

Лабораторная работа №1

Выполнил: Кузнецов Павел М3207

осваиваю *LaTeX*!)

Задача 1. Пусть выборка X_1, \dots, X_n соответствует классу распределений F_θ , $\theta \in E \subset \mathbb{R}$. При каком минимальном объёме выборки n равномерно для $\theta \in E$ выборочное среднее отличается от математического ожидания μ_θ не более чем на $\varepsilon > 0$ с вероятностью, не меньшей $1 - \delta$, $\delta \in (0, 1)$? Сгенерировать 500 выборок найденного объема при $\varepsilon = 0.01$ и $\delta = 0.05$ из указанного распределения F_θ при конкретном параметре θ и посчитать, сколько раз выборочное среднее отличается от математического ожидания μ_θ более чем на ε .

$F_\theta = \text{Bern}(p)$, $p \in (0, 1)$, $p = 2/3$

Решение:

sum-up: у нас есть некоторая выборка X_1, \dots, X_n , которая соответствует классу распределений $F_\theta = \text{Bern}(p)$, $p \in (0, 1)$, нужно найти $\min n : P(|\bar{X}_n - \mu| \leq \varepsilon) \geq 1 - \delta$

$P(|\bar{X}_n - \mu| \leq \varepsilon) = P\left(\frac{\sqrt{n}|\bar{X}_n - \mu|}{\sigma} \leq \frac{\sqrt{n}\varepsilon}{\sigma}\right) \approx 2\Phi\left(\frac{\sqrt{n}\varepsilon}{\sigma}\right)$ (из центральной предельной теоремы, плотность ст.

нормального распределения симметрична)

$$2\Phi\left(\frac{\sqrt{n}\varepsilon}{\sigma}\right) \geq 1 - \delta$$

$$\Phi\left(\frac{\sqrt{n}\varepsilon}{\sigma}\right) \geq \frac{1-\delta}{2}$$

(так как функция Лапласа нормального распределения возрастает от аргумента, минимальное n будет при минимальном значении Φ)

$$\begin{aligned}\Phi\left(\frac{\sqrt{n}\varepsilon}{\sigma}\right) &= \frac{1-\delta}{2} \\ \frac{\sqrt{n}\varepsilon}{\sigma} &= \Phi^{-1}\left(\frac{1-\delta}{2}\right)\end{aligned}$$

$$n_{min} = \left(\frac{\sigma \Phi^{-1}(\frac{1-\delta}{2})}{\varepsilon} \right)^2$$

Найдём конкретные значения n_{min} при заданных значениях $\theta, \varepsilon, \delta$:

```

1  import numpy
2  from scipy.stats import norm
3
4  #mathematical expectation
5  p = 2/3
6
7  eps = 0.01
8  delta = 0.05
9
10 #standard deviation from the Bernoulli distribution
11 sigma = numpy.sqrt(p * (1 - p))
12
13 n = pow(sigma*norm.ppf(1 - delta/2)/eps,2)
14 n_min = round(n)
15 print(n_min)

```

Получившееся значение $n_{min} = 8537$.

Сгенерируем 500 выборок найденного объема при $\varepsilon = 0.01$ и $\delta = 0.05$ из указанного распределения F_p при конкретном параметре p и посчитаем, сколько раз выборочное среднее отличается от математического ожидания μ_p более чем на ε :

```

1  for i in range(500):
2      X = numpy.random.binomial(n, p)/n
3      if (abs(X - p) > eps):
4          count += 1

```

Вывод:

Выборочное среднее отличается от мат. ожидания более чем на epsilon 22 раз. Это соответствует 4.4 % от 500 выборок.

Задача 2-4. В файле `mobile_phones.csv` приведены данные о мобильных телефонах. В сколько моделей можно вставить 2 сим-карты, сколько поддерживают 3G, каково наибольшее число ядер у процессора? Рассчитайте выборочное среднее, выборочную дисперсию, выборочную медиану и выборочную квантиль порядка 2/5, построить график эмпирической функции распределения, гистограмму и box-plot для емкости аккумулятора для всей совокупности и в отдельности для поддерживающих/не поддерживающих Wi-Fi.

Решение: (код)

```
1 import numpy
2 from scipy.stats import norm
3 import pandas
4 import matplotlib.pyplot as plt
5 import statsmodels.api as statsm
```

Листинг 1: Imports

```
1
2 data = pandas.read_csv("/Users/frogge/proggs/mathstatLabs/lab
3 # how many models are available to install 2 sim cards
4 print(data['dual_sim'].sum())
5
6 # how many models support 3-G
7 print(data['three_g'].sum())
8
9 # maximum number of cores
10 print(data['n_cores'].max())
```

Получившиеся значения:

available to install 2 sim cards = 1019

support 3-G = 1523

maximum numb of cores = 8

Далее считаем значения (указанные в задании), строим график, гистограмму и box-plot для емкостей аккумулятора при разных выборках:

```
1 new_data = data['battery_power']
2
3 # selective mean
4 print(round(new_data.mean(), 4))
5
6 # dispersion
7 print(round(new_data.var(), 4))
8
9 # median
10 print(round(new_data.median(), 4))
11
12 # quantile 2/5
13 print(round(new_data.quantile(q = 0.4), 4))
14
15 x = numpy.linspace(min(new_data), max(new_data))
16 y = statsm.distributions.ECDF(new_data)(x)
17 plt.step(x, y)
18 plt.title("Graph of ecdf")
19 plt.show()
20 #plt.hist(new_data, histtype='step', cumulative=True, bins=len(
    sample)) - another way to plot a graph
21
22 plt.hist(new_data)
23 plt.title("Histogram")
24 plt.show()
25
26 plt.boxplot(new_data)
27 plt.title("Box-plot")
28 plt.show()
```

Листинг 2: for all data

Получившиеся значения:

Выборочное среднее = 1238.5185

Выборочная дисперсия = 193088.3598

Выборочная медиана = 1226.0

Выборочная квантиль порядка $2/5 = 1076.0$

Графики в приложении

Далее считаем для выборки из даты, только для моделей поддерживающих WI-FI и наоборот, код остаётся прежним меняем только значение *new_data*

```
1 new_data = data[data['wifi'] == 1]['battery_power'] #with wf
2 new_data = data[data['wifi'] == 0]['battery_power'] #without wf
```

Получившиеся значения:

Выборочное среднее with wf = 1234.9043

Выборочная дисперсия with wf= 190296.4005

Выборочная медиана with wf= 1233.0

Выборочная квантиль порядка 2/5 with wf = 1077.8

Выборочное среднее without= 1242.2353

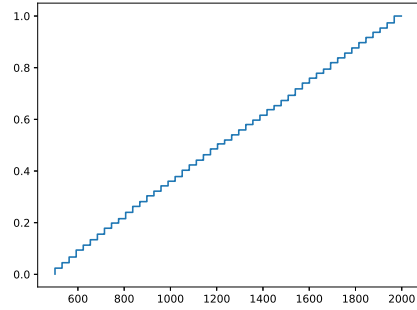
Выборочная дисперсия without= 196128.438

Выборочная медиана without= 1222.0

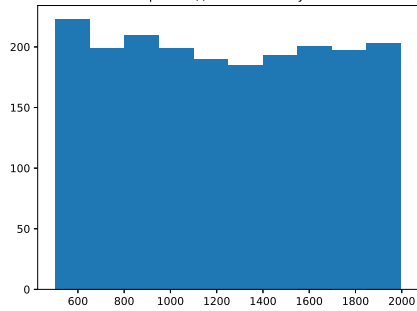
Выборочная квантиль порядка 2/5 without= 1076.8

Приложение

График эмпирической функции распределения для всей совокупности



Гистограмма для всей совокупности



Box-plot для всей совокупности

