

In [1]:

```
1 import numpy as np
2 import scipy.stats as sps
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from statsmodels.distributions.empirical_distribution import ECDF
6 from statsmodels.sandbox.stats.multicomp import multipletests
7 from collections import Counter
8 from tqdm import tqdm_notebook
9
10 red = '#FF3300'
11 blue = '#0099CC'
12 green = '#00CC66'
13
14 %matplotlib inline
15 sns.set(style='ticks', font_scale=1.7)
```

Критерии согласия

Критерий согласия Пирсона (хи-квадрат)

H_0 : Выборка из некоторого класса распределений

H_1 : H_0 не верна.

`chisquare`

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

(`f_obs`, `f_exp=None`, `ddof=0`)

- `f_obs` --- число элементов выборки, попавших в каждый из интервалов
- `f_exp` --- ожидаемое число (по умолчанию равномерное)
- `ddof` --- поправка на число степеней свободы. Статистика асимптотически будет иметь распределение $k - 1 - ddof$, где k --- число интервалов.

In [2]:

```
1 sps.chisquare([16, 18, 16, 14, 12, 12])
```

Out[2]:

Power_divergenceResult(statistic=2.0, pvalue=0.8491450360846096)

In [3]:

```
1 sps.chisquare([16, 18, 16, 14, 12, 12], f_exp=[16, 16, 16, 16, 20, 4])
```

Out[3]:

Power_divergenceResult(statistic=19.7, pvalue=0.0014224993317060594)

In [4]:

```
1 | sps.chisquare([16, 18, 16, 14, 12, 12], f_exp=[16, 16, 16, 16, 24, 0])
```

```
/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:5048: RuntimeWarning: divide by zero encountered in true_divide  
  terms = (f_obs - f_exp)**2 / f_exp
```

Out[4]:

```
Power_divergenceResult(statistic=inf, pvalue=0.0)
```

Эксперимент с группами крови

In [5]:

```
1 | sps.chisquare([121, 120, 79, 33], f_exp=np.array([0.343, 0.340, 0.224, 0.093]))
```

Out[5]:

```
Power_divergenceResult(statistic=0.001011144526478114, pvalue=0.9746327438096731)
```

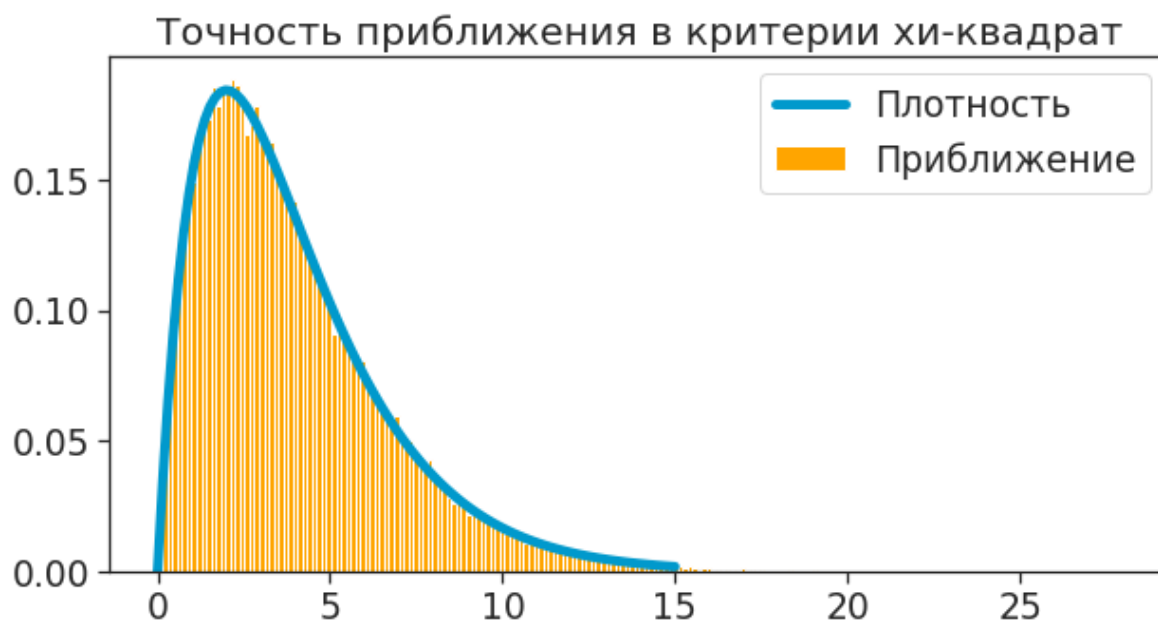
Посмотрим, насколько хорошо выполняется асимптотика при справедливости нулевой гипотезы для разных распределений.

1. Равномерное на 5 элементах

In [6]:

```
1 sample_size = 1000
2 sample_count = 100000
3 chisq_values = np.zeros(sample_count)
4
5 for i in tqdm_notebook(range(sample_count)):
6     sample = sps.randint(low=0, high=5).rvs(size=sample_size)
7     f_obs = np.array(list(Counter(sample).values()))
8     chisq_values[i] = sps.chisquare(f_obs)[0]
9
10 plt.figure(figsize=(10, 5))
11 grid = np.linspace(0, 15, 100)
12 plt.plot(grid, sps.chi2(df=4).pdf(grid),
13          lw=5, color=blue, label='Плотность')
14 plt.hist(chisq_values, bins=200,
15          color='orange', density=True, label='Приближение')
16 plt.title('Точность приближения в критерии хи-квадрат')
17 plt.legend();
```

100% 100000/100000 [11:21<00:00, 146.68it/s]

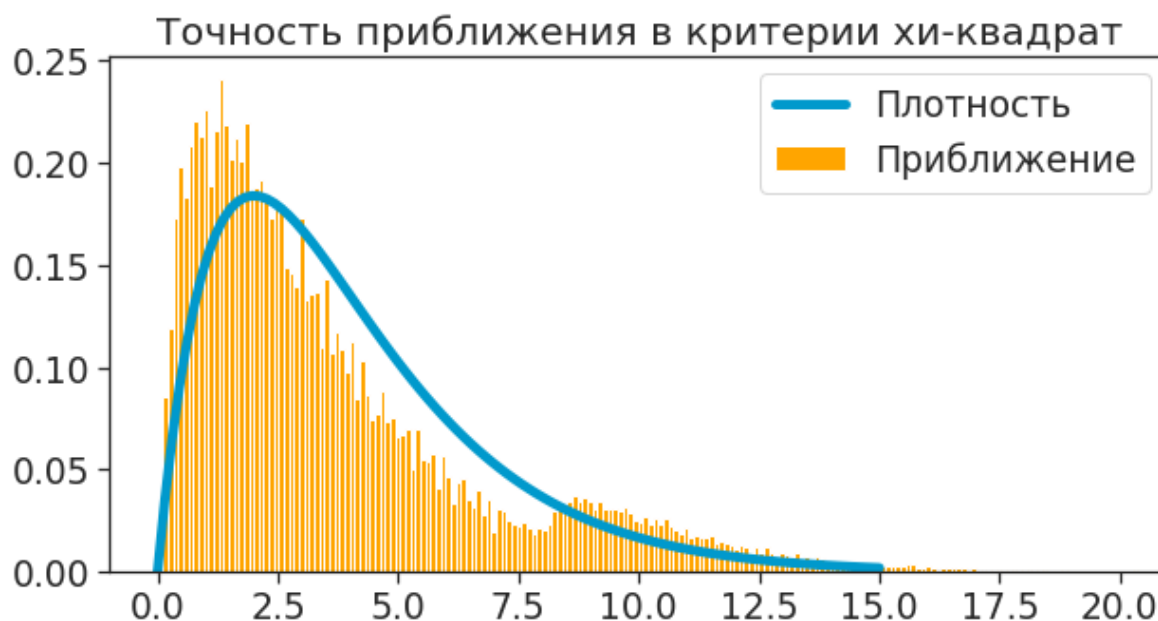


2. Равномерное на 4 элементах, а 5-й с очень маленькой вероятностью.

In [7]:

```
1 chisq_values = np.zeros(sample_count)
2
3 epsilon = 0.0001
4 pk = np.array([0.25 - epsilon / 4] * 4 + [epsilon])
5 some_distribution = sps.rv_discrete(name='1', values=(np.arange(5), pk))
6
7 for i in tqdm_notebook(range(sample_count)):
8     sample = some_distribution.rvs(size=sample_size)
9     f_obs = np.array(list(Counter(sample).values()))
10    if len(f_obs) == 4: f_obs = np.append(f_obs, [0])
11    f_exp = pk * sample_size
12    chisq_values[i] = sps.chisquare(f_obs, f_exp=f_exp)[0]
13
14 plt.figure(figsize=(10, 5))
15 grid = np.linspace(0, 15, 100)
16 plt.plot(grid, sps.chi2(df=4).pdf(grid),
17          lw=5, color=blue, label='Плотность')
18 plt.hist(chisq_values, bins=200,
19         color='orange', density=True, label='Приближение', range=(0, 21))
20 plt.title('Точность приближения в критерии хи-квадрат')
21 plt.xlim((-1, 21));
22 plt.legend();
```

100% 100000/100000 [00:44<00:00, 2248.06it/s]



Вывод: асимптотика ломается, если есть интервалы с малой вероятностью или с малым числом элементов в выборке.

Критерий согласия Колмогорова (-Смирнова) {для нормальных Лиллиефорса}

H_0 : Выборка из некоторого класса распределений

H_1 : H_0 не верна.

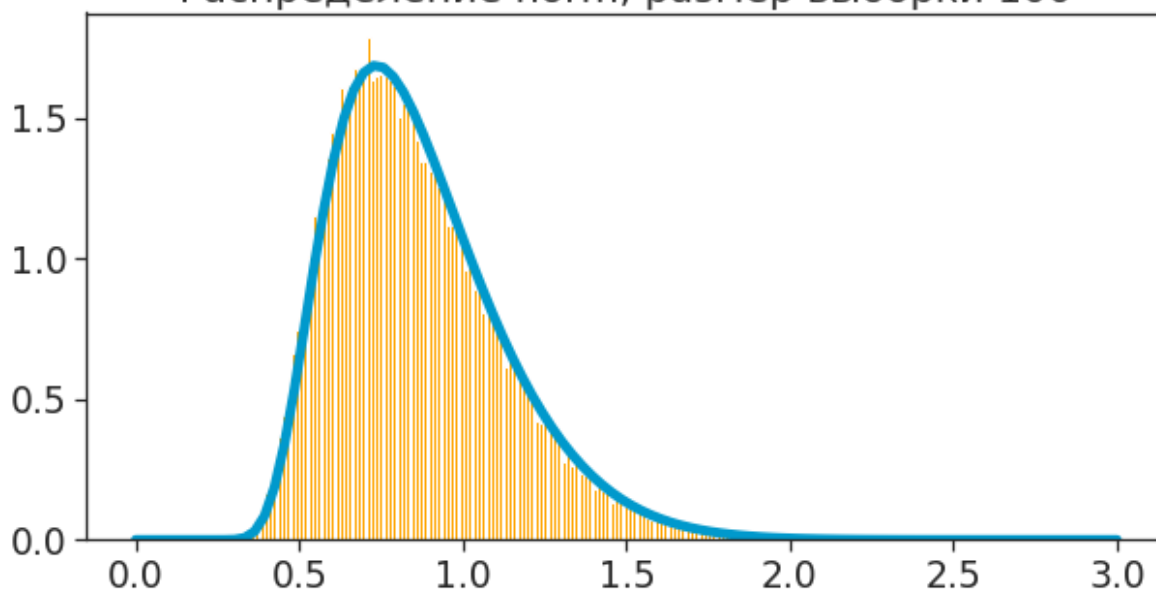
`kstest` (<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.kstest.html#scipy.stats.kstest>)
(`rvs`, `cdf`, `args=()`)

- `rvs` --- выборка
- `cdf` --- функция распределения (сама функция или ее название)
- `args` --- параметры распределения

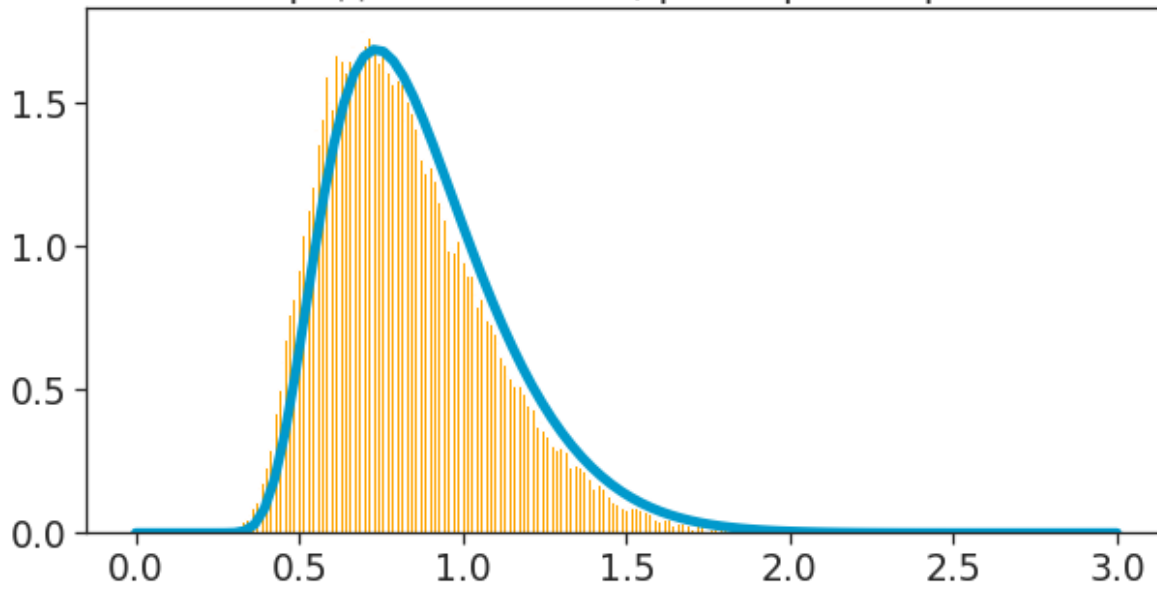
In [8]:

```
1 grid = np.linspace(0, 3, 100)
2 n_iter = 100000
3
4 for size, distr in zip([100, 10, 5, 10], [sps.norm] * 3 + [sps.expon]):
5     D = []
6
7     samples = distr.rvs(size=(n_iter, size))
8     for i in tqdm_notebook(range(n_iter), leave=False):
9         D.append(sps.kstest(samples[i], cdf=distr.cdf)[0])
10
11 plt.figure(figsize=(10, 5))
12 plt.plot(grid, sps.kstwobign.pdf(grid),
13          lw=5, label='Kolmogorov pdf', color=blue)
14 plt.hist(np.array(D) * np.sqrt(size), bins=200, density=True,
15          label='hist  $\sqrt{n} D_n$ ', color='orange')
16 plt.title('Точность приближения в критерии Колмогорова\n' \
17           + 'Распределение {}, размер выборки {}'.format(distr.name, size))
18 plt.show()
```

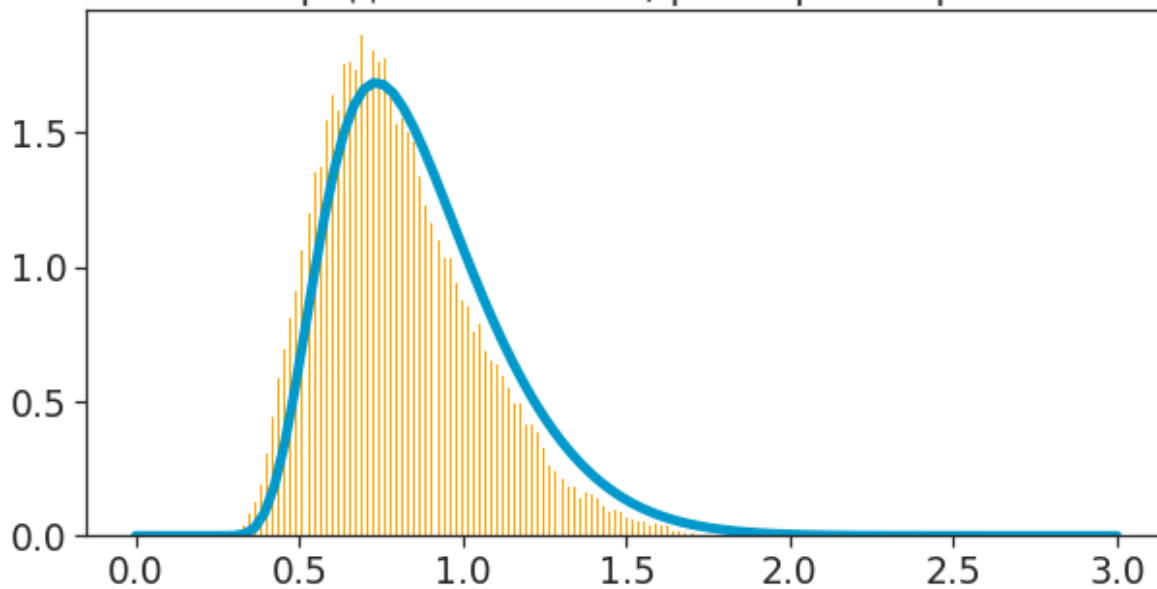
Точность приближения в критерии Колмогорова
Распределение норм, размер выборки 100



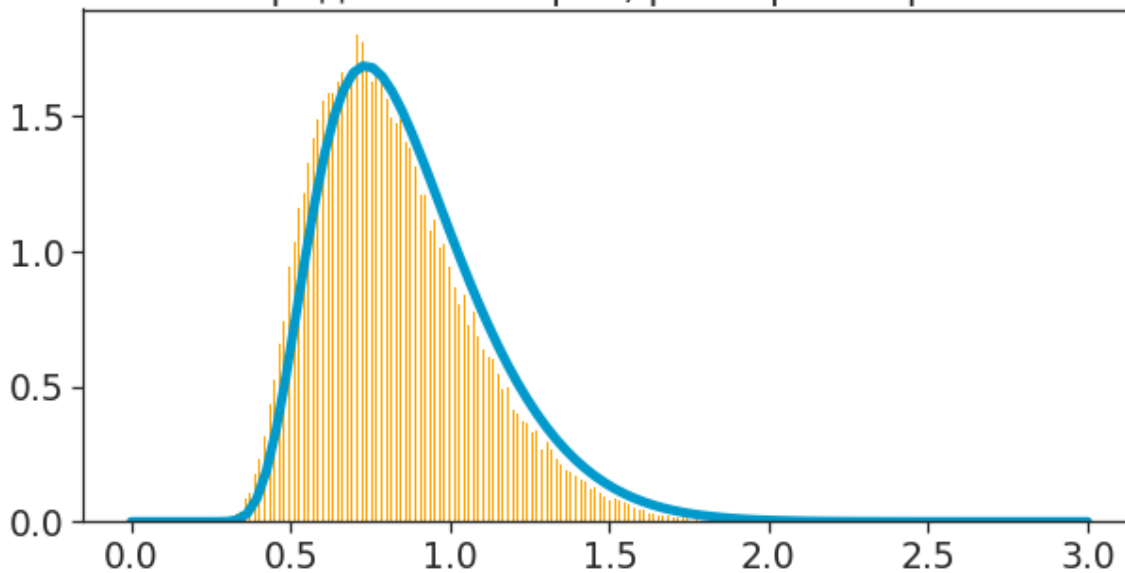
Точность приближения в критерии Колмогорова
Распределение норм, размер выборки 10



Точность приближения в критерии Колмогорова
Распределение норм, размер выборки 5



Точность приближения в критерии Колмогорова Распределение exron, размер выборки 10



Рассмотрим случай $H_0: \mathcal{N}(0, 1)$

In [9]:

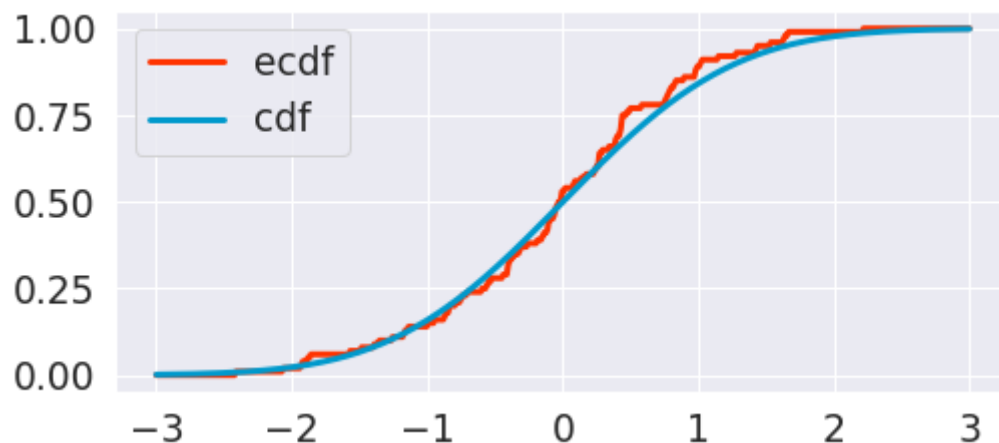
```
1 sns.set(style='darkgrid', font_scale=1.7)
2
3 def apply_kstest(sample, cdf=sps.norm.cdf):
4     print(sps.kstest(sample, sps.norm.cdf))
5
6     ecdf = ECDF(sample)
7     grid = np.linspace(-3, 3, 500)
8
9     plt.figure(figsize=(8, 3.5))
10    plt.plot(grid, ecdf(grid), color=red, label='ecdf', lw=3)
11    plt.plot(grid, cdf(grid), color=blue, label='cdf', lw=3)
12    plt.legend()
13    plt.show()
```

Нулевая гипотеза верна:

In [10]:

```
1 apply_kstest(sps.norm.rvs(size=100))
```

KstestResult(statistic=0.08237600150211866, pvalue=0.4877223292527129)

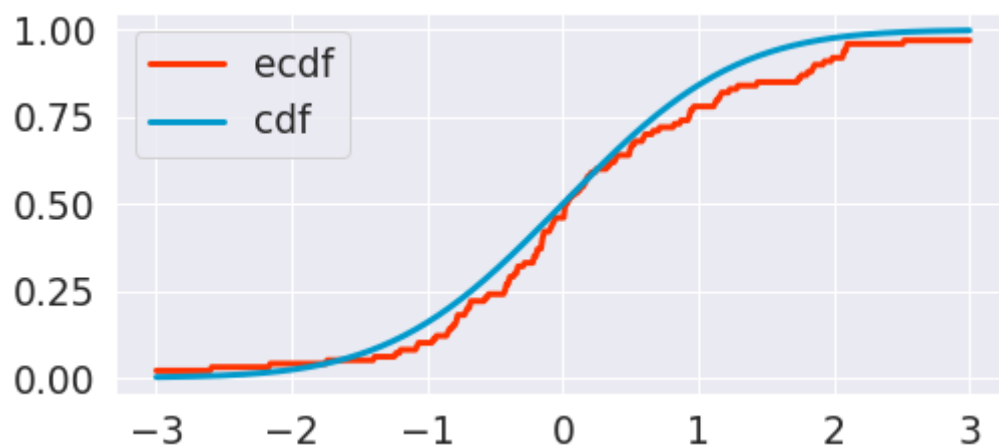


Истинное распределение - Лаплас

In [11]:

```
1 apply_kstest(sps.laplace.rvs(size=100))
```

KstestResult(statistic=0.10815556124535541, pvalue=0.17914691700639876)

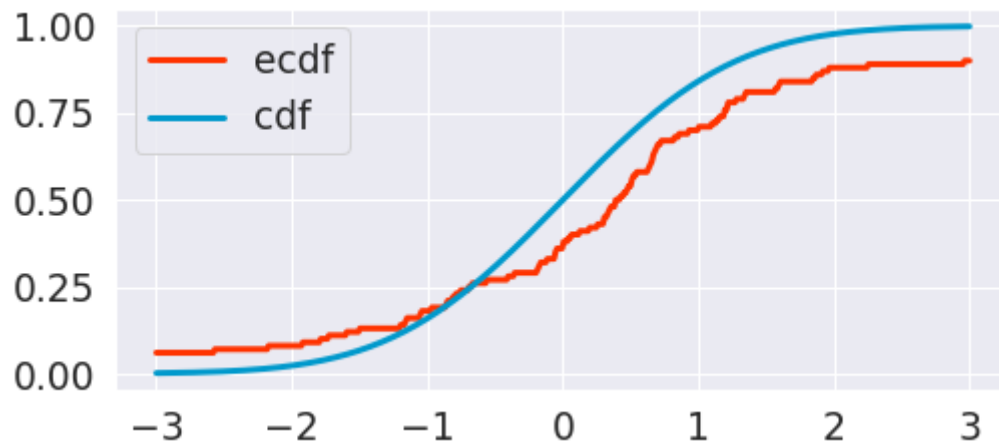


Истинное распределение - Коши

In [12]:

```
1 apply_kstest(sps.cauchy.rvs(size=100))
```

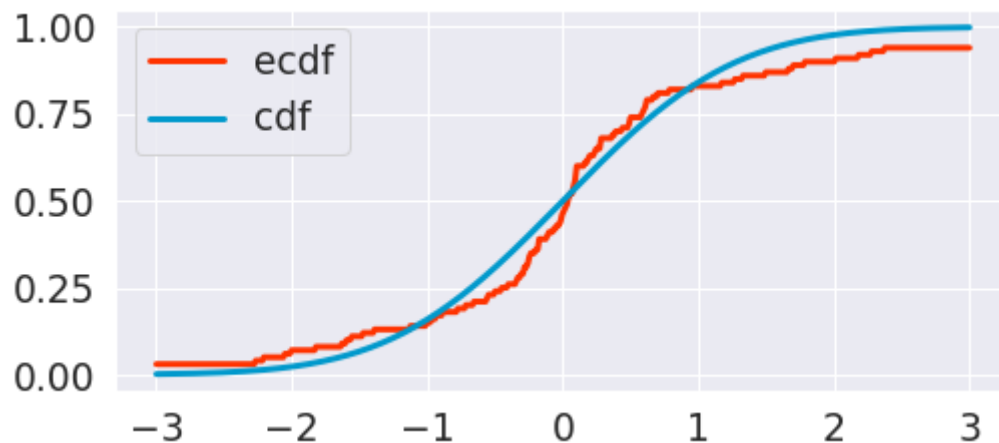
KstestResult(statistic=0.18756670386283075, pvalue=0.001489212227936959)



In [13]:

```
1 apply_kstest(sps.cauchy(scale=0.5).rvs(size=100))
```

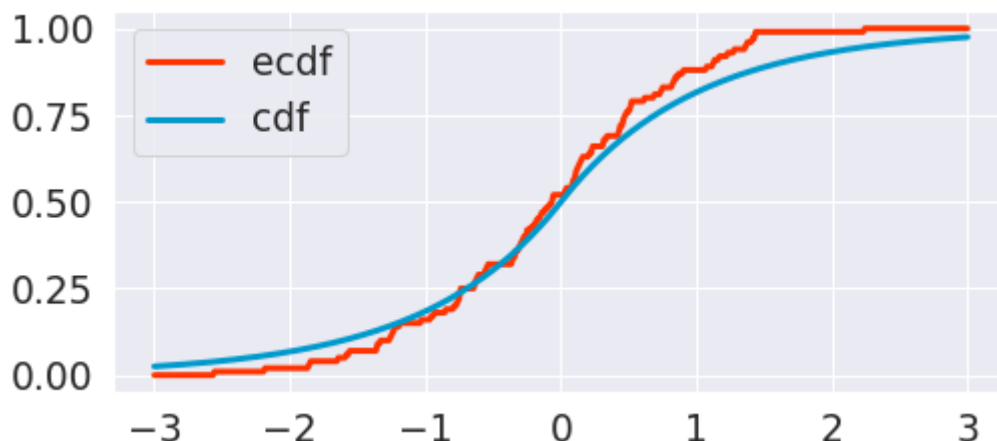
KstestResult(statistic=0.10255129072305458, pvalue=0.2278742984580734)



In [14]:

```
1 apply_kstest(sps.norm.rvs(size=100), cdf=sps.laplace.cdf)
```

```
KstestResult(statistic=0.09056612392589924, pvalue=0.3652768891115752)
```



Вывод: Критерий не чувствителен к хвостам распределений.

Критерий Андерсона-Дарлинга

H_0 : Выборка из некоторого класса распределений

H_1 : H_0 не верна.

[anderson](#)

<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.anderson.html#scipy.stats.anderson>) (x, dist='norm')

- x --- выборка
- dist : {'norm', 'expon', 'logistic', 'gumbel'} --- распределение

Возвращает (statistic, critical_values, significance_level). В critical_values записаны пороги для критериев, которым соответствуют уровни значимости из significance_level в процентах. То есть критерий {statistic > critical_values[i]} имеет уровень значимости significance_level[i].

In [15]:

```
1 sample = sps.norm.rvs(size=100)
2 sps.anderson(sample)
```

Out[15]:

```
AndersonResult(statistic=0.3551151814259441, critical_values=array([0.555, 0.632, 0.759, 0.885, 1.053]), significance_level=array([15. , 10. , 5. , 2.5, 1. ]))
```

In [16]:

```
1 sample = sps.laplace.rvs(size=100)
2 sps.anderson(sample)
```

Out[16]:

```
AndersonResult(statistic=2.1492631844960783, critical_values=array([0.555, 0.632, 0.759, 0.885, 1.053]), significance_level=array([15. , 10. , 5. , 2.5, 1. ]))
```

In [17]:

```
1 sample = sps.cauchy.rvs(size=100)
2 sps.anderson(sample)
```

Out[17]:

```
AndersonResult(statistic=30.650605026187066, critical_values=array([0.555, 0.632, 0.759, 0.885, 1.053]), significance_level=array([15. , 10. , 5. , 2.5, 1. ]))
```

Q-Q plot

[probplot](#)

(<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.probplot.html#scipy.stats.probplot>) (x, sparams=(), dist='norm', fit=True, plot=None)

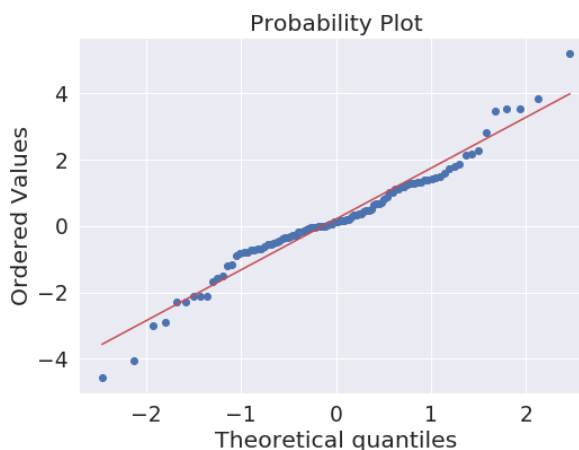
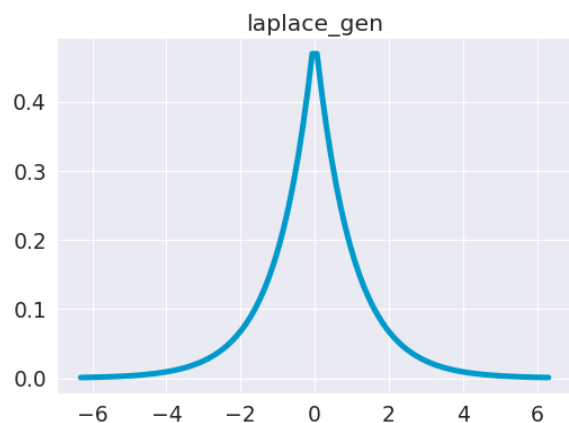
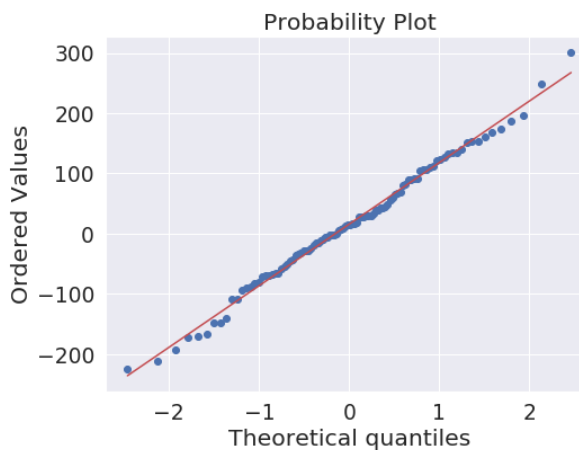
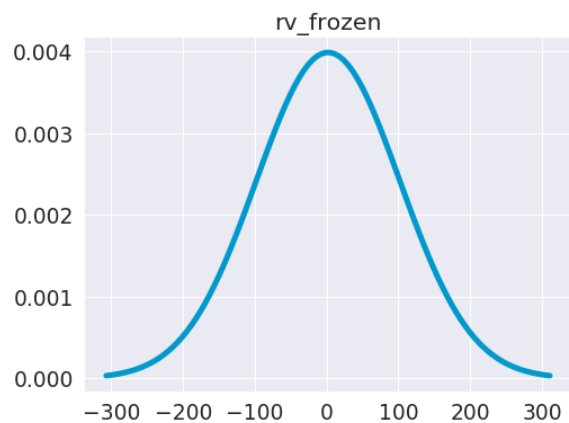
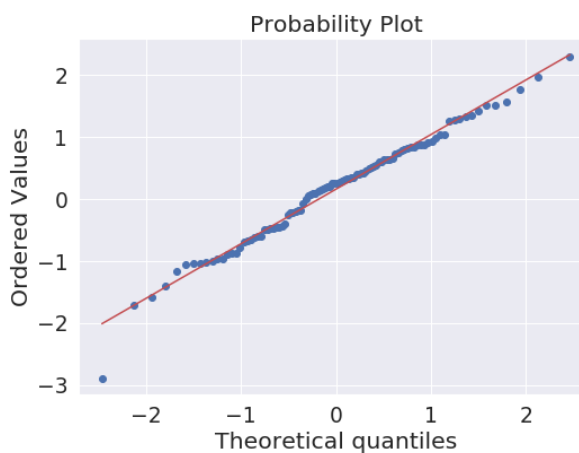
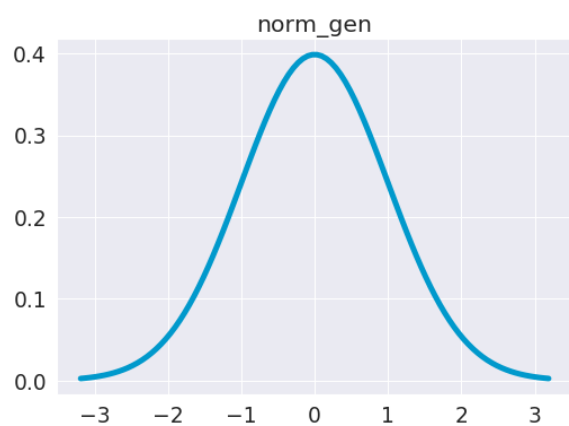
- x --- выборка
- sparams --- параметры распределения
- dist --- распределение
- fit --- подгонять ли линейную регрессию под точки
- plot --- рисовать на уже существующей фигуре

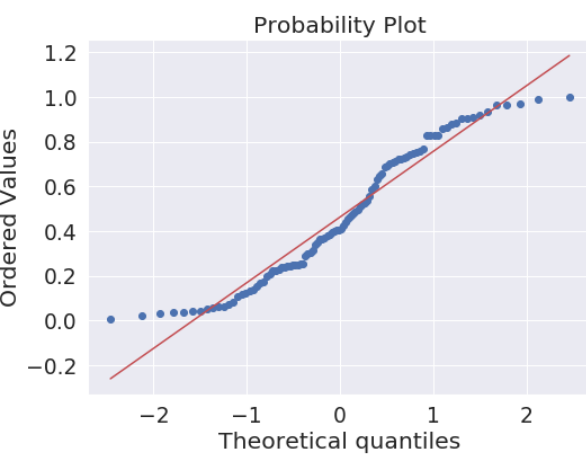
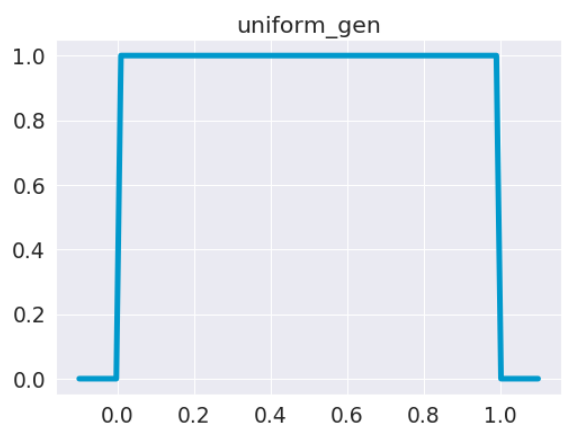
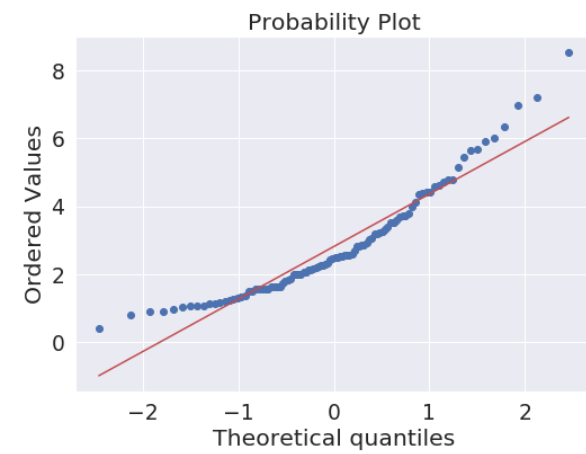
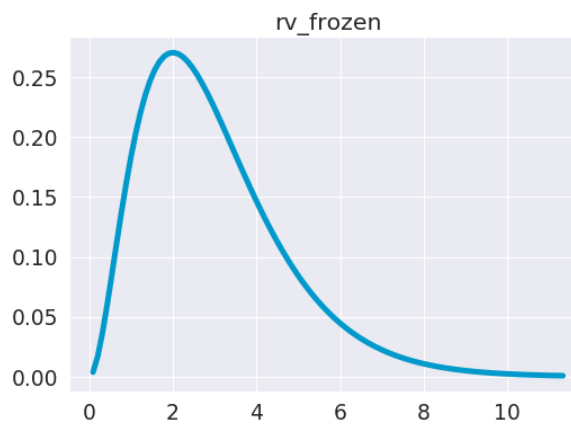
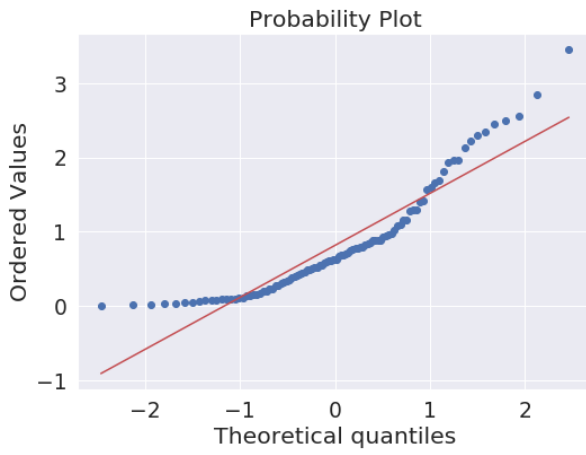
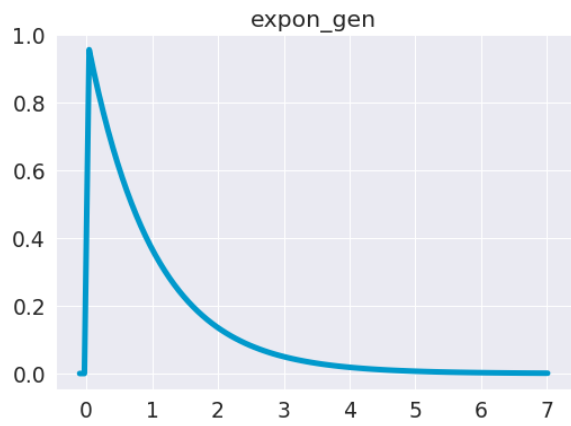
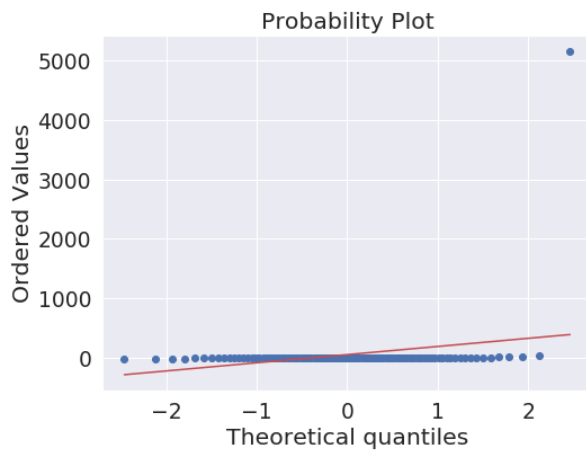
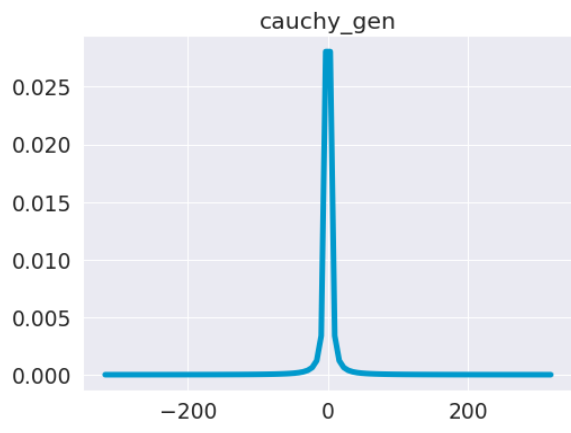
Возвращает:

- (osm, osr) --- теоретические и выборочные квантили
- (slope, intercept, r) --- не возвращается в случае fit=False и plot=None. Значения: коэффициенты линейной регрессии и RSS

In [18]:

```
1 for distr in [sps.norm, sps.norm(2, 100), sps.laplace, sps.cauchy,
2               sps.expon, sps.gamma(a=3), sps.uniform]:
3     sample = distr.rvs(size=100)
4
5     plt.figure(figsize=(15, 6))
6
7     ax = plt.subplot(1, 2, 1)
8     grid = np.linspace(distr.ppf(0.001)-0.1, distr.ppf(0.999)+0.1, 100)
9     plt.plot(grid, distr.pdf(grid), lw=5, color=blue)
10    plt.title(str(distr).split(' ')[0].split('.')[0])
11
12    ax = plt.subplot(1, 2, 2)
13    sps.probplot(sample, plot=ax)
14
15    plt.tight_layout()
16    plt.show()
```





Критерий Шапиро-Уилка

H_0 : Выборка из нормального распределения

H_1 : H_0 : не верна.

`shapiro_` (<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.shapiro.html#scipy.stats.shapiro>)

(x): W, p-value

In [19]:

```
1 sps.shapiro(sps.norm.rvs(size=100))
```

Out[19]:

```
(0.9901418089866638, 0.6760627031326294)
```

In [20]:

```
1 sps.shapiro(sps.norm(20, 100).rvs(size=100))
```

Out[20]:

```
(0.991076648235321, 0.7510301470756531)
```

In [21]:

```
1 sps.shapiro(sps.laplace.rvs(size=100))
```

Out[21]:

```
(0.9878672957420349, 0.498473197221756)
```

In [22]:

```
1 sps.shapiro(sps.cauchy.rvs(size=100))
```

Out[22]:

```
(0.2799574136734009, 6.384199354729874e-20)
```

Критерий Жарка-Бера

H_0 : Выборка из нормального распределения

H_1 : H_0 не верна.

`jarque_bera`

(https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.jarque_bera.html#scipy.stats.jarque_bera)

(x): jb_value, p-value

In [23]:

```
1 sps.jarque_bera(sps.norm.rvs(size=100))
```

Out[23]:

```
(7.107976892375481, 0.028610301142788175)
```

In [24]:

```
1 sps.jarque_bera(sps.norm(20, 100).rvs(size=100))
```

Out[24]:

```
(0.9921562789026134, 0.608914059003439)
```

In [25]:

```
1 sps.jarque_bera(sps.laplace.rvs(size=100))
```

Out[25]:

```
(21.054147137236328, 2.6800940745586033e-05)
```

In [26]:

```
1 sps.jarque_bera(sps.cauchy.rvs(size=100))
```

Out[26]:

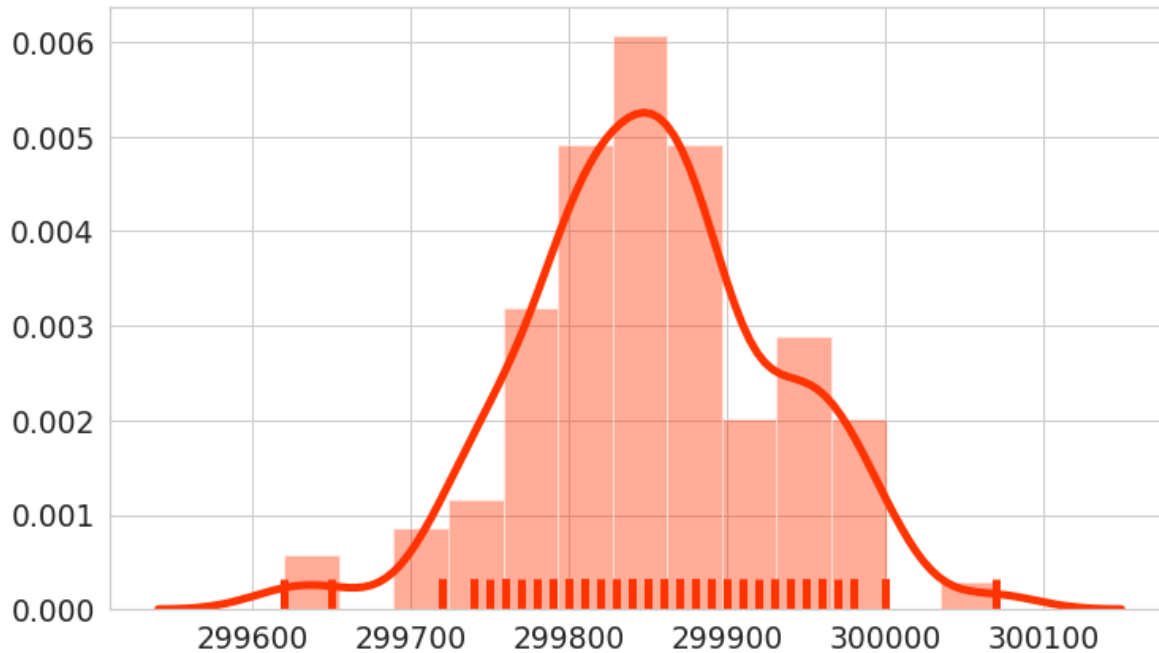
```
(28338.26644038969, 0.0)
```

Пример: эксперимент Майкельсона

Данные классического эксперимента Майкельсона по измерению скорости света с помощью вращающегося зеркала, 100 наблюдений:

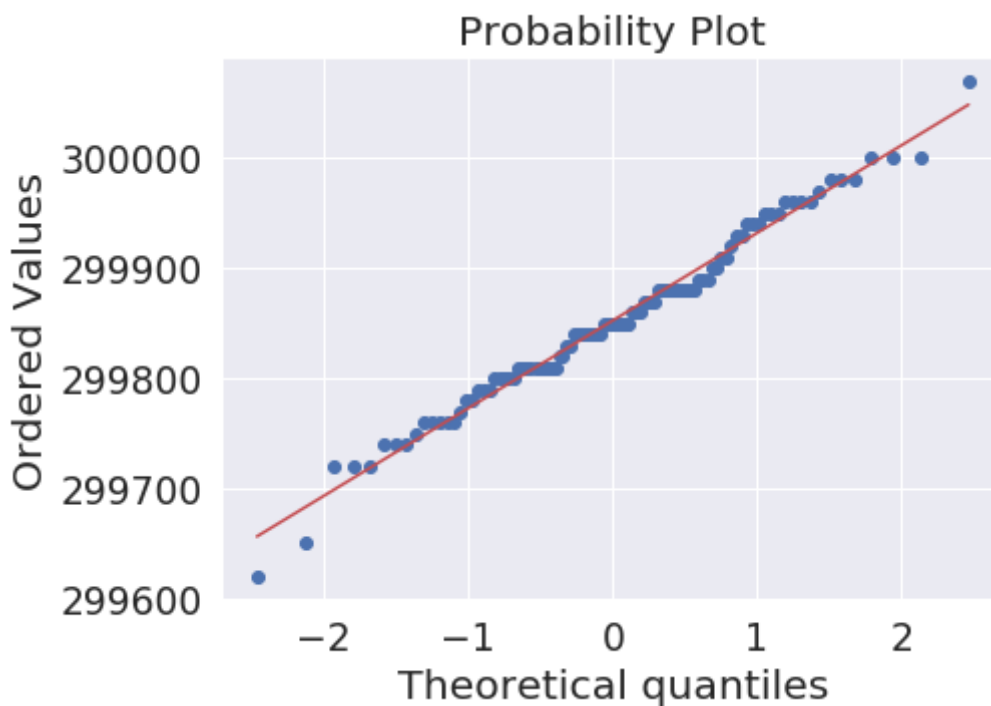
In [27]:

```
1 speed = np.loadtxt('speed.txt')
2
3 sns.set(style='whitegrid', font_scale=1.7)
4 plt.figure(figsize=(12, 7))
5 sns.distplot(speed, rug=True,
6               color=red, kde_kws={'lw': 5}, rug_kws={'lw': 5});
```



In [28]:

```
1 sns.set(style='darkgrid', font_scale=1.7)
2 plt.figure(figsize=(7, 5))
3 sps.probplot(speed, plot=plt.subplot(111));
```



In [29]:

```
1 pvalues = [sps.kstest(speed, sps.norm(*sps.norm.fit(speed)).cdf)[1],
2           sps.shapiro(speed)[1],
3           sps.jarque_bera(speed)[1]]
4
5 pvalues
```

Out[29]:

```
[0.4813808419520875, 0.5140784382820129, 0.8628948044152577]
```

In [30]:

```
1 multipletests(pvalues)
```

Out[30]:

```
(array([False, False, False]),
 array([0.86050917, 0.86050917, 0.8628948 ]),
 0.016952427508441503,
 0.016666666666666666)
```

Статистика, прикладной поток 2019

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