

# Статистика, прикладной поток

## Практическое задание 5

В данном задании вы исследуете некоторые свойства доверительных интервалов и байесовских оценок.

### Правила:

- Дедлайн **25 ноября 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 5". Квадратные скобки обязательны. Если письмо дошло, придет ответ от автоответчика.
- Прислать нужно ноутбук и его pdf-версию (без архивов). Названия файлов должны быть такими: `5.N.ipynb` и `5.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 - 5 баллов **O2**
- Задача 2 - 5 баллов **O2**
- Задача 3 - 5 баллов **O2**
- Задача 4 - 5 баллов **O2**
- Задача 5 - 7 баллов **O2**
- Задача 6 - 7 баллов **O2**
- Задача 7 - 7 баллов **O2**
- Задача 8 - 10 баллов **O3**
- Задача 9 - 6 баллов **O2**

In [2]:

```
1 import numpy as np
2 import pandas as pd
3 import scipy.stats as sps
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import warnings
7
8 sns.set()
9 warnings.filterwarnings("ignore")
10
11 %matplotlib inline
```

## Доверительные интервалы

### Задача 1.

В этой задаче нужно визуализировать доверительные интервалы для выборок из различных распределений.

Чтобы не плодить код, напишите следующую функцию (см. ниже). Пример построения есть в ноутбуке по `matplotlib`.

In [471]:

```
1 def draw_confidence_interval(
2     left, # левая граница интервала
3     right, # правая граница интервала
4     title,
5     estimation=None, # если задана, то рисуется график оценки
6     sample=None, # если задано, то рисуются точки выборки
7     ylim=(None, None), # ограничение по оси y
8 ):
9     time = np.arange(1, len(left) + 1)
10    plt.figure(figsize=(8, 8))
11    plt.title(title)
12    plt.ylim(ylim)
13    if (estimation is not None):
14        plt.plot(time, estimation, label='estimation')
15    if (sample is not None):
16        plt.scatter(time, sample, alpha =0.2, label='sample')
17    plt.fill_between(time, left, right, alpha=0.5, label='interval')
18    plt.legend()
19    plt.show()
```

Рассмотрим следующие ситуации:

1. Выборка из распределения  $\mathcal{N}(0, 1)$ ; точный доверительный интервал минимальной длины в параметрической модели  $\mathcal{N}(\theta, 1)$ .
2. Выборка из распределения  $U[0, 1]$ ; точный доверительный интервал минимальной длины в параметрической модели  $U[0, \theta]$  на основе статистики  $X_{(n)}$ .

Для каждой ситуации из перечисленных выше сгенерируйте выборку  $X_1, \dots, X_{100}$  и постройте график доверительных интервалов уровня доверия 0.95, вычисленных для всех подвыборок вида  $X_1, \dots, X_i$ ,  $1 \leq i \leq 100$ .

Постройте графики зависимости верхних и нижних границ интервала от размера выборки, используя написанную функцию. Нужно нанести на график точки выборки. Для вычисления квантилей у каждого распределения из `scipy.stats` используйте функцию `ppf`.

## Решение

Точный доверительный интервал минимальной длины в параметрической модели  $\mathcal{N}(\theta, 1)$ :

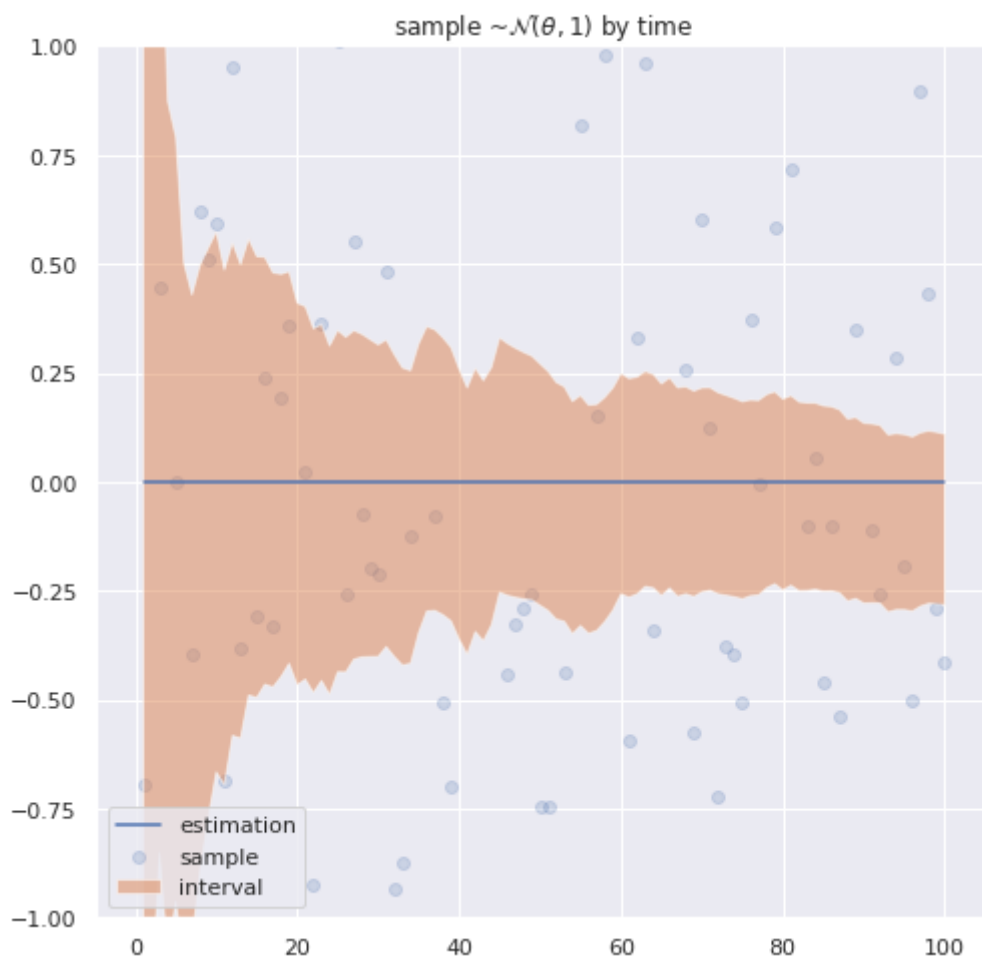
$$\left( \bar{X} - z_{\frac{1+\alpha}{2}} \cdot \frac{1}{\sqrt{n}}, \bar{X} + z_{\frac{1+\alpha}{2}} \cdot \frac{1}{\sqrt{n}} \right)$$

In [472]:

```

1 grid = np.arange(1, 101)
2 alpha = 0.95
3
4 sample = sps.norm().rvs(size = 100)
5
6 cum_mean = sample.cumsum()/grid
7
8 left = cum_mean - sps.norm.ppf((1 + alpha) / 2)/np.sqrt(grid)
9 right = cum_mean + sps.norm.ppf((1 + alpha) / 2)/np.sqrt(grid)
10
11 draw_confidence_interval(left, right,
12                           r'sample ~$\mathcal{N}(\theta, 1)$ by time',
13                           estimation=np.zeros(len(left)),
14                           sample=sample, ylim=(-1, 1))

```

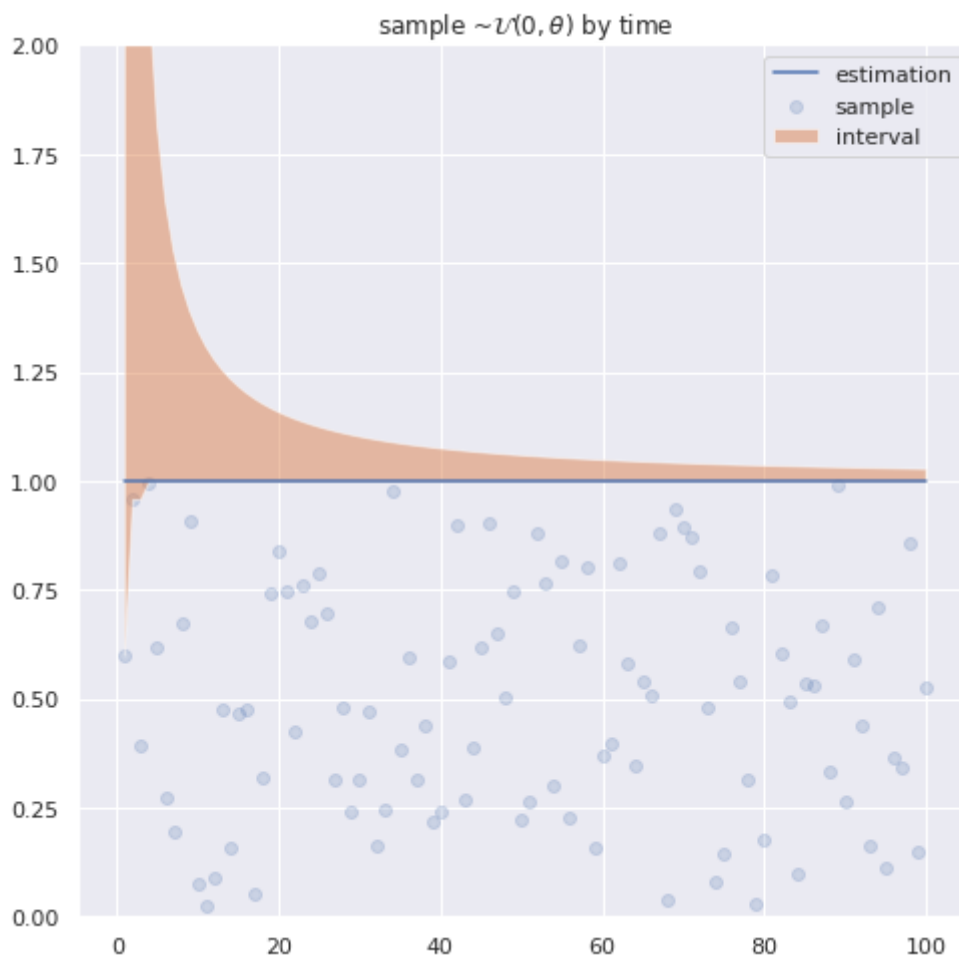


In [473]:

```

1 sample = sps.uniform.rvs(size=100)
2
3 cum_max = np.maximum.accumulate(sample)
4 left = cum_max
5 right = cum_max/(1-alpha)**(1/grid)
6
7 draw_confidence_interval(left, right,
8                           r'sample ~ $\mathcal{U}(0, \theta)$ by time',
9                           estimation=np.ones(len(left)),
10                          sample=sample, ylim=(0,2))
11

```



**Вывод:** Мы видим, что для обеих распределений их точные доверительные интервалы с ростом выборки сходятся к значению оцениваемого параметра. Причем интервал для  $U[0, \theta]$  сходится гораздо быстрее, так как скорость сходимости у него степенная, в отличие от нормального, сходимости у которого  $\frac{1}{n^{(1/2)}}$

## Задача 2.

Аналогично задаче 1 сгенерируйте выборку  $X_1, \dots, X_{100}$  из распределения  $\Gamma(3, 2)$  и постройте доверительные интервалы для следующих случаев:

- точный асимптотический доверительный интервал в параметрической модели  $\Gamma(\theta, 2)$ ; точки выборки наносить на график не нужно;

- точный асимптотический доверительный интервал для  $\theta$  в параметрической модели  $\Gamma(\theta, \beta)$ , причем  $\beta$  неизвестно.

Изобразите интервалы *на одном* графике полупрозрачными цветами. Точки выборки наносить на график не нужно.

### Решение:

- точный асимптотический доверительный интервал в парам. модели  $\Gamma(\theta, \beta)$  при известном  $\beta$  :

$$\left( \frac{\beta}{\bar{X}} - \sqrt{\frac{\beta}{n}} \cdot \frac{z_{\frac{1+\alpha}{2}}}{\bar{X}}, \frac{\beta}{\bar{X}} + \sqrt{\frac{\beta}{n}} \cdot \frac{z_{\frac{1+\alpha}{2}}}{\bar{X}} \right)$$

- точный асимптотический доверительный интервал для  $\theta$  в парам. модели  $\Gamma(\theta, \beta)$ ,  $\beta$  неизвестно.

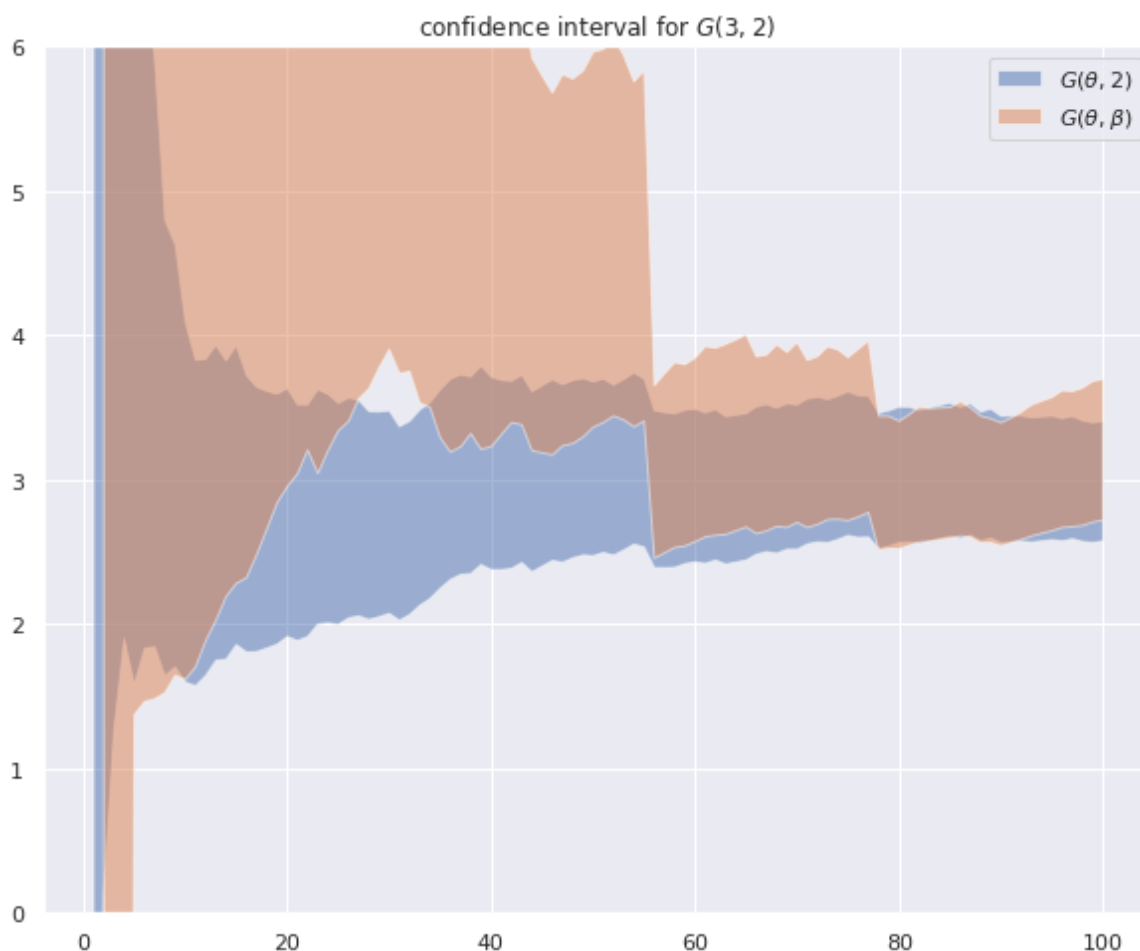
$$\left( \frac{\bar{X}}{S^2} - \frac{z_{\frac{1+\alpha}{2}} \sigma \left( \frac{\bar{X}}{S^2}, \frac{\bar{X}^2}{S^2} \right)}{\sqrt{n}}, \frac{\bar{X}}{S^2} + \frac{z_{\frac{1+\alpha}{2}} \sigma \left( \frac{\bar{X}}{S^2}, \frac{\bar{X}^2}{S^2} \right)}{\sqrt{n}} \right)$$

In [478]:

```

1 sample = sps.gamma(a=2, scale = 1/3).rvs(size=100)
2
3 cum_mean = sample.cumsum()/grid
4
5 Z = sps.norm.ppf((1 + alpha) / 2)
6
7 left1 = (2 - Z*np.sqrt(2/grid)) / cum_mean
8 right1 = (2 + Z*np.sqrt(2/grid)) / cum_mean
9
10 estimation = cum_mean / ((sample**2).cumsum()/grid - cum_mean**2)
11
12 variance = estimation * np.sqrt(2 - 3/(cum_mean*estimation))
13
14 left2 = estimation - Z*variance / np.sqrt(grid)
15 right2 = estimation + Z*variance / np.sqrt(grid)
16
17 plt.figure(figsize=(10, 8))
18
19 plt.title(r"confidence interval for $G(3, 2)$")
20 plt.ylim((0,6))
21
22 plt.fill_between(grid, left1, right1, alpha=0.5,
23                 label = r'$G(\theta, 2)$')
24
25 plt.fill_between(grid, left2, right2, alpha=0.5,
26                 label = r'$G(\theta, \beta)$')
27
28 plt.legend()
29 plt.show()

```



Сравните полученные интервалы.

**Вывод:** Здесь мы видим, что оба наших интервала с ростом выборки лучше и лучше начинают оценивать параметр, что подтверждает теорические свойства этих асимптотических дов. интервалов. Но при этом мы можем видеть, что интервал для  $G(\theta, 2)$ , то есть при известном втором параметре гораздо устойчивей, и гораздо "равномерней" сходится, нежели интервал, где второй параметр неизвестен.

---

### Задача 3.

Аналогично задаче 1 сгенерируйте выборку  $X_1, \dots, X_{100}$  из стандартного распределения Коши и постройте доверительные интервалы для следующих случаев

- точный доверительный интервал минимальной длины в параметрической модели  $\mathcal{N}(\theta, 1)$ ;
- точный асимптотический доверительный интервал в параметрической модели распределения Коши со сдвигом, используя выборочную медиану;
- точный асимптотический доверительный интервал в параметрической модели распределения Коши со сдвигом, используя асимптотически эффективную оценку.

Изобразите интервалы *на одном* графике полупрозрачными цветами. Точки выборки нужно нанести на график.

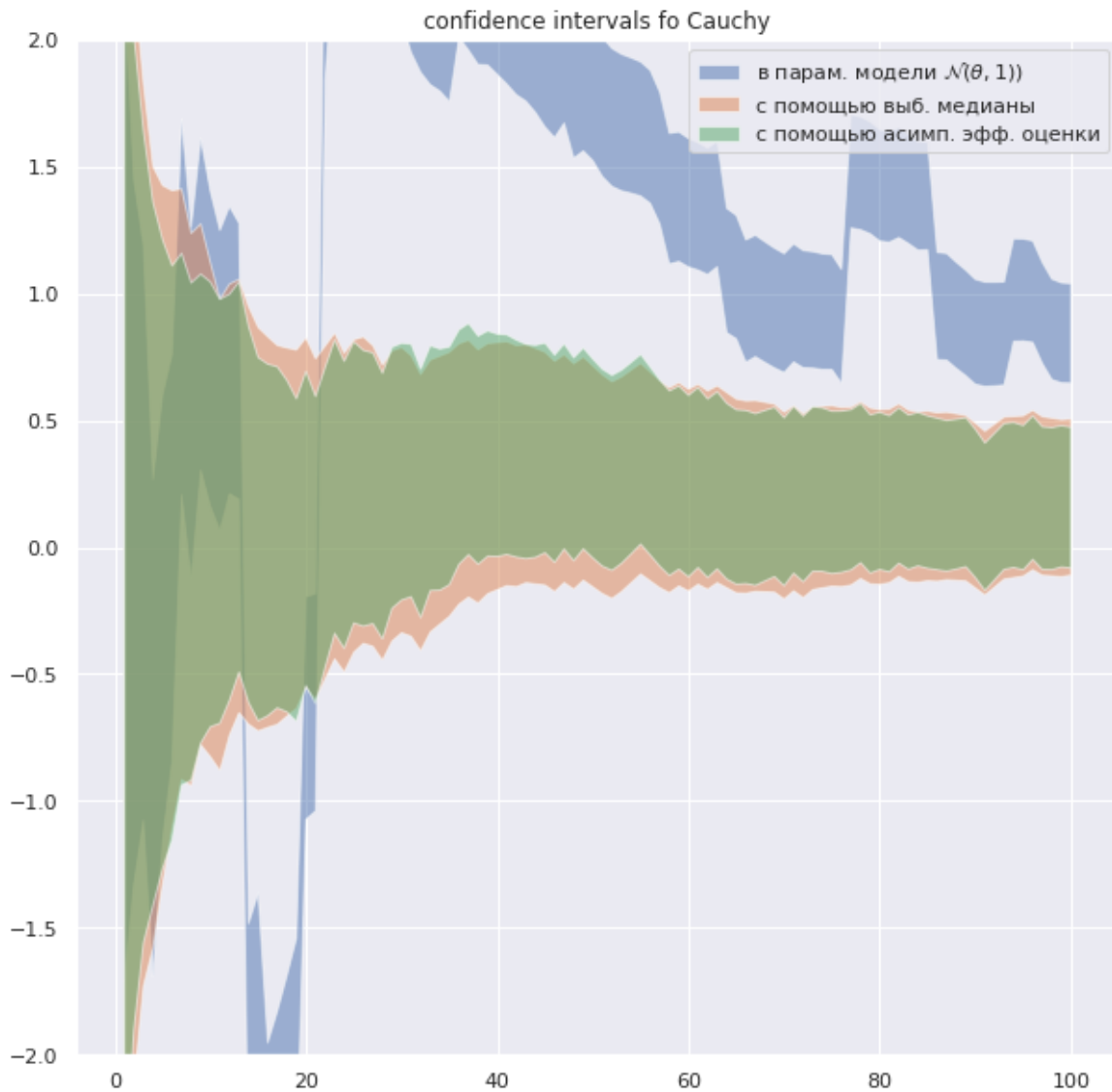
In [489]:

```

1  grid = np.arange(1, 101)
2  alpha = 0.95
3
4  sample = sps.cauchy.rvs(size=100)
5
6  cum_mean = sample.cumsum()/grid
7  cum_median = np.array([np.median(sample[:i+1]) for i in range(len(sample))])
8
9
10 quantile = sps.norm.ppf((1 + alpha) / 2)
11
12 #1ый доверительный интервал
13 left1 = cum_mean - quantile/np.sqrt(grid)
14 right1 = cum_mean + quantile/np.sqrt(grid)
15
16 #2й доверительный интервал
17 left2 = cum_median - np.pi*quantile/(2*np.sqrt(grid))
18 right2 = cum_median + np.pi*quantile/(2*np.sqrt(grid))
19
20 #3ий доверительный интервал
21 newtor_iter = list()
22 for i in range(100):
23     med = cum_median[i]
24     numerator = ((sample[:i+1] - med)/(1 + (sample[:i+1] - med)**2)).sum()
25     denominator=((1-(sample[:i+1]-med)**2)/(1+(sample[:i+1]-med)**2)**2).sum()
26     newtor_iter.append(numerator/denominator)
27
28 newtor_iter = np.array(newtor_iter)
29
30 left3 = cum_median - newtor_iter - quantile*np.sqrt(2/grid)
31 right3 = cum_median - newtor_iter + quantile*np.sqrt(2/grid)
32
33
34 plt.figure(figsize=(10, 10))
35
36 plt.title("confidence intervals fo Cauchy")
37 plt.ylim((-2, 2))
38
39 plt.fill_between(grid, left1, right1, alpha=0.5,
40                 label=r'в парам. модели  $\mathcal{N}(\theta, 1)$ ')
41 plt.fill_between(grid, left2, right2, alpha=0.5,
42                 label='с помощью выб. медианы')
43 plt.fill_between(grid, left3, right3, alpha=0.5,
44                 label='с помощью асимп. эфф. оценки')
45
46
47 plt.legend()
48 plt.show()

```





Сравните полученные интервалы.

**Вывод:** Мы можем видеть, что любой произвольный интервал, какой нам вздумается, не подойдет под оценку, например как в нашем случае с распределением Коши мы можем видеть, что интервал, используемый для модели  $\mathcal{N}(\theta, 1)$  нам совершенно не подходит - он ни к чему не сходится, хотя бы потому что у него нет матожидания. Кроме того мы сравнили интервалы, полученные из медианы, и из асимпт. эфф. оценки, которую мы нашли с помощью итерации в методе Ньютона из оценки для медианы. Можем видеть, что в принципе ничего особенного не поменялось, видимо существенные различия появляются на больших выборках, но тем не менее иногда даже на нашей выборке асимптотически эфф. дов. интервал все-таки немного более узкий, чем интервал, полученный с помощью медианы.

#### Задача 4.

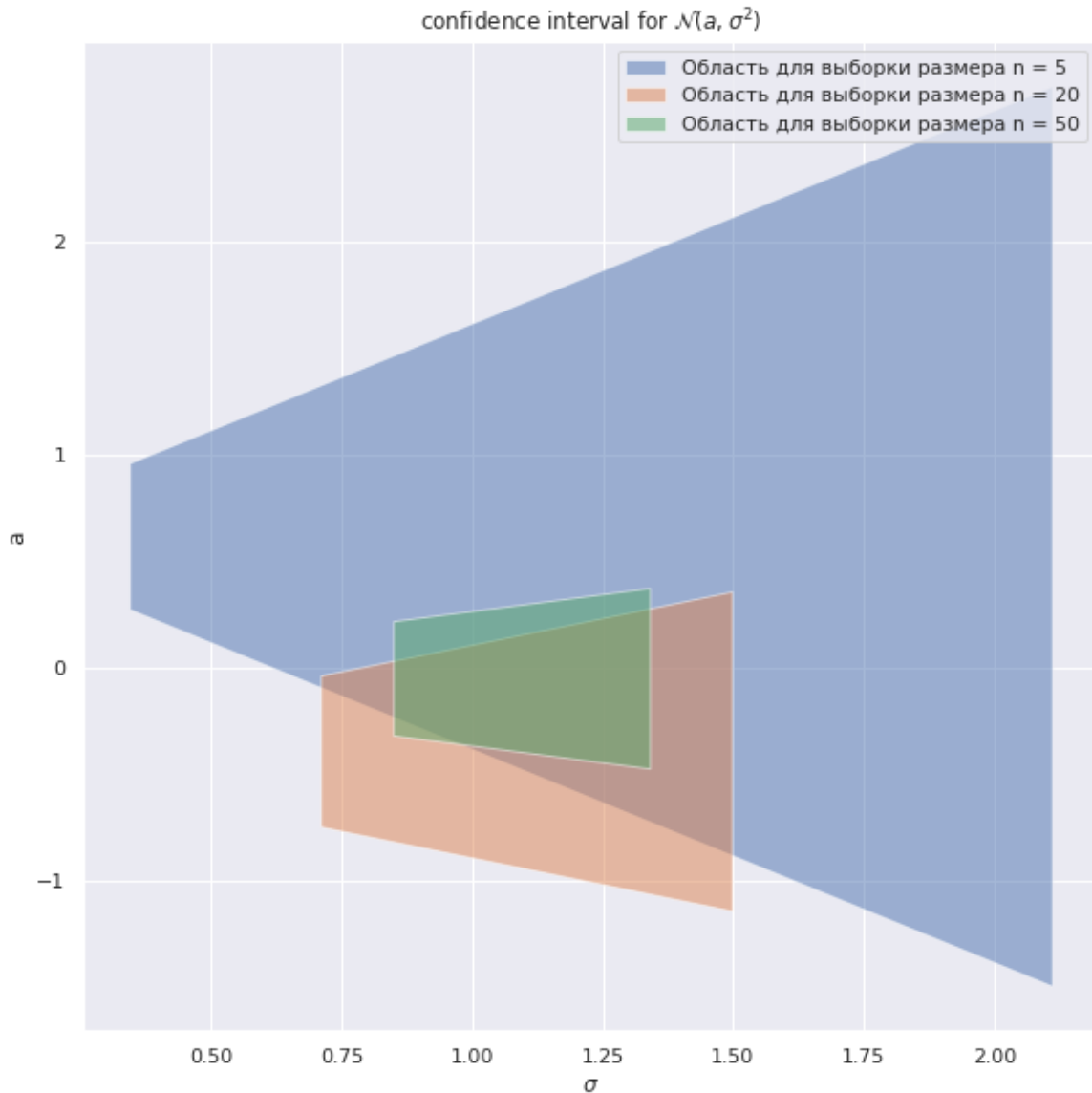
Пусть  $X_1, \dots, X_n$  --- выборка из распределения  $\mathcal{N}(a, \sigma^2)$ . Постройте точную доверительную область для параметра  $\theta = (a, \sigma)$  уровня доверия  $\alpha = 0.95$  для сгенерированной выборки размера  $n \in \{5, 20, 50\}$  из стандартного нормального распределения.

In [96]:

```

1  # функция доверительной области для выборки из нормального
2  def confidence_region(sample):
3      S2 = (sample**2).mean() - (sample.mean())**2
4      left_sigma = ((n * S2) / sps.chi2(n-1).ppf((1+alpha**0.5)/2)) ** 0.5
5      right_sigma = ((n * S2) / sps.chi2(n-1).ppf((1-alpha**0.5)/2)) ** 0.5
6
7      sigmas = np.linspace(left_sigma, right_sigma, 500)
8
9      quantile = sps.norm().ppf((1 + alpha**0.5) / 2)
10
11     left_a = sample.mean() - sigmas * quantile / (n ** 0.5)
12     right_a = sample.mean() + sigmas * quantile / (n ** 0.5)
13
14     return (sigmas, left_a, right_a)
15
16
17 ns = [5, 20, 50]
18
19
20 plt.figure(figsize=(10, 10))
21 plt.title(r"confidence interval for $\mathcal{N}(a, \sigma^2)$")
22
23 plt.xlabel(r"$\sigma$")
24 plt.ylabel(r"a")
25
26 for n in ns:
27     sample = sps.norm.rvs(size=n)
28
29     sigmas, left_a, right_a = confidence_region(sample)
30
31     plt.fill_between(sigmas, left_a, right_a, alpha=0.5,
32                     label="Область для выборки размера n = {}".format((n)))
33
34 plt.legend()
35 plt.show()

```



**Вывод:** Здесь мы можем видеть, что площадь доверительной области уменьшается с ростом размера выборки, и тем не менее в любом случае получается, что точка, соответствующая истинным значением параметра (то есть  $(0, 1)$ ) принадлежит доверительным областям.

## Задача 5.

В данном задании вам нужно изучить доверительные интервалы для параметра сдвига в нормальной модели в случае неизвестной дисперсии. Требуется построить асимптотический доверительный интервал (через центральную предельную теорему и лемму Slutsky) и точный неасимптотический (через распределения хи-квадрат и Стьюдента).

Вывод этих интервалов был разобран на семинарах. Выпишите только ответы.

Асимптотический доверительный интервал:  $(\bar{X} - \frac{S}{\sqrt{n}} Z_{\frac{1+\alpha}{2}}; \bar{X} + \frac{S}{\sqrt{n}} Z_{\frac{1+\alpha}{2}})$

Точный доверительный интервал:  $(\bar{X} - \frac{S}{\sqrt{n-1}} T_{n-1, \frac{1+\alpha}{2}}; \bar{X} + \frac{S}{\sqrt{n-1}} T_{n-1, \frac{1+\alpha}{2}})$

Реализуйте функции построения этих интервалов по выборке.

In [51]:

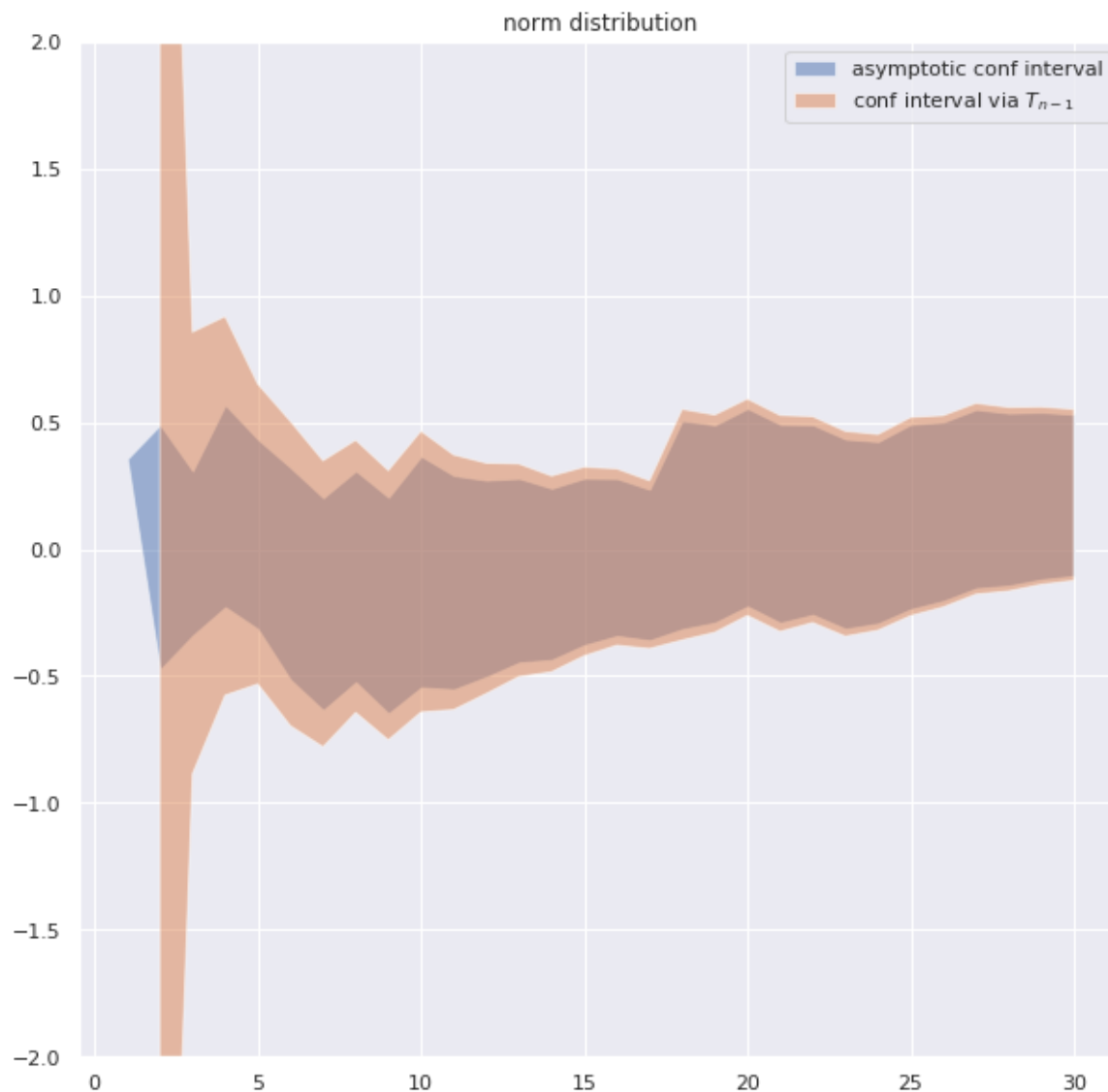
```
1 def calc_asympt_conf_interval(sample, alpha=0.95):
2     n = len(sample)
3     grid = np.arange(1, n + 1)
4     sample_mean = sample.mean()
5     Z = sps.norm().ppf((1 + alpha) / 2)
6     S = np.sqrt((sample**2).mean() - sample_mean**2)
7     return (sample_mean - S/np.sqrt(n)*Z,
8            sample_mean + S/np.sqrt(n)*Z)
9
10 def calc_conf_interval(sample, alpha=0.95):
11     n = len(sample)
12     grid = np.arange(1, n + 1)
13     sample_mean = sample.mean()
14     quantile1 = sps.t(df = n-1).ppf((1 + alpha) / 2)
15     quantile2 = sps.t(df = n-1).ppf((1 - alpha) / 2)
16     S = np.sqrt((sample**2).mean() - sample_mean**2)
17
18     return (sample_mean - S/np.sqrt(n-1)*quantile1,
19            sample_mean - S/np.sqrt(n-1)*quantile2)
```

Сгенерируйте выборку из нормального распределения и сравните два доверительных интервала в зависимости от размера выборки. Для сравнения отобразите оба интервала на одном графике. Поясните теоретическую причину такого поведения доверительных интервалов.

Указание: рассматривайте длину выборки около 20-30.

In [60]:

```
1 grid = np.arange(1, 31)
2 sample = sps.norm().rvs(size=30)
3
4 cum_mean = sample.cumsum() / grid
5
6 vect_asympt = np.vectorize(lambda x:alc_asympt_conf_interval(x),
7                             signature='(n)->(2)')
8
9 Left1 = []
10 Right1 = []
11 Left2 = []
12 Right2 = []
13
14 for i in grid:
15     left, right = calc_asympt_conf_interval(sample[:i])
16     Left1.append(left)
17     Right1.append(right)
18
19     left, right = calc_conf_interval(sample[:i])
20     Left2.append(left)
21     Right2.append(right)
22
23
24
25 plt.figure(figsize=(10, 10))
26
27
28 plt.ylim((-2, 2))
29 plt.title("norm distribution")
30
31 plt.fill_between(grid, Left1, Right1, alpha=0.5,
32                  label='asymptotic conf interval')
33 plt.fill_between(grid, Left2, Right2, alpha=0.5,
34                  label='conf interval via  $T_{n-1}$ ')
35
36 plt.legend()
37 plt.show()
38
39
```



**Вывод:** Здесь мы можем видеть сравнение точного доверительного интервала, и асимптотич. довер. интервала. Если посмотреть на график, то получается странный результат - асимптотический доверительный интервал получается более узким, чем обычный дов. интервал. Видимо это связано с тем, что у нас неизвестная дисперсия, и мы ее должны оценить

Скачайте данные [wine dataset](http://archive.ics.uci.edu/ml/datasets/wine) (<http://archive.ics.uci.edu/ml/datasets/wine>) и рассмотрите столбцы Alkalinity of ash, Nonflavanoid phenols, Proanthocyanins и Hue для вина первого типа (за тип вина отвечает первый столбец).

Постройте доверительные интервалы для параметров сдвига каждого из столбцов, предполагая, что столбцы имеют нормальное распределение. Нужно построить доверительные интервалы обоих рассмотренных выше типов. Запишите их в виде таблицы.

In [116]:

```
1 data = pd.read_csv("wine.data")
2 columns = ["Alkalinity of ash", "Nonflavanoid phenols",
3           "Proanthocyanins", "Hue"]
4
5 data = data[data["Class of wine"] == 1]
6 data = data.loc[:, columns]
```

In [117]:

```

1 array =[]
2 for column in columns:
3     left1, right1 = calc_asympt_conf_interval(data[column])
4     left1 = round(left1, 4)
5     right1 = round(right1, 4)
6
7     left2, right2 = calc_conf_interval(data[column])
8     left2 = round(left2, 4)
9     right2 = round(right2, 4)
10
11     array.append([(left1, right1),(left2, right2)])
12
13 df = pd.DataFrame(array,
14                   index=columns,
15                   columns = ['asympt. conf. interval', 'conf. interval'])
16
17 df

```

Out[117]:

	asympt. conf. interval	conf. interval
<b>Alcalinity of ash</b>	(16.3931, 17.6815)	(16.3737, 17.7009)
<b>Nonflavanoid phenols</b>	(0.2723, 0.3077)	(0.2717, 0.3083)
<b>Proanthocyanins</b>	(1.7951, 2.0036)	(1.7919, 2.0067)
<b>Hue</b>	(1.0326, 1.0915)	(1.0317, 1.0924)

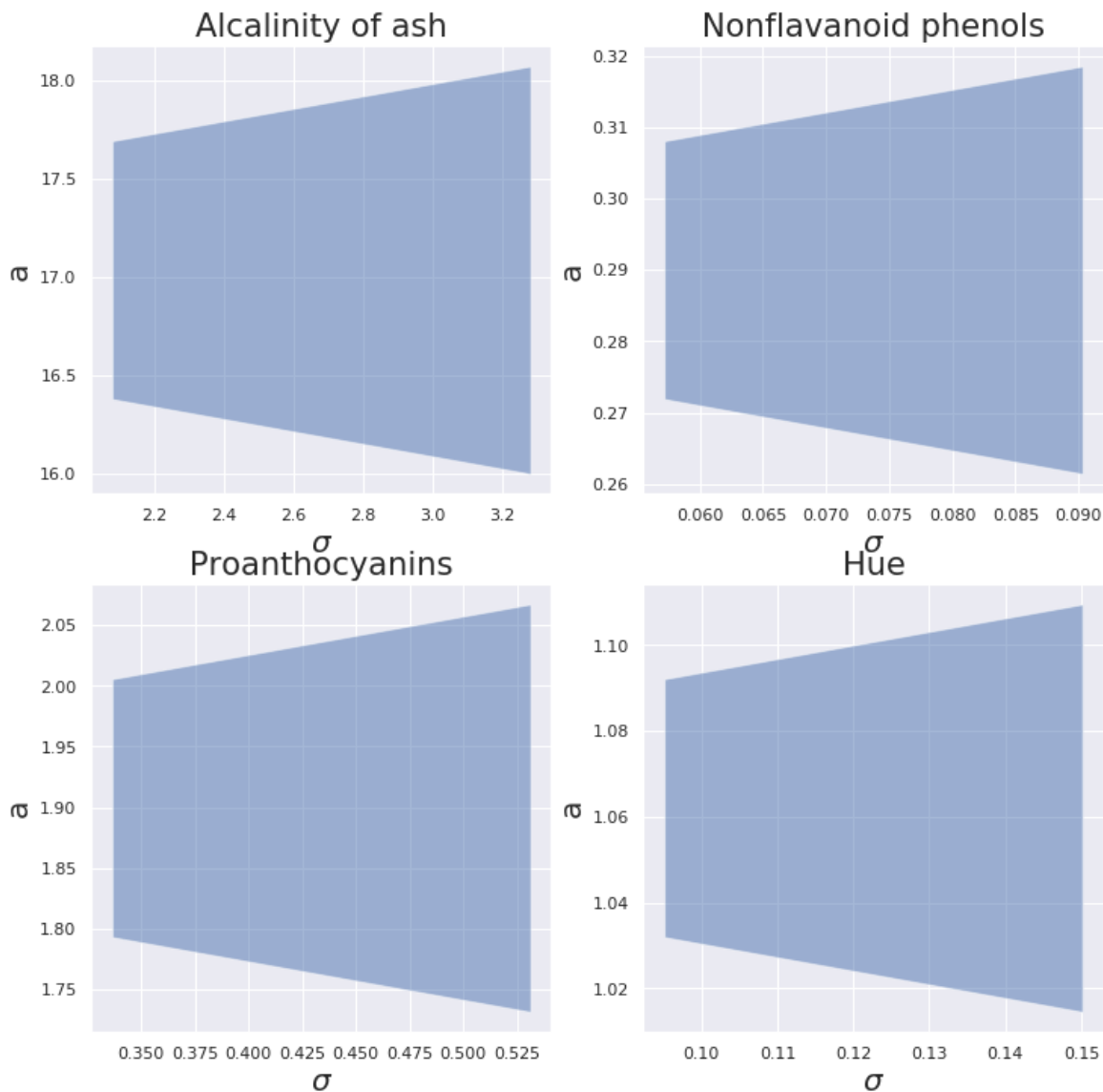
Наконец, постройте точную доверительную область для параметров сдвига и масштаба для каждого из рассматриваемых столбцов. Для экономии места стройте графики в два столбца.

In [124]:

```

1 plt.figure(figsize=(12, 12))
2
3 for ind, column in enumerate(columns):
4     plt.subplot(2,2, ind+1)
5     plt.title(column, fontsize=21)
6
7     plt.ylabel("a", fontsize=20)
8     plt.xlabel(r"$\sigma$", fontsize=20)
9
10    sample = data[column]
11
12    sigmas, left_a, right_a = confidence_region(sample)
13
14    plt.fill_between(sigmas, left_a, right_a, alpha=0.5,
15                    label="Область для выборки {}".format(column))
16
17 plt.show()

```



**Вывод:** Из рисунка и из таблицы видим, что не при всех значениях сигма интервалы, показанные в таблице, будут вложены в дов. области, но уже при небольшом сигма они будут внутри этих областей. Видимо, уже на основе этого мы можем судить о значениях сигмы в распределениях для наших выборок.



## Байесовский подход

### Задача 6.

Пусть  $X_1, \dots, X_n$  --- выборка из распределения  $\mathcal{N}(\theta, \sigma^2)$ , а параметр  $\theta$  в качестве априорного распределения имеет стандартное распределение Коши со сдвигом.

Апостериорное распределение вычисляется по формуле:  $q(\theta | x) = \frac{q(\theta) p_t(x_1) \dots p_t(x_n)}{\int_{\Theta} q(\theta) p_t(x_1) \dots p_t(x_n) d\theta}$ , где  $p_t(x)$  --- плотность распределения  $\mathcal{N}(\theta, \sigma^2)$ , а  $q(\theta)$  --- плотность распределения Коши. Как было сказано на лекции, аналитически интеграл в знаменателе посчитать не удастся. Однако, этот интеграл можно вычислить численно, например, с помощью метода Монте-Карло.

В данном случае, интеграл  $\int_{\mathbb{R}} f(x) p(x) dx$ , где  $p(x)$  -- некоторая плотность, можно оценить как  $\frac{1}{k} \sum_{i=1}^k f(\xi_i)$ , где  $\xi_1, \dots, \xi_k$  -- сгенерированная выборка из распределения, имеющего плотность  $p(x)$ .

Рассмотрим столбец Alkalinity of ash датасета [о вине](http://archive.ics.uci.edu/ml/datasets/Wine) (<http://archive.ics.uci.edu/ml/datasets/Wine>).

Для выборки, образованной этим столбцом посчитайте  $\sigma$  -- знаменатель в формуле Байеса. Параметры априорного распределения выберите некоторым разумным способом, не опираясь на данные. Какой размер вспомогательной выборки в методе приближенного интегрирования необходим, чтобы с большой точностью посчитать значение  $\sigma$ ?

In [208]:

```
1 sigma = 2.7
2 loc_ash = 20
3
4 ash = data["Alkalinity of ash"]
5
6 def func(sample, t):
7     n = len(sample)
8     e = np.exp(-1/2*((sample-t)**2).sum()/(2.7**2))
9     return 1/(2*np.pi*2.7**2)**(n/2)*e
10
11 vect_func = np.vectorize(func, signature="(n),()->()")
12
13 #функция, по Методу Монте Карло вычисляющая матожидание функции func
14 def getApproxMonteCarlo(func, sample, k):
15     ts = sps.cauchy(loc=loc_ash).rvs(size=k)
16     return func(sample, ts).mean()
17
18
19 denominator = getApproxMonteCarlo(vect_func, ash, 100000)
```

Для апостериорного распределения:

- Нарисуйте график плотности;
- Посчитайте математическое ожидание;
- Найдите симметричный 95%-ый доверительный интервал.

In [209]:

```
1 #здесь нашли функцию усл. плотности
2 ts = np.linspace(15, 25, 1000)
3
4 q = sps.cauchy(loc=loc_ash).pdf
5
6
7 condit_q = lambda ts : vect_func(ash, ts)*q(ts)/denominator
8
```

In [210]:

```
1 #посчитаем матож методом монтекарло для функции  $p_{\{t\}}(x)*t$ 
2 numen_func = lambda sample, ts: vect_func(sample, ts)*ts
3
4 #посчитаем матож без знаменателя
5 numerator = getApproxMonteCarlo(numen_func, ash, 1000000)
6
7
8 condit_E = numerator/denominator
9
10 print("условное матожидание : ", condit_E)
```

условное матожидание : 17.06690581137273

In [199]:

```
1 #посчитаем квантили
2 alpha = 0.95
3 eps = 1e-4
4 #посчитаем кумулятивные суммы по сетке
5 ts = np.linspace(15, 20, 10000)
6
7 qs = condit_q(ts).cumsum()/condit_q(ts).sum()
8
9 quantile1 = ts[np.argwhere(np.abs(qs - (1-alpha)/2) < eps).min()]
10 quantile2 = ts[np.argwhere(np.abs(qs - (1+alpha)/2) < eps).min()]
11
12 print("доверительный интервал :", (round(quantile1, 3),
13                                     round(quantile2, 3)))
```

доверительный интервал : (16.418, 17.811)

In [207]:

```

1 #нарисуем график
2
3
4 plt.figure(figsize=(10, 8))
5 plt.title("условная плотность  $q(t|x)$ ", fontsize=20)
6 plt.xlabel("t", fontsize=15)
7 plt.ylabel("conditional pdf", fontsize=13)
8
9 plt.xlim((15, 20))
10
11
12 plt.plot(ts, condit_q(ts), label='density')
13 plt.axvline(x=condit_E, ls='--', color='red',
14             label='condition Expectation')
15 plt.axvspan(quantile1, quantile2, color='brown', alpha=0.2,
16             label='0.95 cond. interval')
17
18 plt.legend(fontsize=15)
19 plt.show()

```



**Вывод:** Здесь мы получили апостериорное распределение параметра  $\theta$ , а также условное матожидание равное  $\approx 17.07$  и дов. интервал, равный (16.418, 17.811). Мы можем видеть, что, сравнивая со значениями из предыдущей задачи, байесовский дов. интервал приблизительно такой же, как интервал из задачи 5. Наш матож тоже лежит в дов. интервале из задачи 5.

Кроме того, мы увидили, что возиться с интегралами было довольно нелегко, поэтому хотелось бы другие методы, которые не потребовали большой вычислительной сложности.

## Задача 7.

Рассмотрим схему испытаний Бернулли (т.е. броски монет) с вероятностью успеха  $p$ .

Постройте несколько графиков сопряженного распределения для разных параметров и охарактеризуйте, как его значения параметров соотносятся с априорными знаниями о монете. Это могут быть, например, знания вида

- монета скорее честная (при таком априорном распределении наиболее вероятны значения  $p$  в окрестности 0.5);
- монета скорее нечестная, перевес неизвестен (наименее вероятны значения  $p$  в окрестности 0.5);
- монета скорее нечестная, перевес в сторону герба (наиболее вероятны значения  $p$  в окрестности 1);
- монета скорее честная, либо с небольшим перекосом вправо (наиболее вероятны значения  $p$  в окрестности ~0.6);
- ничего не известно (все значения равновероятны).

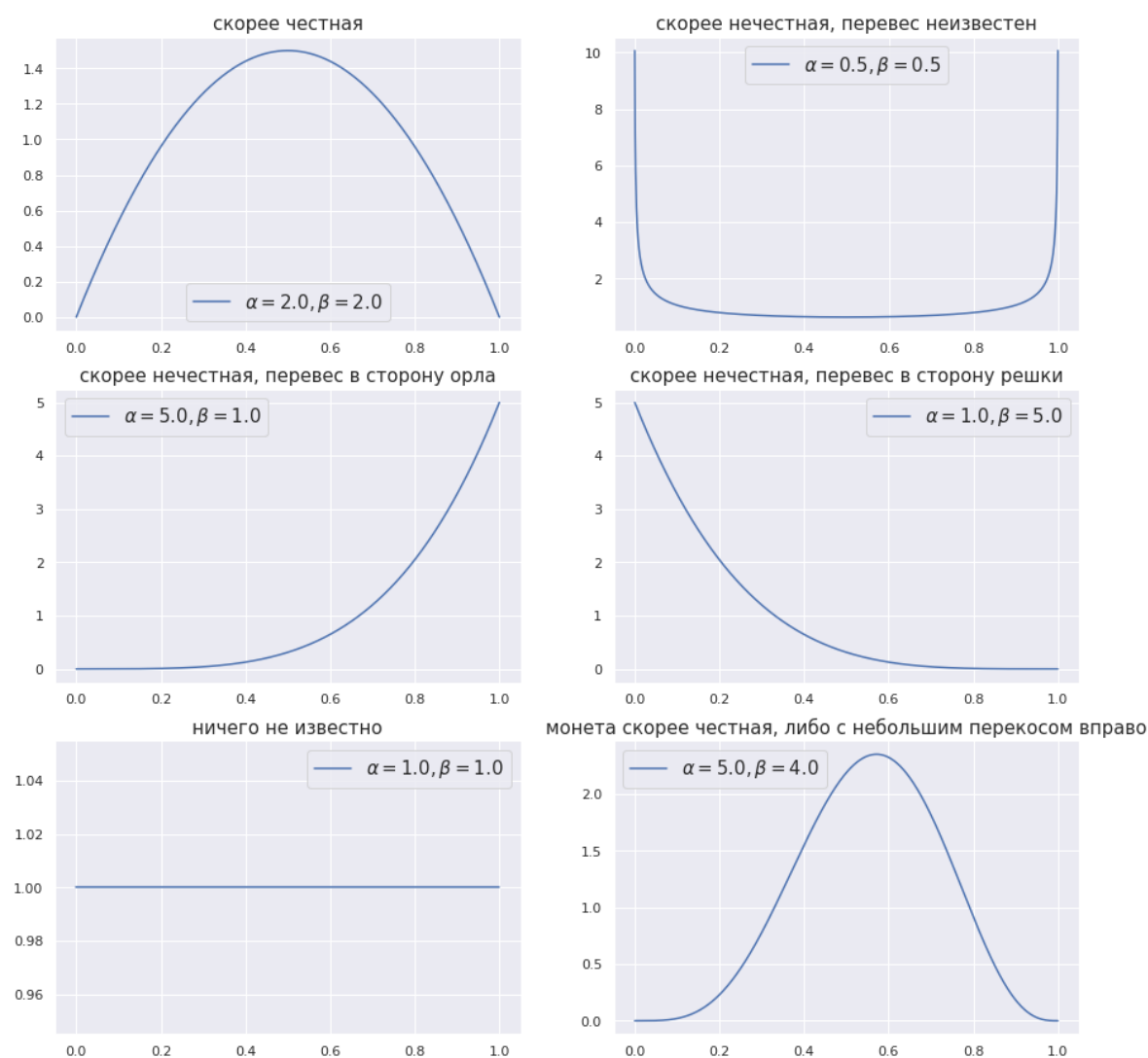
Для каждого случая из перечисленных выше постройте график плотности сопряженного распределения (на одной фигуре).

In [326]:

```

1 grid= np.linspace(0, 1, 1000)
2
3 indices = np.arange(1, 7)
4 alphas = [2, 0.5, 5, 1, 1, 5]
5 betas = [2, 0.5, 1, 5, 1, 4]
6 titles = ["скорее честная",
7           "скорее нечестная, перевес неизвестен",
8           "скорее нечестная, перевес в сторону орла",
9           "скорее нечестная, перевес в сторону решки",
10          "ничего не известно",
11          "монета скорее честная, либо с небольшим перекосом вправо"]
12
13
14 plt.figure(figsize=(15, 15))
15
16 for i, alpha, beta, title in zip(indices, alphas, betas, titles):
17     plt.subplot(3, 2, i)
18     plt.title(title, fontsize=15)
19     plt.plot(grid, sps.beta(a=alpha, b=beta).pdf(grid),
20              label = r"$\alpha = %.1f, \beta = %.1f$" % (alpha, beta))
21     plt.legend(fontsize=15)
22 plt.show()

```



**Вывод:** На примере этих графиков мы видим, что с помощью Бета распределения мы можем довольно хорошо описать различные априорные знания о монетке. Мы можем легко подобрать параметры

распределения так, чтобы данное априорное распределение хорошо описывало наши предварительные знания о монетке.

Ниже приведена реализация некоторых вспомогательных функций.

In [403]:

```

1  def draw_posteriori(grid, distr_class, post_params, ind,
2                      title=None, xlim=None):
3      ''' Рисует серию графиков апостериорных плотностей.
4          grid --- сетка для построения графика
5          distr_class --- класс распределений из scipy.stats
6          post_params --- параметры апостериорных распределений
7                      shape=(размер выборки, кол-во параметров)
8      ...
9
10     size = post_params.shape[0] - 1
11
12     plt.subplot(4, 2, ind)
13     plt.title(title)
14     for n in range(size+1):
15         plt.plot(grid,
16                  distr_class(post_params[n]).pdf(grid) \
17                      if np.isscalar(post_params[n]) \
18                      else distr_class(*post_params[n]).pdf(grid),
19                  label='n=%d: (%.2f,%.2f)' % (n, post_params[n][0],
20                                              post_params[n][-1]),
21                  lw=2.5,
22                  color=(1-n/size, n/size, 0))
23     plt.grid(ls=':')
24     plt.legend(loc=2)
25     plt.xlim(xlim)
26
27
28  def draw_estimations(ml, distr_class, post_params, ind,
29                      title=None, confint=True, ylim=None):
30      ''' Рисует графики байесовской оценки (м.о. и дов. инт.) и ОМП.
31          ml --- Оценка максимального правдоподобия для 1 <= n <= len(sample)
32          distr_class --- класс распределений из scipy.stats
33          post_params --- параметры апостериорных распределений
34                      shape=(размер выборки, кол-во параметров)
35      ...
36
37     size = len(ml)
38     distrs = []
39     for n in range(size+1):
40         distrs.append(distr_class(post_params[n]) \
41                      if np.isscalar(post_params[n]) \
42                      else distr_class(*post_params[n]))
43
44     plt.subplot(4, 2, ind)
45     plt.title(title)
46     plt.plot(np.arange(size+1), [d.mean() for d in distrs],
47             label='Bayes', lw=1.5)
48     plt.fill_between(np.arange(size+1), [d.ppf(0.975) for d in distrs],
49                    [d.ppf(0.025) for d in distrs], alpha=0.1)
50     plt.plot(np.arange(size)+1, ml, label='MLE', lw=1.5)
51     plt.grid(ls=':')
52     plt.ylim(ylim)
53     plt.legend(loc=2)

```

Реализуйте следующую функцию

In [364]:

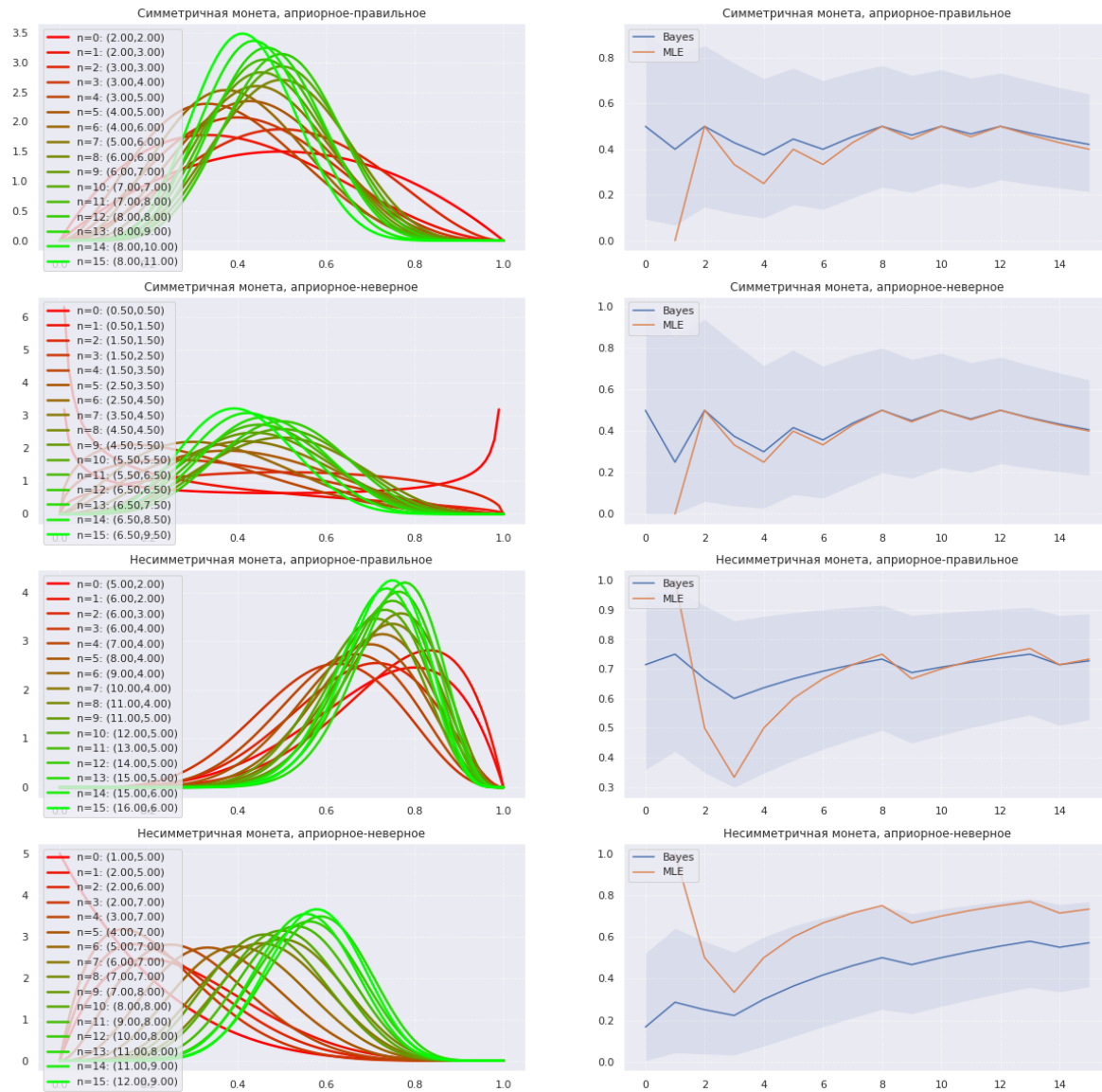
```
1 def bern_posterior_params(sample, a, b):
2     ''' Возвращает параметры апостериорного распределения
3         для всех  $0 \leq n \leq \text{len}(\text{sample})$ .
4         a, b --- параметры априорного распределения.
5     '''
6
7     grid = np.arange(len(sample) + 1)
8     summ = np.zeros(shape=len(sample) + 1)
9     summ[1:] = sample.cumsum()
10    alphas = a + summ
11    betas = b + grid - summ
12
13    params = np.dstack((alphas, betas))
14    return params[0]
```

Проведите по 15 бросков симметричной и несимметричной монет (можно сгенерировать) и рассмотрите для каждой из них два случая --- параметры априорного распределения подобраны правильно или неправильно. Постройте графики, воспользовавшись функциями `draw_posteriori` и `draw_estimations`.



In [385]:

```
1
2 def get_graphics(func, grid, distr_class, coins, a, b, ind,title):
3     draw_posteriori(grid, distr_class,
4                     func(coins, a, b),ind,
5                     title = title)
6
7     ml = coins.cumsum()/np.arange(1, len(coins) + 1)
8     draw_estimations(ml, distr_class,
9                     func(coins, a, b),ind + 1,
10                    title = title)
11
12
13 #генерим выборки
14 grid = np.linspace(0, 1, 100)
15
16 symm_coins = sps.bernoulli(0.5).rvs(15)
17
18 asym_coins = sps.bernoulli(0.8).rvs(15)
19
20 plt.figure(figsize=(20, 20))
21
22 get_graphics(bern_posterior_params, grid, sps.beta, symm_coins, 2, 2, 1,
23             'Симметричная монета, априорное-правильное')
24
25 get_graphics(bern_posterior_params, grid, sps.beta, symm_coins, 0.5, 0.5, 3,
26             'Симметричная монета, априорное-неверное')
27 get_graphics(bern_posterior_params, grid, sps.beta, asym_coins, 5, 2, 5,
28             'Несимметричная монета, априорное-правильное')
29 get_graphics(bern_posterior_params, grid, sps.beta, asym_coins, 1, 5, 7,
30             'Несимметричная монета, априорное-неверное')
31
32 plt.show()
```



Сделайте вывод. Что можно сказать про зависимость от параметров априорного распределения? Сравните байесовские оценки с оценкой максимального правдоподобия.

**Вывод:** Зависимость от параметров априорного распределения конечно есть, но даже если мы выбираем их изначально в корне неправильно, апостериорное распределение с ростом выборки все равно стремится к распределению с более правильными параметрами, хотя стоит отметить, что в нашем случае для выборки из 15 элементов распределение иногда не успевает как следует приблизиться к правильным параметрам.

Байесовские оценки конечно также очень зависят от параметра априорного распределения, и при выборе правильных параметров априорного распределения, мы видим, как хорошо условное Мат.Ожидание и Дов.Интервал описывают истинный параметр распределения бернулли. Мы видим, что в таких случаях оценки MLE и байесовские точечные оценки довольно быстро начинают сходятся к очень близким значениям. Можно даже отметить, что зачастую в случае правильно выбранных параметров байесовская оценка сходится даже быстрее чем MLE.

Но даже при выборе неправильных (даже полностью противоположных истинным) параметров, байесовские оценки, пусть и медленней, чем при выборе правильных параметров, при росте выборки стремится к тому же значению, что и MLE, хотя MLE сходится к правильному значению быстрее. На этапе выборки из 15 элементов, MLE иногда даже не попадает в байесовский Доверительный Интервал. Поэтому можно сказать, что выборка из 15 элементов, по крайней мере для данного опыта с монетой, слишком мала, чтобы с помощью Байесовского метода сказать какие то хорошие результаты.

## Задача 8.

Один экзаменатор на экзамене по математической статистике при выставлении оценки студенту пользуется следующей схемой. В течении экзамена экзаменатор задает студенту некоторое количество вопросов, получая тем самым выборку  $X_1, \dots, X_n \sim \text{Bern}(p)$  --- индикаторы того, что студент на вопрос ответил правильно. При этом сначала он подбирает некоторое априорное распределение на основе его знаний о студенте к моменту начала ответа. После каждого ответа студента экзаменатор вычисляет апостериорное распределение и строит байесовский доверительный интервал для  $p$  уровня доверия 0.95. Если после очередного ответа студента доверительный интервал содержит лишь одно число  $i/10$ , где  $i \in \{0, \dots, 10\}$ , то экзаменатор выставляет студенту оценку  $i+1$ .

Ответьте на следующие вопросы:

- Квантили какого уровня нужно выбирать экзаменатору при построении доверительного интервала, чтобы задавать студенту как можно меньше вопросов? Какие оценки будет выставлять экзаменатор в таком случае?
- Как зависит оценка студента и среднее количество заданных вопросов у различных студентов (по уровню знаний) при различных априорных представлениях экзаменатора?
- Нужно ли дружить с таким экзаменатором?

In [358]:

```

1  #условие окончания экзамена
2  def isOver(q1, q2):
3      if ((q2 - q1) < 0.1):
4          for i in range(11):
5              if (q1 <= i/10 <= q2):
6                  return i + 1
7      return -1
8
9
10 def getMark(p, a, b):
11     """
12     p- параметр бернуллиевского распределения, по
13     которому строится выборка
14     a, b - параметры априорного распределения
15     """
16     params = [(a, b)]
17     sample = list()
18     q1 = 0
19     q2 = 1
20
21     mark = -1
22
23     while (mark == -1):
24         coin = sps.bernoulli(p).rvs()
25         sample.append(coin)
26
27         params.append((params[-1][0] + coin,
28                        params[-1][1] + 1 - coin))
29
30         q1 = sps.beta(*params[-1]).ppf(0.025)
31         q2 = sps.beta(*params[-1]).ppf(0.975)
32
33         mark = isOver(q1, q2)
34     return (mark, len(sample))
35
36
37 getMark(0.8, 1, 1)

```

Out[358]:

(9, 225)

**Вывод:** Слишком много вопросов задает экзаменатор, но оценки ставит на одну больше чем в среднем, что хорошо)

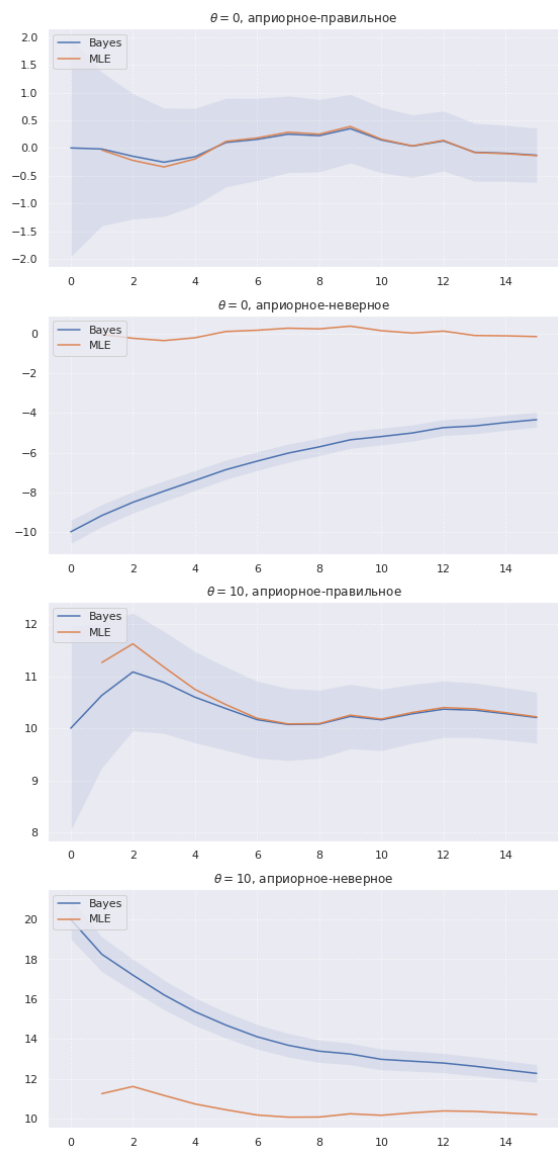
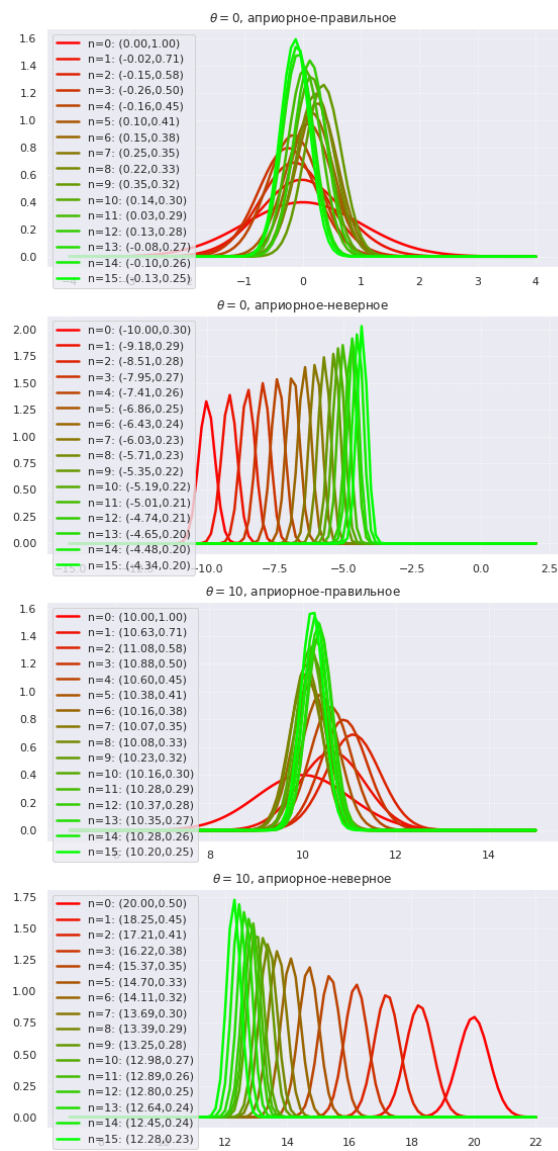
## Задача 9.

Проведите исследование, аналогичное задаче 7 для выборок из распределений

- $\text{mathcal{N}}(\theta, 1)$
- $\text{Exp}(\theta)$

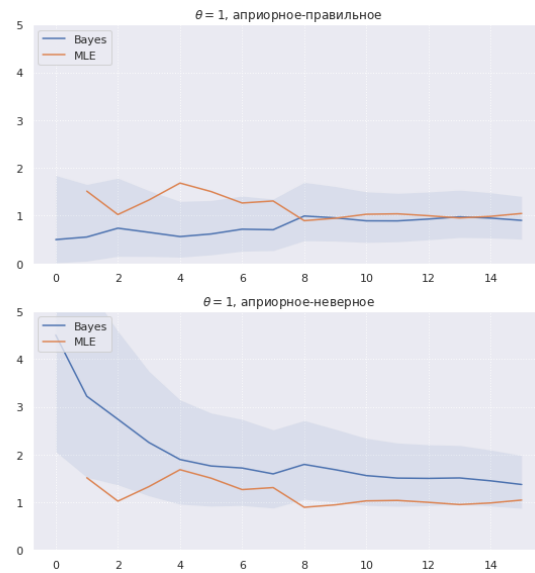
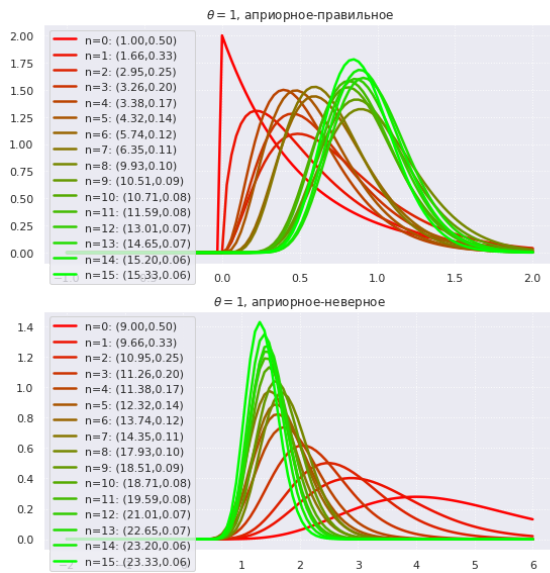
In [395]:

```
1  #для нормального распределения
2
3  def normal_posterior_params(sample, a, sigma):
4      new_sample = np.zeros(len(sample) + 1)
5      new_sample[1:] = sample.cumsum()
6      grid = np.arange(len(sample) + 1)
7
8      mu = (a/sigma**2 + new_sample)/(1/sigma**2 + grid)
9
10     sigmas = 1/np.sqrt(1/sigma**2 + grid)
11
12     params = np.dstack((mu, sigmas))
13     return params[0]
14
15
16
17 #генерим выборки
18 sample1 = sps.norm().rvs(15)
19
20 sample2 = sps.norm(loc=10, scale=1).rvs(15)
21
22 plt.figure(figsize=(20, 20))
23
24 get_graphics(normal_posterior_params, np.linspace(-4, 4, 100),
25             sps.norm, sample1, 0, 1, 1,
26             r'$\theta = 0$, априорное-правильное')
27
28 get_graphics(normal_posterior_params, np.linspace(-15, 2, 100),
29             sps.norm, sample1, -10, 0.3, 3,
30             r'$\theta = 0$, априорное-неверное')
31 get_graphics(normal_posterior_params, np.linspace(5, 15, 100),
32             sps.norm, sample2, 10, 1, 5,
33             r'$\theta = 10$, априорное-правильное')
34 get_graphics(normal_posterior_params, np.linspace(7, 22, 100),
35             sps.norm, sample2, 20, 0.5, 7,
36             r'$\theta = 10$, априорное-неверное')
37
38 plt.show()
```



In [462]:

```
1  # для Exp(theta)
2
3  def expon_posterior_params(sample, a, b):
4      new_sample = np.zeros(len(sample) + 1)
5      new_sample[1:] = sample.cumsum()
6      grid = np.arange(len(sample) + 1)
7
8      mu = a + grid
9
10     bs = b + new_sample
11
12     zer = np.zeros(len(sample) + 1)
13
14     params = np.dstack((bs, zer, 1/mu))
15     return params[0]
16
17
18 def get_exp_graphics(grid, coins, a, b, ind, title):
19     draw_posteriori(grid, sps.gamma,
20                     expon_posterior_params(coins, a, b), ind,
21                     title = title)
22
23     ml = coins.cumsum()/np.arange(1, len(coins) + 1)
24     draw_estimations(1/ml, sps.gamma,
25                     expon_posterior_params(coins, a, b), ind + 1,
26                     title = title, ylim=(0, 5))
27
28
29 #генерим выборки
30
31 sample1 = sps.expon(scale=1).rvs(15)
32
33 plt.figure(figsize=(20, 20))
34
35 get_exp_graphics(np.linspace(-1, 2, 100), sample1, 2, 1, 1,
36                 r'$\theta = 1$, априорное-правильное')
37
38 get_exp_graphics(np.linspace(-2, 6, 100), sample1, 2, 9, 3,
39                 r'$\theta = 1$, априорное-неверное')
40
41
42 plt.show()
```



**Вывод:** Мы видим, что при исследовании  $\mathcal{N}(\theta, 1)$ , и  $\text{Exp}(\theta)$  наблюдаются приблизительно такие же выводы, как и в задаче 7. Но также можно сказать, что при анализе нормального распределения, очень легко подобрать "слишком плохие" параметры сопряженного априорного распределения, и таким образом сходимость байесовских оценок будет очень медленной, а байесовский доверительный интервал будет узким, и даже при относительно больших размерах выборки непригоден. Это будет происходить, если брать слишком маленькую дисперсию у сопряженного априорного распределения, и параметр  $\theta$ , сильно отличающийся от истинного матожидания.

В анализировании  $\text{Exp}(\theta)$  можно подобрать плохие параметры приблизительно похожим образом - чтобы "бугор" сопряженного был далеко от истинного "бугра", и причем был довольно высоким.

In [ ]:

1