

Курсовая работа по курсу математической статистики

Сааков А.С. СКБ182
Версия от 08.12.2020

Содержание

1. [Вероятностные распределения](#)
 - A. [Геометрическое распределение](#)
 - B. [Распределение Максвелла](#)
2. [Основные понятия математической статистики](#)
 - A. [Геометрическое распределение](#)
 - B. [Распределение Максвелла](#)
3. [Оценки](#)
 - A. [Геометрическое распределение](#)
 - B. [Распределение Максвелла](#)
 - C. [Работа с данными](#)
4. [Проверка статистических гипотез](#)
 - A. [Геометрическое распределение](#)
 - B. [Распределение Максвелла](#)
5. [Различение гипотез](#)
 - A. [Геометрическое распределение](#)
 - B. [Распределение Максвелла](#)
6. [Литература](#)

1. Домашнее задание. Вероятностные распределения

```
In [1]: 1 import numpy as np
        2 import matplotlib.pyplot as plt
        3 from scipy.stats import maxwell
        4 import scipy.stats as sts
        5 from scipy.stats import geom
        6 from random import random
        7 from collections import Counter
        8 import copy
        9 import math
       10 from math import *
       11 from random import *
       12 import pandas as pd
       13 import calendar
       14 import statsmodels.api as sm
       15
       16 plt.style.use('ggplot') # Красивые графики
       17 plt.rcParams['figure.figsize'] = (15, 5) # Размер картинок
```

1.1. Геометрическое распределение

1.1.1. Описание основных характеристик распределения

Функция вероятности дискретного распределения: $P_\xi(x) = pq^x, x \in \{0, 1, 2, \dots\}$

Математическое ожидание:

$$M\xi = \sum_{k=1}^{\infty} k p q^{k-1} = p \sum_{k=1}^{\infty} k q^{k-1} = p \sum_{k=1}^{\infty} \frac{d q^k}{d q} = p \frac{d}{d q} \left(\sum_{k=1}^{\infty} q^k \right) = p \frac{d}{d q} \left(\frac{q}{1-q} \right) = p \frac{1}{(1-q)^2} = \frac{1}{p}$$

Дисперсия:

$$\begin{aligned} D\xi &= M(\xi - M\xi)^2 = M\xi^2 - (M\xi)^2 = M(\xi(\xi-1) + \xi) - M\xi^2 = M(\xi(\xi-1)) + M\xi - (M\xi)^2 = \\ M(\xi(\xi-1)) &= p \sum_{k=1}^{\infty} k^2 q^{k-1} = p q \sum_{k=0}^{\infty} \frac{d^2 q^k}{d q^2} = p q \frac{d^2}{d q^2} \left(\sum_{k=0}^{\infty} q^k \right) = p q \frac{d^2}{d q^2} \left(\frac{1}{1-q} \right) = p q \frac{2}{(1-q)^3} = \\ D\xi &= M\xi^2 + M\xi - (M\xi)^2 = \frac{2q}{p^2} + \frac{1}{p} - \frac{1}{p^2} = \frac{2q-1}{p^2} + \frac{1}{p} = \frac{2q-1+p}{p^2} = \frac{2q-1+1-q}{p^2} = \end{aligned}$$

```
In [2]: 1 for p in [0.1, 0.4, 0.6, 0.9]:
2         geom_rv = sts.geom(p)
3         sample = geom_rv.rvs(1000)
4         plt.hist(sample, density = True, label='p = {}'.format(p))
5         plt.legend()
6         plt.show()
7 print('Рис. 1: 1.1.1. Гистограмма вероятностей дискретного распределения')
```

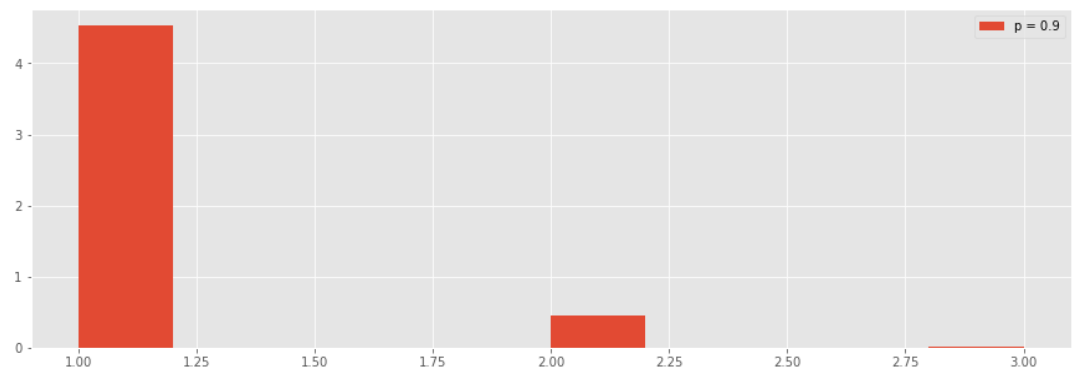
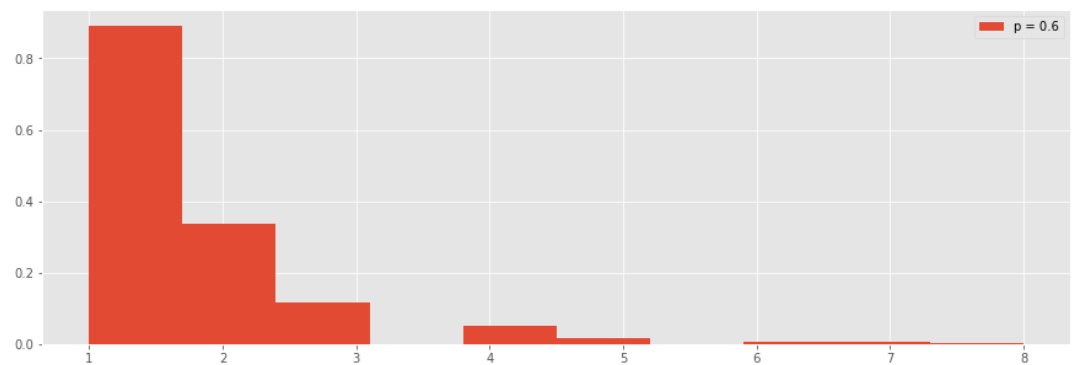
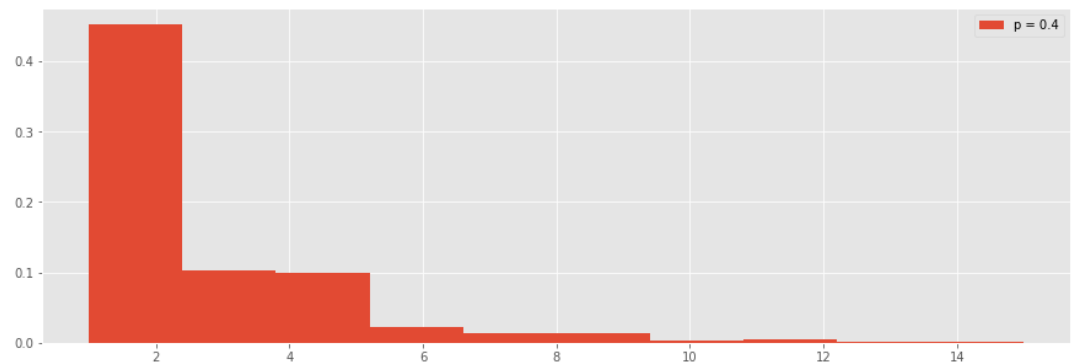
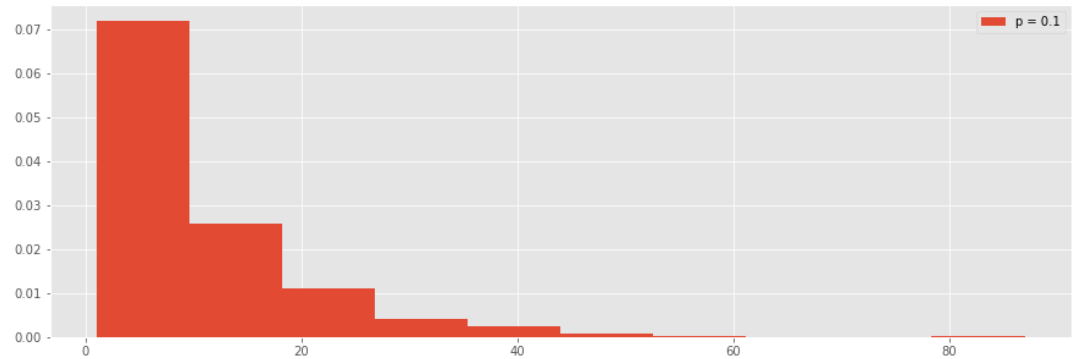


Рис. 1: 1.1.1, Гистограмма вероятностей дискретного распределения

Мода M_0 - значение во множестве наблюдений, которое встречается наиболее часто, для дискретной случайной величины определяется с помощью гистограммы вероятностей. Из гистограмм видно, что $M_0 = 1$

In [3]:

```
1 for p in [0.1, 0.4, 0.6, 0.9]:
2     n = np.arange(0, 8, 1)
3     plt.step(n, 1-(1-p)**(n+1), label='p = {}'.format(p))
4     plt.legend()
5 plt.show()
6 print('Рис. 1: 1.1.1. Гистограмма вероятностей дискретного распределения')
```

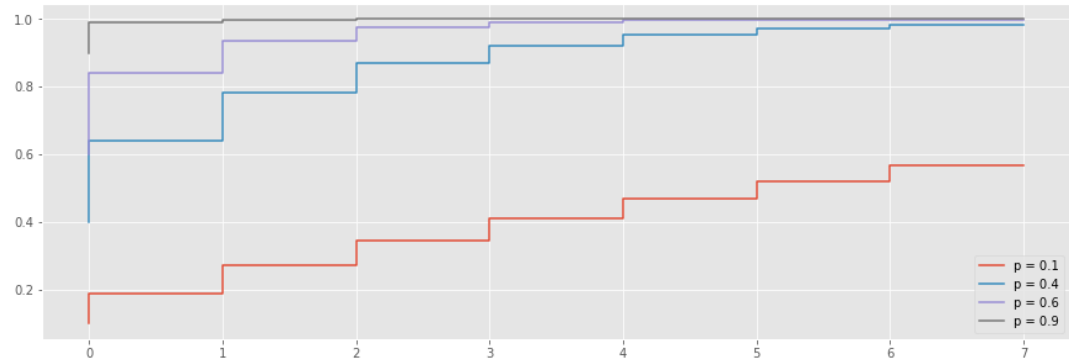


Рис. 1: 1.1.1, Гистограмма вероятностей дискретного распределения

Медиана Me находится из уравнения $P_{\xi}(x) = 0.5$

$$\begin{cases} p + qp + q^2p + \dots + q^{Me-1}p \geq \frac{1}{2} \\ q^{Me-1}p + q^{Me}p + q^{Me+1}p + \dots \geq \frac{1}{2} \end{cases}$$

$$\begin{cases} p \frac{1-q^{Me}}{1-q} \geq \frac{1}{2} \\ q^{Me-1}p \frac{1}{1-q} \geq \frac{1}{2} \end{cases}$$

$$\begin{cases} 1 - q^{Me} \geq 2^{-1} \\ q^{Me-1} \geq 2^{-1} \end{cases}$$

$$\begin{cases} q^{Me} \leq 2^{-1} \\ q^{Me-1} \geq 2^{-1} \end{cases}$$

$$\begin{cases} Me \cdot \log_2 q \leq -1 \\ (Me - 1) \log_2 q \geq -1 \end{cases}$$

Отсюда $-\frac{1}{\log_2 q} \leq Me \leq 1 - \frac{1}{\log_2 q}$

Примеры событий, которые могут быть описаны выбранными случайными величинами

Типичные интерпретации геометрического распределения: описывает количество испытаний n до первого успеха при вероятности наступления успеха в каждом испытании p . Если n подразумевается номер испытания, в котором наступил успех, то геометрическое распределение будет описываться следующей формулой:

$$Geom_p(n) = q^{n-1} p$$

Геометрическое распределение считается дискретной версией экспоненциального распределения. Предположим, что эксперименты Бернулли проводятся через равные промежутки времени. Тогда геометрическая случайная величина X - это время, измеренное в дискретных единицах, которое проходит до того, как мы добьемся первого успеха. Но если мы хотим смоделировать время, прошедшее до того, как данное событие произойдет в непрерывном времени, то подходящим распределением для использования будет экспоненциальное распределение. С математической точки зрения геометрическое распределение обладает тем же свойством без памяти, которым обладает экспоненциальное распределение: в экспоненциальном случае вероятность того, что событие произойдет в течение заданного временного интервала, не зависит от того, сколько времени уже прошло, а событие не произошло; в геометрическом случае вероятность того, что событие произойдет в данный момент (дискретное) времени, не зависит от того, что произошло раньше, потому что эксперимент Бернулли, проведенный в каждый момент времени, не зависит от предыдущих испытаний. Геометрическое распределение полезно для определения вероятности успеха при ограниченном количестве испытаний, что очень применимо к реальному миру, в котором неограниченные испытания редки. Поэтому неудивительно, что различные сценарии хорошо моделируются геометрическими распределениями:

- В спорте, особенно в бейсболе, геометрическое распределение полезно для анализа вероятности того, что отбивающий получит удар, прежде чем он получит три удара; здесь цель - добиться успеха за 3 испытания.
- При анализе затрат и выгод, например, когда компания решает, финансировать ли исследовательские испытания, которые в случае успеха принесут компании некоторую предполагаемую прибыль, цель состоит в том, чтобы достичь успеха до того, как затраты превысят потенциальную выгоду.
- В тайм-менеджменте цель состоит в том, чтобы выполнить задачу за установленный промежуток времени. Другие приложения, подобные вышеупомянутым, также легко создаются. Фактически, геометрическое распределение применяется на интуитивном уровне в повседневной жизни на регулярной основе.

1.1.3 Описание способа моделирования выбранных случайных величин

Существует такой способ реализации метода обратных функций, при котором трудоемкость по крайней мере формально не зависит от p . Действительно, накопленная вероятность $s_{n+1} = p_0 + \dots + p_n$ для геометрического распределения имеет вид

$$s_{n+1} = \sum_{i=0}^n p(1-p)^i = 1 - (1-p)^{n+1}$$

Поэтому событие $\{\xi = n\}$ приобретает вид

$$\begin{aligned} \{\xi = n\} &= \{s_n < \alpha \leq s_{n+1}\} = \{1 - (1-p)^n < \alpha \leq 1 - (1-p)^{n+1}\} = \{(1-p)^{n+1} \leq 1 - \alpha < (1-p)^n\} \\ &= \{(n+1)\ln(1-p) \leq \ln(1-\alpha) < n \cdot \ln(1-p)\} = \{n < \frac{\ln(1-\alpha)}{\ln(1-p)} \leq n+1\}, \end{aligned}$$

и тем самым

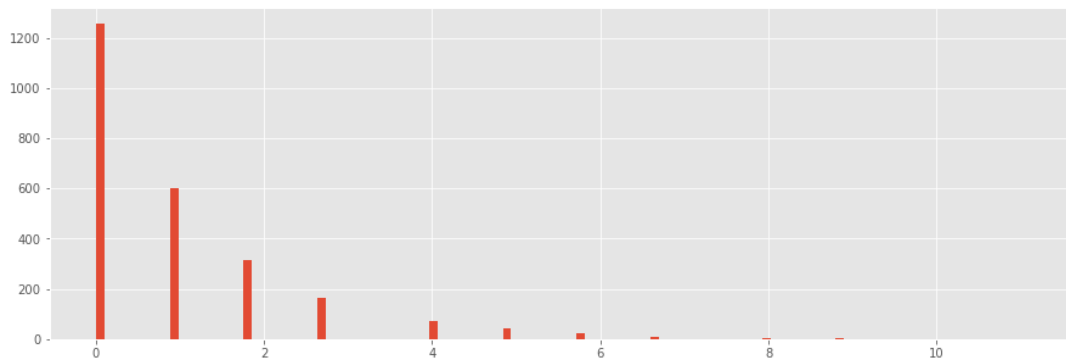
$$\xi = \left\lceil \frac{\ln(1-\alpha)}{\ln(1-p)} \right\rceil$$

Эту же формулу можно получить по-другому. Пусть v - случайная величина, имеющая показательное распределение с параметром λ и $\xi = \lceil n \rceil$. Тогда при $n \geq 0$

$$P(\xi = n) = P(n \leq v < n+1) = e^{-n\lambda} - e^{-(n+1)\lambda} = (1 - e^{-\lambda})e^{-n\lambda}.$$

Поскольку случайная величина $\frac{-\ln(1-\alpha)}{\lambda}$ имеет показательное распределение с параметром λ , то взяв $\lambda = -\ln(1-p)$, приходим к формуле $\xi = \left\lceil \frac{\ln(1-\alpha)}{\ln(1-p)} \right\rceil$

```
In [4]: 1 def sample_(N=2500, scale = 0.5):
2         for x in range(N):
3             je = np.log(random())//np.log(1-scale)#Генерирование случайных чисел
4         return je
5 def Geom(n, p=0.5):
6     x=[sample_(scale=p) for x in range(n)]
7     #print(x)
8     return x
9 plt.hist(Geom(2500,0.5),25, width = 0.1)
10 plt.show()
```



1.2. Распределение Максвелла

1.2.1. Описание основных характеристик распределения

Математическое ожидание:

$$M\xi = \int_0^{\infty} x \sqrt{\frac{2}{\pi}} \frac{x^2}{\lambda^3} e^{-\frac{x^2}{2\lambda^2}} dx = \sqrt{\frac{2}{\pi}} \frac{1}{\lambda^3} \int_0^{\infty} x^3 e^{-\frac{x^2}{2\lambda^2}} dx = 2\lambda^4 \cdot \sqrt{\frac{2}{\pi}} \frac{1}{\lambda^3} = 2\lambda \sqrt{\frac{2}{\pi}}$$

Дисперсия:

$$D\xi = M(\xi - M\xi)^2 = M\xi^2 - (M\xi)^2 = M(\xi(\xi - 1) + \xi) - (M\xi)^2 = M(\xi(\xi - 1)) + M\xi - (M\xi)^2$$

$$M(\xi(\xi - 1)) = \int_0^{\infty} x^2 f(x) dx = \int_0^{\infty} x^2 \sqrt{\frac{2}{\pi}} \frac{x^2}{\lambda^3} e^{-\frac{x^2}{2\lambda^2}} dx = \sqrt{\frac{2}{\pi}} \frac{1}{\lambda^3} \int_0^{\infty} x^4 e^{-\frac{x^2}{2\lambda^2}} dx = \sqrt{\frac{2}{\pi}} \frac{1}{\lambda^3} \cdot 3\lambda^5$$

$$D\xi = M\xi^2 - (M\xi)^2 = 3\lambda^2 - 4\lambda^2 \cdot \frac{2}{\pi} = \frac{3\pi - 8}{\pi} \lambda^2$$

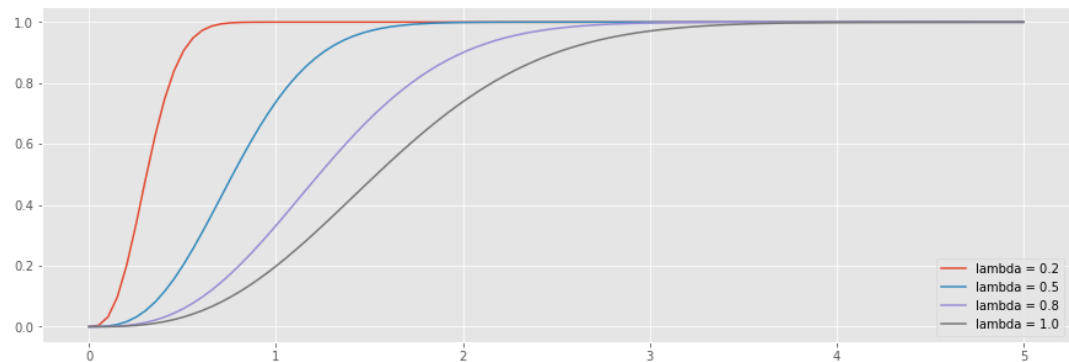
Также использовались известные интегралы, который был взят из курса физики:

$$\int_0^{\infty} x^3 e^{-x^2/2\lambda^2} dx = \frac{1}{2(\frac{1}{\lambda^2})^2} \cdot 4 = 2\lambda^4$$

$$\int_0^{\infty} x^4 e^{-x^2/2\lambda^2} dx = \frac{3}{8} \sqrt{\pi} \left(\frac{1}{2\lambda^2}\right)^{-\frac{5}{2}} = 3\lambda^5 \sqrt{\frac{\pi}{2}}$$

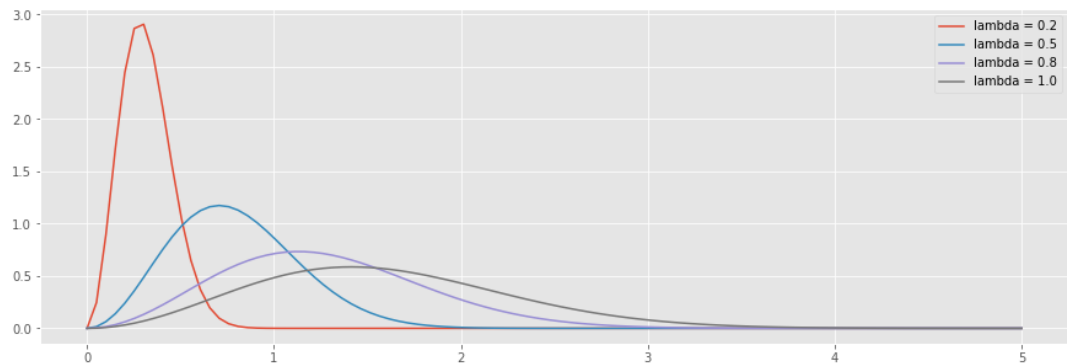
```
In [5]: 1 for lambd in [0.2,0.5,0.8,1.0]:
2         maxwell_rv = sts.maxwell(scale = lambd)
3         x = np.linspace(0,5,100)
4         cdf = maxwell_rv.cdf(x)
5         plt.plot(x, cdf, label = 'lambda = {}'.format(lambd))
6         plt.legend()
7 print('Рис.4: График функции распределения')
```

Рис.4: График функции распределения



```
In [6]: 1 for lambda in [0.2,0.5,0.8,1.0]:
2         maxwell_rv = sts.maxwell(scale = lambda)
3         x = np.linspace(0,5,100)
4         pdf = maxwell_rv.pdf(x)
5         k = max(pdf)
6         plt.plot(x, pdf, label = 'lambda = {}'.format(lambda))
7         plt.legend()
8     print('\n')
9     print('Рис.5: График плотности вероятности распределения')
```

Рис.5: График плотности вероятности распределения



Модой абсолютно непрерывного распределения называют любую точку локального максимума

плотности распределения. $f'(x) = \sqrt{\frac{2}{\pi}} \frac{x^2}{\lambda^3} e^{-\frac{x^2}{2\lambda^2}} = \frac{4x}{\lambda^3 \sqrt{2\pi}} e^{-\frac{x^2}{2\lambda^2}} - \frac{2x^2}{\lambda^3 \sqrt{2\pi}} e^{-\frac{x^2}{2\lambda^2}} \frac{x}{\lambda^2} = 0$

$$\begin{aligned} \frac{4x}{\lambda^3 \sqrt{2\pi}} e^{-\frac{x^2}{2\lambda^2}} &= \frac{2x^2}{\lambda^3 \sqrt{2\pi}} e^{-\frac{x^2}{2\lambda^2}} \frac{x}{\lambda^2} \\ 4 &= 2x \frac{x}{\lambda^2} \\ x^2 &= 2\lambda^2 \\ x = M_0 &= \lambda\sqrt{2} \end{aligned}$$

Медиана

$$\begin{aligned} \int_0^{Me} \sqrt{\frac{2}{\pi}} \frac{1}{\lambda^3} x^2 e^{-\frac{x^2}{2\lambda^2}} dx &= \frac{1}{2} \\ \int_0^{Me} x^2 e^{-\frac{x^2}{2\lambda^2}} dx &= \frac{\lambda^3}{2} \sqrt{\frac{\pi}{2}} \\ (-\lambda^2 e^{-\frac{x^2}{2\lambda^2}} x + \lambda^3 \sqrt{\frac{\pi}{2}}) \Big|_0^{Me} &= \frac{\lambda^3}{2} \sqrt{\frac{\pi}{2}} \\ -Me\lambda^2 e^{-\frac{Me^2}{2\lambda^2}} &= -\frac{\lambda^3}{2} \sqrt{\frac{\pi}{2}} \\ Me \cdot e^{-\frac{Me^2}{2\lambda^2}} &= \frac{\lambda^2}{2} \sqrt{\frac{\pi}{2}} \\ Me &\approx 1,5383\lambda \end{aligned}$$

1.2.2. Примеры событий, которые могут быть описаны выбранными случайными величинами

Впервые распределение было определено и использовалось для описания скоростей частиц в идеализированных газах, где частицы свободноперемещаются внутри стационарного контейнера, не взаимодействуя друг с другом, за исключением очень коротких столкновений, в которых они обмениваются энергией и импульсом друг с другом или со своим тепловым окружением. Термин «частица» в этом контексте относится только к газообразным частицам (атомам или молекулам), и предполагается, что система частиц достигла термодинамического равновесия. Энергии таких частиц следуют так называемой статистике Максвелла – Больцмана, а статистическое распределение скоростей выводится путем приравнивания энергии частиц к кинетической энергии. Распределение Максвелла – Больцмана в основном применяется к скоростям частиц в трех измерениях, но оказывается, что оно зависит только от скорости (величины скорости) частиц. Распределение вероятности скорости частицы указывает, какие скорости более вероятны: частица будет иметь скорость, выбранную случайным образом из распределения, и с большей вероятностью будет находиться в одном диапазоне скоростей, чем в другом.

При тепловом равновесии ($T = \text{const}$) $u_{\text{кв}}$ молекул газа остается постоянной и равной $u = \sqrt{\frac{3kT}{m}}$

Это объясняется тем, что в газе устанавливается стационарное статическое распределение молекул по значениям скоростей, называемое распределением Максвелла:

$$f(u) = \frac{dN(u)}{Ndu} = 4\pi \left(\frac{m}{2\pi kT} \right)^{\frac{3}{2}} \cdot u^2 \cdot e^{-\frac{mu^2}{2kT}}$$

В теории вероятностей рассматривается распределения Максвелла, в котором $x = u$ и $\frac{1}{\lambda^2} = \frac{m}{kT}$

Нетипичной интерпретацией распределения Максвелла будут данные, которые представляют время ремиссии (в месяцах) у пациентов с раком мочевого пузыря и первоначально использовались Lee и Wang.

Ремиссия (лат. remissio «уменьшение, ослабление») — период течения хронической болезни, который проявляется значительным ослаблением (неполная ремиссия) или исчезновением (полная ремиссия) её симптомов (признаков заболевания)

1.2.3. Описание способа моделирования выбранных случайных величин

Способ 1: Существует полярный метод (группа полярных методов предназначена для моделирования распределений, так или иначе связанных с двумерными распределениями, инвариантными относительно вращений), где моделируются две независимые случайные величины ξ_1, ξ_2 , каждая из которых имеет распределение $N(0,1)$.

Полярные координаты. Каждая точка $X = (x, y)^T \in R^2 \setminus \{0\}$ может быть однозначно представлена в виде $X = ||X||\bar{e}$, где $\bar{e} = 1$. Полагая $s = ||X||$ и $\bar{e} = (\cos t, \sin t)^T$, где $t \in [0, 2\pi)$, получаем биекцию $\phi : (x, y)^T \rightarrow (s, t)^T$, действующую из $R^2 \setminus \{0\}$ в $(0, \infty) \times [0, 2\pi)$. Конечно, переменные (s, t) являются полярными координатами вектора X , а обратное отображение $\psi : (0, \infty) \times [0, 2\pi) \rightarrow R^2 \setminus \{0\}$ имеет вид $x = s \cos t, y = s \sin t$ с якобианом $\det \psi'(s, t) = s$.

Если теперь рассмотреть случайный вектор $\bar{\xi} \in R^2$ с плотностью распределения $p_{\bar{\xi}}(x, y)$ и обозначить r, φ (случайные) полярные координаты этого вектора, то, так как в этом случае $n = 1$, мы получим из (7.1.1), что

$$p_{r,\varphi}(s, t) = s p_{\bar{\xi}}(s \cos t, s \sin t) 1_{(0,\infty) \times [0,2\pi)}(s, t). \quad (1.2.3.1)$$

Выражение (1.2.3.1) выглядит особенно просто, если существует такая функция

$f : (0, \infty) \rightarrow (0, \infty)$, что $p_{\bar{\xi}}(x, y) = f(\sqrt{x^2 + y^2})$. В этом случае, очевидно,

$$p_{r,\varphi}(s, t) = s f(s) 1_{(0,\infty) \times [0,2\pi)}(s, t) = 2\pi s f(s) I_{(0,\infty)}(s) \frac{1}{2\pi} I_{[0,2\pi)}(t). \quad (1.2.3.2)$$

Это значит, что случайные величины r и φ независимы, $\varphi \in U(0, 2\pi)$, а r имеет плотность распределения $p_r(s) = 2\pi s f(s), s > 0$

Действительно, поскольку совместная плотность распределения ξ_1, ξ_2 имеет вид

$$p(x, y) = \frac{1}{2\pi} e^{-\frac{x^2+y^2}{2}}, x, y \in R,$$

то, как следовательно из вышенаписанного, полярный радиус r и полярный угол φ случайного вектора (ξ_1, ξ_2) независимы, причем полярный угол равномерно распределен на $[0, 2\pi)$, а

полярный радиус имеет распределение Рэлея ($p(x) = x e^{-\frac{x^2}{2}}, x > 0$).

Отсюда, применяя моделирующую формулу ($\xi = \sqrt{-2 \ln(\alpha)}$), сразу же приходим к представлению

$$\begin{aligned} \xi_1 &= \sqrt{-2 \ln(\alpha_1)} \cos(2\pi\alpha_2), \\ \xi_2 &= \sqrt{-2 \ln(\alpha_1)} \sin(2\pi\alpha_2), \end{aligned}$$

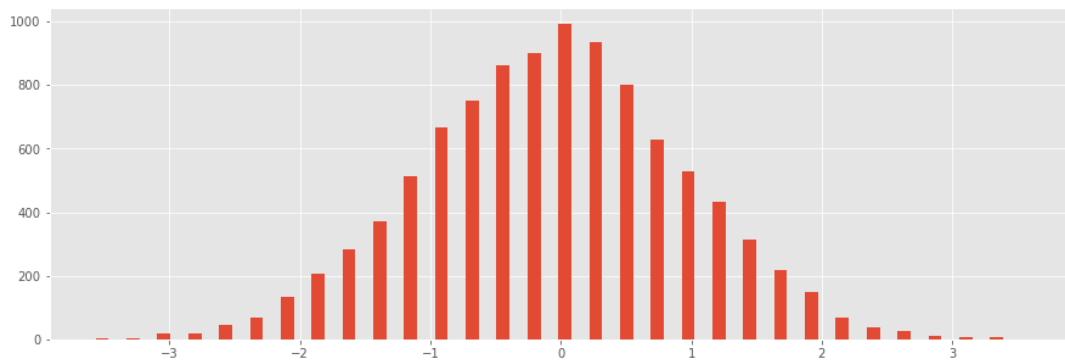
где $\alpha_1, \alpha_2 \in U(0, 1)$

В итоге получаем следующий алгоритм:

```

In [7]: 1 %matplotlib inline
2 import random
3 import pandas as pd
4 import math
5 from scipy import stats
6 import sys
7 plt.style.use('ggplot') # Красивые графики
8 plt.rcParams['figure.figsize'] = (15, 5) # Размер картинок
9 N=10000
10 random.seed(123)
11 epsilon = sys.float_info.epsilon
12
13 def box_muller():
14     u1, u2 = 0.0, 0.0
15     while u1 < epsilon or u2 < epsilon:
16         u1 = random.random()
17         u2 = random.random()
18
19     n1 = math.sqrt(-2 * math.log(u1)) * math.cos(2 * math.pi * u2)
20     n2 = math.sqrt(-2 * math.log(u1)) * math.sin(2 * math.pi * u2)
21     return n1, n2
22
23 # Use KS to test
24 samples = [box_muller()[0] for x in range(N)]
25 test_stat, pvalue = stats.kstest(samples, 'norm', args=(0, 1), N=N)
26
27 # Plot our samples against our reference distribution
28 plt.hist(samples, 30, width = 0.1)
29 plt.show()

```



Стоит заметить, что данный способ достаточно быстр.

Я бы ещё подметил тот факт, что распределение Максвелла с параметром $\lambda = 1$ очень схоже с $N(0,1)$:

$$f(x) = \sqrt{\frac{2}{\pi}} e^{-\frac{x^2}{2}} x^2 - \text{Максвелла}$$

$$f(x) = \sqrt{\frac{2}{\pi}} e^{-\frac{x^2}{2}} \cdot \frac{1}{2} - N(0, 1)$$

Способ 2: По определению, случайный вектор $\bar{\xi} = (\xi_1, \xi_2)^T$ равномерно распределен на единичной окружности S^1 с центром в нуле, если $\xi_1^2 + \xi_2^2 = 1$ с вероятностью 1 и если полярный угол φ вектора $\bar{\xi}$ равномерно распределен на $[0, 2\pi)$. Из этого определения сразу же следует моделирующая формула для равномерного распределения на S^1 :

$$\begin{aligned}\xi_1 &= \cos(2\pi\alpha) \\ \xi_2 &= \sin(2\pi\alpha)\end{aligned}\quad (1.2.3.3)$$

Вычисление двух тригонометрических функций, однако, может оказаться трудоемкой операцией. Стандартной альтернативой формуле (1.2.3.3) является использование метода отбора для моделирования равномерного распределения в единичном круге $B_1(0) = \{(x, y) : x^2 + y^2 < 1\}$ с центром в нуле с последующей нормировкой результата. Формальное обоснование этой процедуры представим ниже.

Аналогично полярным координатам на плоскости, каждый ненулевой вектор $X = (x, y, z)^T \in R^3$ может быть однозначно представлен в виде $X = \|X\| \bar{e}$, где

$$\bar{e} = (\cos(t) \cos(u), \sin(t) \cos(u), \sin(u))^T, \quad t \in [0, 2\pi), u \in [-\pi/2, \pi/2].$$

Это, конечно, соответствует переходу от евклидовой системы координат (x, y, z) к сферической системе (s, t, u) со сферическим радиусом $s = \|X\|$, долготой s и широтой u . Хорошо известно, что якобиан обратного отображения равен $s^2 \cos(u)$.

Поэтому, если случайный вектор $\bar{\xi} = (\xi_1, \xi_2, \xi_3)^T$ имеет плотность распределения $p_\xi(x, y, z)$, то сферические координаты r, φ, θ этого вектора имеют совместную плотность

$$pr, \varphi, \theta(s, t, u) = p_\xi(s \cos(t) \cos(u), s \sin(t) \cos(u), s \sin(u)) s^2 \cos(u), \quad (1.2.3.4)$$

сосредоточенную в области $(0, \infty) \times [0, 2\pi) \times (-\pi/2, \pi/2)$. В случае, когда

$$p_\xi(x, y, z) = f(\sqrt{x^2 + y^2 + z^2}), \quad (1.2.3.5)$$

равенство (1.2.3.4) приобретает вид

$$p_{r, \varphi, \theta}(s, t, u) = 4\pi s^2 f(s^2) I_{(0, \infty)}(s) \frac{1}{2\pi} I_{(0, 2\pi)} \frac{\cos(u)}{2} I_{(-\frac{\pi}{2}, \frac{\pi}{2})}.$$

Таким образом, случайные величины r, φ и θ оказываются независимыми, причем долгота φ равномерно распределена на $(0, 2\pi)$, плотность $p_r(s)$ распределения r равна $4\pi s^2 f(s^2)$, а плотность $p_\theta(u)$ распределения широты θ сосредоточена на $(-\frac{\pi}{2}, \frac{\pi}{2})$ и равна на этом интервале $0.5 \cos(u)$.

Например, если $\bar{\xi} = (\xi^1, \xi^2, \xi^3)^T$ — случайный вектор с независимыми $N(0, 1)$ -распределенными координатами, то его плотность распределения равна $(2\pi)^{-\frac{3}{2}} e^{-\frac{x^2+y^2+z^2}{2}}$, и длина $r = \|\bar{\xi}\|$ этого вектора будет иметь плотность распределения $\sqrt{\frac{2}{\pi}} s^2 e^{-\frac{s^2}{2}}$.

Как уже обсуждалось при первом способе моделирования, если случайный вектор ξ имеет

распределение (1.2.3.5), то вектор $v = \frac{\bar{\xi}}{\|\bar{\xi}\|}$ равномерно распределен на поверхности сферы

$\{(x, y, z) : x^2 + y^2 + z^2 = 1\}$. С другой стороны, координаты v_1, v_2, v_3 этого вектора выражаются через случайные величины φ и θ как

$$v_1 = \cos(\varphi) \cos(\theta), v_2 = \sin(\varphi) \cos(\theta), v_3 = \sin(\theta). \quad (1.2.3.6)$$

Предложение 1:

Пусть $\bar{\xi} \in R^d$ — d -мерный случайный вектор, обладающий распределением P_ξ с плотностью распределения p_ξ , причем $P_\xi(D) = 1$ для некоторого измеримого $D \subset R^d$. Рассмотрим измеримое отображение $\varphi : D \rightarrow R^d$ и предположим, что при $i = 1, \dots, n$ существуют открытые попарно непересекающиеся подмножества $D_i \subset R^d$, удовлетворяющие следующим условиям.

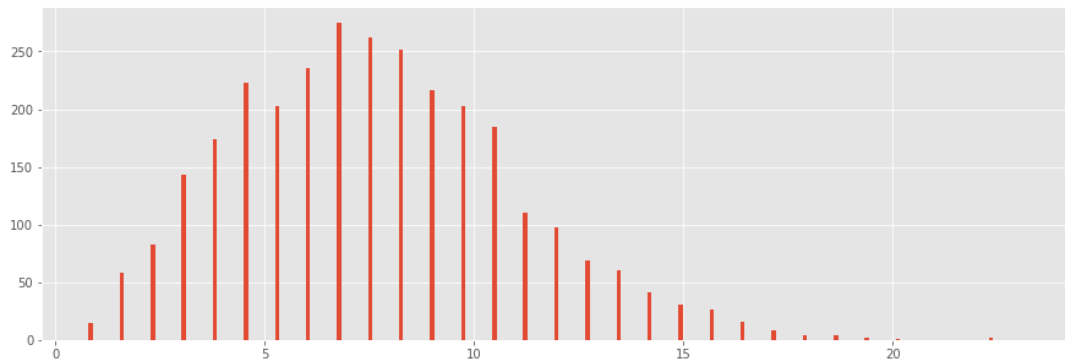
1. Множества D и $\cup_{i=1}^n D_i$ совпадают P_ξ -почти всюду.
2. Обозначим $\phi_i = \phi|_{D_i}$, $\phi(D_i) = \phi_i(D_i) = G_i$ и предположим, что при любом i отображение $\phi_i : D_i \rightarrow G_i$ является биекцией с $\psi_i = \phi_i^{-1}$; непрерывно дифференцируемо с $\det \psi_i' \neq 0$ в G_i .

Тогда случайный вектор $\bar{\eta} = \psi(\bar{\xi})$ обладает плотностью $p_\eta(Y)$ и

$$p_\eta(Y) = \sum_i p_\xi \phi_i(Y) \det \psi_i'(Y) I_{G_i}(Y).$$

—

```
In [8]: 1 def Maxwell(n, lambd = 1):
2         x = [sample_(scale = lambd) for x in range(n)]
3         y=[sample_(scale = lambd) for x in range(n)]
4         z=[sample_(scale = lambd) for x in range(n)]
5         l = []
6         #print(x)
7         for i in range(n):
8             l.append(np.sqrt(x[i]**2+y[i]**2+z[i]**2))
9         return l
10 # Our sample function of N(0,1) using Equation (2)
11 def sample_(N = 3000, scale = 1):
12     return scale*2.0*np.sqrt(N)*(sum(randint(0,1) for x in range(N))/N-0.5)
13 plt.hist(Maxwell(3000,5),30, width = 0.1)
14 plt.show()
```



2 Домашнее задание. Основные понятия математической статистики

2.1 Геометрическое распределение

2.1.1 Моделирование выбранных случайных величин

```
In [66]: 1 # Создание случайной величины с геометрическим распределением, зависящим
2         # от параметра p
3         p = 0.5
4         geom_rv = sts.geom(n)
```

```
In [67]: 1 #Генерация выборки объема n = 5 с выводом
2         for n in [5]:
3             means_5 = []
4             for i in range(5):
5                 sample = geom_rv.rvs(n)
6                 means_5.append(sample)
7             print(sample)
```

```
[2 2 1 3 3]
[2 2 1 2 2]
[3 1 4 1 1]
[1 2 1 1 1]
[5 2 1 1 3]
```

```
In [68]: 1 #Генерация выборки объема n = 10 с выводом
2 for n in [10]:
3     means_10 = []
4     for i in range(5):
5         sample = geom_rv.rvs(n)
6         means_10.append(sample)
7         print(sample)
```

```
[1 1 2 1 4 2 1 1 2 1]
[2 1 3 1 3 1 1 1 3 1]
[1 2 2 2 2 1 4 1 1 1]
[2 2 2 2 1 1 1 2 1 1]
[4 1 1 2 1 1 1 2 3 2]
```

```
In [69]: 1 #Генерация выборки объема n = 100 ,без вывода
2 for n in [100]:
3     means_100 = []
4     for i in range(5):
5         sample = geom_rv.rvs(n)
6         means_100.append(sample)
```

```
In [70]: 1 #Генерация выборки объема n = 1000 ,без вывода
2 for n in [1000]:
3     means_1000 = []
4     for i in range(5):
5         sample = geom_rv.rvs(n)
6         means_1000.append(sample)
```

```
In [71]: 1 #Генерация выборки объема n = 100000 ,без вывода
2 for n in [100000]:
3     means_100000 = []
4     for i in range(5):
5         sample = geom_rv.rvs(n)
6         means_100000.append(sample)
```

2.1.2 Построение эмпирической функции распределения

```
In [31]: 1 #n=5
2 for a in range(5):
3     b=means_5[a]
4     b=sorted(b)
5     print('Empirical distribution function F5(x) for sample',a+1,':')
6     for i in range(4):
7         if(i==0):
8             n=0.
9             g=1
10            print(n,' , x <=',b[i])
11            if(b[i+1]==b[i]):
12                g+=1
13            else:
14                n=round(n+0.2*g,1)
15                g=1
16                print(n,' ,',b[i], '< x <=',b[i+1])
17            if(i==3):
18                n=1.
19                print(n,' , x > ',b[i+1])
```

Empirical distribution function F5(x) for sample 1 :

```
0.0 , x <= 1
0.4 , 1 < x <= 2
0.6 , 2 < x <= 3
0.8 , 3 < x <= 5
1.0 , x > 5
```

Empirical distribution function F5(x) for sample 2 :

```
0.0 , x <= 1
0.6 , 1 < x <= 2
0.8 , 2 < x <= 6
1.0 , x > 6
```

Empirical distribution function F5(x) for sample 3 :

```
0.0 , x <= 1
0.2 , 1 < x <= 2
0.6 , 2 < x <= 4
0.8 , 4 < x <= 5
1.0 , x > 5
```

Empirical distribution function F5(x) for sample 4 :

```
0.0 , x <= 1
0.8 , 1 < x <= 3
1.0 , x > 3
```

Empirical distribution function F5(x) for sample 5 :

```
0.0 , x <= 1
0.2 , 1 < x <= 2
0.4 , 2 < x <= 3
0.8 , 3 < x <= 4
1.0 , x > 4
```

In [32]:

```
1 #n=5
2 for a in range(5):
3     b=means_5[a]
4     b=sorted(b)
5     v=len(b)
6     N=[]
7     for i in range(b[v-1]):
8         N.append(b.count(i))
9         x1=[]
10        y1=[]
11        t=0
12        for i in range(b[v-1]):
13            t+=N[i]
14            x1.append(i)
15            y1.append(t/v)
16            x1.append(i+1)
17            y1.append(t/v)
18        x1.append(b[v-1])
19        y1.append(1)
20        x1.append(b[v-1]+2)
21        y1.append(1)
22        plt.plot(x1,y1,label="ECDF "+str(a+1))
23        plt.legend(loc='lower right')
24 plt.title("Эмпирическая функция выборки объема: "+str(v))
25 n=np.arange(1,8,1)#Построение
26 plt.step(n,1-(1-p)**(n-1),'k-', label='CDF')#теоретической функции
27 plt.legend()#распределения
28 plt.xlabel("numbers")
29 plt.ylabel("probability")
30 plt.show()
31 print("\n")
32 #n=10
33 for a in range(5):
34     b=means_10[a]
35     b=sorted(b)
36     v=len(b)
37     N=[]
38     for i in range(b[v-1]):
39         N.append(b.count(i))
40         x2=[]
41         y2=[]
42         t=0
43         for i in range(b[v-1]):
44             t+=N[i]
45             x2.append(i)
46             y2.append(t/v)
47             x2.append(i+1)
48             y2.append(t/v)
49         x2.append(b[v-1])
50         y2.append(1)
51         x2.append(b[v-1]+2)
52         y2.append(1)
53         plt.plot(x2,y2,label="ECDF "+str(a+1))
54         plt.legend(loc='lower right')
55 plt.title("Эмпирическая функция выборки объема: "+str(v))
56 n=np.arange(1,8,1)#Построение
57 plt.step(n,1-(1-p)**(n-1),'k-', label='CDF')#теоретической функции
58 plt.legend()#распределения
59 plt.xlabel("numbers")
60 plt.ylabel("probability")
61 plt.show()
62 print("\n")
63 #n=100
64 for a in range(5):
65     b=means_100[a]
66     b=sorted(b)
67     v=len(b)
68     N=[]
```


Пусть $X = (X_1, \dots, X_n)$ - выборка из дискретного распределения $\sigma(\xi)$ Величина скачка в точке j есть

$$\Delta \hat{F}_n(j) = \hat{F}_n(j) - \hat{F}_n(j-1) = \frac{v_j}{n},$$

$j = 1, \dots, N$

Здесь $P\{\Delta \hat{F}_n(j) = 0\} = P(v_j = 0) = (1 - p_j)^n$ что мало при больших n , т.е. в большой выборке скачок в точке j наверняка будет иметь место. Более того, так как

$P\{\cup_{j=1}^N \{\Delta \hat{F}_n(j) = 0\}\} \leq \sum_{j=1}^N (1 - p_j)^n \rightarrow 0$, при $n \rightarrow \infty$, то в больших выборках с вероятностью, близкой к 1, скачки э.ф.р. $F_n(x)$ будут иметь место во всех точках $1, 2, \dots, N$, а случайными будут лишь величины этих скачков.

Если же теоретическая функция распределения $F_\xi = F(x)$ непрерывна, то с вероятностью 1 все элементы выборки $X = (X_1, \dots, X_n)$ будут различны, и случайными теперь будут точки скачков, величины же скачков неслучайны и равны $\frac{1}{n}$.

Таким образом, для выборок из дискретных и непрерывных распределений характер соответствующих эмпирических функций распределения будет различным, что можно заметить на получившихся графиках для дискретного и непрерывного распределения. Тем не менее в любом случае э.ф.р. $\hat{F}_n(x)$ с увеличением объема выборки n сближается в каждой точке x с теоретической функцией распределения $F(x)$.

Максимальная точная верхняя граница разности пары эмпирических функций распределения - наибольшая разность между значениями вероятности двух функций в одной точке.

Верхняя граница разности э.ф.р. выборок размера $n = 5$: 0.600

Верхняя граница разности э.ф.р. выборок размера $n = 10$: 0.400

Верхняя граница разности э.ф.р. выборок размера $n = 100$: 0.17

Верхняя граница разности э.ф.р. выборок размера $n = 1000$: 0.061

Верхняя граница разности э.ф.р. выборок размера $n = 100000$: 0.0049

С увеличением объема выборки верхняя граница разности уменьшается, что, впрочем, очевидно.

2.1.3 Построение вариационного ряда выборки

Определение:

Пусть $X = (X_1, \dots, X_n)$ - выборка из некоторого распределения $\sigma(\xi)$

Произвольной реализации $x = (x_1, \dots, x_n)$ этой выборки можно поставить в соответствие упорядоченную последовательность

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

располагая x_1, \dots, x_n в порядке их возрастания, так что $x_{(1)} = \min\{x_1, \dots, x_n\}$, $x_{(2)}$ - второе по величине значение, $x_{(n)} = \max\{x_1, \dots, x_n\}$

Обозначим через $X_{(k)}$ случайную величину, которая для каждой реализации выборки X принимает значение $x_{(k)}$, $k = 1, \dots, n$. Так по выборке X определяют новую последовательность случайных величин $X_{(1)}, \dots, X_{(n)}$, называемых *порядковыми статистиками* выборки. Из определения порядковых статистик следует, что они упорядочены по возрастанию их значений, т.е. они образуют возрастающую последовательность

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)},$$

которая называется *вариационным рядом* выборки X .

```
In [33]: 1 #Вариационный ряд для выборки объема n=5 с выводом
          2 for a in range(5):
          3     b=means_5[a]
          4     b=sorted(b)
          5     print(b)
```

[1, 1, 2, 3, 5]

[1, 1, 1, 2, 6]

[1, 2, 2, 4, 5]

[1, 1, 1, 1, 3]

[1, 2, 3, 3, 4]

```
In [34]: 1 #Вариационный ряд для выборки объема n=10 с выводом
2 for a in range(5):
3     b=means_10[a]
4     b=sorted(b)
5     print(b)
```

```
[1, 1, 1, 1, 1, 2, 2, 3, 5, 5]
[1, 1, 1, 2, 3, 3, 3, 3, 3, 4]
[1, 1, 1, 2, 2, 2, 2, 2, 3, 4]
[1, 1, 1, 1, 1, 2, 2, 2, 3, 3]
[1, 1, 2, 2, 2, 2, 3, 3, 5, 6]
```

```
In [35]: 1 #Вариационный ряд для выборки объема n=100 без вывода
2 for a in range(5):
3     b=means_100[a]
4     b=sorted(b)
```

```
In [36]: 1 #Вариационный ряд для выборки объема n=1000 без вывода
2 for a in range(5):
3     b=means_1000[a]
4     b=sorted(b)
```

```
In [37]: 1 #Вариационный ряд для выборки объема n=100000 без вывода
2 for a in range(5):
3     b=means_100000[a]
4     b=sorted(b)
```

Определение:

α - квантиль случайной величины ξ с функцией распределения $F(x) = P\{\xi < x\}$ — это любое число x_α , удовлетворяющее двум условиям:

1) $F(x_\alpha) \leq \alpha$ 2) $F(x_\alpha + 0) \geq \alpha$.

Исходя из того, что при больших выборках э.ф.р. стремится к теоритической функции распределения, эмпирические квантили так же стремятся к теоритическим по определению.

Пусть $F(x)$ - функция распределения. Тогда квантильная функция:

$F^{-1}(r) = \min\{x \in N_+ : F(x) \geq r\}$ for $r \in (0; 1)$

$F^{-1}(r) = \left[\frac{\ln(1-r)}{\ln(1-p)} \right]$

```
In [38]: 1 k = 1
2 for b in [means_5[a], means_10[a], means_100[a], means_1000[a], means_100000[a]]:
3     if(k==1):
4         print('n = 5')
5     if(k==2):
6         print('n = 10')
7     if(k==3):
8         print('n = 100')
9     if(k==4):
10        print('n = 1000')
11    if(k==5):
12        print('n = 100000')
13    print(np.quantile(b, 0.1))
14    k += 1
```

```
n = 5
1.4
n = 10
1.0
n = 100
1.0
n = 1000
1.0
n = 100000
1.0
```

```
In [39]: 1 #Сравнение
         2 np.log(1-0.1)//np.log(1-n)
```

Out[39]: 0.0

```
In [40]: 1 k = 1
         2 for b in [means_5[a], means_10[a], means_100[a], means_1000[a], means_100000[a]
         3     if(k==1):
         4         print('n = 5')
         5     if(k==2):
         6         print('n = 10')
         7     if(k==3):
         8         print('n = 100')
         9     if(k==4):
        10         print('n = 1000')
        11     if(k==5):
        12         print('n = 100000')
        13     print(np.quantile(b, 0.5))
        14     k += 1
```

```
n = 5
3.0
n = 10
2.0
n = 100
1.0
n = 1000
2.0
n = 100000
1.0
```

```
In [41]: 1 #Сравнение
         2 geom.median(n)
```

Out[41]: 1.0

```
In [42]: 1 #Сравнение
         2 np.log(1-0.5)//np.log(1-n)
```

Out[42]: 1.0

```
In [43]: 1 k = 1
         2 for b in [means_5[a], means_10[a], means_100[a], means_1000[a], means_100000[a]
         3     if(k==1):
         4         print('n = 5')
         5     if(k==2):
         6         print('n = 10')
         7     if(k==3):
         8         print('n = 100')
         9     if(k==4):
        10         print('n = 1000')
        11     if(k==5):
        12         print('n = 100000')
        13     print(np.quantile(b, 0.7))
        14     k += 1
```

```
n = 5
3.0
n = 10
3.0
n = 100
2.0
n = 1000
2.0
n = 100000
2.0
```

```
In [44]: 1 #Сравнение
          2 np.log(1-0.7)/np.log(1-p)
Out[44]: 1.0
```

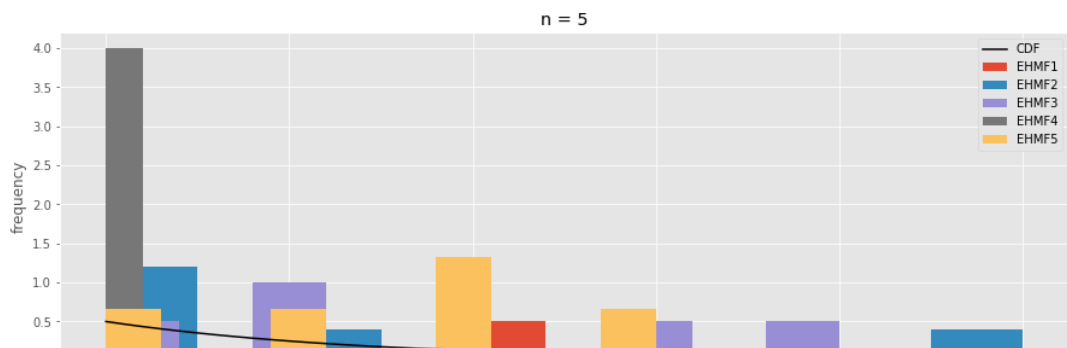
2.1.4 Построение гистограммы и полигона частот

In [47]:

```
1 #n=5
2 for a in range(5):
3     b=means_5[a]
4     b=sorted(b)
5     x=[]
6     y=[]
7     c=Counter(b)
8     for i in c:
9         x.append(i)
10        y.append(b.count(i)/5.0)
11    plt.plot(x,y,label="EPMF "+str(a+1))
12    plt.legend(loc='lower right')
13 plt.title("n = 5")
14 n=np.arange(1,2,0.1)#Построение
15 plt.plot(n,p*(1-p)**(n-1), 'k-',label='CDF')#функции вероятности
16 plt.legend()#распределения
17 plt.xlabel("numbers")
18 plt.ylabel("frequency")
19 plt.show()
20 #n=10
21 for a in range(5):
22     b=means_10[a]
23     b=sorted(b)
24     x=[]
25     y=[]
26     c=Counter(b)
27     for i in c:
28         x.append(i)
29         y.append(b.count(i)/10.0)
30     plt.plot(x,y,label="EPMF "+str(a+1))
31     plt.legend(loc='lower right')
32 plt.title("n = 10")
33 n=np.arange(1,6,0.1)#Построение
34 plt.plot(n,p*(1-p)**(n-1), 'k-',label='CDF')#функции вероятности
35 plt.legend()#распределения
36 plt.xlabel("numbers")
37 plt.ylabel("frequency")
38 plt.show()
39 #n=100
40 for a in range(5):
41     b=means_100[a]
42     b=sorted(b)
43     x=[]
44     y=[]
45     c=Counter(b)
46     for i in c:
47         x.append(i)
48         y.append(b.count(i)/100.0)
49     plt.plot(x,y,label="EPMF "+str(a+1))
50     plt.legend(loc='lower right')
51 plt.title("n = 100")
52 n=np.arange(1,9,0.1)#Построение
53 plt.plot(n,p*(1-p)**(n-1), 'k-',label='CDF')#функции вероятности
54 plt.legend()#распределения
55 plt.xlabel("numbers")
56 plt.ylabel("frequency")
57 plt.show()
58 #n=1000
59 for a in range(5):
60     b=means_1000[a]
61     b=sorted(b)
62     x=[]
63     y=[]
64     c=Counter(b)
65     for i in c:
66         x.append(i)
67         y.append(b.count(i)/1000.0)
68     plt.plot(x,y,label="EPMF "+str(a+1))
```

In [48]:

```
1 #n=5
2 for a in range(5):
3     plt.hist(means_5[a],density=True,label='EHMF{}'.format(a+1))
4     plt.legend()
5 n=np.arange(1,5,0.1)#Построение
6 plt.plot(n,p*(1-p)**(n-1),'k-',label='CDF')#функции вероятности
7 plt.legend()#распределения
8 plt.title("n = 5")
9 plt.xlabel("numbers")
10 plt.ylabel("frequency")
11 plt.show()
12 #n=10
13 for a in range(5):
14     plt.hist(means_10[a],density=True,label='EHMF{}'.format(a+1))
15     plt.legend()
16 n=np.arange(1,6,0.1)#Построение
17 plt.plot(n,p*(1-p)**(n-1),'k-',label='CDF')#функции вероятности
18 plt.legend()#распределения
19 plt.title("n = 10")
20 plt.xlabel("numbers")
21 plt.ylabel("frequency")
22 plt.show()
23 #n=100
24 for a in range(5):
25     plt.hist(means_100[a],density=True,label='EHMF{}'.format(a+1))
26     plt.legend()
27 n=np.arange(1,9,0.1)#Построение
28 plt.plot(n,p*(1-p)**(n-1),'k-',label='CDF')#функции вероятности
29 plt.legend()#распределения
30 plt.title("n = 100")
31 plt.xlabel("numbers")
32 plt.ylabel("frequency")
33 plt.show()
34 #n=1000
35 for a in range(5):
36     plt.hist(means_1000[a],density=True,label='EHMF{}'.format(a+1))
37     plt.legend()
38 n=np.arange(1,15,0.1)#Построение
39 plt.plot(n,p*(1-p)**(n-1),'k-',label='CDF')#функции вероятности
40 plt.legend()#распределения
41 plt.title("n = 1000")
42 plt.xlabel("numbers")
43 plt.ylabel("frequency")
44 plt.show()
45 #n=100000
46 for a in range(5):
47     plt.hist(means_100000[a],density=True,label='EHMF{}'.format(a+1))
48     plt.legend()
49 n=np.arange(1,18,0.1)#Построение
50 plt.plot(n,p*(1-p)**(n-1),'k-',label='CDF')#функции вероятности
51 plt.legend()#распределения
52 plt.title("n = 100000")
53 plt.xlabel("numbers")
54 plt.ylabel("frequency")
55 plt.show()
```



Если наблюдаемая в эксперименте случайная величина ξ дискретна и принимает значения a_1, a_2, \dots , то более наглядное представление о ее законе распределения дадут относительные частоты $v_r^* = \frac{v_r}{n}$, где v_r - число элементов выборки $X = (X_1, \dots, X_n)$, принявших значение a_r : $v_r = \sum_{j=1}^n I(X_j = a_r)$, $r = 1, 2, \dots$, т.е. v_r^* сближается с ростом n с теоретической вероятностью $P\{\xi = a_r\}$, и потому, по крайней мере для больших выборок, относительные частоты v_r^* можно рассматривать в качестве приближенных значений (оценок) для неизвестных вероятностей $P\{\xi = a_r\}$.

Наглядным представлением данных является полигон частот, который представляет собой ломаную с вершинами в точках $(a_r; v_r)$, $r = 1, 2, \dots$.

Можно рассматривать также статистический ряд $\{(a_r; v_r)\}$.

На графиках выше наглядно подтверждаются наши теоретические знания.

2.2 Распределение Максвелла

2.2.1 Моделирование выбранных случайных величин

```
In [42]: 1 # Создание случайной величины с распределением Максвелла, зависящим
2 # от параметра lambda
3 lambda=1.0
4 maxwell_rv=stats.maxwell(scale=lambda)
```

```
In [43]: 1 #Генерация выборки объема n = 5 с выводом
2 for n in[5]:
3     means_5=[]
4     for i in range(5):
5         sample=maxwell_rv.rvs(n)
6         means_5.append(sample)
7     print(sample)

[2.61492501 1.8185356 1.60239099 1.88471966 2.32343524]
[1.37090364 0.81484833 2.95500439 1.00735304 0.77608477]
[1.93052179 0.98916486 0.83490276 1.98262 2.28278868]
[2.21407855 2.51775003 2.5953651 1.29263856 2.0909891 ]
[3.23971372 1.82680786 2.17985755 0.31893926 2.1808736 ]
```

```
In [44]: 1 #Генерация выборки объема n = 10 с выводом
2 for n in[10]:
3     means_10=[]
4     for i in range(5):
5         sample=maxwell_rv.rvs(n)
6         means_10.append(sample)
7     print(sample)

[2.12102194 2.21942606 2.75273494 0.82323637 1.82087582 1.86830689
1.07562858 0.47693664 2.17478789 1.433096 ]
[1.66711706 2.12778881 1.59876399 2.19313351 1.76273771 1.05264685
2.74675272 1.55595092 1.03700015 0.59029844]
[3.33341564 4.23787351 1.52515002 1.63675074 0.71231781 2.58619079
1.33683265 1.02519597 1.97457463 2.09229573]
[1.27631799 2.25908834 1.69106005 2.57793336 0.90366426 1.35694638
1.87914425 0.63669354 1.09633501 1.0690649 ]
[2.14434436 2.28425493 0.67972466 2.71354297 1.8494517 1.18174504
1.4799613 0.52483374 1.61124112 1.55908243]
```

```
In [45]: 1 #Генерация выборки объема n = 100 без вывода
2 for n in[100]:
3     means_100=[]
4     for i in range(5):
5         sample=maxwell_rv.rvs(n)
6         means_100.append(sample)
```

```
In [46]: 1 #Генерация выборки объема n = 1000 без вывода
2 for n in [1000]:
3     means_1000=[]
4     for i in range(5):
5         sample=maxwell_rv.rvs(n)
6         means_1000.append(sample)
```

```
In [47]: 1 #Генерация выборки объема n = 100000 без вывода
2 for n in [100000]:
3     means_100000=[]
4     for i in range(5):
5         sample=maxwell_rv.rvs(n)
6         means_100000.append(sample)
```

```
In [55]: 1 #Вернёмся к медиане и убедимся, что в пункте 1.2.1 она была найдена верно
2 maxwell.median()
```

```
Out[55]: 1.5381722544550522
```

2.2.2 Построение эмпирической функции распределения

In [56]:

```
1 #n=5
2 for a in range(5):
3     b=means__5[a]
4     b=sorted(b)
5     v=len(b)
6     N=[]
7     for i in range((v-1)):
8         N.append(b.count(b[i]))
9         x=[]
10        y=[]
11        t=0
12        x.append(0.0)
13        y.append(0.0)
14        x.append(b[0])
15        y.append(0.0)
16        for i in range((v-1)):
17            t+=N[i]
18            x.append(b[i])
19            y.append(float(t/v))
20            x.append(b[i+1])
21            y.append(float(t/v))
22        x.append(b[v-1])
23        y.append(1)
24        x.append(b[v-1]+2)
25        y.append(1)
26        plt.plot(x,y,label="ECDF "+str(a+1))
27        plt.legend(loc='lower right')
28 plt.title("Эмпирическая функция выборки объема: "+str(v))
29 x=np.linspace(0,5,100)#Построение
30 cdf=maxwell_rv.cdf(x)#теоретической функции
31 plt.plot(x,cdf,label='CDF')#распределения
32 plt.legend()
33 plt.xlabel("numbers")
34 plt.ylabel("probability")
35 plt.show()
36 print("\n")
37 #n=10
38 for a in range(5):
39     b = means__10[a]
40     b = sorted(b)
41     v = len(b)
42     N = []
43     for i in range((v-1)):
44         N.append(b.count(b[i]))
45         x=[]
46         y=[]
47         t=0
48         x.append(0.0)
49         y.append(0.0)
50         x.append(b[0])
51         y.append(0.0)
52         for i in range((v-1)):
53             t+=N[i]
54             x.append(b[i])
55             y.append(float(t/v))
56             x.append(b[i+1])
57             y.append(float(t/v))
58         x.append(b[v-1])
59         y.append(1)
60         x.append(b[v-1]+2)
61         y.append(1)
62         plt.plot(x,y,label="ECDF "+str(a+1))
63         plt.legend(loc='lower right')
64 plt.title("Эмпирическая функция выборки объема: "+str(v))
65 x=np.linspace(0,5,100)#Построение
66 cdf=maxwell_rv.cdf(x)#теоретической функции
67 plt.plot(x,cdf,label='CDF')#распределения
68 plt.legend()
```

2.2.3 Построение вариационного ряда выборки

```
In [39]: 1 #Вариационный ряд для выборки объема n=5 с выводом
2 for a in range(5):
3     b=means__5[a]
4     b=sorted(b)
5     print(b)
```

[0.9048224033555651, 1.0595749548872997, 1.3361146201247887, 1.7684778211524081, 2.064309990285417]
[1.142775449379784, 1.6322198084743014, 1.9269294058020057, 1.9307312276956405, 2.137402279174372]
[1.168850546533639, 1.2873609279696019, 1.4055755577674152, 1.554026224592775, 2.9493391029304523]
[0.981057908261026, 1.200322004547301, 1.8319652264631978, 1.9311087818516284, 2.2702664778364174]
[1.0795029380722423, 1.103674980358039, 1.3308824708067275, 2.6687299642723477, 2.9675918656118716]

```
In [40]: 1 #Вариационный ряд для выборки объема n=10 с выводом
2 for a in range(5):
3     b=means__10[a]
4     b=sorted(b)
5     print(b)
```

[0.4344229664385466, 0.9873280481429386, 1.1514899800296186, 1.1620788799472992, 1.5047110956567034, 1.6532896759587332, 1.9249137266578096, 2.0658538511219224, 2.1461320252051195, 2.324740490202557]
[1.1220497289203735, 1.2401311800502641, 1.2435276204947807, 1.576561337080865, 1.6375956762713444, 1.7697569758530114, 1.8386307894229095, 1.9485280422518474, 1.950915639549961, 1.9767833922678453]
[0.8963933247210342, 1.4705601230263092, 1.4945015272392095, 1.4974428220032772, 1.5122209047811999, 1.6248884551477416, 1.7155281056459615, 1.9907561312818842, 2.3686125510412857, 2.4716462283848437]
[0.39200688435666003, 0.5966409826346033, 0.9698060007821667, 1.0285075860763897, 1.0404367899757856, 1.486851983947381, 2.019114126722199, 2.1664177622078764, 2.641475897208004, 2.84064283839896]
[1.148593485170174, 1.1736791688497565, 1.5230401251658805, 1.553329580159769, 1.5775412015483297, 1.781091079341479, 1.941831168996058, 2.0709499700182676, 2.143143746322394, 2.9585921488894913]

```
In [41]: 1 #Вариационный ряд для выборки объема n=100 без вывода
2 for a in range(5):
3     b=means__100[a]
4     h=sorted(h)
```

```
In [60]: 1 #Вариационный ряд для выборки объема n=1000 без вывода
2 for a in range(5):
3     b=means__1000[a]
4     h=sorted(h)
```

```
In [61]: 1 #Вариационный ряд для выборки объема n=100000 без вывода
2 for a in range(5):
3     b=means__100000[a]
4     h=sorted(h)
```

Возникли сложности при вычислении теоретических значений квантилей, однако был найден справочник: "Справочник по вероятностным распределениям" Р.Н.Вадзинский. В нём была найдена таблица для приближенного решения уравнения $x_\alpha = \lambda m_\alpha$, где $x_\alpha = \lambda m_\alpha$ - квантиль порядка α распределения Максвелла

```
In [62]: 1 k=1
2 for b in [means__5[a],means__10[a],means__100[a],means__1000[a],means__10000[a],means__100000[a]]:
3     if(k==1):
4         print('n = 5')
5     if(k==2):
6         print('n = 10')
7     if(k==3):
8         print('n = 100')
9     if(k==4):
10        print('n = 1000')
11    if(k==5):
12        print('n = 100000')
13    print(np.quantile(b,0.1))
14    k+=1
```

```
n = 5
1.2719859381522172
n = 10
1.1252696205430428
n = 100
0.7938084172369115
n = 1000
0.7374170981472817
n = 100000
0.7603097895100107
```

Сравнение со значением (с теоретическим) из таблицы:
 $\alpha \approx 0.76$

```
In [63]: 1 k=1
2 for b in [means__5[a],means__10[a],means__100[a],means__1000[a],means__10000[a],means__100000[a]]:
3     if(k==1):
4         print('n = 5')
5     if(k==2):
6         print('n = 10')
7     if(k==3):
8         print('n = 100')
9     if(k==4):
10        print('n = 1000')
11    if(k==5):
12        print('n = 100000')
13    print(np.quantile(b,0.5))
14    k+=1
```

```
n = 5
1.5059353635900166
n = 10
1.4534420207082235
n = 100
1.574535047585379
n = 1000
1.4765126474373647
n = 100000
1.537607751024725
```

```
In [64]: 1 #Сравнение
2 maxwell.median()
```

```
Out[64]: 1.5381722544550522
```

```
In [65]: 1 k=1
2 for b in [means__5[a],means__10[a],means__100[a],means__1000[a],means__10000[a],means__100000[a]]:
3     if(k==1):
4         print('n = 5')
5     if(k==2):
6         print('n = 10')
7     if(k==3):
8         print('n = 100')
9     if(k==4):
10        print('n = 1000')
11    if(k==5):
12        print('n = 10000')
13    print(np.quantile(b,0.7))
14    k+=1
```

```
n = 5
1.649424804466466
n = 10
1.7200020124901898
n = 100
1.9686707573406326
n = 1000
1.8574082725897791
n = 10000
1.9163671862147416
```

Сравнение со значением (с теоретическим) из таблицы:
 $\alpha \approx 1.92$

С увеличением объема выборки э.ф.р. стремится к теоритической функции распределения, следовательно, эмпирические квантили так жestreмятся к теоритическим по определению. Что и можно наблюдать выше.

2.2.4 Построение гистограммы и полигона частот

In [66]:

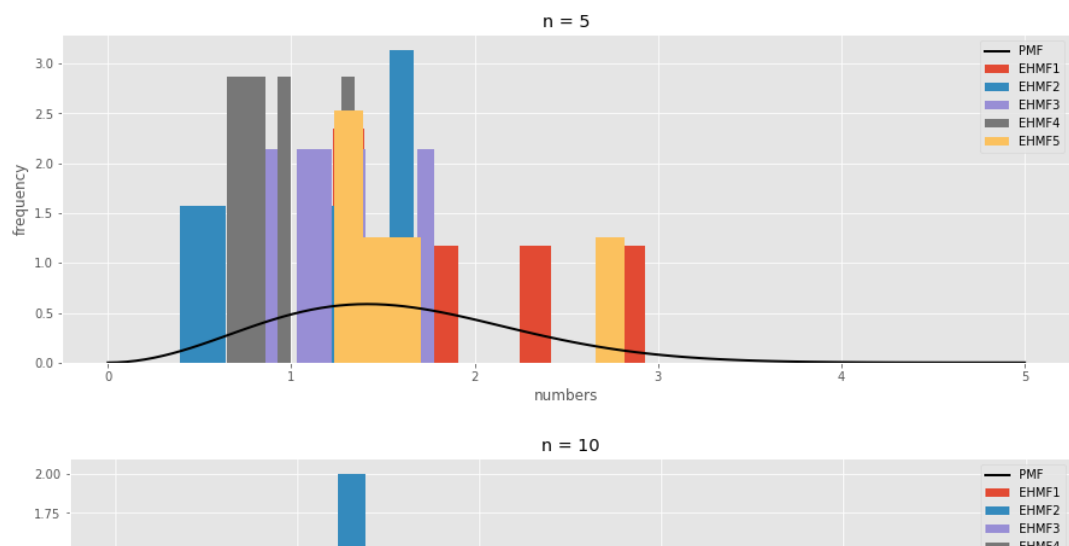
```
1 #n=5
2 for a in range(5):
3     b=means__5[a]
4     mas=list(range(1,6))
5     p=[0,0,0,0,0]
6     for i in range(5):
7         mas[i]=b[i]
8         if mas[i]>0 and mas[i]<1:
9             p[0]=p[0]+1
10        if mas[i]>1 and mas[i]<2:
11            p[1]=p[1]+1
12        if mas[i]>2 and mas[i]<3:
13            p[2]=p[2]+1
14        if mas[i]>3 and mas[i]<4:
15            p[3]=p[3]+1
16        if mas[i]>4 and mas[i]<5:
17            p[4]=p[4]+1
18    print()
19    dob=[]
20    bod=[]
21    keks=0.5
22    for i in range(5):
23        dob.append(keks)
24        bod.append(p[i]/5.0)
25    keks+=1
26    plt.plot(dob,bod,label='EPMF'+str(a+1))
27    rv=maxwell()
28    x=np.linspace(0,5,100)
29    plt.plot(x,rv.pdf(x),'k-',lw=2,label='PMF')
30    plt.legend()
31    plt.title("n = 5")
32    plt.xlabel("numbers")
33    plt.ylabel("frequency")
34    plt.show()
35    #n=10
36    for a in range(5):
37        b=means__10[a]
38        mas=list(range(1,11))
39        p=[0,0,0,0,0]
40        for i in range(10):
41            mas[i]=b[i]
42            if mas[i]>0 and mas[i]<1:
43                p[0]=p[0]+1
44            if mas[i]>1 and mas[i]<2:
45                p[1]=p[1]+1
46            if mas[i]>2 and mas[i]<3:
47                p[2]=p[2]+1
48            if mas[i]>3 and mas[i]<4:
49                p[3]=p[3]+1
50            if mas[i]>4 and mas[i]<5:
51                p[4]=p[4]+1
52        print()
53        dob=[]
54        bod=[]
55        keks=0.5
56        for i in range(5):
57            dob.append(keks)
58            bod.append(p[i]/10.0)
59        keks+=1
60        plt.plot(dob,bod,label='EPMF'+str(a+1))
61    rv=maxwell()
62    x=np.linspace(0,5,100)
63    plt.plot(x,rv.pdf(x),'k-',lw=2,label='PMF')
64    plt.legend()
65    plt.title("n = 10")
66    plt.xlabel("numbers")
67    plt.ylabel("frequency")
68    plt.show()
```

In [67]:

```

1  #n=5
2  for a in range(5):
3      plt.hist(means__5[a],density=True,label='EHMF{}'.format(a+1))
4      plt.legend()
5  rv=maxwell()
6  x=np.linspace(0,5,100)
7  plt.plot(x,rv.pdf(x),'k-',lw=2,label='PMF')
8  plt.legend()
9  plt.title("n = 5")
10 plt.xlabel("numbers")
11 plt.ylabel("frequency")
12 plt.show()
13 #n=10
14 for a in range(5):
15     plt.hist(means__10[a],density=True,label='EHMF{}'.format(a+1))
16     plt.legend()
17 rv=maxwell()
18 x=np.linspace(0,5,100)
19 plt.plot(x,rv.pdf(x),'k-',lw=2,label='PMF')
20 plt.legend()
21 plt.title("n = 10")
22 plt.xlabel("numbers")
23 plt.ylabel("frequency")
24 plt.show()
25 #n=100
26 for a in range(5):
27     plt.hist(means__100[a],density=True,label='EHMF{}'.format(a+1))
28     plt.legend()
29 rv=maxwell()
30 x=np.linspace(0,5,100)
31 plt.plot(x,rv.pdf(x),'k-',lw=2,label='PMF')
32 plt.legend()
33 plt.title("n = 100")
34 plt.xlabel("numbers")
35 plt.ylabel("frequency")
36 plt.show()
37 #n=1000
38 for a in range(5):
39     plt.hist(means__1000[a],density=True,label='EHMF{}'.format(a+1))
40     plt.legend()
41 rv=maxwell()
42 x=np.linspace(0,5,100)
43 plt.plot(x,rv.pdf(x),'k-',lw=2,label='PMF')
44 plt.legend()
45 plt.title("n = 1000")
46 plt.xlabel("numbers")
47 plt.ylabel("frequency")
48 plt.show()

```



Для непрерывной случайной величины ξ , обладающей непрерывной плотностью $f(x)$, также можно построить по соответствующей выборке $X = (X_1, \dots, X_n)$ статистический аналог $\hat{f}_n(x)$ для плотности $f(x)$, который называется гистограммой. Для этого используется метод группировки, в соответствии с которым область Δ возможных значений ξ разбивается на некоторое число N непересекающихся интервалов $\Delta_1, \dots, \Delta_N$ (так что $\Delta = \bigcup_{r=1}^N \Delta_r$, подсчитывают числа v_1, \dots, v_N наблюдений X_1, \dots, X_n , попавших в соответствующие интервалы: $v_r = \sum_{j=1}^n I(X_j \in \Delta_r), r = 1, \dots, N$ (так что $\sum_{r=1}^N v_r = n$, и строят кусочно-постоянную функцию

$$\hat{f}_n(x) = \frac{v_r}{n|\Delta_r|}$$

при $x \in \Delta_r, r = 1, \dots, N$

Здесь $|\Delta_r|$ - длина интервала Δ_r . То, что построенная по такому правилу гистограмма $\hat{f}_n(x)$ действительно "похожа" на теоретическую плотность $f(x)$, следует из закона больших чисел, согласно которому при $n \rightarrow \infty$ относительная частота $\frac{v_r}{n}$ сближается с теоретической вероятностью

$$P\{\xi \in \Delta_r\} = \int_{\Delta_r} f(x)dx$$

Но этот интеграл по теореме о среднем равен $f(a_r)|\Delta_r|$ где a_r - некоторая внутренняя точка интервала Δ_r (при малом Δ_r в качестве a_r можно взять, например, середину интервала). Таким образом, при больших n и достаточно "мелком" разбиении $\{\Delta_r\}$ $\hat{f}_n(x) \approx f(a_r)$ при $x \in \Delta_r$ т.е. гистограмма $\hat{f}_n(x)$ будет достаточно хорошо приближать график плотности $f(x)$, следовательно, $\hat{f}_n(x)$ можно рассматривать в качестве статистического аналога (оценки) для $f(x)$. Наряду с гистограммой, в качестве приближения для неизвестной теоретической плотности $f(x)$ можно использовать кусочно-линейный график называемый полигоном частот. Он также считается статистическим аналогом теоретической плотности. Данные на полигоне частот и гистограммах подтверждают теоретические знания: с увеличением объема выборки полигон частот и гистограммы практически совпадают с теоретической плотностью $f(x)$.

3 Домашнее задание. Оценки

3.1. Геометрическое распределение

3.1.1 Нахождение выборочного среднего и выборочной дисперсии геометрического распределения

Наиболее важными характеристиками случайной величины ξ являются ее моменты $\alpha_k = M\xi^k$, а также центральные моменты $\mu^k = M(\xi - \alpha_1)^k$ (когда они существуют). Их статистическими аналогами, вычисляемыми по соответствующей выборке $X = (X_1, \dots, X_n)$, являются выборочные моменты соответственно обычные:

$$\hat{\alpha}_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

и центральные:

$$\hat{\mu}_k = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\alpha}_1)^k$$

$\hat{\alpha}_1$ (принято обозначать, как \bar{X}) называют выборочным средним, μ^2 - выборочной дисперсией. Таким образом, выборочное среднее и выборочная дисперсия являются статистическими аналогами теоретических среднего (математического ожидания) $M\xi$ и дисперсии $D\xi$, когда они существуют.

Выборочное среднее, относящийся к выборке X , считается как:

$$\bar{X} = \hat{\alpha}_1 = \frac{1}{n} \sum_{i=1}^n X_i$$

Выборочная дисперсия, относящийся к выборке X , подсчитывается как:

$$S^2 = \hat{\mu}_2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

Найдем математическое ожидание и дисперсию выборочного среднего и выборочной дисперсии:

$$M\bar{X} = \frac{1}{n} \sum_{i=1}^n MX_i = M\xi = \alpha_1$$

$$D\bar{X} = \frac{1}{n^2} \sum_{i=1}^n DX_i = \frac{1}{n} D\xi = \frac{\mu_2}{n}$$

Для выборочной дисперсии введем обозначение: $Y_i = X_i - \alpha_1$:

$$S^2 = \hat{\mu}_2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}^2$$

Поскольку $MY_i = 0$, $MY_i^2 = \mu_2$ и $MY_iY_j = MY_jY_i = 0$, ($i \neq j$), то:

$$M\bar{Y}^2 = \frac{1}{n^2} \sum_{i,j=1}^n MY_iY_j = \frac{1}{n^2} \sum_{i=1}^n MY_i^2 = \frac{\mu_2}{n}$$

Отсюда следует, что

$$MS^2 = \frac{n-1}{n} \mu_2$$

Перейдём к вычислению DS^2

$$(S^2)^2 = \frac{1}{n^2} \left(\sum_{i=1}^n Y_i^2 \right)^2 - \frac{2}{n} \bar{Y}^2 \sum_{i=1}^n Y_i^2 + \bar{Y}^4$$

Так как случайные величины Y_1, \dots, Y_n независимы и $MY_i = 0$, то в правой части равенства

$$M\bar{Y}^4 = \frac{1}{n^4} (n\mu_4 + 3n(n-1)\mu_2^2) = \frac{\mu_4 + 3(n-1)\mu_2^2}{n^3}$$

Аналогично находим

$$\frac{1}{n^2} M \left(\sum_{i=1}^n Y_i^2 \right) = \frac{\mu_4 + (n-1)\mu_2^2}{n} = M(\bar{Y}^2 \sum_{i=1}^n Y_i^2)$$

С учётом этих соотношений по формуле

$$DS^2 = M(S^2)^2 - (MS^2)^2$$

получим

$$DS^2 = \frac{\mu_4 - \mu_2^2}{n} - \frac{2(\mu_4 - 2\mu_2^2)}{n^2} + \frac{\mu_4 - 3\mu_2^2}{n^3} = \frac{(n-1)^2}{n^3} \left(\mu_4 - \frac{n-3}{n-1} \mu_2^2 \right)$$

Аналогично можно находить моменты и более высоких порядков, хотя с увеличением порядка вид формул и их вывод усложняются.

Теперь рассмотрим свойства выборочных среднего и дисперсии при неограниченном возрастании объема выборки n , которые дадут нам ответ на вопрос, оценками каких параметров распределений они являются. Чтобы подчеркнуть зависимость моментов $\hat{\alpha}_k, \hat{\mu}_k$ от объема выборки, будем в дальнейшем приписывать дополнительный индекс n : $\hat{\alpha}_{nk}, \hat{\mu}_{nk}$

$$M\hat{\alpha}_{nk} = \frac{1}{n} \sum_{i=1}^n M X_i^k = M\xi^k = \hat{\alpha}_k$$

$$D\hat{\alpha}_{nk} = \frac{1}{n^2} \sum_{i=1}^n D X_i^k = \frac{1}{n} D\xi^k = \frac{1}{n} (M\xi^{2k} - (M\xi^k)^2) = \frac{\alpha_{2k} - \alpha_k^2}{n}$$

На основании неравенства Чебышева, отсюда следует, что для любого $\epsilon > 0$ при $n \rightarrow \infty$

$$P|\hat{\alpha}_{nk} - \alpha_k| < \epsilon \rightarrow 1$$

т.е. выборочный момент $\hat{\alpha}_{nk}$ сходится по вероятности при $n \rightarrow \infty$ к соответствующему теоретическому моменту α_k . Таким образом, $\hat{\alpha}_{nk}$ можно использовать в качестве оценки α_k , когда объем выборки достаточно велик. Аналогичное утверждение справедливо и для центральных моментов:

$$P|\mu_{nk} - \mu_k| < \epsilon \rightarrow 1$$

т.е. μ_{nk} можно использовать в качестве оценки μ_k , когда объем выборки достаточно велик.

1. Оценка $\hat{\theta}(X)$ параметра θ называется несмещенной, если:

$$E(\hat{\theta}(X)) = \theta$$

2. Оценка $\hat{\theta}(X) = \hat{\theta}_n(X_1, \dots, X_n)$ параметра θ называется состоятельной, если при $n \rightarrow \infty$ соблюдается:

$$\hat{\theta}_n(X_1, \dots, X_n) \xrightarrow{p} \theta$$

При этом для проверки состоятельности достаточно убедиться, что соблюдены следующие два условия:

$$\lim_{n \rightarrow \infty} E(\hat{\theta}_n(X_1, \dots, X_n)) = \theta$$

$$\lim_{n \rightarrow \infty} Var(\hat{\theta}_n(X_1, \dots, X_n)) = 0$$

Выборочное среднее является несмещенной оценкой для теоретического математического ожидания.

$$\hat{\alpha}_1 = \frac{1}{n} \sum_{i=1}^n X_i$$

$$M\hat{\alpha}_1 = M\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n} \cdot n \cdot MX = MX$$

$$M\hat{\alpha}_1 = MX$$

Выборочное среднее является состоятельной оценкой для теоретического математического ожидания.

$$\hat{\alpha}_1 = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\lim_{n \rightarrow \infty} \hat{\alpha}_1 = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \lim_{n \rightarrow \infty} \frac{1}{n} \cdot n \cdot MX = MX$$

Выборочная дисперсия S^2 является состоятельной и несмещенной оценкой для теоретической дисперсии.

$$S^2 = \hat{\mu}_2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \frac{1}{n} \sum_{i=1}^n (\hat{\alpha}_1)^2 = MX^2 - \frac{1}{n} n (\hat{\alpha}_1)^2 = \hat{\alpha}_2 - (\hat{\alpha}_1)^2$$

$$MS^2 = M(\hat{\alpha}_2 - (\hat{\alpha}_1)^2) = M\left(\frac{1}{n} \sum_{i=1}^n X_i^2\right) - M\left(\left(\frac{1}{n} \sum_{i=1}^n X_i\right)^2\right) = MX^2 - (MX)^2 = DX$$

Продолжим исследование свойств выборочных моментов для больших выборок и рассмотрим теперь асимптотическое поведение их выборочных распределений.

Если распределение случайной величины v_n сходится при $n \rightarrow \infty$ к распределению случайной величины v и при этом $\zeta(v) = N(\mu, \sigma^2)$, то будем писать $\zeta(v_n) \rightarrow N(\mu, \sigma^2)$. Далее иногда будем говорить, что случайная величина v_n асимптотически нормальна $N(\mu_n, \sigma_n^2)$, и записывать это следующим образом:

$$\zeta(v_n) \approx N(\mu_n, \sigma_n^2), \text{ если } \zeta\left(\frac{v_n - \mu_n}{\sigma_n}\right) \rightarrow N(0, 1).$$

Найдем сначала асимптотические распределения выборочных моментов $\hat{\alpha}_{nk}$. Величина $n\hat{\alpha}_{nk} = \sum_{i=1}^n X_i^k$ является суммой независимых одинаково распределенных случайных величин. Если конечен момент $\alpha_{2k} = M\xi^{2k}$, то к этой сумме можно применить центральную предельную теорему теории вероятностей. Так как $MX_i^k = \alpha^k$, $DX_i^k = \alpha_{2k} - \alpha_k^2$, то величина

$$\frac{n\hat{\alpha}_{nk} - n\alpha_k}{\sqrt{n(\alpha_{2k} - \alpha_k^2)}} = \frac{\hat{\alpha}_{nk} - \alpha_k}{\sqrt{\frac{\alpha_{2k} - \alpha_k^2}{n}}}$$

асимптотически нормальна $N(0, 1)$. Таким образом справедлива следующая теорема:

Если конечен теоретический момент α_{2k} , то при $n \rightarrow \infty$ выборочный момент $\hat{\alpha}_{nk}$ асимптотически нормален $N(\alpha_k, \frac{\alpha_{2k} - \alpha_k^2}{n})$

Из теоремы следует, что если существует теоретическая дисперсия, то выборочное среднее $\hat{\alpha}_{n1}$ асимптотически нормально $N(\alpha_1, \frac{\mu_2}{n})$

Из теоремы об асимптотической нормальности функций от выборочных моментов следует, что асимптотически нормальными являются и центральные выборочные моменты $\hat{\mu}_{nk}$, поскольку они являются непрерывными функциями (многочленами) от обычных выборочных моментов.

Выборочное среднее

In [68]:

```
1 #n=5
2 vs_5 = []
3 for i in range(5):
4     vs_5.append(np.mean(means_5[i]))
5 print(vs_5)
6 #n=10
7 vs_10 = []
8 for i in range(5):
9     vs_10.append(np.mean(means_10[i]))
10 print(vs_10)
11 #n=100
12 vs_100 = []
13 for i in range(5):
14     vs_100.append(np.mean(means_100[i]))
15 print(vs_100)
16 #n=1000
17 vs_1000 = []
18 for i in range(5):
19     vs_1000.append(np.mean(means_1000[i]))
20 print(vs_1000)
21 #n=100000
22 vs_100000 = []
23 for i in range(5):
24     vs_100000.append(np.mean(means_100000[i]))
25 print(vs_100000)
```

[2.4, 2.2, 2.8, 1.4, 2.6]
 [2.2, 2.4, 2.0, 1.7, 2.7]
 [2.02, 2.02, 1.91, 2.17, 1.98]
 [1.978, 2.011, 2.058, 1.959, 1.954]
 [2.0093, 2.0101, 2.01032, 1.99982, 2.00051]

$$M\xi = \frac{1}{p}, \text{ при } p = 0.5 \quad M\xi = 2$$

```
In [69]: 1 #Сравнение
          2 p = 0.5
          3 neom.stats(n.moments = 'm')
```

Out[69]: array(2.)

Выборочная дисперсия

```
In [70]: 1 #n=5
          2 vd_5 = []
          3 for i in range(5):
          4     vd_5.append(round(np.var(means_5[i]),6))
          5 print(vd_5)
          6 #n=10
          7 vd_10 = []
          8 for i in range(5):
          9     vd_10.append(round(np.var(means_10[i]),6))
          10 print(vd_10)
          11 #n=100
          12 vd_100 = []
          13 for i in range(5):
          14     vd_100.append(round(np.var(means_100[i]),6))
          15 print(vd_100)
          16 #n=1000
          17 vd_1000 = []
          18 for i in range(5):
          19     vd_1000.append(round(np.var(means_1000[i]),6))
          20 print(vd_1000)
          21 #n=100000
          22 vd_100000 = []
          23 for i in range(5):
          24     vd_100000.append(round(np.var(means_100000[i]),6))
          25 print(vd_100000)
```

```
[2.24, 3.76, 2.16, 0.64, 1.04]
[2.36, 1.04, 0.8, 0.61, 2.41]
[2.5196, 2.2396, 1.6819, 2.3811, 2.8596]
[1.955516, 2.090879, 2.188636, 1.845319, 1.699884]
[2.033354, 2.020038, 2.033233, 2.02872, 1.99785]
```

$$D\xi = \frac{q}{p^2} = \frac{1-p}{p^2}$$

При $p = 0.5$, $D\xi = 2$

```
In [71]: 1 #Сравнение
          2 neom.stats(n.moments = 'v')
```

Out[71]: array(2.)

Как видно из полученных значений, чем больше объём выборки, тем менее отличаются выборочное среднее от теоретического математического ожидания и выборочная дисперсия от теоретической дисперсии

3.1.2 Построение доверительного интервала для выборочного среднего

Определение: γ - доверительным интервалом для g называется такой случайный интервал $(T_1(X), T_2(X))$, $T_1(X) < T_2(X)$, который содержит внутри себя (накрывает) неизвестное значение g с вероятностью, не меньшей γ :

$$P\{T_1(X) < g < T_2(X)\} \geq \gamma$$

Здесь $T_1(X)$ и $T_2(X)$ - некоторые статистики (функции от выборки), называемые соответственно нижней и верхней доверительными границами, а γ - задаваемый заранее доверительный уровень, который обычно выбирается близким к 1. Длина доверительного интервала характеризует точность локализации оцениваемой характеристики g , а величина γ является показателем надежности доверительного интервала.

В сформулированной ранее теореме [Если конечен теоретический момент α_{2k} , то при $n \rightarrow \infty$ выборочный момент $\hat{\alpha}_{nk}$ асимптотически нормален $N(\alpha_k, \frac{\alpha_{2k} - \alpha_k^2}{n})$] можно заменить асимптотическую дисперсию $\frac{\alpha_{2k} - \alpha_k^2}{n}$ ее оценкой $\frac{\hat{\alpha}_{n,2k} - \hat{\alpha}_{nk}^2}{n}$. Это дает искомый асимптотический γ -доверительный интервал для момента α_k вида:

$$(\hat{\alpha}_{nk} \mp c_\gamma \sqrt{\frac{\hat{\alpha}_{n,2k} - \hat{\alpha}_{nk}^2}{n}})$$

Полагая здесь $k = 1$, получим соответствующий интервал для теоретического среднего $\alpha_1 = M\xi$:

$$(\bar{X} \mp \frac{c_\gamma S}{\sqrt{n}})$$

Чтобы построить асимптотический γ -доверительный интервал для теоретической дисперсии $\mu_2 = D\xi$, надо просто воспользоваться результатом теоремы об асимптотической нормальности выборочной дисперсии $[\zeta(\frac{\sqrt{n}(S^2 - \mu_2)}{\sqrt{\hat{\mu}_{n4} - S^4}}) \rightarrow N(0, 1)]$: искомый интервал есть

$$(S^2 \mp c_\gamma \sqrt{\frac{\hat{\mu}_{n4} - S^4}{n}})$$

Положим $\gamma = 0.95$ и найдем доверительный интервал для выборочного среднего.

$$\Phi(c_\gamma) = \frac{\gamma}{2} = 0.475$$

Из таблицы значений функции Лапласа $c_\gamma \approx 1.96$

In [77]:

```
1 #n=5
2 print('n = 5')
3 for i in range(5):
4     print('(', vs_5[i], '-+ 1.96 *', np.sqrt(vd_5[i]/5), ') = (', vs_5[i],
5 #n=10
6 print('n = 10')
7 for i in range(5):
8     print('(', vs_10[i], '-+ 1.96 *', np.sqrt(vd_10[i]/10), ') = (', vs_10[i],
9 #n=100
10 print('n = 100')
11 for i in range(5):
12     print('(', vs_100[i], '-+ 1.96 *', np.sqrt(vd_100[i]/100), ') = (', vs_
13 #n=1000
14 print('n = 1000')
15 for i in range(5):
16     print('(', vs_1000[i], '-+ 1.96 *', np.sqrt(vd_1000[i]/1000), ') = (',
17 #n=100000
18 print('n = 100000')
19 for i in range(5):
20     print('(', vs_100000[i], '-+ 1.96 *', np.sqrt(vd_100000[i]/100000), ')
```

```
n = 5
( 1.8 -+ 1.96 * 0.33466401061363027 ) = ( 1.8 -+ 0.655941 )
( 1.6 -+ 1.96 * 0.35777087639996635 ) = ( 1.6 -+ 0.701231 )
( 2.6 -+ 1.96 * 0.8294576541331088 ) = ( 2.6 -+ 1.625737 )
( 1.4 -+ 1.96 * 0.35777087639996635 ) = ( 1.4 -+ 0.701231 )
( 2.4 -+ 1.96 * 0.4560701700396552 ) = ( 2.4 -+ 0.893898 )
n = 10
( 2.0 -+ 1.96 * 0.4 ) = ( 2.0 -+ 0.784 )
( 3.2 -+ 1.96 * 0.7720103626247513 ) = ( 3.2 -+ 1.51314 )
( 2.1 -+ 1.96 * 0.7543208866258444 ) = ( 2.1 -+ 1.478469 )
( 1.6 -+ 1.96 * 0.322490309931942 ) = ( 1.6 -+ 0.632081 )
( 1.5 -+ 1.96 * 0.291547594742265 ) = ( 1.5 -+ 0.571433 )
n = 100
( 1.97 -+ 1.96 * 0.13817018491700733 ) = ( 1.97 -+ 0.270814 )
( 1.72 -+ 1.96 * 0.10007996802557444 ) = ( 1.72 -+ 0.196157 )
( 2.45 -+ 1.96 * 0.17342145196024625 ) = ( 2.45 -+ 0.339906 )
( 1.94 -+ 1.96 * 0.1412940196894405 ) = ( 1.94 -+ 0.276936 )
( 1.94 -+ 1.96 * 0.13024592124131948 ) = ( 1.94 -+ 0.255282 )
n = 1000
( 2.027 -+ 1.96 * 0.04624144245154989 ) = ( 2.027 -+ 0.090633 )
( 1.943 -+ 1.96 * 0.04330994112210267 ) = ( 1.943 -+ 0.084887 )
( 1.998 -+ 1.96 * 0.04409077000915271 ) = ( 1.998 -+ 0.086418 )
( 2.001 -+ 1.96 * 0.04506660626228694 ) = ( 2.001 -+ 0.088331 )
( 1.987 -+ 1.96 * 0.04618258329716951 ) = ( 1.987 -+ 0.090518 )
n = 100000
( 1.99964 -+ 1.96 * 0.004461636471072021 ) = ( 1.99964 -+ 0.008745 )
( 1.99584 -+ 1.96 * 0.004471759161672283 ) = ( 1.99584 -+ 0.008765 )
( 2.00301 -+ 1.96 * 0.004486001560409893 ) = ( 2.00301 -+ 0.008793 )
( 1.99926 -+ 1.96 * 0.004447649941261115 ) = ( 1.99926 -+ 0.008717 )
( 2.0108 -+ 1.96 * 0.004515952834120392 ) = ( 2.0108 -+ 0.008851 )
```

*Округлено до 6 знаков после запятой.

3.1.3 Нахождение оптимальности рассматриваемых оценок

Для построения теории оптимального оценивания прежде всего надо договориться о мере точности оценок, т.е. уточнить смысл приближенного равенства $T(X) \approx g$. Если статистика $T(x)$ используется для оценивания g , то одной из разумных мер расхождения между ними является $(T(X) - g)^2$, или квадратичная ошибка. Но так как это величина случайная используется среднеквадратичная ошибка (с. к. о.) $\Delta(T) = M(T(X) - g)^2$.

Определение: Оценка минимизирующая с. к. о. в данном классе оценок T_g называется оптимальной в среднеквадратичном смысле и обозначается T^* :

$$T^* = \operatorname{argmin}_{T \in T_g} \Delta(T)$$

Пусть требуется оценить заданную параметрическую функцию $\tau(\theta)$ в модели $F = F(x; \theta)$, $\theta \in \Theta$ по соответствующей выборке $X = (X_1, \dots, X_n)$. Обозначим τ_τ класс всех несмещенных оценок $T = T(X)$ для $\tau(\theta)$ и предположим, что он не пуст. Дополнительно предположим, что дисперсии всех оценок из класса τ_τ конечны: $D\theta T = M_\theta(T - \tau(\theta))^2 < \infty$, в этом случае мерой точности оценок является их дисперсия.

Утверждение: Для несмещенных оценок среднеквадратичное отклонение совпадает с ее дисперсией, а для смещенной оценки больше ее дисперсии.

Доказательство:

$$M_\theta(T - \tau)^2 = M(T - MT + MT - \tau)^2 = M(T - MT)^2 + M(MT - \tau)^2 + 2M((T - MT)(MT - \tau))$$

$$b^2 = 0 \Leftrightarrow MT = \tau$$

Теорема Рао-Блэкуэлла-Колмогорова: Оптимальная оценка, если она существует, является функцией от достаточной статистики.

По определению достаточная статистика $T = T(X)$ называется полной, если для всякой функции $\varphi(T)$ из того, что

$$M_\theta \varphi(T) = 0, \forall \theta$$

следует $\varphi(t) \equiv 0$ на всем множестве значений статистики T .

Теорема: Если существует полная достаточная статистика, то всякая функция от нее является оптимальной оценкой своего математического ожидания.

Итак, пусть существует полная достаточная статистика $T = T(X)$ и требуется оценить заданную параметрическую функцию $\tau(\theta)$. Тогда:

1) Если существует какая-то несмещенная оценка $\tau(\theta)$, то существует и несмещенная оценка, являющаяся функцией от T ; можно так же сказать, что если нет несмещенных оценок вида $H(T)$, то класс несмещенных оценок τ_τ для $\tau(\theta)$ пуст;

2) оптимальная (н.о.р.м.д.) оценка когда она существует, всегда является функцией от T и она однозначно определяется уравнением $M_\theta H(T) = \tau(\theta)$

3) оптимальную оценку τ^* можно искать по формуле:

$$\tau^* = H(T) = M_\theta(T_1 | T)$$

исходя из любой несмещенной оценки T_1 функции $\tau(\theta)$. Условие регулярности (Рао-Крамера):

Параметрическое семейство $F = \{F_\theta, \theta \in \Theta\}$ называется регулярным, если выполнены следующие условия (условия регулярности):

1) $L(\bar{x}, \theta) > 0$ для всех значений $\bar{x} \in B$ из выборочного пространства и дифференцируема по θ , $\forall \theta \in \Theta$

2) Случайная величина $V(X; \theta)$, называемая функцией вклада выборки и определенная равенством

$$V(X; \theta) = \frac{\partial \ln(L(\bar{x}; \theta))}{\partial \theta} = \sum_{i=1}^n \frac{\partial \ln(f_\theta(\bar{x}_j))}{\partial \theta}$$

имеет ограниченную дисперсию:

$$0 < M_\theta V^2(X; \theta) < \infty$$

при этом значение $\frac{\partial \ln(f_\theta(\bar{x}_j))}{\partial \theta}$ будем называть вкладом i -ого наблюдения выборки.

3) $\forall \theta \in \Theta$ статистики $T(x)$ верно равенство:

$$\frac{\partial}{\partial \theta} \int_{R^n} T(\bar{x}) L(\bar{x}; \theta) d\bar{x} = \int_{R^n} T(\bar{x}) \frac{\partial L(\bar{x}; \theta)}{\partial \theta} d\bar{x}$$

Условия регулярности:

- 1) Так как $0 < \theta < 1$, то $L(\bar{x}, \theta) > 0$. Функция дифференцируема по θ .
- 2) Воспользуемся равенством:

$$i_n(\theta) = M_\theta V^2(X; \theta) = n * i(\theta)$$

Найдём $i(\theta)$:

$$i(\theta) = -M_\theta \frac{\partial^2 (\ln f_\theta(x_1))}{\partial \theta^2} = \frac{\partial^2 (\ln(f(x_1, \theta)))}{\partial \theta^2} = \frac{\partial^2 (x \ln((1 - \theta)\theta))}{\partial \theta^2}$$

Первая производная:

$$\frac{x\theta + \theta - 1}{(\theta - 1)\theta}$$

Вторая производная:

$$-\frac{\theta^2(x + 1) + 2\theta - 1}{(\theta - 1)^2 \theta^2}$$

$$M(x) = \frac{1 - \theta}{\theta}$$

$$M_\theta \left(\frac{\theta^2(x + 1) + 2\theta - 1}{(\theta - 1)^2 \theta^2} \right) = M_\theta \left(\frac{\theta^2 x}{(\theta - 1)^2 \theta^2} \right) + \frac{\theta + 2\theta - 1}{(\theta - 1)^2 \theta^2} = \left(\frac{\theta^2 M(x)}{(\theta - 1)^2 \theta^2} \right) + \frac{3\theta - 1}{(\theta - 1)^2 \theta^2} = \left(-\frac{6}{(\theta - 1)^2 \theta^2} \right) + \frac{3\theta - 1}{(\theta - 1)^2 \theta^2} = -\frac{\theta^2 - 4\theta + 1}{(\theta - 1)^2 \theta^2}$$

В силу того, что $0 < \theta < 1$, то информация Фишера конечна, следовательно, свойство выполняется.

- 3) Выборочное среднее не зависит от параметра θ , то свойство выполняется.

$$f(x) = (1 - p)^x p$$

Запишем логарифмическую функцию правдоподобия:

$$L(\bar{x}; \theta) = \prod_{i=1}^n (1 - \theta)^{x_i} \theta$$

Прологарифмируем обе части:

$$\begin{aligned} \ln L(\bar{x}; \theta) &= \sum_{i=1}^n \ln(1 - \theta)^{x_i} \theta = \sum_{i=1}^n \ln(\theta) + \sum_{i=1}^n \ln((1 - \theta)^{x_i}) \\ \frac{\partial(\ln(L(x, \theta)))}{\partial \theta} &= \sum_{i=1}^n \left(\frac{1}{\theta}\right) + \sum_{i=1}^n \left(-\frac{x_i}{1 - \theta}\right) \end{aligned}$$

Так как

$$\sum_{i=1}^n \left(\frac{x_i}{1 - \theta}\right) = \sum_{i=1}^n \left(\frac{1}{\theta}\right), \theta \neq 0; 1,$$

то

$$\begin{aligned} \frac{1}{1 - \theta} \bar{X} &= \frac{1}{\theta} \\ \bar{X} &= \frac{1 - \theta}{\theta} \\ M\left(\frac{1 - \theta}{\theta}\right) &= M\bar{X} = M\xi \end{aligned}$$

А значит, оценка несмещенная. Также оценка является состоятельной, в силу того, что

$$M\bar{X} \xrightarrow{P} \bar{X}$$

Теперь приступим к оцениванию эффективности данной оценки:

Вспомним, что

$$-M_{\theta} \frac{\partial^2(\ln(f_{\theta}(X_1)))}{\partial \theta} = i_1(\theta)$$

$f(x)$ геометрического распределения можно дважды продифференцировать.

$$\begin{aligned} \frac{\partial^2(\ln(f_{\theta}(X_1)))}{\partial \theta} &= -\frac{1}{\sigma^2} \\ e(T) &= \frac{(r'(\theta))^2}{i_n(\theta) D_{\theta} T} \end{aligned}$$

3.2 Распределения Максвелла

3.2.1 Нахождение выборочного среднего и выборочной дисперсии геометрического распределения

Выборочное среднее


```

In [72]: 1 #n=5
          2 vs_5 = []
          3 for i in range(5):
          4     vs_5.append(round(np.mean(means__5[i]),6))
          5 print(vs_5)
          6 #n=10
          7 vs_10 = []
          8 for i in range(5):
          9     vs_10.append(round(np.mean(means__10[i]),6))
         10 print(vs_10)
         11 #n=100
         12 vs_100 = []
         13 for i in range(5):
         14     vs_100.append(round(np.mean(means__100[i]),6))
         15 print(vs_100)
         16 #n=1000
         17 vs_1000 = []
         18 for i in range(5):
         19     vs_1000.append(round(np.mean(means__1000[i]),6))
         20 print(vs_1000)
         21 #n=100000
         22 vs_100000 = []
         23 for i in range(5):
         24     vs_100000.append(round(np.mean(means__100000[i]),6))
         25 print(vs_100000)

[1.915438, 1.095802, 1.249879, 0.910554, 1.714316]
[1.577273, 1.727918, 1.602861, 1.788271, 1.522487]
[1.664325, 1.623293, 1.442467, 1.652924, 1.605711]
[1.599327, 1.606045, 1.588832, 1.574616, 1.550436]
[1.590289, 1.593489, 1.59538, 1.593763, 1.594673]

```

```

In [73]: 1 #Сравнение
          2 2*np.sqrt(2/np.pi)

```

Out[73]: 1.5957691216057308

$$M\xi = 2\lambda\sqrt{\frac{2}{\pi}}$$

$$\text{При } \lambda = 1.0 \quad M\xi = 2 \cdot \sqrt{\frac{2}{\pi}} \approx 1.5957691216057308$$

Выборочная дисперсия

```
In [74]: 1 #n=5
2 vd_5 = []
3 for i in range(5):
4     vd_5.append(round(np.var(means__5[i]),6))
5 print(vd_5)
6 #n=10
7 vd_10 = []
8 for i in range(5):
9     vd_10.append(round(np.var(means__10[i]),6))
10 print(vd_10)
11 #n=100
12 vd_100 = []
13 for i in range(5):
14     vd_100.append(round(np.var(means__100[i]),6))
15 print(vd_100)
16 #n=1000
17 vd_1000 = []
18 for i in range(5):
19     vd_1000.append(round(np.var(means__1000[i]),6))
20 print(vd_1000)
21 #n=100000
22 vd_100000 = []
23 for i in range(5):
24     vd_100000.append(round(np.var(means__100000[i]),6))
25 print(vd_100000)

[0.406101, 0.287325, 0.097846, 0.057745, 0.327642]
[0.251981, 0.210423, 0.524976, 0.650905, 0.281158]
[0.359419, 0.474391, 0.422773, 0.447515, 0.395477]
[0.457396, 0.431927, 0.44121, 0.44778, 0.477308]
[0.451604, 0.45236, 0.453762, 0.453338, 0.456103]
```

```
In [75]: 1 #Сравнение
2 (3*np.pi-8)/np.pi
```

Out[75]: 0.4535209105296745

$$D_{\xi}^{\lambda} = \frac{3\pi - 8}{\pi} \cdot \lambda$$

При $\lambda = 1.0$ $D_{\xi}^{\lambda} = \frac{3\pi - 8}{\pi} \approx 0.4535209105296745$

Как видно из полученных значений, чем больше объём выборки, тем менее отличаются выборочное среднее от теоретического математического ожидания и выборочная дисперсия от теоретической дисперсии.

3.2.2 Построение доверительного интервала для выборочного среднего

Положим $\gamma = 0.95$ и найдем доверительный интервал для выборочного среднего.

$$\Phi(c_{\gamma}) = \frac{\gamma}{2} = 0.475$$

Из таблицы значений функции Лапласа $c_{\gamma} \approx 1.96$

In [76]:

```
1 #n=5
2 print('n = 5')
3 for i in range(5):
4     print('(', vs_5[i], '-+ 1.96 *', np.sqrt(vd_5[i]/5), ') = (', vs_5[i]
5 #n=10
6 print('n = 10')
7 for i in range(5):
8     print('(', vs_10[i], '-+ 1.96 *', np.sqrt(vd_10[i]/10), ') = (', vs_1
9 #n=100
10 print('n = 100')
11 for i in range(5):
12     print('(', vs_100[i], '-+ 1.96 *', np.sqrt(vd_100[i]/100), ') = (', v
13 #n=1000
14 print('n = 1000')
15 for i in range(5):
16     print('(', vs_1000[i], '-+ 1.96 *', np.sqrt(vd_1000[i]/1000), ') = (
17 #n=100000
18 print('n = 100000')
19 for i in range(5):
20     print('(', vs_100000[i], '-+ 1.96 *', np.sqrt(vd_100000[i]/100000), '
```

```
n = 5
( 1.915438 -+ 1.96 * 0.284991578822954 ) = ( 1.915438 -+ 0.558583 )
( 1.095802 -+ 1.96 * 0.23971858501167573 ) = ( 1.095802 -+ 0.469848 )
( 1.249879 -+ 1.96 * 0.13988995675172683 ) = ( 1.249879 -+ 0.274184 )
( 0.910554 -+ 1.96 * 0.10746627377926528 ) = ( 0.910554 -+ 0.210634 )
( 1.714316 -+ 1.96 * 0.2559851558196295 ) = ( 1.714316 -+ 0.501731 )
n = 10
( 1.577273 -+ 1.96 * 0.15873909411357998 ) = ( 1.577273 -+ 0.311129 )
( 1.727918 -+ 1.96 * 0.14505964290594403 ) = ( 1.727918 -+ 0.284317 )
( 1.602861 -+ 1.96 * 0.22912354745857091 ) = ( 1.602861 -+ 0.449082 )
( 1.788271 -+ 1.96 * 0.2551283990464409 ) = ( 1.788271 -+ 0.500052 )
( 1.522487 -+ 1.96 * 0.16767766696850242 ) = ( 1.522487 -+ 0.328648 )
n = 100
( 1.664325 -+ 1.96 * 0.05995156378277384 ) = ( 1.664325 -+ 0.117505 )
( 1.623293 -+ 1.96 * 0.06887604808639938 ) = ( 1.623293 -+ 0.134997 )
( 1.442467 -+ 1.96 * 0.06502099660878784 ) = ( 1.442467 -+ 0.127441 )
( 1.652924 -+ 1.96 * 0.06689656194454241 ) = ( 1.652924 -+ 0.131117 )
( 1.605711 -+ 1.96 * 0.06288696208277197 ) = ( 1.605711 -+ 0.123258 )
n = 1000
( 1.599327 -+ 1.96 * 0.021386818370201774 ) = ( 1.599327 -+ 0.041918 )
( 1.606045 -+ 1.96 * 0.020782853509564082 ) = ( 1.606045 -+ 0.040734 )
( 1.588832 -+ 1.96 * 0.021004999404903586 ) = ( 1.588832 -+ 0.04117 )
( 1.574616 -+ 1.96 * 0.021160812838830177 ) = ( 1.574616 -+ 0.041475 )
( 1.550436 -+ 1.96 * 0.021847379705584834 ) = ( 1.550436 -+ 0.042821 )
n = 100000
( 1.590289 -+ 1.96 * 0.0021250976448154092 ) = ( 1.590289 -+ 0.004165 )
( 1.593489 -+ 1.96 * 0.002126875642815066 ) = ( 1.593489 -+ 0.004169 )
( 1.59538 -+ 1.96 * 0.0021301690073794615 ) = ( 1.59538 -+ 0.004175 )
( 1.593763 -+ 1.96 * 0.0021291735485863992 ) = ( 1.593763 -+ 0.004173 )
( 1.594673 -+ 1.96 * 0.0021356568076355336 ) = ( 1.594673 -+ 0.004186 )
```

*Округлено до 6 знаков после запятой.

3.2.3 Нахождение оптимальности рассматриваемых оценок

Для простоты предположим, что все частоты $p_i, i = \overline{1, n}$ равны единице.
Запишем функцию максимального правдоподобия для закона Максвелла.

$$L(\lambda) = \left(\frac{2}{\pi}\right)^{\frac{n}{2}} \frac{\prod_{i=1}^n x_i^2}{\lambda^{3n}} e^{-\frac{1}{2\lambda^2} \sum_{i=1}^n x_i^2}$$

$$\ln(L(\lambda)) = 2 \sum_{i=1}^n \ln x_i - 3n \ln \lambda - \frac{n}{2} \ln \frac{2}{\pi} - \frac{1}{2\lambda^2} \sum_{i=1}^n x_i^2$$

$$\frac{\partial(\ln L(\lambda))}{\partial \lambda} = -\frac{3n}{\lambda} + \frac{\sum_{i=1}^n x_i^2}{\lambda^3} = 0$$

Переходя к статистическому ряду (не все p_i равны 1, $i = \overline{1, n}$), получим уравнение для нахождения λ :

$$\lambda = \sqrt{\frac{\sum_{i=1}^m p_i x_i^2}{3n}}$$

Оценка методом моментов:

Поскольку по выборке оценивается лишь один параметр, то для нахождения λ используем оценку математического ожидания.

$$M[X] = \bar{X} = 2\lambda \sqrt{\frac{2}{\pi}},$$

где $\bar{X} = \frac{1}{n} \sum_{i=1}^m x_i p_i, \sum_{i=1}^m p_i = n$

Отсюда $\lambda = \sqrt{\frac{\pi}{8} \bar{X}}$

3.3 Работа с данными

В данном задании проведем анализ реальных данных. Воспользуемся нетипичной интерпретацией геометрического распределения - игра бейсбол. В бейсболе геометрическое распределение полезно для анализа вероятности того, что отбивающий получит удар, прежде чем он получит три удара; здесь цель - добиться успеха за 3 испытания.

Рассмотрим финальную серию чемпионата МЛБ-2020 между Лос Анджелес Доджерс и Тампа-Бэй Рэйс. Проанализируем статистику бэттеров первых 5 матчей, чтобы оценить какие бэттеры должны чаще выходить на поле в последнем матче серии и сравним с тем, как на самом деле это было.

Почему принято решение выбрать именно такой подход? Всё из-за того, что чем выше процент отбивания у бэттера, тем выше вероятность того, что произойдет успех. Докажем это:

Пусть вероятность того, что бэттер отобьёт удар равна p . Тогда рассчитаем вероятность успеха:

$$P(X = 0) + P(X = 1) + P(X = 2) = q^0 p + q^1 p + q^2 p = p + (1 - p)p + (1 - p)^2 p = p + p - p^2$$

Пройдемся по циклу и убедимся, что при увеличении p увеличивается и вероятность успеха.

```
In [77]: 1 for p in [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]:
          2     print(n*(n**2-3*n+3))
```

```
0.271
0.488
0.6569999999999999
0.784
0.875
0.9359999999999999
0.9730000000000000
0.9919999999999999
0.9989999999999999
```

Рассмотрим статистические показатели, которые нам пригодятся:

AB - At Bats = PA - BB - IBB - HBP - CI - SF - SH (На бите): Выходы на биту бэттера, за исключением уоков, ударов мячом, жертвованных ударов, умышленных помех со стороны защиты или других препятствий.

H - Hits (Хиты): общее количество хитов (1B, 2B, 3B и HR). Хит - удар, давший возможность выйти на базу. При этом защита не совершила ошибку.

BA - Batting average = H / AB (он же AVG - средний коэффициент результативности отбивания): число хитов, деленное на число выходов на биту.

```
In [78]: 1 match_1 = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Матч-1.csv', ', ',
2 match_2 = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Матч-2.csv', ', ',
3 match_3 = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Матч-3.csv', ', ',
4 match_4 = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Матч-4.csv', ', ',
5 match_5 = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Матч-5.csv', ', '.
```

```
In [79]: 1 match_1
```

Out[79]:

	Игрок	Команда	AB	H	BA
0	Adames W.	TAM	2	0	0.000
1	Arozarena R.	TAM	3	0	0.000
2	Barnes A.	LAD	4	0	0.000
3	Bellinger C.	LAD	4	1	0.250
4	Betts M.	LAD	4	2	0.500
5	Brosseau M.	TAM	1	1	1.000
6	Choi Ji-Man	TAM	0	0	0.000
7	Diaz Y.	TAM	4	1	0.250
8	Hernandez E.	LAD	2	1	0.500
9	Kiermaier K.	TAM	3	2	0.667
10	Lowe B.	TAM	4	0	0.000
11	Margot M.	TAM	4	1	0.250
12	Meadows A.	TAM	2	0	0.000
13	Muncy M.	LAD	4	2	0.500
14	Pederson J.	LAD	2	0	0.000
15	Renfroe H.	TAM	2	0	0.000
16	Seager C.	LAD	2	0	0.000
17	Taylor C.	LAD	3	2	0.667
18	Turner J.	LAD	4	1	0.250
19	Wendle J.	TAM	4	1	0.250
20	Zunino M.	TAM	3	0	0.000
21	Smit U.	LAD	5	1	0.200

In [80]: 1 match 2

Out[80]:

	Игрок	Команда	AB	H	BA
0	Adames W.	TAM	4	1	0.167
1	Arozarena R.	TAM	3	1	0.167
2	Barnes A.	LAD	1	0	0.000
3	Bellinger C.	LAD	3	0	0.143
4	Betts M.	LAD	3	0	0.286
5	Brosseau M.	TAM	2	0	0.333
6	Choi Ji-Man	TAM	3	1	0.333
7	Diaz Y.	TAM	1	1	0.400
8	Hernandez E.	LAD	1	0	0.333
9	Kiermaier K.	TAM	4	0	0.286
10	Lowe B.	TAM	5	2	0.222
11	Margot M.	TAM	3	2	0.429
12	Meadows A.	TAM	3	1	0.200
13	Muncy M.	LAD	3	0	0.286
14	Pederson J.	LAD	1	0	0.000
15	Phillips B.	TAM	0	0	0.000
16	Pollock A.	LAD	2	0	0.000
17	Renfroe H.	TAM	0	0	0.000
18	Rios E.	LAD	2	0	0.000
19	Seager C.	LAD	4	2	0.333
20	Taylor C.	LAD	4	1	0.429
21	Turner J.	LAD	4	1	0.250
22	Wendle J.	TAM	3	1	0.286
23	Zunino M.	TAM	4	0	0.000
24	Smit U.	LAD	4	1	0.222

In [81]: 1 match 3

Out[81]:

	Игрок	Команда	AB	H	BA
0	Adames W.	TAM	3	1	0.222
1	Arozarena R.	TAM	4	1	0.200
2	Barnes A.	LAD	3	1	0.125
3	Bellinger C.	LAD	4	1	0.182
4	Betts M.	LAD	5	2	0.333
5	Choi Ji-Man	TAM	4	0	0.143
6	Hernandez E.	LAD	1	0	0.250
7	Kiermaier K.	TAM	2	0	0.222
8	Lowe B.	TAM	4	0	0.154
9	Margot M.	TAM	3	1	0.400
10	Meadows A.	TAM	4	1	0.222
11	Muncy M.	LAD	4	2	0.364
12	Pederson J.	LAD	3	1	0.167
13	Seager C.	LAD	3	1	0.333
14	Taylor C.	LAD	4	0	0.273
15	Tsutsugo Y.	TAM	1	0	0.000
16	Turner J.	LAD	5	2	0.308
17	Wendle J.	TAM	3	0	0.200
18	Zunino M.	TAM	2	0	0.000
19	Perez M.	TAM	0	0	0.000
20	Smit U.	LAD	4	0	0.154

In [82]: 1 match 4

Out[82]:

	Игрок	Команда	AB	H	BA
0	Adames W.	TAM	4	1	0.231
1	Arozarena R.	TAM	4	3	0.357
2	Bellinger C.	LAD	4	0	0.133
3	Betts M.	LAD	5	0	0.235
4	Brosseau M.	TAM	2	1	0.400
5	Choi Ji-Man	TAM	0	0	0.143
6	Diaz Y.	TAM	3	0	0.250
7	Hernandez E.	LAD	4	1	0.250
8	Kiermaier K.	TAM	4	2	0.308
9	Lowe B.	TAM	4	1	0.176
10	Margot M.	TAM	2	0	0.333
11	Meadows A.	TAM	2	0	0.182
12	Muncy M.	LAD	4	1	0.333
13	Pederson J.	LAD	2	2	0.375
14	Phillips B.	TAM	1	1	1.000
15	Pollock A.	LAD	2	1	0.250
16	Renfroe H.	TAM	4	1	0.167
17	Seager C.	LAD	5	4	0.500
18	Taylor C.	LAD	5	1	0.250
19	Tsutsugo Y.	TAM	1	0	0.000
20	Turner J.	LAD	5	4	0.444
21	Wendle J.	TAM	1	0	0.182
22	Zunino M.	TAM	2	0	0.000
23	Smit U.	LAD	4	1	0.176

In [83]: 1 match 5

Out[83]:

	Игрок	Команда	AB	H	BA
0	Adames W.	TAM	4	0	0.176
1	Arozarena R.	TAM	4	1	0.333
2	Barnes A.	LAD	2	0	0.100
3	Bellinger C.	LAD	4	1	0.158
4	Betts M.	LAD	5	1	0.227
5	Brosseau M.	TAM	0	0	0.400
6	Choi Ji-Man	TAM	0	0	0.143
7	Diaz Y.	TAM	3	2	0.364
8	Hernandez E.	LAD	1	0	0.222
9	Kiermaier K.	TAM	3	2	0.375
10	Lowe B.	TAM	4	0	0.143
11	Margot M.	TAM	3	2	0.400
12	Meadows A.	TAM	2	0	0.154
13	Muncy M.	LAD	3	2	0.389
14	Pederson J.	LAD	2	1	0.400
15	Renfroe H.	TAM	1	0	0.143
16	Seager C.	LAD	3	1	0.471
17	Taylor C.	LAD	4	0	0.200
18	Tsutsugo Y.	TAM	1	0	0.000
19	Turner J.	LAD	4	0	0.364
20	Wendle J.	TAM	4	0	0.133
21	Zunino M.	TAM	2	0	0.000
22	Perez M.	TAM	0	0	0.000
23	Smit U.	LAD	4	0	0.143

```
In [84]: 1 itog = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Итого.csv', ',', 'pa
2 iton
```

Out[84]:

	Игрок	Команда	AB	H	BA
0	Adames W.	TAM	17	3	0,176470588235294
1	Arozarena R.	TAM	18	6	0,333333333333333
2	Barnes A.	LAD	10	1	0,1
3	Bellinger C.	LAD	19	3	0,157894736842105
4	Betts M.	LAD	22	5	0,227272727272727
5	Brosseau M.	TAM	5	2	0,4
6	Choi Ji-Man	TAM	7	1	0,142857142857143
7	Diaz Y.	TAM	11	4	0,363636363636364
8	Hernandez E.	LAD	9	2	0,222222222222222
9	Kiermaier K.	TAM	16	6	0,375
10	Lowe B.	TAM	21	3	0,142857142857143
11	Margot M.	TAM	15	6	0,4
12	Meadows A.	TAM	13	2	0,153846153846154
13	Muncy M.	LAD	18	7	0,388888888888889
14	Pederson J.	LAD	10	4	0,4
15	Phillips B.	TAM	1	1	1
16	Pollock A.	LAD	4	1	0,25
17	Rios E.	LAD	2	0	0
18	Seager C.	LAD	17	8	0,470588235294118
19	Taylor C.	LAD	20	4	0,2
20	Turner J.	LAD	22	8	0,363636363636364
21	Smit U.	LAD	21	3	0,142857142857143
22	Renfroe H.	TAM	7	1	0,142857142857143
23	Wendle J.	TAM	15	2	0,133333333333333
24	Zunino M.	TAM	13	0	0
25	Tsutsugo Y.	TAM	3	0	0

А теперь попробуем спрогнозировать кто из бэттеров будет подходить к бите больше в своей команде, анализируя приведенные выше статистические данные по итогам первых 5 матчей финальной серии. Выведем данные в порядке убывания.

```
In [85]: 1 prognoz_TAM = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Итого_TAM.
2 prognoz_TAM
```

Out[85]:

	Игрок	Команда	AB	H	BA	Комментарий
0	Margot M.	TAM	15	6	0,4	Лучший BA, один из лучших AB
1	Arozarena R.	TAM	18	6	0,3333333333333333	Один из лидеров по AB и отличный BA
2	Adames W.	TAM	17	3	0,176470588235294	Один из лидеров по AB и средний BA
3	Lowe B.	TAM	21	3	0,142857142857143	Лучший AB, но значимо хуже BA
4	Kiermaier K.	TAM	16	6	0,375	Очень хороший BA и неплохой AB
5	Wendle J.	TAM	15	2	0,1333333333333333	Примерно равные средние показатели
6	Meadows A.	TAM	13	2	0,153846153846154	Примерно равные средние показатели
7	Diaz Y.	TAM	11	4	0,363636363636364	Отличный BA, но все же низкий AB
8	Zunino M.	TAM	13	0	0	Приличный AB, но худший BA
9	Choi Ji-Man	TAM	7	1	0,142857142857143	Низкие AB и BA
10	Renfroe H.	TAM	7	1	0,142857142857143	Низкие AB и BA
11	Brosseau M.	TAM	5	2	0,4	Один из худших AB, но отличный BA
12	Tsutsugo Y.	TAM	3	0	0	Очень низкие AB и BA
13	Phillips B.	TAM	1	1	1	Подходил к бите всего лишь раз

А теперь посмотрим, сколько в итоге было подходов к бите у игроков TAM. Выведем данные в порядке убывания.

```
In [86]: 1 match_6_TAM = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Матч-6_TAM
2 match_6_TAM
```

Out[86]:

	Игрок	Команда	AB	H
0	Margot M.	TAM	4	0
1	Arozarena R.	TAM	4	2
2	Adames W.	TAM	4	0
3	Lowe B.	TAM	3	0
4	Kiermaier K.	TAM	3	1
5	Wendle J.	TAM	3	0
6	Meadows A.	TAM	3	1
7	Zunino M.	TAM	3	1
8	Choi Ji-Man	TAM	2	0
9	Diaz Y.	TAM	1	0
10	Renfroe H.	TAM	1	0
11	Brosseau M.	TAM	1	0

```
In [87]: 1 prognos_LAD = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Итого_LAD.
2 prognos_LAD
```

Out[87]:

	Игрок	Команда	AB	H	BA	Комментарий
0	Turner J.	LAD	22	8	0,363636363636364	Лучший AB, один из лучших BA
1	Betts M.	LAD	22	5	0,227272727272727	Лучший AB, высокий BA
2	Seager C.	LAD	17	8	0,470588235294118	Высокий AB, лучший BA
3	Muncy M.	LAD	18	7	0,388888888888889	Высокий AB, один из лучших BA
4	Taylor C.	LAD	20	4	0,2	Высокий AB, неплохой BA
5	Smit U.	LAD	21	3	0,142857142857143	Высокий AB, неплохой BA
6	Bellinger C.	LAD	19	3	0,157894736842105	Высокий AB, неплохой BA
7	Pederson J.	LAD	10	4	0,4	Средний AB, высокий BA
8	Hernandez E.	LAD	9	2	0,222222222222222	Средние показатели
9	Barnes A.	LAD	10	1	0,1	Средние показатели
10	Pollock A.	LAD	4	1	0,25	Низкий AB, неплохой BA
11	Rios E.	LAD	2	0	0	Всего лишь 2 подхода к бите

А теперь посмотрим, сколько в итоге было подходов к бите у игроков LAD. Выведем данные в порядке убывания.

```
In [88]: 1 match_6_LAD = pd.read_csv('/home/alexander/Рабочий стол/Курсовая/Матч-6_LAD
2 match_6_LAD
```

Out[88]:

	Игрок	Команда	AB	H
0	Betts M.	LAD	4	2
1	Muncy M.	LAD	4	0
2	Turner J.	LAD	3	0
3	Seager C.	LAD	3	0
4	Taylor C.	LAD	3	1
5	Smit U.	LAD	3	1
6	Bellinger C.	LAD	3	0
7	Barnes A.	LAD	3	1
8	Pollock A.	LAD	2	0
9	Hernandez E.	LAD	1	0
10	Pederson J.	LAD	0	0
11	Rios E.	LAD	0	0

Оценка того, кто больше из игроков будет подходить к бите в своей команде оказалось достаточно точной, несмотря на то, что данные были взяты лишь по 5 последним играм, что подтверждает важность этих показателей у игроков. Конечно же, по данным за такой короткий период корректировку могли внести такие ситуации как, например, травмы игроков, из-за которых они не смогли бы участвовать в 6 финальном матче или же участвовать в нём не в полной мере. Обычно в бейсболе статистику игроков оценивают по всему прошедшему сезону, что даст более высокую точность.

Как можно заметить, рассмотренные статистические показатели у игроков команды LAD лучше, чем у игроков команды TAM. Отсюда неудивительно, что LAD выиграли не только последний матч у TAD со счетом 3:1, но и всю финальную серию со счетом 4:2 и стали чемпионами MLB-2020.

Ссылка на данные: <https://www.scoreboard.com/ru/baseball/usa/mlb-2020/> (<https://www.scoreboard.com/ru/baseball/usa/mlb-2020/>)

4 Домашнее задание. Проверка статистических гипотез

1. Критерий согласия Колмагорова

Пусть дана выборка $X = (X_1, \dots, X_n)$ из распределения $L(\xi)$ и $F\xi$ - неизвестное распределение.

- $H_0 : F\xi = F(x)$ - простая гипотеза
- $H_1 : \text{не } F(x)$

Критерий Колмогорова основан на теореме Колмогорова:

$$D_n = D_n(x) = \sup |\hat{F}_n(x) - F(x)|_{x \in R}$$

где D_n - это отклонение эмпирической функции распределения от теоретической функции распределения.

\hat{F}_n - оптимальная несмещенная состоятельная оценка для $F(x)$.

Замечание: D_n не должно сильно отклоняться от 0.

По т. Колмогорова:

$$P(nD_n \geq \lambda_\alpha | H_0) = 1 - K(\lambda_\alpha) = \alpha$$

по $\alpha \rightarrow \lambda_\alpha$

Проверяем, выполняется ли неравенство: $nD_n \geq \lambda_\alpha$

Известно, что

$$X_1 = \{x : D_n(x)\sqrt{n} \geq \lambda_\alpha\}$$

Следовательно, H_0 отвергается $\Leftrightarrow nD_n \geq \lambda_\alpha$

По Долошеву: $\frac{6nD_n}{6\sqrt{n}}$ сходится к распределению Колмогорова, причем $\sqrt{n}D_n \in \frac{1}{6\sqrt{n}}$

Способ вычисления $D_n = \sup |F_n(x) - F(x)|_{x \in R}$. Вычисление супремума функции не является

тривиальной задачей. Однако в данном случае $\hat{F}_n(x)$ принимает конечное число значений:

$\{\frac{1}{n}, \frac{2}{n}, \dots, \frac{n}{n}\}$, что значительно упрощает задачу. Пусть у нас есть вариационный ряд выборки:

$X_{(1)}, X_{(2)}, \dots, X_{(n)}$. Определим следующие две функции:

$$D_n^+ = \max_{1 \leq k \leq n} \left| \frac{k}{n} - F(x_{(k)}) \right|$$

$$D_n^- = \max_{1 \leq k \leq n} \left| F(x_{(k)}) - \frac{k-1}{n} \right|$$

Тогда вычислить D_n можно следующим образом:

$$D_n = \max\{D_n^-, D_n^+\}$$

Однако критерий Колмагорова обладает рядом минусов: 1. Функция $D_n = \sup |F_n(x) - F(x)|_{x \in R}$ не зависит от вида функции распределения $F(x)$, только в случае если $F(x)$ непрерывная. Встает вопрос, что делать если $F(x)$ имеет точки разрыва.

Пусть Y_1, \dots, Y_n - н.о.р.сл.в. $Y_i \approx R[0, 1]$. А X_1, \dots, X_n - выборка из некоторого распределения, функция которого имеет точки разрыва. Построим следующую случайную величину:

$U_i = F(x_i^-) + Y_i[F(x_i) - F(x_i^-)]$, где $F(x_i^-) = \lim_{z \downarrow x_i} F(z)$. Доказывается, что случайная величина $U_i \approx R[0, 1]$

2. В случае сложных гипотез распределение $D_n(\theta)$, зависит как от вида априорных распределений, так и от способа получения оценок, размера выборки n , вида θ .

$$H_0 : F_\xi(x) \in F_0 = \{F_\theta(x), \theta \in \Theta\}$$

придется исследовать следующую статистику

$$b_n = \sup_{x \in R} \hat{F}_n(x) - F_{\hat{\theta}(x)},$$

где $\hat{\theta}$ есть зависимость от θ . Это плохо потому что, по одной выборке считается статистический параметр $\hat{\theta}$ и вычисляется критерий. То есть, есть зависимость. В таком случае, при наличии достаточно большой выборки длины m следует выбрать границу n и использовать первую часть для оценки параметра, а вторую для вычисления критерия.

Определение:

Пусть случайные величины $\xi_1, \xi_2, \dots, \xi_n$ имеют стандартное нормальное распределение, тогда случайная величина:

$$\chi_n^2 = \sum_{i=1}^n \xi_i^2$$

имеет распределение, которое называется хи квадрат с n - степенями свободы.

2. Критерий Хи-квадрат

Пусть ξ - случайный вектор $\xi = (\xi_1, \dots, \xi_n)$, и $\xi_i \approx N(0, 1)$. И, вектор ξ имеет единичную матрицу ковариаций. Пусть также $c = (c_1, c_2, \dots, c_n) \in R^n$, такой что $|c| = 1$.

Рассмотрим проекцию $\xi^{(c)}$ на гиперплоскость $L_{\bar{c}} = x \in R^n : (\bar{x}, \bar{c}) = 0$, которая ортогональна вектору \bar{c} . Тогда вектор ξ имеет математическое ожидание равное $\theta = (0, \dots, 0)$ и матрицу ковариации

$$C(\xi^{(c)}) = E - ||c_i c_j||_{i,j=1}^n$$

4.1 Геометрическое распределение

4.1.1 Проверка гипотез о виде распределения

Сформулируем гипотезу H_0 и H_1 : Пусть дана выборка $X = (X_1, \dots, X_n)$ из распределения $L(\xi)$ и F_ξ - неизвестное распределение.

- $H_0 : F_\xi = \text{Geom}(x, 0.5)$ - простая гипотеза, Geom - выбранное дискретное распределение
- H_1 : не $\text{Geom}(x, 0.5)$

Для проверки гипотезы H_0 воспользуемся критерием Пирсона (хи-квадрат). Для каждой выборки объемов $n = 10, 100, 1000, 100000$ найдем значение критерия, границу критического множества для уровня значимости $\alpha = 0.1$ и $\alpha = 0.05$. В каждом случае возьмем $N = 15$. Разделим каждую выборку на равновероятностные интервалы $\text{pr.arange}(0.01, 1, \frac{1}{15})$, в случае если взятый интервал меньше, чем 1, он склеивается со следующим.

ЗДЕСЬ ДОЛЖЕН БЫТЬ КОД, ВЫЧИСЛЯЮЩИЙ ЗНАЧЕНИЕ КРИТЕРИЯ, ГРАНИЦЫ ДЛЯ ЗАДАННЫХ УРОВНЕЙ И ОПРЕДЕЛЯЮЩИЙ ДЛЯ КАЖДОЙ ВЫБОРКИ ПРИНИМАЕТСЯ ЛИ ГИПОТЕЗА.

```
In [99]: 1 means 10
```

```
Out[99]: [array([1, 1, 2, 1, 4, 2, 1, 1, 2, 1]),
          array([2, 1, 3, 1, 3, 1, 1, 1, 3, 1]),
          array([1, 2, 2, 2, 2, 1, 4, 1, 1, 1]),
          array([2, 2, 2, 2, 1, 1, 1, 2, 1, 1]),
          array([4, 1, 1, 2, 1, 1, 1, 2, 3, 2])]
```

```
In [94]: 1 def intervals_v2(X):
2         vr = []
3         xr = [0]
4         v = len(X)
5         r = 0
6         Xs = sorted(X)
7         a1 = Xs[0]
8         an = Xs[len(X) - 1]
9         d = (an - a1)/15 #длина интервала
10        s = 0
11        x0 = -d
12        x1 = 0
13        k = 0
14        while s!=v: #подсчёт значений vr(попаданий в интервал)
15            r += 1
16            x0 += d
17            x1 += d
18            for z in range(v):
19                if x0 <= X[z] < x1:
20                    k += 1
21            if k >=4: #проверка, в каждый интервал не меньше 4 значений
22                vr.append(k)
23                xr.append(x1)
24                s += k
25                k = 0
26            elif s+k == v:
27                vr.append(k)
28                xr.append(x1)
29                break
30        return vr, xr
```

```
In [95]: 1 def pearson_criterion(X):
2         print("Выборки объема: ", len(X[0]))
3         v = len(X[0])
4         for t in range(5):
5             print("Выборка: ", t+1)
6             chi = 0
7             vr, xr = intervals_v2(X[t])
8             for i in range(len(vr)):
9                 v_i = vr[i]
10                a = xr[i]
11                b = xr[i+1]
12                p_i = e**(-a*a/8)-e**(-b*b/8)
13                q = ((v_i - v*p_i)**2/(v*p_i))
14                chi += q
15            print("Статистика критерия Пирсона Chi2:", chi)
16            print("Число степеней свободы: ", len(vr)-1)
```

```
In [96]: 1 intervals_v2(means 1000[0])
```

```
Out[96]: ([483, 280, 112, 58, 31, 21, 8, 4, 3],
[0,
1.4666666666666666,
2.1999999999999997,
3.6666666666666665,
4.3999999999999995,
5.133333333333333,
6.6,
7.333333333333333,
9.533333333333331,
12.466666666666661])
```

```
In [97]: 1 y = [24, 27, 21, 24, 20, 28, 21, 25, 24, 24], [23, 28, 29, 17, 23, 28, 24,
```



```
In [98]: 1 pearson_criterion(means_10)
2 print()
3 pearson_criterion(means_100)
4 print()
5 pearson_criterion(means_1000)
6 print()
7 pearson_criterion(means_100000)
```

Выборки объема: 10
Выборка: 1
Статистика критерия Пирсона Chi2: 12.786544815008364
Число степеней свободы: 1
Выборка: 2
Статистика критерия Пирсона Chi2: 16.932053428952944
Число степеней свободы: 1
Выборка: 3
Статистика критерия Пирсона Chi2: 9.599063445717778
Число степеней свободы: 2
Выборка: 4
Статистика критерия Пирсона Chi2: 11.887917289322628
Число степеней свободы: 1
Выборка: 5
Статистика критерия Пирсона Chi2: 7.4846893492343
Число степеней свободы: 2

Выборки объема: 100
Выборка: 1
Статистика критерия Пирсона Chi2: 117.5990433572617
Число степеней свободы: 5
Выборка: 2
Статистика критерия Пирсона Chi2: 34.79772027043223
Число степеней свободы: 5
Выборка: 3
Статистика критерия Пирсона Chi2: 101.44634548906481
Число степеней свободы: 4
Выборка: 4
Статистика критерия Пирсона Chi2: 152.3897880295792
Число степеней свободы: 5
Выборка: 5
Статистика критерия Пирсона Chi2: 71.69093973912976
Число степеней свободы: 4

Выборки объема: 1000
Выборка: 1
Статистика критерия Пирсона Chi2: 1257.2838385174853
Число степеней свободы: 8
Выборка: 2
Статистика критерия Пирсона Chi2: 103515.17036948926
Число степеней свободы: 9
Выборка: 3
Статистика критерия Пирсона Chi2: 217.5278124946671
Число степеней свободы: 8
Выборка: 4
Статистика критерия Пирсона Chi2: 513.2244278989731
Число степеней свободы: 7
Выборка: 5
Статистика критерия Пирсона Chi2: 1483.5149838733266
Число степеней свободы: 9

Выборки объема: 100000
Выборка: 1
Статистика критерия Пирсона Chi2: 176254935.18994218
Число степеней свободы: 12
Выборка: 2
Статистика критерия Пирсона Chi2: 82414828589691.83
Число степеней свободы: 13
Выборка: 3
Статистика критерия Пирсона Chi2: 120776582.70235875
Число степеней свободы: 14

In []: 1

Теперь проверим заведомо ложную гипотезу:

- $H_0 : F_\xi = \text{Geom}(x, 0.8)$ - простая гипотеза, Geom - выбранное дискретное распределение
- H_1 : не Geom(x, 0.8)

In []: 1

Видно, что для близких по значениям параметров гипотез критерий уверенно отклоняет неверную гипотезу только при больших объемах выборки. Это говорит о том, что при малых объемах выборки критерий говорит лишь о приближенных значениях параметров.

3. Критерий хи-квадрат для сложных гипотез В общем случае сложные для полиномиального распределения, используемого в критерии χ^2 , гипотезы будут принимать следующий вид.

$$H_0 : p = p(\theta), \theta = (\theta_1, \dots, \theta_r), \theta \in \Theta, r < N - 1$$

Тогда по аналогии с предыдущим случаем можем получить статистику $\hat{X}_n^2(\theta) = \sum_{j=1}^N \frac{(v_j - np_j(\theta))^2}{np_j(\theta)}$

Эта статистика зависит от неизвестного параметра, поэтому использовать непосредственно её нельзя. Для этого параметр θ заменяют некоторой оценкой $\hat{\theta}$ и получают в итоге статистику

$\hat{X}_n^2 = X_n^2(\hat{\theta})$. Однако узнать распределение \hat{X}_n^2 при гипотезе H_0 представляется трудной задачей.

Кроме того, величины $p_j(\hat{\theta})$ представляют собой функции от наблюдений.

Простая гипотеза H_0 заключается в том, что $p = \dot{p} = (\dot{p}_1, \dot{p}_2, \dots, \dot{p}_N)$ - заданный вероятностный вектор ($0 < \dot{p}_j < 1, j = 1, \dots, N; \dot{p}_1 + \dots + \dot{p}_N = 0$) Р. Фишер в 1924 г. получил, предельное

распределение статистики \hat{X}_n^2 , использующая оценку максимального правдоподобия

$$\hat{\theta} = \operatorname{argmax}_{\theta} \prod_{j=1}^N (p_j(\theta))^{v_j}$$

или иначе оценка по видоизменённому методу минимизации χ^2 , является распределением $\chi^2(N - 1 - r)$

$$L(\hat{X}_n^2 | H_0) \rightarrow \chi^2(N - 1 - r), n \rightarrow \infty$$

Тогда можно заключить, что критерий имеет вид:

H_0 отвергается $\leftrightarrow \hat{X}_n^2 > 1 - \alpha$ -квантиль распределения $\chi^2(N - 1 - r)$

Для упрощения задачи будем пользоваться оценкой параметра M , найденной методом максимального правдоподобия: $\hat{\theta}_1 = -\frac{N}{M}$

Type *Markdown* and LaTeX: α^2

4.1.2 Проверка параметрических гипотез

Данное задание посвящено вопросу различения двух простых параметрических гипотез, а также закреплению основных понятий, и практическому применению методов математической статистики для решения поставленной задачи.

4.1.2.1 Выбор данных

По причине конечности данного на решение задачи времени для ускорения завершения данного домашнего задания выберем геометрическое распределение и рассмотрим две выборки с разными (но известными) параметрами. А именно, к уже сконфигурированным выборкам из распределения с параметром $p = 0.5$, добавим к рассмотрению выборку из распределения с параметром $p = 0.8$. Сгенерируем новые выборки размеров 10, 100, 1000, 10000 и запишем их в файлы по уже отработанной схеме моделирования.

4.1.2.2 Постановка задачи

Критерий однородности Смирнова используется для проверки гипотезы о принадлежности двух независимых выборок одному закону распределения, то есть о том, что два эмпирических распределения соответствуют одному и тому же закону.

Мы сопоставляем сначала частоты по первому разряду, потом по сумме первого и второго разрядов, потом по сумме первого, второго и третьего разрядов и т. д. Таким образом, мы сопоставляем всякий раз накопленные к данному разряду частоты.

Если различия между двумя распределениями существенны, то в какой-то момент разность накопленных частот достигнет критического значения, и мы сможем признать различия статистически достоверными. В формулу критерия λ включается эта разность. Чем больше эмпирическое значение λ , тем более существенны различия.

Обозначим за нулевую гипотезу H_0 , что две исследуемые выборки подчиняются одному распределению случайной величины ξ . Соответственно, H_1 - исследуемые выборки не однородны.

Обозначим как X_0 - часть пространства наблюдений такая, что если $x \in X_0$, то следует принять H_0 , и как X_1 - часть пространства наблюдений такая, что если $x \in X_1$, то следует принять H_1 . Простым языком, если $x \in X_1$, а на самом деле истинна гипотеза H_0 , то говорится, что допущена ошибка первого рода. Если с точностью до наоборот - это ошибка второго рода. Вероятность $P_1(X_1)$ отвергнуть гипотезу H_0 , когда она действительно является ложной, называется мощностью критерия.

$P(x \in X_1 | H_0) = \alpha$ - ошибка 1 рода.

$P(x \in X_0 | H_1) = \beta$ - ошибка 2 рода.

Функция мощности критерия - функционал на множестве допустимых распределений F и выборке X .

$$W(F) = W(F; X_{1,\alpha}) = P(x \in X_{1,\alpha} | F),$$

где $P(x \in X_{1,\alpha} | F)$ - вероятность попасть в $X_{1,\alpha}$, если F - истинная гипотеза. Также

$$\alpha = \sup_{F \in F_0} W(F)$$

$$\beta = \sup_{F \in F_1} 1 - W(F)$$

4.1.3 Вычисление функции отношения правдоподобия

Рассмотрим подробнее теорию по функции отношения правдоподобия и ее применению.

$$W(\theta_0, X_{1\alpha}) = P_{\theta_0}(\bar{X} \in X_{1\alpha}) \leq \alpha$$

$$W(\theta_1, X_{1\alpha}) = P_{\theta_1}(\bar{X} \in X_{1\alpha}) \geq \alpha$$

$$W(\theta_0, X_{1\alpha}) = \sum_{X_{1\alpha}} L(\bar{X}, \theta_0) \leq \alpha$$

$$W(\theta_1, X_{1\alpha}) = \sum_{X_{1\alpha}} L(\bar{X}, \theta_1) \geq \alpha$$

$$l(\bar{X}) = \frac{L(\bar{X}, \theta_0)}{L(\bar{X}, \theta_1)} = \frac{\prod_{i=1}^n f_1(x_i)}{\prod_{i=1}^n f_2(x_i)}$$

Если $l(\bar{X}) \geq c$, то принимаем гипотезу H_0 , иначе принимаем гипотезу H_1 . Функция $l(\bar{X})$ называется статистикой отношения правдоподобия.

Рассмотрим функцию $\psi(c) = \sum_{\bar{X}: l(\bar{X}) \geq c} L(\bar{X}, \theta_0)$, $\psi(0) = 1$. Заметим, что чем больше c , тем меньше значение $\psi(c)$, то есть заданная функция убывающая.

Теорема Леймана-Пирсона:

Пусть $\forall \alpha \in (0, 1)$: $\psi(c_\alpha) = \alpha$, тогда критическая область $X_{1\alpha}$ наиболее мощный критерий для гипотезы H_0 с уровнем значимости α относительно альтернативы H_1 , среди всех критериев с уровнем значимости α .

Предположим, что нам известно, что обе выборки подчиняются закону распределения $Geom(x, \theta)$, но θ - неизвестный параметр.

- $H_0 : F_1 = F_2 = Geom(x, \theta_1)$ - сложная гипотеза, $Geom$ - выбранное дискретное распределение
- $H_1 : F_2 = Geom(x, \theta_2)$

Тогда $l(\bar{X}) =$

Type *Markdown* and *LaTeX*: α^2

Так как распределение дискретное, для равенства стоит выбрать критерий следующего вида (рандомизированный критерий - наиболее мощный для дискретных распределений):

$$\phi^* = 0, \sum x_i < t_{1-\alpha}$$

$$\phi^* = \frac{\alpha - P(\xi_0 < t_{1-\alpha})}{P(\xi_0 < t_{1-\alpha})}, \sum x_i = t_{1-\alpha}$$

$$\phi^* = 1, \sum x_i > t_{1-\alpha}$$

Или для удобства вычисления:

$$\phi^* = 0, \frac{1}{m} \sum x_i < \frac{1}{m} t_{1-\alpha}$$

$$\phi^* = \frac{\alpha - P(\xi_0 < t_{1-\alpha})}{P(\xi_0 < t_{1-\alpha})}, \frac{1}{m} \sum x_i = \frac{1}{m} t_{1-\alpha} \quad \phi^* = 1, \frac{1}{m} \sum x_i > \frac{1}{m} t_{1-\alpha}$$

Для вычисления β для каждой выборки можно воспользоваться следующей формулой:

$$\beta = P(\xi_1 \leq t_{1-\alpha}), \xi_1 \approx Geom(\dots)$$

Применим критерий к выборкам каждого распределения объемов 10, 100 и 1000, и найдем для каждой значение ошибки второго рода при $\alpha = 0.1, 0.5$.

Как видно из результатов обработки сгенерированных выборок, критерий и подобранное значение β оказываются очень точными, в следствии чего достижение равенства в нашем случае не происходит и дискретное распределение ведет себя, как непрерывное с точки зрения построения оптимального критерия.

In []: 1

Минимальный необходимый объем выборки можно определить из условия $\beta \approx \alpha$

In []: 1

Беря разный объем выборки, смотрим значение ошибки второго рода и получаем, что для $\alpha = 0.1$ $m = \dots$, а для $\alpha = 0.05$ $m = \dots$. Такой объем выборки крайне мал и то, что критерий для них работает, говорит в пользу его оптимальности.

4.2 Распределение Максвелла

4.2.1 Проверка гипотез о виде распределения

Сформулируем гипотезу H_0 и H_1 : Пусть дана выборка $X = (X_1, \dots, X_n)$ из распределения $L(\xi)$ и F_ξ - неизвестное распределение.

- $H_0 : F_\xi = \text{Maxwell}(x, 1.0)$ - простая гипотеза, Maxwell - выбранное непрерывное распределение
- H_1 : не $\text{Maxwell}(x, 1.0)$

Для проверки гипотезы H_0 воспользуемся критерием Пирсона (хи-квадрат). Для каждой выборки объемов $n = 10, 100, 1000, 100000$ найдем значение критерия, границу критического множества для уровня значимости $\alpha = 0.1$ и $\alpha = 0.05$. В каждом случае возьмем $N = 15$. Разделим каждую выборку на равновероятностные интервалы $\text{pr.arange}(0.01, 1, \frac{1}{15})$, в случае если взятый интервал меньше, чем 1, он склеивается со следующим.

ЗДЕСЬ ДОЛЖЕН БЫТЬ КОД, ВЫЧИСЛЯЮЩИЙ ЗНАЧЕНИЕ КРИТЕРИЯ, ГРАНИЦЫ ДЛЯ ЗАДАННЫХ УРОВНЕЙ И ОПРЕДЕЛЯЮЩИЙ ДЛЯ КАЖДОЙ ВЫБОРКИ ПРИНИМАЕТСЯ ЛИ ГИПОТЕЗА.

```
In [58]: 1 pearson_criterion(means__10)
2 print()
3 pearson_criterion(means__100)
4 print()
5 pearson_criterion(means__1000)
6 print()
7 pearson_criterion(means__100000)
```

Выборки объема: 10
Выборка: 1
Статистика критерия Пирсона Chi2: 12.948205533236099
Число степеней свободы: 2
Выборка: 2
Статистика критерия Пирсона Chi2: 11.360567891920532
Число степеней свободы: 2
Выборка: 3
Статистика критерия Пирсона Chi2: 7.094003975477621
Число степеней свободы: 2
Выборка: 4
Статистика критерия Пирсона Chi2: 16.794661342915962
Число степеней свободы: 2
Выборка: 5
Статистика критерия Пирсона Chi2: 11.78496970099767
Число степеней свободы: 2

Выборки объема: 100
Выборка: 1
Статистика критерия Пирсона Chi2: 128.36468124889564
Число степеней свободы: 11
Выборка: 2
Статистика критерия Пирсона Chi2: 189.05692009135595
Число степеней свободы: 11
Выборка: 3
Статистика критерия Пирсона Chi2: 178.93265353183912
Число степеней свободы: 10
Выборка: 4
Статистика критерия Пирсона Chi2: 158.92587831338133
Число степеней свободы: 10
Выборка: 5
Статистика критерия Пирсона Chi2: 145.56230697737487
Число степеней свободы: 12

Выборки объема: 1000
Выборка: 1
Статистика критерия Пирсона Chi2: 1365.186737434531
Число степеней свободы: 14
Выборка: 2
Статистика критерия Пирсона Chi2: 1424.212155401219
Число степеней свободы: 15
Выборка: 3
Статистика критерия Пирсона Chi2: 1491.3398510264155
Число степеней свободы: 14
Выборка: 4
Статистика критерия Пирсона Chi2: 1465.854198192775
Число степеней свободы: 14
Выборка: 5
Статистика критерия Пирсона Chi2: 1562.4093704485408
Число степеней свободы: 12

Выборки объема: 100000
Выборка: 1
Статистика критерия Пирсона Chi2: 159925.1065809577
Число степеней свободы: 14
Выборка: 2
Статистика критерия Пирсона Chi2: 156714.40761985185
Число степеней свободы: 15
Выборка: 3
Статистика критерия Пирсона Chi2: 159240.2061818955
Число степеней свободы: 14

4.2.2 Проверка параметрических гипотез

4.2.2.1 Выбор данных

По причине конечности данного на решение задачи времени для ускорения завершения данного домашнего задания выберем распределение Максвелла и рассмотрим две выборки с разными (но известными) параметрами. А именно, к уже сконфигурированным выборкам из распределения с параметром $\lambda = 1.0$, добавим к рассмотрению выборку из распределения с параметром $\lambda = 1.5$. Сгенерируем новые выборки размеров 10, 100, 1000, 10000 и запишем их в файлы по уже отработанной схеме моделирования.

4.2.3 Вычисление функции отношения правдоподобия

$$l(\bar{X}) = \frac{L(\bar{X}, \theta_0)}{L(\bar{X}, \theta_1)} = \frac{\prod_{i=1}^n f_1(x_i)}{\prod_{i=1}^n f_2(x_i)}$$

5 Домашнее задание. Различение гипотез

5.1 Геометрическое распределение

5.1.1 Выбор данных

Необходимо найти два набора данных, соответствующих разным распределениям. Это могут быть оценки к фильмам, статистики игр разных команд и т.п.

В случае наличия сложностей можно выбрать одно из выбранных на первом домашнем задании распределений и рассмотреть две выборки с разными (но известными) параметрами.

В крайнем случае можно рассмотреть распределение Бернулли с разными вероятностями успеха p .

5.1.2 Постановка задачи

Необходимо описать постановку задачи. Что является гипотезой H_0 , что H_1 ? Что такое ошибка первого и второго рода, функция мощности?

5.1.3 Вычисление функции отношения правдоподобия.

Необходимо описать вид функции $l(X)$ отношения правдоподобия, описать способ ее вычисления.

5.1.4 Вычисление критической области/количества материала

Рассмотрим один из самых сложных вопросов данной контрольной работы — вычисление критической области.

Для оценки ошибок первого и второго рода по материалу или вычислению необходимого материала при фиксированных ошибках необходимо знать распределение статистики в случае верности гипотезы $H_0 — I(X|H_0)$ и в случае верности гипотезы $H_1 — I(X|H_1)$. Для большинства распределений это сделать достаточно сложно.

В случае, если не удастся вычислить распределение статистики $I(\bar{X})$ в случае верности разных гипотез, предлагается рассмотреть асимптотический подход к различению гипотез.

Прологарифмировав функцию отношения правдоподобия получим сумму одинаково распределенных независимых случайных величин вида

$$z_i = \ln \frac{f_1(X_i)}{f_2(X_i)}$$

Используя Ц.П.Т. можно легко получить распределение статистики $\ln I(\bar{X})$ в случае верности каждой из гипотез.

Имея две нормально распределенные случайные величины с разными параметрами задача вычисления ошибок первого/второго рода решается легко, как и вычисление минимально необходимого количества материала для достижения нужных ошибок первого и второго рода.

5.2 Распределение Максвелла

5.2.1 Выбор данных

5.2.2 Постановка задачи

5.2.3 Вычисление функции отношения правдоподобия.

5.2.4 Вычисление критической области/количества материала

6 Литература

- [1] "Справочник по вероятностным распределениям" Р.Н.Вадзинский
https://fileskachat.com/view/10838_b741e0be3370efed892ccfe2b6c1358f.html (https://fileskachat.com/view/10838_b741e0be3370efed892ccfe2b6c1358f.html)
- [2] "Введение в математическую статистику" (Ивченко Г.И., Медведев Ю.И.)
<http://bookre.org/reader?file=1221378&pg=101> (<http://bookre.org/reader?file=1221378&pg=101>)
- [3] Power Maxwell distribution:
<https://arxiv.org/pdf/1807.01200.pdf> (<https://arxiv.org/pdf/1807.01200.pdf>)
- [4] Geometric Distribution:
<https://brilliant.org/wiki/geometric-distribution/> (<https://brilliant.org/wiki/geometric-distribution/>)
- [5] The Maxwell Distribution:
<https://randomservices.org/random/special/Maxwell.html> (<https://randomservices.org/random/special/Maxwell.html>)
- [6] The Geometric Distribution:
<https://randomservices.org/random/bernoulli/Geometric.html> (<https://randomservices.org/random/bernoulli/Geometric.html>)
- [7] Sampling from a Normal Distribution
<http://bjlkeng.github.io/posts/sampling-from-a-normal-distribution/> (<http://bjlkeng.github.io/posts/sampling-from-a-normal-distribution/>)
- [8] "Моделирование распределений В.В.Некруткин"
<https://clck.ru/RHRdy> (<https://clck.ru/RHRdy>)
- [9] Сайт со статистикой бейсбола
<https://www.scoreboard.com/ru/baseball/usa/mlb-2020/> (<https://www.scoreboard.com/ru/baseball/usa/mlb-2020/>)

In []:

1