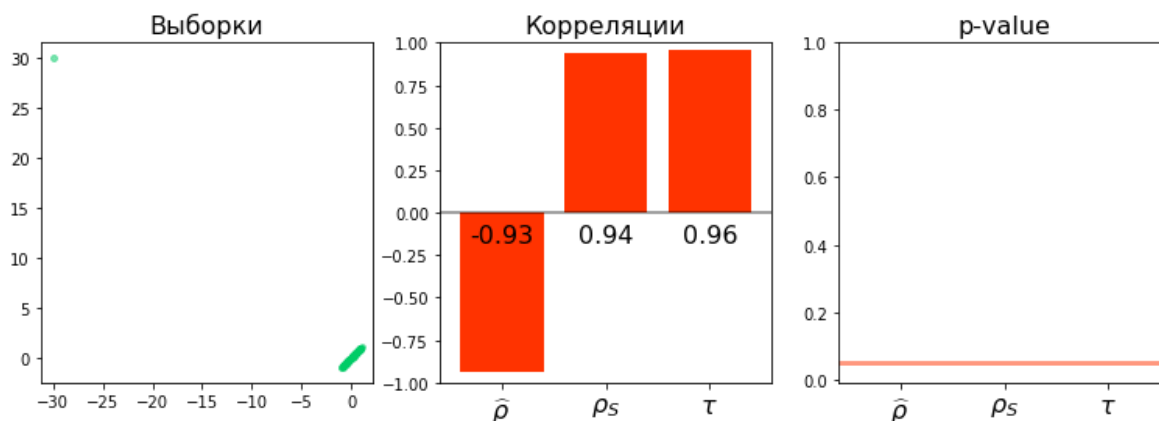
```
1 x1 = sps.uniform(loc=-1, scale=2).rvs(size=100)
2 x2 = np.array(x1)
3 x1[-1] = -6
4 x2[-1] = 6
5 draw_graphics(x1, x2)
```



Если выбор "слишком большой", то коэффициент корреляции Пирсона может быть близок к -1, что означает отрицательную линейную зависимость, несмотря на то, что на самом деле она положительна. Другие два коэффициента корреляции практически не меняются.

In [21]:

```
1 x1 = sps.uniform(loc=-1, scale=2).rvs(size=100)
2 x2 = np.array(x1)
3 x1[-1] = -30
4 x2[-1] = 30
5 draw_graphics(x1, x2)
```



[Еще несколько примеров с Википедии](#)

(https://upload.wikimedia.org/wikipedia/commons/0/02/Correlation_examples.png)

[Игра - отгадайте значение коэффициента корреляции Пирсона по выборке.](#)

(<http://guessthecorrelation.com>)

Никита Волков

<https://mipt-stats.gitlab.io/> (<https://mipt-stats.gitlab.io/>).