


```
In [23]: dots_num = 500

y = np.linspace(data["citric acid"].min(), data["citric acid"].max(), dots_num, endpoint=True)
x = np.linspace(data["alcohol"].min(), data["alcohol"].max(), dots_num, endpoint=True)
grid = np.stack(np.meshgrid(x, y), axis=2)
plt.figure(figsize=(12, 8))

# high quality
plt.pcolormesh(
    x,
    y,
    sps.multivariate_normal.pdf(grid, mean=mu_high_quality, cov=sigma_high_quality),
    cmap=get_density_cmap("blues"),
    label="High quality wines"
)

CS = plt.contour(
    x,
    y,
    sps.multivariate_normal.pdf(grid, mean=mu_high_quality, cov=sigma_high_quality)
)

plt.clabel(
    CS,
    inline=1
)

plt.scatter(
    data=data[data["quality"] == 8],
    x="alcohol",
    y="citric acid",
    label="High quality wines",
    c="b"
)

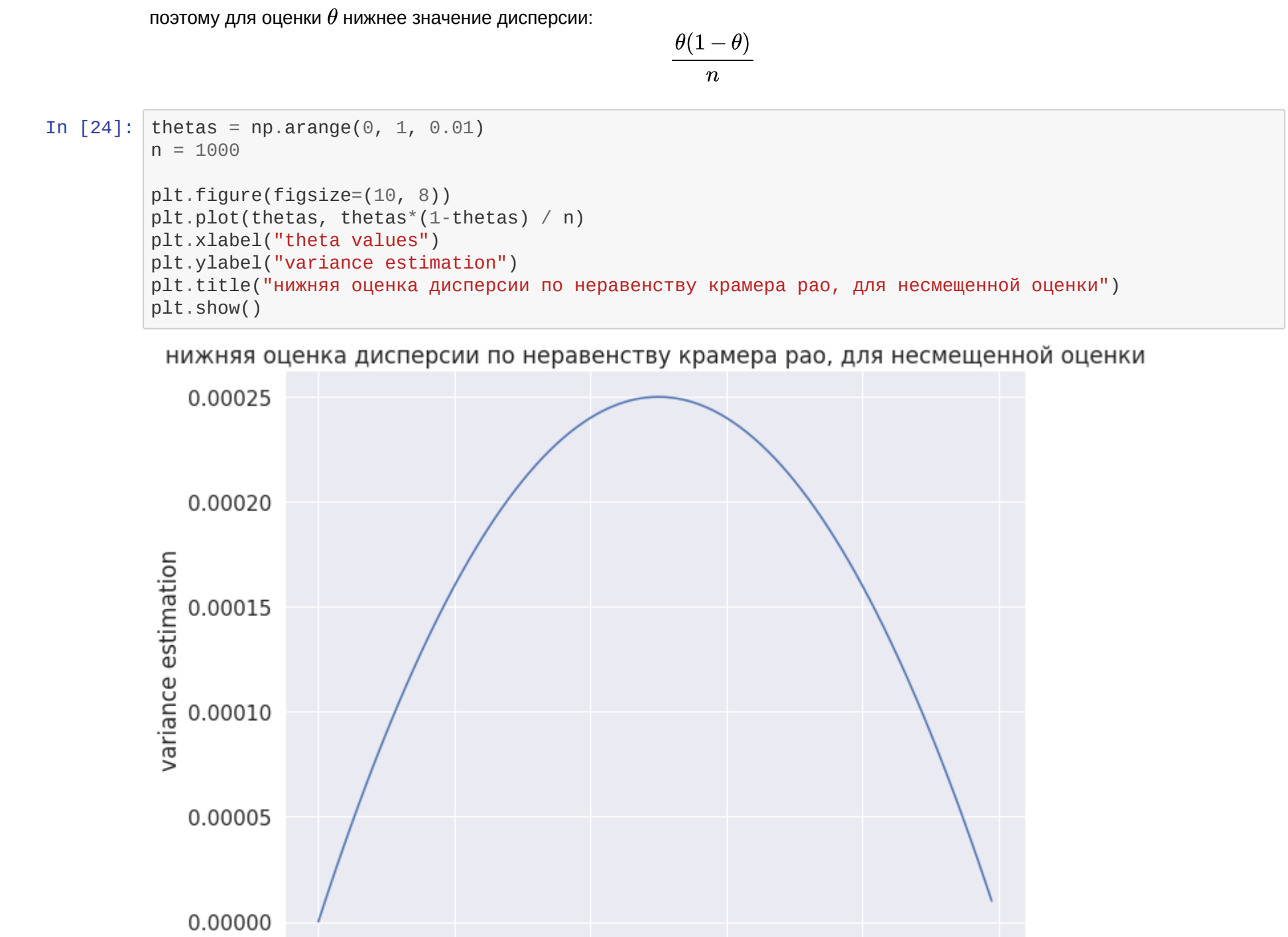
# low quality
plt.pcolormesh(
    x,
    y,
    sps.multivariate_normal.pdf(grid, mean=mu_low_quality, cov=sigma_low_quality),
    cmap=get_density_cmap("greens"),
    label="Low quality wines"
)

CS = plt.contour(
    x,
    y,
    sps.multivariate_normal.pdf(grid, mean=mu_low_quality, cov=sigma_low_quality)
)

plt.clabel(
    CS,
    inline=1
)

plt.scatter(
    data=data[data["quality"] == 3],
    x="alcohol",
    y="citric acid",
    label="Low quality wines",
    c="g"
)

plt.title("Distribution density")
plt.xlabel("alcohol")
plt.ylabel("citric acid")
plt.legend()
plt.show()
```



Что можно сказать о вине, которому сомелье дала наивысший балл по сравнению с вином, которому дала наименьший балл, основываясь на график выше?

- Ответ:**
- В среднем алкогольность и содержание лимонной кислоты у вина лучшего качества выше
 - Дисперсия по алкогольности у вина высокого качества больше
 - Дисперсия по содержанию лимонной кислоты у вина низкого качества выше

Задача 3

Рассмотрим $X_1, \dots, X_n \sim \text{Bern}(\theta)$. По сетке значений $\theta \in [0, 1]$ с шагом 0.01 постройте график зависимости нижней оценки дисперсии произвольной несмещенной оценки из неравенства Рао-Крамера от θ .

Нер-во Крамера-Рао:

$$D_\theta \theta^* \geq \frac{(\tau'(\theta))^2}{ni(\theta)}$$

Из критерия эффективности (расписывая вклад) можно получить, что $i(\theta) = \frac{1}{\theta(1-\theta)}$

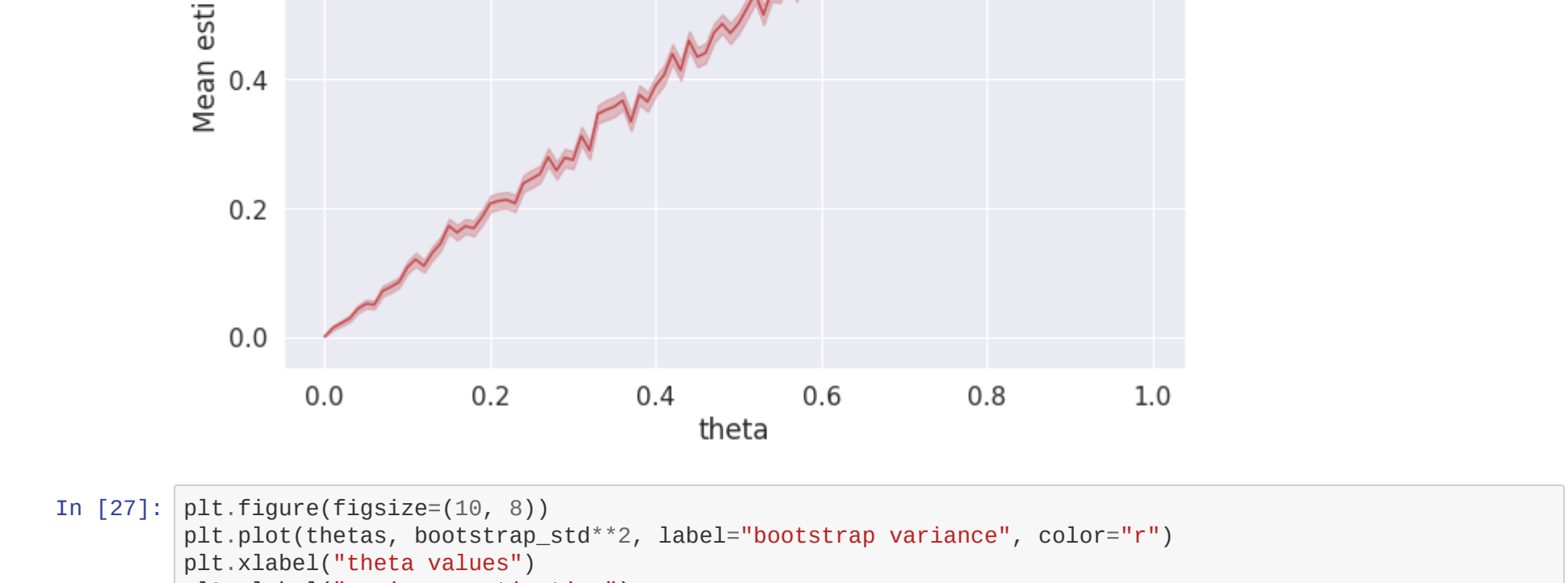
поэтому для оценки θ нижнее значение дисперсии:

$$\frac{\theta(1-\theta)}{n}$$

```
In [24]: thetas = np.arange(0, 1, 0.01)
n = 1000

plt.figure(figsize=(10, 8))
plt.plot(thetas, thetas*(1-thetas) / n)
plt.xlabel("thetas values")
plt.ylabel("variance estimation")
plt.title("нижняя оценка дисперсии по неравенству крамера рао, для несмещенной оценки")
plt.show()
```

нижняя оценка дисперсии по неравенству крамера рао, для несмещенной оценки



Какой можно сделать вывод (напишите в комментариях)?

- Вывод**
- График симметричен отн. 0.5 тк $D_\theta \theta^* = D_\theta(1-\theta^*)$
 - Максимальное значение дисперсии достигается при $\theta = \frac{1}{2}$ тк при данном значении θ информация фишера выборки минимальна

Для каждого значения θ (для той же сетки) сгенерируйте выборку размера $n = 1000$ для параметра θ , посчитайте эффективную оценку θ и бутстрепную оценку дисперсии (количество бутстрепных выборок равно 1000) этой эффективной оценки θ .

Эффективная оценка θ , \bar{X} (из все же критерия эффективности)

```
In [25]: theta_estimates = np.zeros(len(thetas))
bootstrap_estimates = np.zeros(len(thetas))
bootstrap_std = np.zeros(len(thetas))
n = 1000

for i, theta in enumerate(thetas):
    sample = sps.bernoulli.rvs(theta, size=n)
    estimates = np.array([np.random.choice(sample, size=n).mean() for _ in range(n)])

    theta_estimates[i] = sample.mean()
    bootstrap_estimates[i] = estimates.mean()
    bootstrap_std[i] = estimates.std()
```

Нарисуйте график зависимости полученных бутстрепных оценок от θ .

```
In [26]: plt.figure(figsize=(10, 8))
plt.plot(thetas, bootstrap_estimates, label="bootstrap  $\overline{X}$ ", color="r")
plt.fill_between(
    x=thetas,
    y1=bootstrap_estimates + bootstrap_std,
    y2=bootstrap_estimates - bootstrap_std,
    alpha=0.5,
    label="bootstrap std",
    color="r"
)

plt.xlabel("thetas")
plt.ylabel("Mean estimation")
plt.title("effective theta estimation")
plt.legend()
plt.show()
```



```
In [27]: plt.figure(figsize=(10, 8))
plt.plot(thetas, bootstrap_std**2, label="bootstrap variance", color="r")
plt.ylabel("variance estimation")
plt.xlabel("thetas values")
plt.title("variance estimation by bootstrap")
plt.legend()
plt.show()
```



Вывод

- График зависимости эффективной оценки $\theta^* = \bar{X}$ от θ ожидаемо похож на прямую
- Величина стандартного отклонения увеличивается до $\theta = \frac{1}{2}$, в ней достигает максимума, затем убывает
- График дисперсии оценки, полученной бутстрепом, ожидаемо похож на график нижней оценки дисперсии из неравенства Крамера-Рао, тк оценка $\theta^* = \bar{X}$ является эффективной