

```
In [ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import random as rnd
import re
import sklearn
from sklearn.tree import DecisionTreeClassifier
from collections import defaultdict

104

In [ ]: df = pd.read_csv("C:\Users\Мой ПК (LAPTOP-TOID7NG7)\Downloads\122-23_new_probability.csv", sep = ",")
df = df[['Isamen', 'Last name']]
df.fillna(0, inplace=True)
df

Out[ ]:
   Isamen  Last name
0      4.0  Переноска
1      0.0  Роднярина
2      5.0  Сафина
3      9.0  Садырна
4      6.0  Соколов
...      ...      ...
332     0.0  Нахичевань
333     0.0  Нахичевань
334     0.0  Герасимов
335     0.0  Табадекурузы
336     0.0  Ханомов
...      ...      ...
337 rows x 2 columns

In [ ]: group1 = df[df['Last name'].str[0].isin(['А', 'В', 'Г', 'Д', 'Е', 'Ж', 'З', 'И', 'Й'])['Isamen']
group2 = df[df['Last name'].str[0].isin(['К', 'Л', 'М', 'Н', 'О', 'П', 'Р', 'С', 'Т', 'У', 'Ф', 'Х', 'Ц', 'Ч', 'Ш', 'Щ', 'Ъ', 'Ы', 'Э', 'Ю', 'Я'])['Isamen']

a) [5] Итоговая таблица тест-Узна

In [ ]: import pandas as pd
import numpy as np
from scipy import stats

mean1 = group1.mean()
mean2 = group2.mean()
print("Средняя оценка для группы 1:", mean1)
print("Средняя оценка для группы 2:", mean2)

t_stat, p_value = stats.ttest_ind(group1, group2)
print("p-значение:", p_value)
if p_value < 0.05:
    print("Отвергаем нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, не равны.")
else:
    print("Мы не можем отвергнуть нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, равны.")

Средняя оценка для группы 1: 4.611111111111111
Средняя оценка для группы 2: 4.989722222222222
p-значение: 0.5180855951057487
Мы не можем отвергнуть нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, равны.

b) [5] Итоговая таблица тест-Бутстреп

In [ ]: n_bootstrap = np.mean(group1)
n_bootstrap = np.mean(group2)
obs_diff = n_bootstrap1 - n_bootstrap2
n_iter = 1000
boot_diffs = []
for i in range(n_iter):
    np.random.seed(i)
    boot_bootstrap1 = np.random.choice(group1, size=len(group1), replace=True)
    boot_bootstrap2 = np.random.choice(group2, size=len(group2), replace=True)
    boot_diff = np.mean(boot_bootstrap1) - np.mean(boot_bootstrap2)
    boot_diffs.append(boot_diff)

p_value = np.sum(np.abs(boot_diffs) >= np.abs(obs_diff)) / n_iter

print("Среднее значение ожидаемых результатов экзамена у тех, у кого фамилия начинается с гласной буквы:", n_bootstrap1)
print("Среднее значение ожидаемых результатов экзамена у тех, у кого фамилия начинается с согласной буквы:", n_bootstrap2)
print("Разница средних значений:", obs_diff)
if p_value < 0.05:
    print("Отвергаем нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, не равны.")
else:
    print("Мы не можем отвергнуть нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, равны.")

Среднее значение ожидаемых результатов экзамена у тех, у кого фамилия начинается с гласной буквы: 4.611111111111111
Среднее значение ожидаемых результатов экзамена у тех, у кого фамилия начинается с согласной буквы: 4.989722222222222
Разница средних значений: 0.2986111111111111
p-значение: 0.627
Мы не можем отвергнуть нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, равны.

c) [5] Итоговая таблица тест-Т-статистики

In [ ]: from scipy.stats import ttest_ind
n_bootstrap1 = np.mean(group1)
n_bootstrap2 = np.mean(group2)
# вычисление разницы средних значений между группами
diff_mean = n_bootstrap1 - n_bootstrap2
# создание массива для бутстрапирования
bootstrap_diff = np.empty(1000)
# бутстрапирование
for i in range(1000):
    bootstrap_group1 = np.random.choice(group1, size=len(group1), replace=True)
    bootstrap_group2 = np.random.choice(group2, size=len(group2), replace=True)
    bootstrap_diff[i] = np.mean(bootstrap_group1) - np.mean(bootstrap_group2)
# вычисление p-значения
p_value = np.sum(bootstrap_diff >= diff_mean) / len(bootstrap_diff)
print("p-значение бутстрапа:", p_value)
# вывод результата
if p_value < 0.05:
    print("Отвергаем нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, не равны.")
else:
    print("Мы не можем отвергнуть нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, равны.")

p-значение бутстрапа: 0.687
Мы не можем отвергнуть нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, равны.

d) [5] Итоговая таблица перестановочный тест

In [ ]: vowel_results = group1.to_list()
consonant_results = group2.to_list()
def mean(lst):
    return sum(lst)/len(lst)
vowel_mean = mean(vowel_results)
consonant_mean = mean(consonant_results)
# определение функции для вычисления статистической значимости
def permutation_test(lst1, lst2, nperm):
    obs_diff = mean(lst1) - mean(lst2)
    pool_diff = lst1 - lst2
    count = 0
    for i in range(nperm):
        random_shuffle(pool_diff)
        perm_diff = mean(random_shuffle(pool_diff)) - mean(random_shuffle(pool_diff))
        if perm_diff >= obs_diff:
            count += 1
    return count/nperm
# вычисление статистической значимости
p_value = permutation_test(vowel_results, consonant_results, 10000)
print("p-значение перестановочного теста:", p_value)
# вывод результата
if p_value < 0.05:
    print("Отвергаем нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, не равны.")
else:
    print("Мы не можем отвергнуть нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, равны.")

p-значение перестановочного теста: 0.7425
Мы не можем отвергнуть нулевую гипотезу: ожидаемые результаты экзамена по теории вероятностей у студентов, чья фамилия начинается с гласной и с согласной буквы, равны.

105

In [ ]: gr1 = df[df['Last name'].str[0].isin(['А', 'В', 'Г', 'Д', 'Е', 'Ж', 'З', 'И', 'Й'])['Last name'].to_list()]
gr2 = df[df['Last name'].str[0].isin(['К', 'Л', 'М', 'Н', 'О', 'П', 'Р', 'С', 'Т', 'У', 'Ф', 'Х', 'Ц', 'Ч', 'Ш', 'Щ', 'Ъ', 'Ы', 'Э', 'Ю', 'Я'])['Last name'].to_list()]
median = df.median()

g1_bot = 0
g1_top = 0
sig_bot = 0
sig_top = 0

for index, row in df.iterrows():
    if row['Last name'] in gr1:
        if row['Isamen'] >= median.bot():
            g1_bot += 1
        else:
            g1_top += 1
    if row['Isamen'] >= median.bot():
        sig_bot += 1
    else:
        sig_top += 1

print("Фамилии на согласной и меньше медианы", sig_bot)
print("Фамилии на согласной и больше медианы", sig_top)
print("Фамилии на гласной и меньше медианы", g1_bot)
print("Фамилии на гласной и больше медианы", g1_top)

Фамилии на согласной и меньше медианы 324
Фамилии на согласной и больше медианы 377
Фамилии на гласной и меньше медианы 10
Фамилии на гласной и больше медианы 336
C:\Users\Мой ПК\AppData\Local\Temp\ipykernel_22580\3507167334.py:4: FutureWarning: The default value of numeric_only in DataFrame.median is deprecated. In a future version, it will default to False. In addition, specifying "numeric_only=True" is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.
    median = df.median()

a) [5] Постройте 95% асимптотический интервал для отношения шансов хорошо написать экзамен («несогласные» к «согласным»). Проверьте гипотезу о том, что отношение шансов равно 1 и укажите P-значение

In [ ]: import numpy as np
from scipy.stats import norm

n1, n2 = sig_bot, sig_top
m1, m2 = g1_bot, g1_top

p1 = n1 / n1
p2 = n2 / n2

OR = (p2 / (1 - p2)) / (p1 / (1 - p1))

SE_logOR = np.sqrt(1 / n2 + 1 / (n2 - n2) + 1 / n1 + 1 / (n1 - n1))
CI_logOR = np.log(OR) - norm.ppf(0.975) * SE_logOR, np.log(OR) + norm.ppf(0.975) * SE_logOR
CI_OR = np.exp(CI_logOR)

print("95% interval for log(OR):", CI_logOR)
print("95% interval for OR:", CI_OR)
print("OR =", OR)

p_value = 2 * (1 - norm.cdf(np.abs(np.log(OR)) / SE_logOR))

print("P-value =", p_value)
if p_value < 0.05:
    print("Отвергаем нулевую гипотезу: То есть, можно сказать, что отношение шансов не равно 1, и студенты, чья фамилия начинается на гласную, имеют значительно больше шансов написать экзамен хорошо, чем те, у кого фамилия начинается на согласную.")
else:
    print("Мы не можем отвергнуть нулевую гипотезу: То есть, можно сказать, что отношение шансов равно 1, и студенты, чья фамилия начинается на гласную, имеют одинаковые шансы написать экзамен хорошо, как и те, у кого фамилия начинается на согласную.")
```

```
95% interval for ln(OR): (-8.4382984064620175, 8.57737151295964)
OR = 1.1128819570556
p-value = 8.43462581528952
We cannot start with a null hypothesis: To eat, it is better to say that the probability of getting a good grade is 1, than to say that the probability of getting a bad grade is 1.
6) [5] Постройте 95% доверительный интервал для отношения вероятностей хорошо написать экзамен. Проверьте гипотезу о том, что отношение вероятностей равно 1 и укажите P-значение
```

```
In [ ]: from statsmodels.stats import proportion
# because these arguments, maximum number of iterations, 8 and 0.001
p_consistent = log_bol / (log_bol + log_mom)
p_wilson = (g1_bol) / (g1_bol + g1_mom)
# using the Wilson method for the proportion
lb, ub = proportion.proportion_confint(log_bol, log_mom, alpha=0.05, method='wilson')
# the Wilson method
print(f'proportion consistent: {p_consistent:.2f}')
print(f'proportion wilson: {p_wilson:.2f}')
print(f'CI for proportion difference (using wilson method):', lb, ub)
proportion consistent: 0.50
proportion wilson: 0.53
CI for proportion difference (using wilson method): [0.5166885 0.57885938] [0.54219281 0.68813958]
```

```
In [ ]: from statsmodels.stats.proportion import proportions_test
count = np.array([g1_bol, log_mom])
nobs = np.array([log_bol + log_mom, log_bol + log_mom])
stat, pval = proportions_test(count, nobs, value=0, alternative='two-sided', prop_var=proportion)
print(f'z-statistic: {stat:.2f}')
print(f'p-value: {pval:.2f}')
z-statistic: -0.09
p-value: 0.49
```

```
In [ ]: print("Разница между долями не является статистически значимой, т.е. p-value > 0.05")
# the difference between the two proportions is not statistically significant, i.e. p-value > 0.05
[5] Постройте 95% интервал для отношения шансов хорошо написать экзамен с помощью названного бутстрапа. Проверьте гипотезу о том, что отношение шансов равно 1 и укажите P-значение.
```

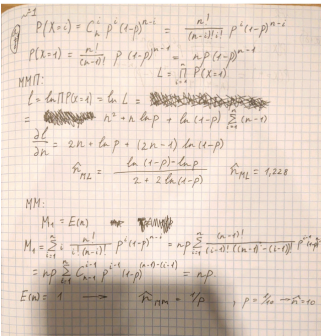
```
In [ ]: n = 1000
ix_above_median = [random.choice([0, 1], weights=[0.4, 0.6]) for _ in range(n)]
initials = [random.choice('a', 'z'), random.choice('a', 'z')] for _ in range(n)
exam_results = [random.choice([0, 1], weights=[0.5, 0.5]) for _ in range(n)]
In [ ]: data = pd.DataFrame({'ix_above_median': ix_above_median, 'initials': initials, 'exam_results': exam_results})
table = pd.crosstab(data['ix_above_median'], data['initials'], normalize='all')
table.iloc[0, 0] = 177
table.iloc[0, 1] = 124
table.iloc[1, 0] = 39
table.iloc[1, 1] = 17
```

```
print(table)
initials
ix_above_median
a      177.0  124.0
z       39.0   17.0
OR = (table.iloc[0, 0] * table.iloc[1, 1]) / (table.iloc[0, 1] * table.iloc[1, 0])
print(f'Отношение шансов: ', OR)
Отношение шансов: 1.277464850881102
```

```
In [ ]: n_bootstrap = 1000
OB_bootstrap = np.zeros(n_bootstrap)
for i in range(n_bootstrap):
    random.seed(i)
    sample = data.sample(n=len(data), replace=True)
    table_bootstrap = pd.crosstab(sample['ix_above_median'], sample['initials'], normalize='all')
    OB_bootstrap[i] = (table_bootstrap.iloc[0, 0] * table_bootstrap.iloc[1, 1]) / (table_bootstrap.iloc[0, 1] * table_bootstrap.iloc[1, 0])
CI = np.percentile(OB_bootstrap, q=[2.5, 97.5])
print(f'Доверительный интервал для отношения шансов (95%): {CI[0]:.2f}, {CI[1]:.2f}')
Доверительный интервал для отношения шансов (95%): [0.925789 1.5495815]
```

```
In [ ]: SE = np.sqrt(np.var(OB_bootstrap))
z = (OR - 1) / SE
print(f'z-критерий: ', z)
z-критерий: 1.686880235143467
```

```
In [ ]: from scipy.stats import norm
p_value = 2 * norm.sf(abs(z))
print(f'P-значение: ', p_value)
P-значение: 0.092849128565641
Таким образом, мы можем отвергнуть нулевую гипотезу о том, что отношение шансов равно 1, на уровне значимости 0.05.
NOT
[5] [6]
```



4) [15] Предположим, что частоты p равно 100. Проведем 10000 симуляций вызовов такси до первого повторного, рассчитаем 10000 оценок методом моментов и 10000 оценок методом максимального правдоподобия. Постройте гистограммы для оценок двух методов. Оцените смещение, дисперсию и среднеквадратичную ошибку двух методов

```
In [ ]: # generating the 10000 simulations
n = 100
simulations = []
for i in range(10000):
    taxi_arrived = []
    found = False
    while not found:
        taxi = random.randint(1, n)
        if taxi in taxi_arrived:
            found = True
        taxi_arrived.append(taxi)
    simulations.append(len(taxi_arrived))

# method of moments
moment_estimates = [n*(np.mean(simulations)/np.mean(np.arange(1, n+1)))**2]*10000

# maximum likelihood method
def log_likelihood(p, n, n1):
    return n*log(1-p) + np.sum(np.log(np.arange(n+1, n+1)))
def maximize_likelihood(n, n1):
    return n*log(1-p) + np.sum(np.log(np.arange(n+1, n+1)))
MLE_estimates = [maximize_likelihood(n, n1) for n1 in simulations]

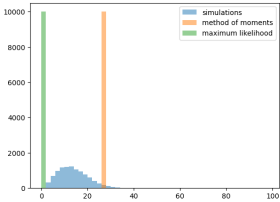
# histogram
bins = np.arange(0, 100, 2)
plt.figure(figsize=(10, 5))
plt.hist(simulations, bins=bins, alpha=0.5, label='simulations')
plt.hist(moment_estimates, bins=bins, alpha=0.5, label='method of moments')
plt.hist(MLE_estimates, bins=bins, alpha=0.5, label='maximum likelihood')
plt.legend()
plt.show()

# bias
bias_moments = np.mean(moment_estimates) - n
bias_MLE = np.mean(MLE_estimates) - n

# variance
var_moments = np.var(moment_estimates)
var_MLE = np.var(MLE_estimates)

# mean squared error
MSE_moments = (np.mean(moment_estimates) - n)**2
MSE_MLE = (np.mean(MLE_estimates) - n)**2

print('Method of moments:')
print('Bias:', abs(bias_moments))
print('Variance:', var_moments)
print('Mean Squared Error:', MSE_moments)
print('Maximum Likelihood:')
print('Bias:', abs(bias_MLE))
print('Variance:', var_MLE)
print('Mean Squared Error:', MSE_MLE)
```



```

import random
import numpy as np
import matplotlib.pyplot as plt

# true parameter
n = 20

# generate the true population
population = [random.choice(['A', 'B', 'C', 'D', 'E', 'F']) for _ in range(n)]

# method of moments estimation
def mom_estimation(sample):
    return len(set(sample))

# maximum likelihood estimation
def mla_estimation(sample):
    counts = dict.fromkeys(['A', 'B', 'C', 'D', 'E', 'F'], 0)
    for name in sample:
        counts[name] += 1
    return max(counts.values())

# simulate 10000 experiments
n_experiments = 10000
n_calls = 10
mom_estimates = np.zeros(n_experiments)
mla_estimates = np.zeros(n_experiments)
for i in range(n_experiments):
    # generate the sample
    sample = [random.choice(population) for _ in range(n_calls)]
    # calculate the estimates
    mom_est = mom_estimation(sample)
    mla_est = mla_estimation(sample)
    # store the estimates
    mom_estimates[i] = mom_est
    mla_estimates[i] = mla_est

# plot the histograms of the estimates
fig, axs = plt.subplots(1, 2, figsize=(10, 4))
axs[0].hist(mom_estimates, bins=range(1, 7))
axs[0].set_title('Method of Moments')
axs[0].set_xlabel('Number of unique names')
axs[0].set_ylabel('Frequency')
axs[1].hist(mla_estimates, bins=range(1, 7))
axs[1].set_title('Maximum Likelihood')
axs[1].set_xlabel('Number of most frequent name')
axs[1].set_ylabel('Frequency')
plt.show()

# calculate the bias, variance, and mean squared error of the estimates
mom_bias = np.mean(mom_estimates) - n
mla_bias = np.mean(mla_estimates) - n
mom_var = np.var(mom_estimates)
mla_var = np.var(mla_estimates)
mom_mse = mom_bias**2 + mom_var
mla_mse = mla_bias**2 + mla_var

print('Method of Moments:')
print('Bias:', mom_bias)
print('Variance:', mom_var)
print('MSE:', mom_mse)

print('Maximum Likelihood:')
print('Bias:', mla_bias)
print('Variance:', mla_var)
print('MSE:', mla_mse)

```



Method of Moments	Maximum Likelihood
Bias: -15.4862	Bias: -15.4862
Variance: 8.434895000000001	Variance: 8.434895000000001
MSE: 240.1108	MSE: 240.1108

```

# функция для построения доверительного интервала на основе бутстреп t-статистики
def conf_int_t_bootstrap(data, boot_reps=1000, alpha=0.05):
    n = len(data)
    boot_tstats = []
    for rep in range(boot_reps):
        boot_sample = np.random.choice(data, size=n, replace=True)
        boot_mean = np.mean(boot_sample)
        boot_std = np.std(boot_sample, ddof=1)
        boot_se = boot_std / np.sqrt(n)
        boot_tstat = (boot_mean - intensity) / boot_se
        boot_tstats.append(boot_tstat)
    left = intensity - np.percentile(boot_tstats, 100 * (1 - alpha) / 2) * np.std(data, ddof=1) / np.sqrt(n)
    right = intensity + np.percentile(boot_tstats, 100 * alpha / 2) * np.std(data, ddof=1) / np.sqrt(n)
    return left, right

# параметры
n = 20
alpha = 0.05
intensity = 1
reps = 10000

# генерация данных и оценка бутстрепом
np.random.seed(42)
prob_cover = 0
for rep in range(reps):
    data = np.random.exponential(scale=intensity, size=n)
    left, right = conf_int_t_bootstrap(data, alpha=alpha)
    if left < intensity < right:
        prob_cover += 1
print('Процент попаданий в доверительный интервал для бутстрапа t-статистики:', prob_cover)
Вероятность попадания в доверительный интервал для бутстрапа t-статистики: 0.982
0

```

```

import numpy as np
from scipy.stats import t
import bootstrap.bootstrap as bs
import bootstrap.bootstrap_functions as bs_funcs

# исходные данные
np.random.seed(42)
n = 20
data = np.random.normal(loc=15, scale=5, size=n)

# классический доверительный интервал
alpha = 0.05
mean = np.mean(data)
std = np.std(data, ddof=1)
n_ad = t.ppf(1 - alpha / 2, n - 1)
ci_low = mean - n_ad * std / np.sqrt(n), mean + n_ad * std / np.sqrt(n)
print('Классический доверительный интервал: [(ci_low, -2F), (ci_high, -2F)]')
Классический доверительный интервал: [12.83, 15.97]

# асимптотический нормальный доверительный интервал
alpha = t.ppf(alpha / 2, n - 1)
t = (np.mean(data) - 15) / (np.std(data, ddof=1) / np.sqrt(n))
ci_low asympt = np.mean(data) - alpha * np.std(data, ddof=1) / np.sqrt(n)
ci_high asympt = np.mean(data) + alpha * np.std(data, ddof=1) / np.sqrt(n)
print('Асимптотический нормальный доверительный интервал: [(ci_low asympt, -2F), (ci_high asympt, -2F)]')
Асимптотический нормальный доверительный интервал: [12.83, 15.97]

# бутстреп
resampled = np.random.choice(data, size=(n, n))
means = np.mean(resampled, axis=1)
ci_low bs = bs.bootstrap(means, stat_funcs, stats.mean, alpha=alpha, num_iterations=10000)
print('Доверительный интервал с помощью бутстрапа t-статистики: [(ci_low bs, lower_bound, -2F), (ci_high bs, upper_bound, -2F)]')
Доверительный интервал с помощью бутстрапа t-статистики: [14.28, 14.85]

# бутстреп t-статистики
resampled = np.random.choice(data, size=(n, n))
means = np.mean(resampled, axis=1)
t_values = (means - np.mean(data)) / (np.std(data, ddof=1) / np.sqrt(n))
ci_low t = bs.bootstrap(t_values, stat_funcs, stats.mean, alpha=alpha, num_iterations=10000)
ci_high t = np.percentile(t_values, 100 * alpha / 2)
ci_low t_corrected = np.percentile(t_values, 100 * (1 - alpha) / 2)
ci_high t_corrected = np.percentile(t_values, 100 * alpha / 2)
print('Доверительный интервал с помощью бутстрапа t-статистики: [(ci_low t_corrected, -2F), (ci_high t_corrected, -2F)]')
Доверительный интервал с помощью бутстрапа t-статистики: [14.29, 14.85]

# бутстреп t-статистики
resampled = np.random.choice(data, size=(n, n))
means = np.mean(resampled, axis=1)

```

```
[t_values, (mean - np.mean(data)) / np.std(data, ddof=1)
ci_ba_t = ba.bootstrap(t_values, stat_func=ba_stat.mean, alpha=alpha, num_iterations=10000)
se_ba_t = np.std(t_values, ddof=1)
ci_ba_t_corrected = np.mean(data) - ci_ba_t_upper_bound + se_ba_t, np.mean(data) - ci_ba_t_lower_bound + se_ba_t
print(f"Доверительный интервал с помощью бутстрапа t-статистики: [{ci_ba_t_corrected[0]:.2f}, {ci_ba_t_corrected[1]:.2f}]")

Доверительный интервал с помощью бутстрапа t-статистики: [14.18, 14.42]

In [ ]: def coverage_probability(ci, mu):
    alpha = 0.05
    n = 20
    df = n - 1
    t_alpha2 = 1.ppf(1 - alpha / 2, df)
    t_mu = (mu - np.mean(data)) / (np.std(data, ddof=1) / np.sqrt(n))
    t_1 = (ci[0] - mu) / (np.std(data, ddof=1) / np.sqrt(n)), (ci[1] - mu) / (np.std(data, ddof=1) / np.sqrt(n))
    p = 1 - df(t_1, df) + 1.cdf(t_1, df) / (1 - alpha)
    return p

# классический доверительный интервал
p_classic = coverage_probability(ci, 15)
print(f"Вероятность накрытия классического доверительного интервала: {p_classic:.4f}")

# асимптотический нормальный доверительный интервал
p_asympt = coverage_probability(ci_asympt, 15)
print(f"Вероятность накрытия асимптотического нормального доверительного интервала: {p_asympt:.4f}")

# модный бутстрап
p_ba = coverage_probability(ci_ba_lower_bound, ci_ba_upper_bound, 15)
print(f"Вероятность накрытия доверительного интервала с помощью модного бутстрапа: {p_ba:.4f}")

# бутстрап t-статистики
p_ba_t = coverage_probability(ci_ba_t_corrected, 15)
print(f"Вероятность накрытия доверительного интервала с помощью бутстрапа t-статистики: {p_ba_t:.4f}")

Вероятность накрытия классического доверительного интервала: 0.9367
Вероятность накрытия асимптотического нормального доверительного интервала: 0.9367
Вероятность накрытия доверительного интервала с помощью модного бутстрапа: 0.9859
Вероятность накрытия доверительного интервала с помощью бутстрапа t-статистики: 0.9164

Самый лучший интервал с помощью асимптотического нормального доверительного интервала

NPG

0)

In [ ]: df = pd.read_csv("C:/Users/Мой ПК/Dropbox/Мой ПК (LAPTOP-TOID7NG7)/Downloads/22-23_hse_probability.csv", sep = ";")
df = df[["Знамен", "Last name"]]
df.fillna(0, inplace=True)
df

Out[ ]:
Знамен    Last name
0         4.0    Речников
1         0.0    Родионова
2         5.0    Сафина
3         9.0    Сахаров
4         6.0    Соколов
...
332        0.0    Наумовича
333        0.0    Хайкина
334        0.0    Герасимов
335        0.0    Тахирбекулов
336        0.0    Кичинкина
337 rows x 2 columns

In [ ]: # импортируем библиотеку pandas для работы с данными в таблице
import pandas as pd

df = df.copy()
# считаем количество дней в каждой фамилии
df["last name"] = df["last name"].apply(len)
df

Out[ ]:
Знамен    Last name
0         4.0         9
1         0.0         9
2         5.0         6
3         9.0         7
4         6.0         8
...
332        0.0        11
333        0.0         7
334        0.0         9
335        0.0        14
336        0.0         9
337 rows x 2 columns

In [ ]: EF = df["last name"].sum() / 337
EV = df[["Знамен", "Last name"]].sum() / 337
print(EF, EV)
7.891738074183977 4.838888514224629
p = EV/EF

In [ ]: b = EV/EF
b

Out[ ]: 0.6130388751287959
cov(X, Y) = E[(Y - EV)(X - EF)]

In [ ]: df["cov"] = (df["Знамен"] - EV) * (df["Last name"] - EF)
sum_cov = df["cov"].sum()

Out[ ]: -1.88982872513347
x_bar^2 = E[(Y - EV)^2 | h]
x_bar^2 = E[(X - EF)^2 | h]

In [ ]: df["x_bar^2"] = (df["Знамен"] - EV) ** 2
df["x_bar^2"] = (df["Last name"] - EF) ** 2
x_bar = df["x_bar^2"].sum() / 336
x_bar = df["x_bar^2"].sum() / 336
print(x_bar, x_bar)
7.259957173943762 3.3337925678728937
r = cov(X, Y) / (x_bar * x_bar)

In [ ]: r = sum_cov / (x_bar * x_bar)
r

Out[ ]: -0.8469467031678548
0)

In [ ]: from scipy.stats import pearsonr

a = df["Знамен"].to_list()
b = df["Last name"].to_list()
# вычисление фактической корреляции
corr, p = pearsonr(a, b)

# создание массива для сохранения корреляций при перестановках
perms = np.zeros(10000)

# вычисление перестановочного теста
for i in range(10000):
    # случайная перестановка для функций
    permuted_lengths = np.random.permutation(b)
    # вычисление корреляции для переставленных для функций
    permuted_corr = pearsonr(a, permuted_lengths)
    # сохранение корреляции
    perms[i] = permuted_corr

# вычисление p-значения
p_value = np.sum(perms > corr) / len(perms)

# вывод результатов
print("Fact correlation:", corr)
print("p-value:", p_value)
if p_value < 0.05:
    print("Случайная нулевая гипотеза: здесь есть корреляция")
else:
    print("Мы не можем отвергнуть нулевую гипотезу: здесь нету корреляции")

Fact correlation: -0.88885888877282936
p-value: 0.3818
Мы не можем отвергнуть нулевую гипотезу: здесь нету корреляции

NPG

Указание задачи:

В первой урне 7 белых и 3 черных шара, во второй урне 8 белых и 4 черных шара, в третьей урне 2 белых и 13 черных шаров. Из этих урн наугад выбирается одна урна. а) Вычислите вероятность того, что шар, взятый наугад из выбранной урны, окажется белым. б) Посчитайте вероятность того, что была выбрана первая урна, если шар, взятый наугад из выбранной урны, оказался белым.
```

```
Решение птг-4 чзгг:
а) Обозначим событие выбора первой, второй и третьей урны соответственно через А, В и С. Тогда общая вероятность каждого события равна 1/3, так как у нас есть 3 урны, из которых нужно выбрать одну.
Вероятность того, что белый шар будет выбран из первой урны:
P(белый шар из 1 урны) = 7/10
Вероятность того, что белый шар будет выбран из второй урны:
P(белый шар из 2 урны) = 8/12 = 2/3
Вероятность того, что белый шар будет выбран из третьей урны:
P(белый шар из 3 урны) = 2/15
Тогда общая вероятность выбора белого шара
P(белый шар) = P(A)*P(белый шар из 1 урны) + P(B)*P(белый шар из 2 урны) + P(C)*P(белый шар из 3 урны) = (1/3)*(7/10) + (1/3)*(2/3) + (1/3)*(2/15) = 0.5 = 50%
б) Определим вероятность выбора первой урны при условии, что был выбран белый шар:
Событие А – выбор первой урны. Событие В – выбран белый шар.
Общая вероятность выбора белого шара:
P(B) = P(A)*P(белый шар из 1 урны) + P(B)*P(белый шар из 2 урны) + P(C)*P(белый шар из 3 урны) = (1/3)*(7/10) + (1/3)*(2/3) + (1/3)*(2/15) = 0.5 = 50%
Вероятность выбора первой урны при условии наступления события В:
P(A|B) = P(A*В) / P(B)
P(A*В) – вероятность того, что выбрана первая урна и шар из нее оказался белым.
Из вероятности условия В известно, что был выбран белый шар. Значит, интересующие нас событие – это выбор первой урны и выбор белого шара из нее.
P(A*В) = P(белый шар из 1 урны) * P(A) = (7/10)*(1/3) = 7/30
P(B) известна и равна 0.5.
Тогда искомая вероятность:
P(A|B) = (7/30) / 0.5 = 0.4667 = 46.67%
Таким образом, вероятность того, что была выбрана первая урна, если шар, взятый наугад из выбранной урны, оказался белым, составляет 46.67%.
а) Вероятность выбора первой урны равна 1/3. Вероятность выбора белого шара из первой урны равна 7/10. Вероятность выбора второй урны равна также 1/3. Вероятность выбора белого шара из второй урны равна 8/12, что упрощается до 2/3. Вероятность выбора третьей урны также равна 1/3. Вероятность выбора белого шара из третьей урны равна 2/15. Тогда общая вероятность выбора белого шара будет равна:
(1/3)*(7/10) + (1/3)*(2/3) + (1/3)*(2/15) =
б) Вероятность выбора первой урны и выбора белого шара из нее равна (1/3)*(7/10) = 7/30 (это мы уже посчитали в пункте а). Вероятность выбора белого шара из любой урны равна 11/30 (это мы также посчитали в пункте а). Тогда по формуле Байеса.
P(выбрана 1-я урна | белый шар) = P(белый шар | выбрана 1-я урна) * P(выбрана 1-я урна) / P(белый шар)
P(белый шар) = (1/3)*(7/10) + (1/3)*(2/3) + (1/3)*(2/15) =
P(выбрана 1-я урна | белый шар) = (7/10)*(1/3) / (11/30) =
NB
Часто всего я не могу справиться на лекциях, поэтому в дополнение к просмотру семинаров, я смотрю в записки видео прошлых лет на канале https://www.youtube.com/@mathmeth1995. Где видео снимал Тильман Николай Петрович. Я также пользовался материалами, которые предоставляли преподаватели по анализу данных по информатике для понимания того, как применять пробранный материал в сфере программирования (https://github.com/hee-econ-data-science/indat\_2023).
Иногда смотрю некоторые пояснения в открытом учебном пособии "Математическая статистика", написанным Н.И.Черновой (https://phs.msu.ru/~chernova/mst/mst\_mio07.pdf).
И наконец, игрался с визуализированными данными на примере простых задач через сайт (https://seeing-theory.brown.edu/). Это было полезно, так как, например, я понял, как по графику определить корреляцию и её силу (помогло во время прорешивания МСР по майнору "математико-статистический анализ?")
```