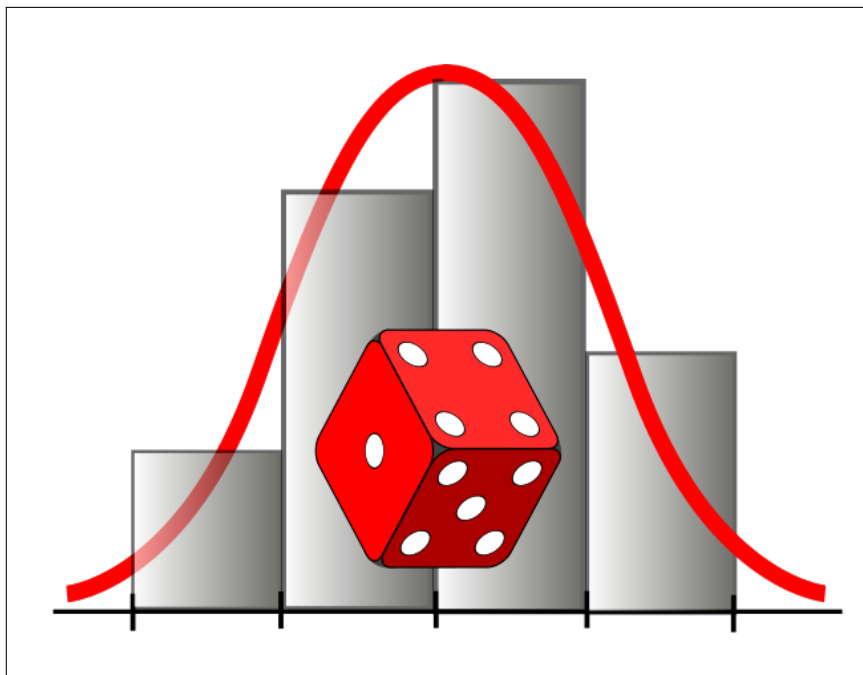


”Лабораторная” работа 2 по Теории вероятностей

Мочекон Семён



1 Задача 1.1.

Пусть такое возможно и $X = X_1 + X_2$, где X_1, X_2 независимы и имеют невырожденное распределение.

В силу независимости, вероятность того, что X примет значение равное k^2 , представляется по формуле свёртки:

$$P(X = k^2) = P(X_1 + X_2 = k^2) = \sum_j P(X_1 = j) \cdot P(X_2 = k^2 - j)$$

В связи с этим встаёт вопрос: Какие j могут быть в сумме?

Из условия имеем, что носителем X является множество квадратов натуральных чисел с нулём. То есть:

$$\text{supp}(X) = \{n^2 : n \in \mathbb{N}_0\}$$

Но тогда:

$$\text{supp}(X) = \text{supp}(X_1 + X_2) = \{n^2 : n \in \mathbb{N}_0\}$$

То есть все возможные суммы случайных величин X_1, X_2 должны давать в точности множество квадратов \mathbb{N}_0 .

Теперь обратим внимание на второе свойство X_1, X_2 . Их распределение не вырождено, то есть каждая из случайных величин принимает больше чем одно значение. Ну и дополнительно отметим, что носитель хотя бы одной должен быть счётен, иначе не выполнена счётность носителя X .

Так как эта пара случайных величин задаёт своей суммой все квадраты \mathbb{N}_0 , то каждая из сумм свёртки не пуста.

Не умаляя общности будем отталкиваться от того, что некоторое значение принимает именно X_1 , а дополняющее до квадрата X_2 .

Тогда для некоторого m и \hat{j} :

$$P(X = m^2) = P(X_1 = \hat{j}) \cdot P(X_2 = m^2 - \hat{j}) + \sum_{j, j \neq \hat{j}} P(X_1 = j) \cdot P(X_2 = m^2 - j)$$

Аналогично для некоторого d и \dot{j} :

$$P(X = d^2) = P(X_1 = \dot{j}) \cdot P(X_2 = d^2 - \dot{j}) + \sum_{j, j \neq \dot{j}} P(X_1 = j) \cdot P(X_2 = d^2 - j)$$

При том \hat{j}, \dot{j} таковы, что $\hat{j} \neq \dot{j}$ так как X_1 принимает хотя бы два значения.

Тогда, возвращаясь к носителю случайной величины, мы приходим к тому, что на самом деле $\text{supp}(X_1 + X_2)$ не будет совпадать с квадратами натуральных чисел с нулём "почти наверное". Так как, то же $m^2 - \hat{j} + \dot{j}$ является квадратом в исключительных случаях.

Предположение привело нас к противоречию.

Но если допустить, что у одной величины вырожденное распределение, то есть $X_1 \in \{c\}$, тогда можно представить X как: $X = X_1 + X_2$, $X_1 \in \{c\}$, $X_2 \in \{k^2 - c\}_{k=0}^{+\infty}$

2 Задача 2.3.

Определение стандартного гауссовского вектора:

Вектор $X = (X_1, \dots, X_n)$ будем называть **стандартным гауссовским**, если его плотность имеет вид

$$p(x_1, \dots, x_n) = \frac{1}{(2\pi)^{n/2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n x_i^2 \right\} = \frac{1}{(2\pi)^{n/2}} \exp \left\{ -\frac{x^T x}{2} \right\}.$$

Покажем, что совместная плотность X, Y удовлетворяет этому определению.

$$X = \sqrt{-2 \ln V} \cos(2\pi U), \quad Y = \sqrt{-2 \ln V} \sin(2\pi U)$$

Введём:

$$\rho = \sqrt{-2 \ln V}, \quad \varphi = 2\pi U$$

Тогда:

$$V = e^{-\frac{\rho^2}{2}}, \quad U = \frac{\varphi}{2\pi}$$

Что в принципе даёт выражение для полярных координат, так как $\varphi \in [0; 2\pi]$, $\rho \in [0; +\infty)$. По области определения U, V Так как старые случайные величины имеющие равномерное распределение на $[0; 1]$, то плотность каждой равна $p_U = p_V = 1$, а их совместная плотность, в силу независимости $p_{UV} = 1 \cdot 1 = 1$.

Тогда по теореме о преобразовании плотности:

$$p_{XY} = p_{UV} \cdot \left| \frac{1}{\det(J)} \right|$$

$$J = \begin{pmatrix} \frac{\partial X}{\partial U} & \frac{\partial X}{\partial V} \\ \frac{\partial Y}{\partial U} & \frac{\partial Y}{\partial V} \end{pmatrix}$$

$$\frac{\partial X}{\partial U} = -2\pi\rho \sin(2\pi U), \quad \frac{\partial X}{\partial V} = -\frac{\cos(2\pi U)}{V\rho}$$

$$\frac{\partial Y}{\partial U} = 2\pi\rho \cos(2\pi U), \quad \frac{\partial Y}{\partial V} = -\frac{\sin(2\pi U)}{V\rho}$$

Тогда:

$$\left| \frac{1}{\det(J)} \right| = \frac{V}{2\pi} = \frac{1}{2\pi} e^{-\frac{\rho^2}{2}}$$

Ну и тогда переводя полярные координаты в обычные декартовы:

$$p_{XY} = 1 \cdot \frac{1}{2\pi} e^{-\frac{\rho^2}{2}} = \frac{1}{2\pi} e^{-\frac{\rho^2(\cos^2(\varphi) + \sin^2(\varphi))}{2}} = \frac{1}{2\pi} e^{-\frac{x^2 + y^2}{2}}$$

Что в точности соответствует определению *стандартного гауссовского* вектора.

3 Задача 3.4

$$p(x) = \frac{e^{-|x|}}{2} \cdot (1 - e^{-1})^{-1} \cdot \mathbb{I}(|x| \leq 1)$$

```

1 from scipy import stats
2 import numpy as np
3 from time import time
4
5 class CustomDistribution(stats.rv_continuous):
6     def _pdf(self, x):
7         return np.exp(-np.abs(x)) / (2 * (1 - np.exp(-1))) * (np.abs(x) <= 1)
8
9     def _cdf(self, x):
10        result = np.zeros_like(x, dtype=float)
11
12        cond_first = x < -1
13        result[cond_first] = 0
14
15        cond_second = (x >= -1) & (x <= 0)
16        result[cond_second] = (np.exp(x[cond_second])) / (2 * (1 - np.exp(-1)))
17
18        cond_third = (x > 0) & (x <= 1)
19        result[cond_third] = -(np.exp(-x[cond_third])) / (2 * (1 - np.exp(-1)))
20
21        cond_fourth = x > 1
22        result[cond_fourth] = 1
23
24        return result
25
26 def inverse_cdf_sampling(n):
27     """делаем семплы по обратной функции распределения"""
28     u = np.random.uniform(0, 1, size=n)
29     result = np.zeros_like(u)
30
31     cond1 = u <= 0.5
32     result[cond1] = np.log(2 * u[cond1] * (1 - np.exp(-1)) + np.exp(-1))
33
34     cond2 = u > 0.5
35     result[cond2] = -np.log(1 - 2 * (u[cond2] - 0.5) * (1 - np.exp(-1)))
36
37     return result
38
39 def rejecting_sampling(n):
40     """делаем семплы rejecting sampling"""
41     samples = []
42     pdf_bound = 1 / (1 - np.exp(-1))
43
44     while len(samples) < n:
45         batch_size = min(2 * (n - len(samples)), n)
46         x_proposals = np.random.uniform(0, 1, size=batch_size)
47         pdf_values = np.exp(-np.abs(x_proposals)) / (2 * (1 - np.exp(-1)))
48         u = np.random.uniform(0, pdf_bound, size=batch_size)
49         accepted = x_proposals[u <= pdf_values]
50
51         samples.extend(accepted[:n - len(samples)])
52
53     return np.array(samples)
54
55
56 custom_dist = CustomDistribution()

```

```

57 def method1_test(n_values):
58     """проверяем ООП"""
59     results = {}
60
61     for n in n_values:
62         start_time = time()
63         samples = custom_dist.rvs(size=n)
64         elapsed = time() - start_time
65
66         results[n] = {'samples': samples,
67                     'time': elapsed}
68
69     return results
70
71 def method2_test(n_values):
72     """проверяем обратную функцию распределения"""
73     results = {}
74
75     for n in n_values:
76         start_time = time()
77         samples = inverse_cdf_sampling(n)
78         elapsed = time() - start_time
79
80         results[n] = {'samples': samples,
81                     'time': elapsed}
82
83     return results
84
85 def method3_test(n_values):
86     """проверяем rejecting sampling"""
87     results = {}
88
89     for n in n_values:
90         start_time = time()
91         samples = rejecting_sampling(n)
92         elapsed = time() - start_time
93
94         results[n] = {'samples': samples,
95                     'time': elapsed}
96
97     return results
98
99 def eherements():
100     """эксперементы"""
101     n_values_method1 = [10, 500, 1000]
102     n_values_method23 = [10, 500, 1000, 5000, 10000, 100000]
103
104     results_method1 = method1_test(n_values_method1)
105     results_method2 = method2_test(n_values_method23)
106     results_method3 = method3_test(n_values_method23)
107
108     makefile(results_method1, results_method2, results_method3,
109             n_values_method1, n_values_method23)
110
111 def makefile(results_method1, results_method2, results_method3,
112             n_values_method1, n_values_method23, filename='teorver_stat.txt'):
113     with open(filename, 'w') as f:
114         f.write('{:<10} {:<15} {:<12}\n'.format(
115             "Method", "Sample Size", "Time (s)"))
116         f.write('_' * 40 + '\n')
117
118     for n in n_values_method1:

```

```

119         f.write('{:<10} {:<15d} {:<12.6f}\n'.format(
120             'OOP', n, results_method1[n]['time'])
121         ))
122     f.write('-' * 40 + '\n')
123
124     for n in n_values_method23:
125         for method_name, results in [('Inverse CDF', results_method2),
126                                     ('Rejection', results_method3)]:
127             f.write('{:<10} {:<15d} {:<12.6f}\n'.format(
128                 method_name, n, results[n]['time'])
129             ))
130         f.write('-' * 40 + '\n')
131
132 if __name__ == '__main__':
133     eherements()

```

Method	Sample Size	Time (s)

OOP	10	0.080208
OOP	500	4.034261
OOP	1000	6.558655

Inverse CDF	10	0.000000
Rejection	10	0.000000

Inverse CDF	500	0.000000
Rejection	500	0.000000

Inverse CDF	1000	0.000000
Rejection	1000	0.000000

Inverse CDF	5000	0.001000
Rejection	5000	0.001034

Inverse CDF	10000	0.000000
Rejection	10000	0.001196

Inverse CDF	100000	0.004676
Rejection	100000	0.015585

Rejecting sampling

Метод отбраковки используется для генерации случайных чисел из целевого распределения $p(x)$, если напрямую сгенерировать из него сложно. Для этого выбирается вспомогательное распределение $q(x)$, из которого выборка генерируется легко, и константа M , удовлетворяющая условию $p(x) \leq Mq(x)$ для всех x . Это гарантирует, что график $Mq(x)$ полностью покрывает график $p(x)$.

Алгоритм

1. Сгенерировать точку x из распределения $q(x)$
2. Сгенерировать число u из равномерного распределения $U[0; 1]$
3. Если $u \leq \frac{p(x)}{Mq(x)}$, принять x как образец из $p(x)$. В противном случае — отбраковать x и повторить шаги

Обоснование

Вероятность принятия точки x пропорциональна отношению $\frac{p(x)}{Mq(x)}$. Совместная плотность принятых точек равна $q(x) \cdot \frac{p(x)}{Mq(x)} = \frac{p(x)}{M}$.

После нормировки на вероятность принятия $\int \frac{p(x)}{M} dx = \frac{1}{M}$, получаем, что принятые точки имеют распределение $p(x)$.

Вывод

Как видно из результатов, ООП проигрывает абсолютно. А наиболее эффективным себя показывает метод обратной функции распределения.

4 Задача 4.1

Имеем:

$$X_i \in \text{Bern}(\theta \leftrightarrow p), \theta \leftrightarrow p \in (0; 1)$$

$$P(|\bar{X}_{n_{\varepsilon\delta}} - \mu_\theta| \leq \varepsilon) \geq 1 - \delta \implies P(|\bar{X}_{n_{\varepsilon\delta}} - \mu_\theta| \geq \varepsilon) \leq \delta$$

Наибольшая дисперсия $\sigma^2 = \text{Var}(X) = p(1-p) \rightarrow_{\max} 0.25$. Это следует из анализа производной $(p(1-p))'_p = 1 - 2p$.

Чтобы оценки были верны для всех p , будем пользоваться этим далее.

- Тогда неравенство Чебышёва даёт:

$$P(|\bar{X}_{n_{\varepsilon\delta}} - \mu_\theta| \geq \varepsilon) \leq \frac{\text{Var}(\bar{X}_{n_{\varepsilon\delta}})}{\varepsilon^2} = \frac{\sigma^2}{n_{\varepsilon\delta}\varepsilon}$$

Хотим:

$$\frac{\sigma^2}{n_{\varepsilon\delta}\varepsilon} \leq \delta \implies \frac{\sigma^2}{\varepsilon\delta} \leq n_{\varepsilon\delta}$$

Тогда получим оценку:

$$\frac{\sigma^2}{\varepsilon\delta} \leq n_{\varepsilon\delta} \implies \frac{0.25}{\varepsilon\delta} \leq n_{\varepsilon\delta}$$

- Неравенство Хёфдинга, от части нр-в типа Чернова, даст нам:

$$X_i \in \{a_i; b_i\} = \{0; 1\}$$

$$P(|\bar{X}_n - \mu_\theta| \geq \varepsilon) \leq 2\exp\left(-\frac{2\varepsilon^2 n^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

$$P(|\bar{X}_{n_{\varepsilon\delta}} - \mu_\theta| \geq \varepsilon) \leq 2\exp(-2\varepsilon^2 n_{\varepsilon\delta})$$

Тогда чтобы получить, то что требуется:

$$\exp(-2\varepsilon^2 n_{\varepsilon\delta}) \leq \delta \implies \frac{\ln(\frac{2}{\delta})}{2\varepsilon^2} \leq n_{\varepsilon\delta}$$

- Чтобы получить оценку по ЦПТ:

$$\frac{n(\bar{X}_{n_{\varepsilon\delta}} - \mu_\theta)}{\sqrt{n}\sigma} \xrightarrow{n \rightarrow +\infty} \mathcal{N}(0, 1)$$

Что даст нам:

$$\bar{X}_{n_{\varepsilon\delta}} \approx \mathcal{N}(\mu_\theta, \frac{\sigma^2}{n})$$

Не уходя далеко от стандартного нормального закона, применяя $T = \bar{X}_{n_{\varepsilon\delta}} - \mu_\theta$:

$$P(|\bar{X}_{n_{\varepsilon\delta}} - \mu_\theta| \geq \varepsilon) \approx P(|T| \geq \frac{\varepsilon\sqrt{n_{\varepsilon\delta}}}{\sigma})$$

$$P(|T| \geq \frac{\varepsilon\sqrt{n_{\varepsilon\delta}}}{\sigma}) = 1 - P(|T| \leq \frac{\varepsilon\sqrt{n_{\varepsilon\delta}}}{\sigma}) = 2\Phi\left(-\frac{\varepsilon\sqrt{n_{\varepsilon\delta}}}{\sigma}\right) \leq \delta \implies \frac{\varepsilon\sqrt{n_{\varepsilon\delta}}}{\sigma} \geq Q_{\mathcal{N}}\left(1 - \frac{\delta}{2}\right)$$

Выражая через квантили ($Q_{\mathcal{N}}$) стандартного нормального закона, получаем в итоге оценку на $n_{\varepsilon\delta}$:

$$\frac{\sigma^2 Q_{\mathcal{N}}^2(1 - \frac{\delta}{2})}{\varepsilon^2} \leq \frac{0.25 \cdot Q_{\mathcal{N}}^2(1 - \frac{\delta}{2})}{\varepsilon^2} \leq n_{\varepsilon\delta}$$

Теперь программно проверим наши оценки.

```

1  import numpy as np
2  from scipy.stats import norm
3
4  def compute_n_chebyshev(eps, delta, var=0.5):
5      """Вычисляем n из неравенства Чебышёва"""
6      return int(np.ceil(var**2 / (delta * eps**2)))
7
8  def compute_n_chernoff(eps, delta):
9      """Вычисляем n из неравенства типа Чернова"""
10     return int(np.ceil((1 / (2 * eps**2)) * np.log(2 / delta)))
11
12 def compute_n_cpt(eps, delta, var=0.5):
13     """Вычисляем n по ЦПТ"""
14     Q = norm.ppf(1 - delta / 2)
15     return int(np.ceil((Q * var / eps)**2))
16
17 def generate_sample(n, p):
18     """Генерируем выборку из Bern(p)."""
19     return np.random.binomial(1, p, size=n)
20
21 def is_successful(sample_mean, true_mean, eps):
22     """Проверяем удовлетворены ли условия"""
23     return abs(sample_mean - true_mean) <= eps
24
25 def run_experiment(n, p, epsilon, num_samples=100):
26     successful_count = 0
27     for i in range(num_samples):
28         sample = generate_sample(n, p)
29         sample_mean = np.mean(sample)
30         if is_successful(sample_mean, p, epsilon):
31             successful_count += 1
32     success_rate = successful_count / num_samples
33     return successful_count, success_rate
34
35 if __name__ == '__main__':
36     eps = 0.01
37     delta = 0.05
38     p = 0.5 # Точка с наибольшей дисперсией для Bern(p)
39     N = 100
40
41     n_cheb = compute_n_chebyshev(eps, delta)
42     succ_cheb, rate_cheb = run_experiment(n_cheb, p, eps, N)
43
44     n_cher = compute_n_chernoff(eps, delta)
45     succ_cher, rate_cher = run_experiment(n_cher, p, eps, N)
46
47     n_cpt = compute_n_cpt(eps, delta)
48     succ_cpt, rate_cpt = run_experiment(n_cpt, p, eps, N)
49
50     with open('sammmples.txt', 'w', encoding='utf-8') as f:
51         f.write('Итог:\n')
52         f.write(f'Чебышёв:  n = {n_cheb}, Успешные выборки: {succ_cheb}/{N}, Доля:
53         ↪ {rate_cheb:.2%}\n')
54         f.write(f'Чернов:    n = {n_cher}, Успешные выборки: {succ_cher}/{N}, Доля:
55         ↪ {rate_cher:.2%}\n')
56         f.write(f'ЦПТ:      n = {n_cpt}, Успешные выборки: {succ_cpt}/{N}, Доля:
57         ↪ {rate_cpt:.2%}\n')

```

Итог:

Чебышёв: $n = 50000$, Успешные выборки: $100/100$, Доля: 100.00%

Чернов: $n = 18445$, Успешные выборки: $99/100$, Доля: 99.00%

ЦПТ: $n = 9604$, Успешные выборки: $95/100$, Доля: 95.00%

Оценка Чебышёва в силу грубости даёт наибольший успех в выборках, но требует очень большой объём.

Оценка по неравенству типа Чернова сильно улучшает границу, сохраняя хороший процент доли успеха.

Но всё же ЦПТ справляется наиболее лучшим образом. Теряя лишь пять процентов в точности.