

Задача 3

Вариант 1

от Татаринова Никиты Алексеевича

к 17.03.2021

Условие

В былые времена возле станции метро "Площадь Революции" проводились различные обследования (опросы и дегустации), за участие в которых можно было получить коробку конфет, шоколадку, бутылку пива или даже торт. Представьте себя на месте добрых людей, раздающих эти продукты.

Производитель газированной воды, планирующий продвижение нового товара (вода A), заказал вам малое обследование потенциальных потребителей. Вы уже собрали нужные данные и записали их в файл "Данные к задаче 3.ods":

- a – общая оценка респондентом воды A по семибалльной шкале (1 – совсем не понравилось, 7 – превосходно);
- b – оценка респондентом воды B (предполагаемого конкурента);
- sex – пол респондента (0 – мужской, 1 – женский).

Требуется ответить на 2 вопроса.

1. Есть ли основание считать, что предполагаемый потребитель предпочитает воду A воде B ?
2. Связано ли отношение к воде A с полом потребителя? От этого зависит стратегия продвижения товара.

Для ответа на первый вопрос решено использовать критерий знаков, для ответа на второй – критерий ранговых сумм Уилкоксона. Выбран уровень значимости 10%.

Дайте ответы на поставленные вопросы. В каждом случае сформулируйте основную и альтернативную гипотезы, рассчитайте статистику, приведите критическое значение (или значения) и сделайте вывод.

Решение

В качестве языка программирования для решения данной задачи используется C++ (исходные коды представлены в виде файлов "question1.cpp" (для первого вопроса) и "question2.cpp" (для второго вопроса); скриншоты исходных кодов в приложении).

Вопрос 1

Введём обозначения:

- A_i – оценка воды A i -м человеком
- B_i – оценка воды B i -м человеком
- n^+ – число наблюдений, в которых $A_i > B_i$
- n^- – число наблюдений, в которых $B_i > A_i$

Основная гипотеза – у предполагаемого потребителя нет предпочтений между водой A и водой B . Альтернативная гипотеза – предполагаемый потребитель предпочитает воду A воде B .

$$\begin{cases} H_0: P\{A_i > B_i\} = P\{A_i < B_i\} \\ H_A: P\{A_i > B_i\} > P\{A_i < B_i\} \end{cases}$$

Воспользуемся критерием знаков.

Для этого сначала отбросим все наблюдения, в которых $A_i = B_i$. Во-первых, они бесполезны, так как разница в оценках нам не важна – важен только факт превосходства одной воды над другой. Во-вторых, они нарушают распределение n^+ и n^- .

[В методе `get_data` считываем данные из файла "data.txt", игнорируя наблюдения с одинаковыми оценками для воды A и воды B . Отсеянные наблюдения сохраняем в файл "data_without_equal_appraisals.txt".]

Удалив вышеописанные наблюдения, получаем:

$$\begin{cases} n^+ \overset{H_0}{\sim} Bi(n^+ + n^-; 0.5) \\ n^- \overset{H_0}{\sim} Bi(n^+ + n^-; 0.5) \end{cases}$$

В качестве статистики используем $T = n^-$.

[В методе `calculate_n_plus_and_n_minus` считаем n^+ и n^- , сохраняя их в файл "n_plus_and_n_minus.txt". Из файла получаем $n^+ = 26$ и $n^- = 14$.]

Критическая область:

$$T \leq T_{n^+ + n^-; \alpha} = \max\{k: P_{H_0}\{n^- \leq k\} \leq \alpha\} = \max\left\{k: \sum_{i=0}^k C_{n^+ + n^-}^i \cdot 0.5^{n^+ + n^-} \leq \alpha\right\}$$

[В методе `calculate_critical_value` вычисляем критическое значение для области, сохраняя его в файл "critical_value.txt". Из файла получаем $T_{n^+ + n^-; \alpha} = 15$.]

В таком случае, мы попадаем в критическую область, так как $T = 14 \leq 15 = T_{40; 0.1}$. Значит, мы отвергаем основную гипотезу в пользу альтернативной, то есть **есть основание считать**, что предполагаемый потребитель **предпочитает воду A воде B** .

Вопрос 2

Введём обозначения:

- n – исходное количество наблюдений
- n_0 – количество наблюдений, где пол человека – мужской
- n_1 – количество наблюдений, где пол человека – женский
- sex_i – пол i -го человека
- A_i – оценка воды A i -м человеком
- $rank_i$ – ранг i -го наблюдения

Основная гипотеза – отношение к воде A не связано с полом (выборки по полу однородны). Альтернативная гипотеза – отношение к воде связано с полом (выборки неоднородны – двусторонняя альтернатива).

Воспользуемся критерием Уилкоксона.

Для этого сначала отбросим все оценки для воды B (так как мы анализируем не предпочтение между водой A и водой B в зависимости от пола, отношение к воде A в зависимости от пола), отсортируем наблюдения в порядке возрастания по оценке воды A и посчитаем ранги наблюдений.

[В методе "get_data" считываем данные из файла "data.txt", игнорируя оценки для воды B . В этом же методе происходит сортировка наблюдений по возрастанию оценки воды A и присвоение рангов этим наблюдениям. Полученные данные сохраняю в файл "ranked_data.txt".]

Для выбора статистики вычислим n_0 и n_1 .

[В методе calculate_n0_and_n1 вычисляем n_0 и n_1 , сохраняя их в файл "n0_and_n1.txt". Из файла получаем $n_0=20$ и $n_1=24$.]

Рассмотрим ранги наблюдений, где пол человека - мужской, и используем статистику Уилкоксона $W_{\text{набл}} = \sum_{i=0}^{n_0} \text{rank}_i$.

[В методе calculate_Wilcoxon_statistics вычисляем значение W -статистики, сохраняя его в файл "Wilcoxon_statistics.txt". Из файла получаем $W_{\text{набл}}=318.5$.]

Критическая область:

$$\begin{cases} W_{\text{набл}} \geq W(\frac{\alpha}{2}, n_1, n_0) \\ W_{\text{набл}} \leq n_0 \cdot (n_1 + n_0 + 1) - W(\frac{\alpha}{2}, n_1, n_0) \end{cases} \Leftrightarrow \begin{cases} W_{\text{набл}} \geq n_0 \cdot (n_1 + n_0 + 1) - W(\frac{\alpha}{2}, n_0, n_1) \\ W_{\text{набл}} \leq W(\frac{\alpha}{2}, n_0, n_1) \end{cases}$$

Воспользовавшись таблицей критических значений критерия Уилкоксона, получаем $W(\frac{\alpha}{2}, n_0, n_1) = W(0.05, 20, 24) = 379$. Тогда, критическая область имеет вид:

$$\begin{cases} W_{\text{набл}} \geq 521 \\ W_{\text{набл}} \leq 379 \end{cases}$$

В таком случае, мы попадаем в критическую область, так как $W_{\text{набл}}=318.5 \leq 379 = W(\frac{\alpha}{2}, n_0, n_1)$. Значит, мы отвергаем основную гипотезу в пользу альтернативной, то есть **отношение к воде A связано с полом потребителя.**

Приложение

Код программы

question1.cpp

```
1 //на основе get_path получает наблюдения, первые значения которых сохраняют
2 //в массив a (оценки воды A), а также в массив b (оценки воды B). Если вышло
3 //чтобы наблюдение воды A совпало с наблюдением воды B, то наблюдения аннулируются.
4 //На основе get_path получает наблюдения, первые значения которых сохраняют
5 //в массив a (оценки воды A), а также в массив b (оценки воды B). Если вышло
6 //чтобы наблюдение воды A совпало с наблюдением воды B, то наблюдения аннулируются.
7 void get_data(const std::string &path, const std::string &save_path,
8               std::vector<int> &a, std::vector<int> &b) {
9     //создаем массив для хранения наблюдений
10    a.erase(a.begin(), a.end());
11    b.erase(b.begin(), b.end());
12
13    std::ifstream fin;
14    fin.open(get_path);
15    std::ofstream fout;
16    fout.open(save_path);
17
18    //проходим по наблюдениям
19    int cur_a, cur_b, cur_n;
20    while (fin >> cur_a && cur_b && cur_n) {
21        //если оценки не равны, то добавляем данные наблюдений
22        //в массивы a и b
23        if (cur_a != cur_b) {
24            a.push_back(cur_a);
25            b.push_back(cur_b);
26            fout << cur_n << " " << cur_a << " " << cur_b << "\n";
27        }
28    }
29    fin.close();
30    fout.close();
31 }
```

```
1 //Вычисляет и сохраняет в файл path количество наблюдений, в которых оценки воды A больше
2 //оценок воды B (n_plus), количество наблюдений, в которых оценки воды B больше
3 //оценок воды A (n_minus), суммарное количество наблюдений (n_plus + n_minus).
4 //Если вышло, что наблюдение аннулировалось, то вычисляем сумму, вычисляем std::log10_error.
5 //Вычисляем сумму, тогда вычисляем n_plus = n_plus - n_minus, n_minus = n_minus.
6 int calculate_n_plus_and_n_minus(const std::string &path, const std::vector<int> &a,
7                                const std::vector<int> &b) {
8     std::ifstream fin;
9     fin.open(path);
10
11     int n_plus = 0, n_minus = 0;
12     for (int i = 0; i < a.size(); i++) {
13         if (a[i] > b[i]) {
14             n_plus++;
15         } else if (a[i] < b[i]) {
16             n_minus++;
17         } else {
18             std::log10_error("There should be no equal appraisals.");
19         }
20     }
21     fin.close();
22     n_plus = n_plus - n_minus;
23     n_minus = n_minus - n_plus;
24     return n_plus;
25 }
```

```
1 //Вычисляет и возвращает число сочетаний из n по k.
2 long long calculate_number_of_combinations(int k, int n) {
3     if (k == 0) {
4         return 1;
5     }
6     long long res = 1;
7     for (int i = n - k + 1; i <= n; i++) {
8         res *= i;
9     }
10    return res;
11 }
```

Вычисление C_n^k

Получение и редактирование
данных

Вычисление n^+ и n^-

```
78 //Вычисляет, сохраняет в файл path и возвращает критическое значение для n наблюдений
79 //на уровне значимости alpha.
80 int calculate_critical_value(const std::string &path, const int n, const double alpha) {
81     std::ofstream fout;
82     fout.open(path);
83
84     //Так как во всех слагаемых присутствует постоянный коэффициент 0.5^n,
85     //его можно вынести. Поделив обе части неравенства на этот коэффициент,
86     //в правой части получаем alpha / 0.5^n = 2^n * alpha.
87     long long helper = 1;
88     for (int i = 0; i < n; i++) {
89         helper *= 2;
90     }
91     //В таком случае, в левой части останутся только целые числа. Тогда, не имеет
92     //смысла хранить вещественное число - округлив в меньшую сторону, мы не повлияем
93     //на ответ.
94     helper = static_cast<long long>(static_cast<double>(helper) * alpha);
95
96     int k;
97     long long p = 0;
98     for (k = -1; p + calculate_number_of_combinations(k + 1, n) <= helper; k++) {
99         p += calculate_number_of_combinations(k + 1, n);
100     }
101     fout << "T_{" << n << " alpha} ";
102     if (k == -1) {
103         fout << " does not exist";
104     } else {
105         fout << " = " << k;
106     }
107     fout.close();
108     return k;
109 }
```

```
110 int main() {
111     //Массивы для хранения наблюдений. a - оценки воды A; b - оценки воды B. Пол в данном
112     //вопросе не фигурирует, поэтому мы не будем сохранять его программно.
113     std::vector<int> a, b;
114
115     get_data(get_path, save_path, "data_without_equal_appraisals.txt", &a, &b);
116     int n = calculate_n_plus_and_n_minus(path, "n_plus_and_n_minus.txt", a, b);
117     calculate_critical_value(path, "calculate_critical_value.txt", n, n[0] + n[1], alpha: 0.1);
118
119     delete[] n;
120     return 0;
121 }
```

Запуск всех операций

Вычисление критического значения

question2.cpp

```

1 #include <iostream>
2 #include <fstream>
3 #include <vector>
4 #include <string>
5
6 //Компилятор для загрузки на Raspberry Pi
7
8 bool cmp(pair<int>> pair<int, int>, int> i, int> j) {
9     if (pair<int>> pair<int, int>, int> i, int> j) {
10         return i.first.second < j.first.second;
11     }
12 }

```

Метод сравнения наблюдений (для сортировки)

[illegible]

```

54 //Высчитать и вывести в файл путь, количество найденных, и каталог пап
55 //Наименование - первый файл, количество найденных, и каталог пап целиком
56 //Второй файл, суммарное количество найденных (nb + n1)
57 //Выводимого массива, первый элемент которого - nb, второй - n1
58 int calculate_nb_and_n1(const std::string path,
59                       const std::vector<std::pair<std::pair<int, int>, double>> &box_a_rank) {
60     std::ofstream fout;
61     fout.open(path);
62
63     int nb = 0, n1 = 0;
64     for (auto k1 : box_a_rank) {
65         if (k1.first.first() <
66             0)
67             n1++;
68         else {
69             nb++;
70         }
71     }
72     fout << "nb = " << nb << "n1 = " << n1 << "nb + n1 = " << nb + n1;
73
74     int n = box_a_rank.size();
75     n[0] = nb;
76     n[1] = n1;
77     return n;
78 }

```

Вычисление n_0 и n_1

Получение и редактирование данных

```

98 //вычисляет, сохраняет в path_path и возвращает значение W-статистики для мабера
99 //наблюдений sex_a_rank.
100 double calculate_Wilcoxon_statistics(const std::string &path,
101                                     const std::vector<
102                                     std::pair<std::pair<int, int>, double>> &sex_a_rank)
103 {
104     std::ofstream fout;
105     fout.open(path);
106
107     double res = 0;
108     //Во избежание проверки пола можно допознать на (1 - пол): так как женский пол
109     //закодирован нулем, а мужской единицей, то при допознании каждого ранга на
110     //(1 - пол) получаем сумму рангов наблюдений, где пол человека - мужской, то есть
111     //получаем искомого статистику.
112     for (auto &i : sex_a_rank) {
113         res += (1 - i.first.first) * i.second;
114     }
115     fout << "W = " << res;
116
117     fout.close();
118     return res;
119 }

```

```

100
101 int main() {
102     std::vector<std::pair<std::pair<int, int>, double>> sex_a_rank;
103
104     get_data(get_path<"data.txt">, save_path<"ranked_data.txt">, &sex_a_rank);
105     delete[] calculate_n0_and_n1(path<"n0_and_n1.txt">, sex_a_rank);
106     calculate_wilcoxon_statistics(path<"wilcoxon_statistics.txt">, sex_a_rank);
107
108     return 0;
109 }

```

Запуск всех операций

Вычисление W -статистики

Названия и содержимое файлов, связанных с программой

- 1) "question1.txt" – исходный код программы для 1-го вопроса
- 2) "question2.txt" – исходный код программы для 2-го вопроса
- 3) "data.txt" – входной файл с наблюдениями (каждое наблюдение – строка из 3-х чисел, разделённых пробелами: 1-е число – пол человека (0 или 1); 2-е число – оценка воды A (от 1 до 7 включительно); 3-е число – оценка воды B (от 1 до 7 включительно))
- 4) "data_without_equal_appraisals.txt" (1-й вопрос) – файл с наблюдениями, в которых оценка для воды A не равна оценке для воды B
- 5) "n_plus_and_n_minus.txt" (1-й вопрос) – файл, содержащий: количество наблюдений, к которых оценка воды A больше оценки воды B ; количество наблюдений, к которых оценка воды B больше оценки воды A ; суммарное количество наблюдений
- 6) "critical_value.txt" (1-й вопрос) – файл, содержащий критическое значение статистики
- 7) "ranked_data.txt" (1-й вопрос) (2-й вопрос) – файл содержащий наблюдения в виде строки из 3-х чисел: 1-е – пол человека; 2-е – оценка воды A ; 3-е – ранг наблюдения
- 8) "n0_and_n1.txt" (2-й вопрос) – файл, содержащий: количество наблюдений, в которых пол человека - мужской; количество наблюдений, в которых пол человека - женский; суммарное количество наблюдений
- 9) "Wilcoxon_statistics.txt" (2-й вопрос) – файл, содержащий значение W -статистики