

Правительство Российской Федерации

---

Федеральное государственное автономное образовательное учреждение  
высшего профессионального образования «Национальный исследовательский  
университет «Высшая школа экономики»

Московский институт электроники и математики Национального исследовательского  
университета "Высшая школа экономики"



ЭКЗАМИНАЦИОННЫЙ ПРОЕКТ

по дисциплине

"Программирование алгоритмов защиты  
информации

«Корреляционные тесты»

Выполнили

Кузьмина Ксения  
Монахова Мария  
Родионов Филипп  
Япаров Матвей

---

Проверил

Нестеренко А.Ю.  
Преподаватель дисциплины

---

Москва, 2024 г.

# Содержание

<b>1</b>	<b>Постановка задачи</b>	<b>2</b>
<b>2</b>	<b>Теоретическая часть</b>	<b>2</b>
2.1	Знаковый корреляционный критерий Нелсона . . . . .	2
2.1.1	Базовая информация . . . . .	2
2.1.2	Подробное обоснование . . . . .	2
2.2	Сериальный критерий Шведа-Эйзенхарта . . . . .	4
2.2.1	Базовая информация . . . . .	4
2.2.2	Подробное обоснование . . . . .	4
<b>3</b>	<b>Практическое применение</b>	<b>6</b>

# 1 Постановка задачи

Провести корреляционный анализ алгоритма выработки производных ключей, в основе которого лежит функция хеширования Стрибог512. Для данной работы были выбраны два теста – Знаковый корреляционный критерий Нелсона и Серийный критерий Шведа-Эйзенхарта.

Далее будет рассмотрена сначала теоретическая часть, необходимая для применения этих тестов. Далее приведён код на языке C++ для анализа выборок различного размера из файлов.

## 2 Теоретическая часть

### 2.1 Знаковый корреляционный критерий Нелсона

#### 2.1.1 Базовая информация

Критерий, предложенный Нелсоном позволяет установить наличие корреляции, непрерывно анализируя совместное поведение пар  $(x_i, y_i)$  по мере их появления в эксперименте. Критерий основан на числе знаков последовательного изменения величин пар  $(x_i, y_i)$ .

Если  $x_i > x_{i-1}$ ,  $y_i > y_{i-1}$  или  $x_i < x_{i-1}$ ,  $y_i < y_{i-1}$ , то паре  $(x_i, y_i)$  приписывается знак  $+$ , в ином случае знак  $-$ . Другими словами, если значения пар  $(x_i, y_i)$  изменилось в одном направлении, то это отображается знаком  $+$ , в разных направлениях – знаком  $-$ . Если в паре одно или оба значения не изменилось, то ей приписывается значение 0. Статистикой критерия является наименьшее количество  $S$  знаков одного вида ( $+$  или  $-$ ). Корреляция признаётся значимой при  $S > S_\alpha$  ( $S_\alpha$  – критическое значение из справочной таблицы).

При  $n > 90$  сумма  $S$  распределена асимптотически нормально и  $S_\alpha \equiv \frac{n}{2} + u_{1-\alpha} \cdot \sqrt{\frac{11n-2}{36}}$ .

#### 2.1.2 Подробное обоснование

1. Гипотеза  $H_0$  (независимость):

- Предполагается, что  $x_i$  и  $y_i$  независимы.
- Следовательно, направления изменений  $x_i > x_{i-1}$ ,  $x_i < x_{i-1}$  и  $y_i > y_{i-1}$ ,  $y_i < y_{i-1}$  также независимы.

2. Распределение знаков при  $H_0$ :

- Для каждой пары  $(x_i, y_i)$  вероятность получить:
  - $+$ :  $\Pr(+|H_0) = 0.25$ : оба числа растут или убывают;
  - $-$ :  $\Pr(-|H_0) = 0.25$ : одно число растёт, второе убывает;
  - $0$ :  $\Pr(0|H_0) = 0.25$ : любое неизменение равно нулю.

3. Пусть  $n$  – общее число значимых наблюдений (опускаем 0, рассматриваем только  $+$  и  $-$ ).  $S_+$  – число знаков  $+$ ,  $S_-$  – число знаков  $-$ , где  $S_+ + S_- = n$ . Предполагается, что знаки  $+$  и  $-$  равновероятны при  $H_0$ , то есть вероятность получить знак  $+$  или  $-$  равна  $p = 0.5$ .

- **Математическое ожидание.**  $S_+$  следует биномиальному распределению

$$S_+ \sim \text{Binomial}(n, p = 0.5)$$

- **Среднее значение** для  $S_+$ :

$$\mathbb{E}[S_+] = n \cdot p = \frac{n}{2}$$

$$\mathbb{E}[S_-] = n \cdot (p - 1) = \frac{n}{2}$$

- **Дисперсия.** Для биномиального распределения:

$$\text{Var}[S_+] = n \cdot p \cdot (1 - p) = n \cdot 0.5 \cdot 0.5 = \frac{n}{4}$$

$$\text{Var}[S_-] = n \cdot p \cdot (1 - p) = n \cdot 0.5 \cdot 0.5 = \frac{n}{4}$$

4. Статистика критерия  $S = \min(S_+, S_-)$

- Теперь вводим статистику  $S = \min(S_+, S_-)$ , которая используется для оценки минимального количества знаков одного типа. Так как было установлено, что  $S_+ + S_- = n$ , то статистика  $S$  зависит от распределения  $S_+$ .

Так как  $S_+ \sim \text{Binomial}(n, p = 0.5)$  при достаточно большом  $n$  биномиальное распределение можно приближать нормальным:

$$S_+ \sim N\left(\frac{n}{2}, \frac{n}{4}\right)$$

Среднее значение  $\mu = \frac{n}{2}$ , дисперсия  $\sigma^2 = \frac{n}{4}$ .

Соответственно, стандартное отклонение будет равно:

$$\sigma = \sqrt{\frac{n}{4}} = \frac{\sqrt{n}}{2}$$

5. Поскольку  $S = \min(S_+, S_-)$  и нам известно, что  $S$  не может превышать  $\frac{n}{2}$ . При этом  $S$  меньше половины выборки только в случае, когда  $S_+$  отклоняется от  $\frac{n}{2}$ . Чтобы определить критическое значение  $S_\alpha$  нужно использовать нормальное приближение. На уровне значимости  $\alpha$ , вероятность отклонения определяется через квантиль стандартного нормального распределения  $u_{1-\alpha}$ .

6. Вычисление  $S_\alpha$  Для определения  $S_\alpha$  используем теоретическое распределение знаков:

$$S_\alpha = \mathbb{E}[S] + u_{1-\alpha} \cdot \sigma_S$$

- **Математическое ожидание.** Среднее значение  $S$ :

$$\mathbb{E}[S] = \frac{n}{2}$$

- **Стандартное отклонение.** По теоретической модели Нелсона:

$$\text{Var}[S] = \frac{11n - 2}{36}$$

Откуда стандартное отклонение

$$\sigma_S = \sqrt{\frac{11n - 2}{36}}$$

- **Учёт квантиля.**

– Квантиль  $u_{1-\alpha}$  зависит от уровня значимости  $\alpha$ .

\* Для  $\alpha = 0.05$ ,  $u_{1-\alpha} = 1.645$ ;

\* Для  $\alpha = 0.01$ ,  $u_{1-\alpha} = 2.33$ ;

7. Проверка гипотезы.

- Если  $S \leq S_\alpha$ , принимаем нулевую гипотезу  $H_0$ , предполагающую независимость.
- Если  $S \geq S_\alpha$ , отвергаем нулевую гипотезу  $H_0$ , что свидетельствует о значимой корреляции между  $x_i$  и  $y_i$ .

## 2.2 Сериальный критерий Шведа-Эйзенхарта

### 2.2.1 Базовая информация

Совокупность пар  $(x_i, y_i)$  разбивается на две равные совокупности, отвечающие условиям  $y_i > \tilde{y}$  и  $y_i < \tilde{y}$  (где  $\tilde{y}$  – медиана ряда  $y_i$ , при нечётном  $n$  значение  $y_i = \tilde{y}$  исключается). Затем наблюдения ранжируются по возрастающим значениям  $y_i$ . Для последовательных пар значений  $(x_i, y_i)$  с  $y_i > \tilde{y}$  будем применять символ  $a$ , для последовательных пар  $(x_i, y_i)$  с  $y_i < \tilde{y}$  – символ  $b$ . В результате получим последовательность вида  $a, b, a, a, b$ . Последовательность элементов одного вида, ограниченная с двух сторон элементами другого вида (замыкающие интервал последовательности одного вида ограничены с одной стороны последовательностями другого вида), называется *серией*. Количество  $m$  серий является статистикой рассматриваемого критерия.

Корреляция признаётся значимой, если  $m \leq m_\alpha$  (критические значения приведены в таблицах). При чётных  $n > 40$  можно использовать приближение:

$$m_{0.95} = \left\lfloor \frac{n+1}{2} - 0.82 \cdot \sqrt{n-1} \right\rfloor,$$
$$m_{0.99} = \left\lfloor \frac{n+1}{2} - 1.16 \cdot \sqrt{n-1} \right\rfloor,$$

### 2.2.2 Подробное обоснование

Сериальный критерий Шведа-Эйзенхарта используется для проверки гипотезы о независимости последовательности знаков  $\{+, -\}$ . Основная идея – анализ числа серий в последовательности. Серия определяется как непрерывный участок одинаковых знаков. Например:

- Для последовательности  $+++--++$  количество серий  $m = 3$ .
- Для последовательности  $+++++$  количество серий  $m = 1$ .

1. Гипотеза  $H_0$  (независимость):

- Знаки  $+$  и  $-$  независимы и равновероятны.

Гипотеза  $H_1$  (независимость):

- Знаки  $+$  и  $-$  не независимы (наблюдается корреляция).

Если  $H_0$  верна, число серий  $m$  должно быть близко к математическому ожиданию  $\mathbb{E}[m]$ , которое вычисляется на основе теории вероятностей. Отклонение от этого значения позволяет судить о корреляции.

2. Математическое ожидание и дисперсия числа серий  $m$ .

Пусть  $\{s_1, s_2, \dots, s_n\}$  – последовательность знаков  $+$  и  $-$ . Число серий  $m$  определяется как:

$$m = 1 + \sum_{i=1}^{n-1} I(s_i \neq s_{i+1})$$

, где  $I(s_i \neq s_{i+1})$  – индикатор смены знака:

$$I(s_i \neq s_{i+1}) = \begin{cases} 1, & \text{если } s_i \neq s_{i+1}, \\ 0, & \text{если } s_i = s_{i+1} \end{cases}. \quad (1)$$

(а) Математическое ожидание  $\mathbb{E}[m]$ .

При гипотезе независимости  $H_0$  вероятность смены знака между  $s_i$  и  $s_{i+1}$  равна 0.5, так как  $\Pr(s_i = +) = \Pr(s_i = -) = 0.5$ . Тогда математическое ожидание смены знаков:

$$\mathbb{E} \left[ \sum_{i=1}^{n-1} I(s_i \neq s_{i+1}) \right] = (n-1) \cdot \Pr(s_i \neq s_{i+1}) = \frac{n-1}{2}$$

Добавим первую серию, которая есть всегда, для получения математического ожидания числа серий:

$$\mathbb{E}[m] = 1 + \frac{n-1}{2} = \frac{n+1}{2}.$$

(b) **Дисперсия**  $\text{Var}[m]$ .

Для индикаторной функции  $I(s_i \neq s_{i+1})$  дисперсия равна:

$$\text{Var}[I(s_i \neq s_{i+1})] = \Pr(s_i \neq s_{i+1}) \cdot (1 - \Pr(s_i \neq s_{i+1})) = 0.5 \cot 0.5 = 0.25.$$

Так как индикаторы независимы при  $H_0$  дисперсия суммы:

$$\text{Var}\left[\sum_{i=1}^{n-1} I(s_i \neq s_{i+1})\right] = (n-1) \cdot 0.25 = \frac{n-1}{4}.$$

Добавляем первую серию, которая не вносит дисперсию, и получаем дисперсию числа серий:

$$\text{Var}[m] = \frac{n-1}{4}$$

3. **Стандартное отклонение**  $\sigma_m$ .

Стандартное отклонение числа серий:

$$\sigma_m = \sqrt{\text{Var}[m]} = \sqrt{\frac{n-1}{4}} = \frac{\sqrt{n-1}}{2}.$$

4. **Приближённое распределение числа серий**  $m$ .

При  $n > 40$  распределение числа серий  $m$  приближается нормальным распределением:

$$m \sim N\left(\frac{n+1}{2}, \frac{\sqrt{n-1}}{2}\right)$$

5. **Критическое значение**  $m_\alpha$ .

Для проверки гипотезы  $H_0$  используется критическое значение  $m_\alpha$ , которое зависит от уровня значимости  $\alpha$ . При  $n > 40$  значение  $m_\alpha$  вычисляется как:

$$m_\alpha = \left\lceil \frac{n+1}{2} - z_\alpha, \frac{\sqrt{n-1}}{2} \right\rceil$$

, где  $Z_\alpha$  – квантиль стандартного нормального распределения:

- $z_{0.05} = 1.645$  (для уровня значимости  $\alpha = 0.05$ );
- $z_{0.01} = 2.33$  (для уровня значимости  $\alpha = 0.01$ );

6. **Формулы для**  $m_\alpha$

- Для  $\alpha = 0.05$ :

$$m_{0.95} = \left\lceil \frac{n+1}{2} - 0.82 \cdot \sqrt{n-1} \right\rceil$$

.

- Для  $\alpha = 0.01$ :

$$m_{0.99} = \left\lceil \frac{n+1}{2} - 1.16 \cdot \sqrt{n-1} \right\rceil$$

.

7. **Валидность критерия**

- При независимости  $H_0$ :
  - Знаки  $+$  и  $-$  распределены случайно, и число серий  $m$  будет подчиняться теоретическому распределению  $\mathbb{E} = \frac{n+1}{2}$  и  $\text{Var}[m] = \frac{n-1}{4}$ .
  - При больших  $n$  это распределение приближается к нормальному.
- При зависимости  $H_1$ :
  - Если знаки зависимы, то  $m$  отклоняется от ожидаемого значения:
    - \* Кластеризация знаков  $(+++--)$  уменьшает  $m$ .
    - \* Слишком частое чередование  $(+-+-)$  уменьшает  $m$ .

### 3 Практическое применение

```

1 #include <iostream>
2 #include <fstream>
3 #include <vector>
4 #include <string>
5 #include <cmath>
6 #include <algorithm>
7
8 #if defined(_WIN32) || defined(_WIN64)
9 #include <direct.h> // For Windows (_getcwd)
10 #define GetCWD _getcwd
11 #else
12 #include <unistd.h> // For Linux/macOS (getcwd)
13 #define GetCWD getcwd
14 #endif
15
16 using namespace std;
17
18 // Function to read nhmac values from the file
19 vector<string> readNHMACValues(const string &filename) {
20     vector<string> nhmacValues;
21     ifstream file(filename);
22     if (!file.is_open()) {
23         cerr << "Error: Failed to open file " << filename << endl;
24
25         char current_path[FILENAME_MAX];
26         if (GetCWD(current_path, sizeof(current_path))) {
27             cout << "Current working directory: " << current_path << endl;
28         } else {
29             cerr << "Error: Unable to determine the current working directory." << endl;
30         }
31         return nhmacValues;
32     }
33
34     string line;
35     while (getline(file, line)) {
36         // Look for the line starting with "nhmac:"
37         if (line.rfind("nhmac:", 0) == 0) { // Check if the line starts with "nhmac:"
38             string nhmac = line.substr(6); // Extract the part after "nhmac:"
39             nhmac.erase(remove_if(nhmac.begin(), nhmac.end(), ::isspace), nhmac.end());
40             // Remove any whitespace
41             nhmacValues.push_back(nhmac);
42         }
43     }
44
45     file.close();
46     return nhmacValues;
47 }
48
49 // Function to calculate median using lexicographical order
50 string calculateMedian(vector<string> &yValues) {
51     sort(yValues.begin(), yValues.end()); // Sort lexicographically

```

```

52     size_t n = yValues.size();
53     if (n % 2 == 0) {
54         return yValues[n / 2 - 1]; // For even n, return the lower middle value
55     } else {
56         return yValues[n / 2]; // For odd n, return the middle value
57     }
58 }
59
60 // Function to generate the symbol sequence (a, b) based on the median
61 vector<char> generateSymbolSequence(vector<pair<string, string>> &keyPairs, const string
    &median) {
62     vector<char> symbols;
63
64     for (const auto &pair : keyPairs) {
65         const string &yValue = pair.second;
66         if (yValue < median) {
67             symbols.push_back('a');
68         } else if (yValue > median) {
69             symbols.push_back('b');
70         }
71     }
72
73     return symbols;
74 }
75
76 // Function to sort pairs by x-values
77 void sortPairsByX(vector<pair<string, string>> &keyPairs) {
78     sort(keyPairs.begin(), keyPairs.end(), [](const pair<string, string> &a, const pair<
        string, string> &b) {
79         return a.first < b.first;
80     });
81 }
82
83 // Function to count the number of series in the symbol sequence
84 int countSeries(const vector<char> &symbols) {
85     if (symbols.empty()) return 0;
86
87     int seriesCount = 1; // Start with one series
88     for (size_t i = 1; i < symbols.size(); ++i) {
89         if (symbols[i] != symbols[i - 1]) {
90             seriesCount++;
91         }
92     }
93     return seriesCount;
94 }
95
96 // Calculation of critical series value m_alpha
97 int calculateCriticalSeries(int n, double alpha) {
98     double criticalValue;
99     if (alpha == 0.05) {
100         criticalValue = (n + 1) / 2.0 - 0.82 * sqrt(n - 1);
101     } else if (alpha == 0.01) {
102         criticalValue = (n + 1) / 2.0 - 1.16 * sqrt(n - 1);
103     } else {
104         cerr << "Unsupported alpha = " << alpha << ". Defaulting to 0.05." << endl;
105         criticalValue = (n + 1) / 2.0 - 0.82 * sqrt(n - 1);
106     }
107     return static_cast<int>(round(criticalValue)); // Round to nearest integer
108 }
109
110 // Nelson's Criterion
111 int nelsonTest(const vector<pair<string, string>> &keyPairs, int &n) {
112     vector<int> signs;
113     for (size_t i = 1; i < keyPairs.size(); ++i) {
114         int deltaX = keyPairs[i].first.compare(keyPairs[i - 1].first);
115         int deltaY = keyPairs[i].second.compare(keyPairs[i - 1].second);
116
117         // Determine the sign of the pair
118         if ((deltaX > 0 && deltaY > 0) || (deltaX < 0 && deltaY < 0)) {
119             signs.push_back(1); // "+"

```



```

120     } else if ((deltaX > 0 && deltaY < 0) || (deltaX < 0 && deltaY > 0)) {
121         signs.push_back(-1); // "-"
122     } else {
123         signs.push_back(0); // "0"
124     }
125 }
126
127 // Count the number of "+" and "-"
128 int plusCount = count(signs.begin(), signs.end(), 1);
129 int minusCount = count(signs.begin(), signs.end(), -1);
130
131 // n - total number of "+" and "-"
132 n = plusCount + minusCount;
133
134 // Return the minimum of "+" and "-"
135 return min(plusCount, minusCount);
136 }
137
138 // Function to write results to a CSV file
139 void writeResultsToCSV(const string &filename, const string &testFile, int n, int
    nelsonS, double criticalSAlpha,
140                      int series, int criticalM, bool independenceS, bool independenceM
141 ) {
142     ofstream file(filename, ios::app); // Append mode
143     if (!file.is_open()) {
144         cerr << "Error: Unable to open file " << filename << " for writing results." <<
            endl;
145         return;
146     }
147
148     // Write headers if the file is empty
149     static bool headerWritten = false;
150     if (!headerWritten) {
151         file << "TestFile,n,NelsonS,CriticalSAlpha,Series,CriticalM,NelsonIndependence,
            SeriesIndependence\n";
152         headerWritten = true;
153     }
154
155     // Write test data
156     file << testFile << "," << n << "," << nelsonS << "," << criticalSAlpha << ","
157         << series << "," << criticalM << "," << (independenceS ? "1" : "0") << ","
158         << (independenceM ? "1" : "0") << "\n";
159
160     file.close();
161 }
162
163 // Main function to perform the tests
164 void performTests(const string &filename, double alpha) {
165     auto nhmacValues = readNHMACValues(filename);
166     if (nhmacValues.size() < 2) {
167         cerr << "Insufficient data for analysis in file: " << filename << endl;
168         return;
169     }
170
171     // Convert nhmacValues to key pairs
172     vector<pair<string, string>> keyPairs;
173     for (size_t i = 0; i < nhmacValues.size() - 1; ++i) {
174         keyPairs.emplace_back(nhmacValues[i], nhmacValues[i + 1]);
175     }
176
177     // Nelson's Test
178     int n = 0;
179     int nelsonS = nelsonTest(keyPairs, n);
180     double criticalSAlpha = n / 2.0 + 1.645 * sqrt((11.0 * n - 2.0) / 36.0);
181     bool independenceS = nelsonS <= criticalSAlpha;
182
183     // Swed-Eisenhart Serial Criterion
184     vector<string> yValues;
185     for (const auto &pair : keyPairs) {

```

```

186     yValues.push_back(pair.second);
187 }
188
189 string median = calculateMedian(yValues);
190 sortPairsByX(keyPairs);
191 auto symbols = generateSymbolSequence(keyPairs, median);
192 int series = countSeries(symbols);
193 int criticalM = calculateCriticalSeries(n, alpha);
194 bool independenceM = series <= criticalM ;
195
196 // Output results
197 cout << "File: " << filename << endl;
198 cout << "Sample size: " << n << " data points" << endl;
199 cout << "Nelson's Criterion (S-statistic): " << nelsonS
200     << " (critical value S_alpha = " << criticalSAlpha << ")" << endl;
201 cout << "Nelson's result: " << (independenceS ? "Independence" : "Dependence") <<
202     endl;
203 cout << "Swed-Eisenhart Serial Criterion (m): " << series
204     << " (critical value m for alpha = " << alpha << ": " << criticalM << ")" <<
205     endl;
206 cout << "Swed-Eisenhart result: " << (independenceM ? "Dependence" : "Independence")
207     << endl;
208 cout << endl;
209
210 // Write results to CSV
211 writeResultsToCSV("results.csv", filename, n, nelsonS, criticalSAlpha, series,
212     criticalM, independenceS, independenceM);
213 }
214
215 int main() {
216     vector<string> files = {
217         "test_results_100_0.txt",
218         "test_results_200_0.txt",
219         "test_results_300_0.txt",
220         "test_results_400_0.txt",
221         "test_results_500_0.txt",
222         "test_results_600_0.txt",
223         "test_results_700_0.txt",
224         "test_results_800_0.txt",
225         "test_results_900_0.txt",
226     };
227
228     double alpha = 0.05; // Significance level
229
230     for (const auto &file : files) {
231         performTests(file, alpha);
232     }
233
234     return 0;
235 }

```

Листинг 1: Реализация обработки поданных на вход файлов и корреляционных тестов

Теперь рассмотрим вывод в **result.csv**:

TestFile	n	NelsonS	CriticalSAlpha	Series	CriticalM	Nelson Independence	SeriesIndependence
test_results_100_0.txt	98	32	57.9933	54	41	1	1
test_results_200_0.txt	198	71	111.789	106	88	1	1
test_results_300_0.txt	298	94	164.692	154	135	1	1
test_results_400_0.txt	398	120	217.136	198	183	1	1
test_results_500_0.txt	498	168	269.288	250	231	1	1
test_results_600_0.txt	598	182	321.233	316	279	1	1
test_results_700_0.txt	696	231	371.986	350	327	1	1
test_results_800_0.txt	798	278	424.684	392	376	1	1
test_results_900_0.txt	898	304	476.246	464	425	1	1

Таблица 1: Таблица result.csv

Ниже приведены графики, показывающие для каждого критерия зависимости  $S$ ,  $S_\alpha$  и  $m$  от размеров выборки:

