

Санкт-Петербургский государственный университет

ПОТЕШКИН Егор Павлович

Выпускная квалификационная работа

**МЕТОД МОНТЕ-КАРЛО SSA для одномерных и многомерных
ВРЕМЕННЫХ РЯДОВ**

Уровень образования: бакалавриат

Направление 01.03.02 «Прикладная математика и информатика»

Основная образовательная программа СВ.5004.2020 «Прикладная математика и
информатика»

Научный руководитель:

Доцент, кафедра статистического
моделирования

д. ф.-м. н., доцент Н. Э. Голяндина

Рецензент:

Программист, Майкрософт

А. Ю. Шлемов

Санкт-Петербург

2024

Saint Petersburg State University
Applied Mathematics and Computer Science

POTESHKIN Egor Pavlovich

Graduation Project

**MONTÉ CARLO SSA METHOD FOR ONE-DIMENSIONAL AND
MULTIVARIATE TIME SERIES**

Scientific Supervisor:

Associate Professor, Department of
Statistical Modelling N. Golyandina

Reviewer:

Software developer, Microsoft R&D

A. Shlemov

Saint Petersburg

2024

Оглавление

Введение	5
Глава 1. Метод MSSA и его модификации	7
1.1. Вспомогательные определения	7
1.2. Метод MSSA	8
1.2.1. Вложение	8
1.2.2. Разложение	8
1.2.3. Группировка	9
1.2.4. Диагональное усреднение	9
1.3. Этап разложения	9
1.3.1. Basic MSSA	9
1.3.2. Toeplitz Block MSSA	10
1.3.3. Toeplitz Sum MSSA	11
1.4. Сравнение методов MSSA	11
1.4.1. Теоретическое сравнение методов	11
1.4.2. Численное сравнение методов	12
Глава 2. Метод Monte Carlo (M)SSA	14
2.1. Проверка статистических гипотез	14
2.1.1. Проблема критериев с неизвестным распределением статистики	15
2.1.2. Поправка неточных критериев	16
2.2. Monte Carlo SSA	17
2.2.1. Постановка задачи	17
2.2.2. Одиночный тест	17
2.2.3. Множественный тест	18
2.2.4. Используемый вариант MC-SSA	20
2.2.5. Зависимость мощности от параметров сигнала и шума	20
2.2.6. Численное сравнение MC-SSA с другими критериями	22
2.3. Monte Carlo MSSA	24
2.3.1. Отличия от одномерного случая	24
2.3.2. Численное сравнение модификаций MC-MSSA	25

Глава 3. Применение метода Monte Carlo SSA на практике	31
3.1. Зависимость радикальности и мощности от параметра L	31
3.2. Оценка параметров красного шума	35
3.2.1. Искажение критерия при использовании оценок	35
3.3. Наличие мешающего сигнала	36
3.3.1. Периодическая компонента	38
3.3.2. Тренд	40
3.4. Применение к реальным данным	42
Заключение	47
Список литературы	49
Приложение А. Графики	51
А.1. Численное сравнение модификаций MC-MSSA	51
А.2. Искажение критерия при использовании оценок	53
А.3. Наличие мешающего сигнала	54
А.3.1. Периодическая компонента	54
А.3.2. Тренд	56
Приложение Б. Таблицы	60
Б.1. Численное сравнение модификаций MC-MSSA	60
Б.2. Сравнение двух алгоритмов MC-SSA с мешающим сигналом	62

Введение

Метод Singular Spectrum Analysis (SSA) [1, 2] является мощным инструментом для анализа временных рядов. Он позволяет разложить ряд на интерпретируемые компоненты, такие как тренд, периодические колебания и шум, что значительно упрощает процесс анализа.

Метод Monte Carlo SSA [3], в свою очередь, решает задачу обнаружения сигнала в красном шуме. Методы MSSA и Monte Carlo MSSA являются обобщением SSA и Monte Carlo SSA на многомерный случай, когда временной ряд представляет собой коллекцию временных рядов.

Наиболее недавними работами, посвященными Monte Carlo SSA, являются [4, 5, 6, 7]. Тем не менее, до сих пор следующие проблемы не были исследованы. Наиболее часто используемый вариант критерия Monte Carlo SSA для проверки гипотезы об отсутствии сигнала является радикальным [8, 9, 6], то есть вероятность ложно отвергнуть гипотезу больше заданного уровня значимости. Поправка на множественное тестирование, предлагаемая в [6], только частично решает эту проблему, так как второй причиной является используемая в критерии подгонка под ряд. Поэтому оказывается необходимой поправка критерия, получаемой с помощью моделирования, и чем больше радикальность критерия, тем более трудоемкая становится эта задача. Поэтому на практике проблема слишком сильной радикальности приводит к необходимости поиска компромисса.

В данной работе сделана попытка поиска этого практического компромисса, при этом для уменьшения радикальности предлагается использовать метод Toeplitz SSA, причем в случае MSSA в двух вариантах, один из которых до этого, по-видимому, не рассматривался в литературе. Проведено сравнение этих вариантов. Проведено численное исследование алгоритмов Monte Carlo (M)SSA по выбору модификации метода и его параметров исходя из подхода, основанного на том, что среди не слишком радикальных критериев необходимо выбирать наиболее мощные. При этом можно пожертвовать мощностью для уменьшения вычислительных затрат.

Кроме этих исследований была поставлена задача доведения алгоритма до возможности его применения на практике, в частности, выработки рекомендаций по выбору параметров, включения оценки параметров красного шума в алгоритм и предположения

о существовании в исходном ряде сигнала, которые не нужно обнаруживать, например, тренда или сезонности.

Опишем структуру работы.

В главе 1 приведено описание метода MSSA и двух его модификаций, и их численное сравнение. В главе 2 представлен метод Monte Carlo SSA и его численное сравнение с другими критериями проверяющими гипотезу об отсутствии сигнала в красном шуме. В этой же главе проведено численное сравнение модификаций Monte Carlo MSSA на различных примерах. В главе 3 проведено исследование зависимости радикальности и мощности Monte Carlo SSA от параметра L , была исследована степень искажения критерия, если параметры красного шума неизвестны. Также был рассмотрен случай Monte Carlo SSA, когда во временном ряде присутствует мешающий сигнал, и продемонстрирован пример применения метода на реальных данных.

В приложение помещены графики и таблицы, на основе которых делаются выводы в основном тексте, но которые затрудняют его чтение.

Для реализации рассматриваемых алгоритмов использовался язык программирования R и пакет Rssa [10].

Глава 1

Метод MSSA и его модификации

В этой главе рассматривается метод Multivariate Singular Spectrum Analysis (сокращенно MSSA) [11] и его модификации. В разделе 1.1 представлены вспомогательные определения, нужные в дальнейшем. В разделе 1.2 представлен алгоритм метода MSSA, а в разделах 1.3.1 и 1.3.2 — его базовый и стандартный теплицев [12] варианты. В разделе 1.3.3 предлагается другая теплицева модификация MSSA, и в разделе 1.4 происходит сравнение методов MSSA: как теоретическое, так и численное.

1.1. Вспомогательные определения

Определение 1. Пусть $\mathbf{X} = (x_1, \dots, x_N)$ — одномерный временной ряд длины N . Выберем параметр L , называемый *длиной окна*, $1 < L < N$. Рассмотрим $K = N - L + 1$ векторов вложения $X_i = (x_i, \dots, x_{i+L-1})^T$, $1 \leq i \leq K$. Определим оператор вложения \mathcal{T} следующим образом:

$$\mathcal{T}(\mathbf{X}) = \mathbf{X} = [X_1 : \dots : X_K] = \begin{pmatrix} x_1 & x_2 & \cdots & x_K \\ x_2 & x_3 & \cdots & x_{K+1} \\ \vdots & \vdots & \ddots & \vdots \\ x_L & x_{L+1} & \cdots & x_N \end{pmatrix}. \quad (1.1)$$

Определение 2. Матрицу \mathbf{X} из (1.1) называют траекторной матрицей.

Заметим, что матрица \mathbf{X} является *ганкелевой*, т.е. на всех ее побочных диагоналях стоят одинаковые элементы, а оператор \mathcal{T} задает взаимно-однозначное соответствие между множеством временных рядов длины N и множеством ганкелевых матриц $L \times K$.

Определение 3. Пусть $\mathbf{Y} = \{y_{ij}\}_{i,j=1}^{L,K}$ — некоторая матрица. Определим оператор ганкелизации \mathcal{H} :

$$(\mathcal{H}(\mathbf{Y}))_{ij} = \sum_{(l,k) \in A_s} y_{lk} / w_s, \quad (1.2)$$

где $s = i + j - 1$, $A_s = \{(l, k) : l + k = s + 1, 1 \leq l \leq L, 1 \leq k \leq K\}$ и $w_s = |A_s|$ — количество элементов в множестве A_s . Это соответствует усреднению элементов матрицы \mathbf{Y} по побочным диагоналям.

Определение 4. Случайный процесс $\xi = (\xi_1, \dots, \xi_n, \dots)$ называется стационарным, если $\forall k \geq 1 \ E\xi_k = \text{const}$ и $\forall k, l \geq 1$

$$K(k, l) \stackrel{\text{def}}{=} \text{cov}(\xi_k, \xi_l) = \tilde{K}(k - l).$$

Определение 5. Белый гауссовский шум $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n, \dots)$ — стационарный случайный процесс с $\varepsilon_n \sim N(0, \sigma^2) \ \forall n$ и $\tilde{K}(n) = 0$, при $n > 0$.

Определение 6. Детерминированный временной ряд $\mathbf{X} = (x_1, \dots, x_n, \dots)$ называют стационарным, если существует функция $R_x(k)$ ($k \in \mathbb{Z}$) такая, что $\forall k, l \geq 0$

$$R_x^{(N)}(k, l) \stackrel{\text{def}}{=} \frac{1}{N} \sum_{m=1}^N x_{k+m} x_{l+m} \xrightarrow{N \rightarrow \infty} R_x(k - l).$$

1.2. Метод MSSA

Метод MSSA состоит из четырех этапов: *вложения, разложения, группировки и диагонального усреднения*. Введем понятие многоканального временного ряда.

Определение 7. Рассмотрим вещественнозначные одномерные временные ряды $\mathbf{X}^{(d)}$ длины N_d , $d = 1, \dots, D$. Тогда составленный из этих рядов $\mathbf{X} = \{\mathbf{X}^{(d)}\}_{d=1}^D$ — D -канальный временной ряд с длинами N_d .

Замечание 1. Понятие стационарности естественным образом обобщается на многомерный случай. В частности, $\tilde{K}(k)$ и $R_x(k)$ из определений 4, 5 будут матрицами размера $D \times D$, учитывающими кросс-ковариации между рядами.

1.2.1. Вложение

Зафиксируем L , $1 < L < \min(N_1, \dots, N_D)$. Для каждого ряда $\mathbf{X}^{(d)}$ составим траекторную матрицу $\mathbf{X}^{(d)}$. Обозначим $K = \sum_{d=1}^D K_d$. Результатом этапа вложения является траекторная матрица многоканального временного ряда

$$\mathbf{X} = \mathcal{T}_{\text{MSSA}}(\mathbf{X}) = [\mathcal{T}(\mathbf{X}^{(1)}) : \dots : \mathcal{T}(\mathbf{X}^{(D)})] = [\mathbf{X}^{(1)} : \dots : \mathbf{X}^{(D)}]. \quad (1.3)$$

1.2.2. Разложение

Задача этапа разложения — разбить траекторную матрицу \mathbf{X} в сумму матриц единичного ранга: $\mathbf{X} = \mathbf{X}_1 + \dots + \mathbf{X}_p$.

Определение 8. Пусть $\mathbf{X} = \sum_i \sigma_i P_i Q_i^T$ — любое разложение \mathbf{X} в сумму матриц ранга 1. Будем называть P_i *левыми*, а Q_i — *правыми векторами* матрицы \mathbf{X} .

1.2.3. Группировка

На этом шаге множество индексов $I = \{1, \dots, p\}$ разбивается на m непересекающихся множеств I_1, \dots, I_m и матрица \mathbf{X} представляется в виде суммы

$$\mathbf{X} = \sum_{k=1}^m \mathbf{X}_{I_k},$$

где $\mathbf{X}_{I_k} = \sum_{i \in I_k} \mathbf{X}_i$.

1.2.4. Диагональное усреднение

Пусть $\mathbf{Y} = [\mathbf{Y}^{(1)} : \dots : \mathbf{Y}^{(M)}]$ — некоторая составная матрица, тогда оператор ганкелизации для составной матрицы

$$\mathcal{H}_{\text{stacked}}(\mathbf{Y}) = [\mathcal{H}(\mathbf{Y}^{(1)}) : \dots : \mathcal{H}(\mathbf{Y}^{(M)})]. \quad (1.4)$$

Финальным шагом MSSA является преобразование каждой матрицы \mathbf{X}_{I_k} , составленной в разделе 1.2.3, в D -канальный временной ряд:

$$\tilde{\mathbf{X}}_{I_k} = \mathcal{T}_{\text{MSSA}}^{-1} \circ \mathcal{H}_{\text{stacked}}(\mathbf{X}_{I_k}), \quad (1.5)$$

где $\mathcal{T}_{\text{MSSA}}$ — оператор вложения (1.3), $\mathcal{H}_{\text{stacked}}$ — оператор ганкелизации (1.4).

Замечание 2. При $D = 1$ \mathbf{X} — одномерный временной ряд, и приведенный выше алгоритм совпадает с алгоритмом Basic SSA, описанный в [2].

1.3. Этап разложения

Модификации MSSA отличаются только этапом разложения, остальные этапы остаются неизменными.

1.3.1. Basic MSSA

Базовый вариант MSSA использует сингулярное разложение (SVD) матрицы \mathbf{X} . Положим $\mathbf{S} = \mathbf{X}\mathbf{X}^T$. Пусть λ_i — собственные числа, а U_i — ортонормированная система векторов матрицы \mathbf{S} . Упорядочим λ_i по убыванию и найдем p такое, что $\lambda_p > 0$, а $\lambda_{p+1} = 0$. Тогда

$$\mathbf{X} = \sum_{i=1}^p \sqrt{\lambda_i} U_i V_i^T = \sum_{i=1}^p \mathbf{X}_i, \quad (1.6)$$

где $V_i = \mathbf{X}^T U_i / \sqrt{\lambda_i}$. Тройку $(\sqrt{\lambda_i}, U_i, V_i)$ принято называть i -й собственной тройкой сингулярного разложения, $\sqrt{\lambda_i}$ — сингулярным числом, U_i — левым сингулярным вектором, а V_i — правым сингулярным вектором.

Замечание 3. Левые сингулярные векторы U_i являются собственными векторами матрицы $\mathbf{X}\mathbf{X}^T$, а правые V_i в свою очередь — матрицы $\mathbf{X}^T\mathbf{X}$. В одномерном случае U_i и V_i равносильны с точностью до замены L на $N - L + 1$. Но при $D > 1$ это не так по построению матрицы \mathbf{X} . Для наглядности рассмотрим случай $D = 2$, тогда

$$\mathbf{X}\mathbf{X}^T = \mathbf{X}^{(1)}(\mathbf{X}^{(1)})^T + \mathbf{X}^{(2)}(\mathbf{X}^{(2)})^T, \quad \mathbf{X}^T\mathbf{X} = \begin{pmatrix} (\mathbf{X}^{(1)})^T\mathbf{X}^{(1)} & (\mathbf{X}^{(1)})^T\mathbf{X}^{(2)} \\ (\mathbf{X}^{(2)})^T\mathbf{X}^{(1)} & (\mathbf{X}^{(2)})^T\mathbf{X}^{(2)} \end{pmatrix}.$$

1.3.2. Toeplitz Block MSSA

Если предполагается, что ряд \mathbf{X} является отрезком стационарного ряда, то имеет смысл заменить матрицу $\mathbf{X}\mathbf{X}^T$ (или $\mathbf{X}^T\mathbf{X}$ в силу замечания 3) на некоторую другую. Для этого введем следующее обозначение.

Определение 9. Пусть $\mathbf{X} = \{\mathbf{X}^{(d)}\}_{d=1}^D$ — D -канальный временной ряд с $N_d = N$. Зафиксируем $1 < M < N$. Обозначим за $\mathbf{T}_{l,k}^{(M)} \in \mathbb{R}^{M \times M}$ матрицу с элементами

$$\left(\mathbf{T}_{l,k}^{(M)}\right)_{ij} = \frac{1}{N - |i - j|} \sum_{n=1}^{N-|i-j|} x_n^{(l)} x_{n+|i-j|}^{(k)}, \quad 1 \leq i, j \leq M,$$

Замечание 4. Если ряд \mathbf{X} — отрезок стационарного ряда, матрица $\mathbf{T}_{l,k}^{(M)}$ является оценкой кросс-ковариационной матрицы l -го и k -го каналов.

В работе [12] предложен способ разложения \mathbf{X} , который мы назовем Toeplitz Block MSSA. Этот способ использует вместо матрицы $\mathbf{X}^T\mathbf{X}$ матрицу

$$\mathbf{T}_{\text{Block}} = \begin{pmatrix} \mathbf{T}_{1,1}^{(K)} & \mathbf{T}_{1,2}^{(K)} & \cdots & \mathbf{T}_{1,D}^{(K)} \\ \mathbf{T}_{2,1}^{(K)} & \mathbf{T}_{2,2}^{(K)} & \cdots & \mathbf{T}_{2,D}^{(K)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{T}_{D,1}^{(K)} & \mathbf{T}_{D,2}^{(K)} & \cdots & \mathbf{T}_{D,D}^{(K)} \end{pmatrix} \in \mathbb{R}^{DK \times DK},$$

где $K = N - L + 1$. Найдя ортонормированные собственные векторы Q_1, \dots, Q_{DK} матрицы $\mathbf{T}_{\text{Block}}$, получаем разложение траекторной матрицы \mathbf{X} :

$$\mathbf{X} = \sum_{i=1}^{DK} \sigma_i P_i Q_i^T = \mathbf{X}_1 + \dots + \mathbf{X}_{DK}, \quad (1.7)$$

где $Z_i = \mathbf{X}Q_i$, $P_i = Z_i / \|Z_i\|$, $\sigma_i = \|Z_i\|$.

1.3.3. Toeplitz Sum MSSA

Вместе с методом Block рассмотрим другой вариант разложения, который назовем Sum. Рассмотрим вместо матрицы $\mathbf{X}\mathbf{X}^T$ матрицу $\mathbf{T}_{\text{Sum}} = \sum_{d=1}^D \mathbf{T}_{d,d}^{(L)} \in \mathbb{R}^{L \times L}$. Найдем ортонормированные собственные векторы P_1, \dots, P_L матрицы \mathbf{T}_{Sum} и разложим траекторную матрицу \mathbf{X} следующим образом:

$$\mathbf{X} = \sum_{i=1}^L \sigma_i P_i Q_i^T = \mathbf{X}_1 + \dots + \mathbf{X}_L, \quad (1.8)$$

где $S_i = \mathbf{X}^T P_i$, $Q_i = S_i / \|S_i\|$, $\sigma_i = \|S_i\|$.

1.4. Сравнение методов MSSA

1.4.1. Теоретическое сравнение методов

Базовый вариант MSSA можно использовать для временных рядов с разными длинами каналов. Toeplitz Sum MSSA также позволяет это делать, в отличие от Toeplitz Block MSSA. Связано это с вычислением матриц $\mathbf{T}_{l,k}^{(K)}$, для которых при $l \neq k$ требуется выполнение условия $N_l = N_k$. В методе Sum такой проблемы не возникает, поскольку $l = k$ всегда.

Сравним обе модификации по трудоемкости. Самым трудоемким является нахождение собственных векторов матриц $\mathbf{T}_{\text{Block}}$ и \mathbf{T}_{Sum} соответственно. Если использовать спектральное разложение, то трудоемкость составляет $\mathcal{O}(D^3 K^3)$ для метода Block и $\mathcal{O}(L^3)$ для метода Sum. Ее можно уменьшить, используя численные методы нахождения SVD типа метода Ланцоша, до $\mathcal{O}(r D^2 K^2)$ и $\mathcal{O}(r L^2)$ соответственно, где r — желаемое количество ненулевых собственных чисел.

Замечание 5. Для теплицевых матриц размера $L \times L$ можно еще сильнее уменьшить трудоемкость, используя быстрое преобразование Фурье для умножения векторов, до $\mathcal{O}(r L \log L)$, взяв за основу подход из [13, 11]. Но матрица $\mathbf{T}_{\text{Block}}$, в отличие от \mathbf{T}_{Sum} , не является теплицевой, поэтому реализовать более быстрый вариант разложения матрицы метода Block на данный момент не представляется возможным.

Таким образом, можно сравнивать методы Toeplitz Sum MSSA и Toeplitz Block MSSA по трудоемкости, сравнивая L и DK .

1.4.2. Численное сравнение методов

Посмотрим на точность базового и модифицированных методов MSSA для разных значений параметра L . Рассмотрим следующий двухканальный временной ряд длины $N = 71$:

$$\{F^{(1)}, F^{(2)}\} = \{S^{(1)}, S^{(2)}\} + \{N^{(1)}, N^{(2)}\},$$

где $S^{(1)}, S^{(2)}$ — некоторые сигналы, а $N^{(1)}, N^{(2)}$ — независимые реализации гауссовского белого шума с $\sigma^2 = 25$.

Рассмотрим 3 случая, первые два из которых рассматривались ранее в [11]:

1. Косинусы с одинаковыми частотами:

$$s_n^{(1)} = 30 \cos(2\pi n/12), \quad s_n^{(2)} = 20 \cos(2\pi n/12), \quad n = 1, \dots, N.$$

2. Косинусы с разными частотами:

$$s_n^{(1)} = 30 \cos(2\pi n/12), \quad s_n^{(2)} = 20 \cos(2\pi n/8), \quad n = 1, \dots, N.$$

3. Полиномы первой степени (нестационарные ряды):

$$s_n^{(1)} = 1.2n, \quad s_n^{(2)} = 0.8n, \quad n = 1, \dots, N.$$

В качестве оценки точности восстановления сигнала было взято среднеквадратичное отклонение от истинного значения. В таблице 1.1 представлены результаты на основе 10000 реализаций шума. Наиболее точные результаты для каждого метода были выделены жирным шрифтом. Лучший результат для каждого случая выделен отдельно синим.

Как видно из таблицы 1.1, в первом случае метод Block лучше всего выделял сигнал. В случае разных частот каналы имеют разную структуру, поэтому наиболее оптимальным является использовать Toeplitz SSA для каждого канала по отдельности. В третьем случае мы имеем дело с нестационарными рядами одинаковой структуры, поэтому стандартный MSSA справляется лучше всего.

Заметим, что преимущество Block перед Sum в первом случае не очень большое. Также, если сравнивать методы по трудоемкости, для оптимальной длины окна метод Sum численно эффективнее Block: в случае Sum для оптимального $L \approx 2N/3$ матрица размера $2N/3 \times 2N/3$, в случае Block оптимальное $L \approx N/2$ и матрица размера $N \times N$. Учитывая замечание 5 и то, что Sum можно применять к временным рядам разной длины, рекомендуется использовать именно Toeplitz Sum MSSA.

Таблица 1.1. MSE восстановления сигнала.

Случай 1 ($\omega_1 = \omega_2$)	$L = 12$	$L = 24$	$L = 36$	$L = 48$	$L = 60$
SSA	3.25	2.01	2.00	2.01	3.25
Toeplitz SSA	3.2	1.87	1.63	1.59	1.67
MSSA	3.18	1.83	1.59	1.47	2.00
Toeplitz Sum MSSA	3.17	1.75	1.44	1.32	1.33
Toeplitz Block MSSA	1.39	1.26	1.25	1.33	1.97
Случай 2 ($\omega_1 \neq \omega_2$)	$L = 12$	$L = 24$	$L = 36$	$L = 48$	$L = 60$
SSA	3.25	2.01	2.00	2.01	3.25
Toeplitz SSA	3.2	1.87	1.63	1.59	1.67
MSSA	6.91	3.77	3.07	2.88	3.84
Toeplitz Sum MSSA	6.88	3.65	2.64	2.37	2.27
Toeplitz Block MSSA	4.47	3.67	3.22	3.23	3.8
Случай 3 (тренд)	$L = 12$	$L = 24$	$L = 36$	$L = 48$	$L = 60$
SSA	3.65	2.08	1.96	2.08	3.65
Toeplitz SSA	3.33	2.43	3.74	7.84	16.29
MSSA	3.42	1.94	1.63	1.57	2.27
Toeplitz Sum MSSA	3.32	2.24	3.04	5.91	11.95
Toeplitz Block MSSA	12.55	6.18	2.97	1.78	1.97

Глава 2

Метод Monte Carlo (M)SSA

Эта глава делится на две части. В разделе 2.1 приведены нужные понятия из теории проверки статистических гипотез. В разделе 2.2 рассматривается метод Monte Carlo SSA (сокращенно MC-SSA), а в разделе 2.3 его многомерное обобщение (MC-MSSA). В названии методов присутствует (M)SSA, поскольку при построении критериев используется траекторная матрица временного ряда, а самый распространенный вариант использует разложение этой траекторной матрицы. В связи с тем, что получить разложение траекторной матрицы можно разными способами, возникают разные варианты MC-(M)SSA.

В разделе 2.1.1 описана проблема критериев с неизвестным распределением статистики при верной нулевой гипотезе, а также алгоритм поправки неточных критериев в разделе 2.1.2. В разделах 2.2.2 и 2.2.3 приведены алгоритмы метода MC-SSA и его модификации с поправкой на множественные тестирования [6], в которых используется траекторная матрица ряда, но не используется ее разложение. В разделе 2.2.4 рассмотрен выбор параметров метода на основе разложения траекторной матрицы исходного ряда и описаны его особенности. В дальнейшем под MC-SSA будем подразумевать вариант из раздела 2.2.4. В разделе 2.2.6 проведено численное сравнение MC-SSA с другими статистическими критериями. Завершает эту главу раздел 2.3.2, где происходит численное сравнение модификаций метода MC-MSSA на нескольких примерах.

2.1. Проверка статистических гипотез

Рассмотрим правосторонний критерий со статистикой T . Введем некоторые обозначения.

Определение 10. p -value — такое значение p , что при значениях уровня значимости α , больших p , H_0 отвергается (статистика попадает в критическую область $A_{\text{крит}}(\alpha)$), а при меньших — не отвергается.

Определение 11. Ошибка первого рода — вероятность отвергнуть нулевую гипотезу, если она верна: $\alpha_I(\alpha) = P_{H_0}(T \in A_{\text{крит}}(\alpha))$.

Определение 12. Мощность критерия — вероятность отвергнуть нулевую гипотезу, если верна альтернативная: $\beta(\alpha) = P_{H_1}(T \in A_{\text{крит}}(\alpha))$.

Замечание 6. По определению p-value $\alpha_I(\alpha) = P_{H_0}(p < \alpha)$ и $\beta(\alpha) = P_{H_1}(p < \alpha)$.

Определение 13. ROC-кривая — это кривая, задаваемая параметрически

$$\begin{cases} x = \alpha_I(\alpha) \\ y = \beta(\alpha) \end{cases}, \quad \alpha \in [0, 1]$$

Замечание 7. С помощью ROC-кривых можно сравнивать по мощности неточные (в частности, радикальные) критерии. Отметим, что для точного критерия ROC-кривая совпадает с графиком мощности, так как $\alpha_I(\alpha) = \alpha$.

2.1.1. Проблема критериев с неизвестным распределением статистики

Если распределение T при верной H_0 неизвестно, оно оценивается с помощью G суррогатных выборок. Приведем алгоритм для правостороннего критерия.

Алгоритм 1. Критерий, использующий суррогатные выборки

1. По выборке $X = (X_1, \dots, X_n)$ строится статистика критерия $t = T(X)$.
2. Моделируется G выборок R_1, \dots, R_G объема n при верной H_0 , и строятся статистики $t_i = T(R_i)$, $i = 1, \dots, G$.
3. Находится оценка t_α критического значения \hat{t}_α как выборочный $(1 - \alpha)$ -квантиль (t_1, \dots, t_G) .
4. Если $t > \hat{t}_\alpha$, то H_0 отвергается.

Замечание 8. Пусть T имеет функцию распределения F при верной H_0 . Тогда p-value вычисляется как $p = 1 - F(t)$. Если распределение T неизвестно, вместо F используют ее оценку

$$\hat{F}_G(x) = \frac{1}{G} \sum_{i=1}^G \mathbb{1}_{t_i \leq x},$$

и $\hat{p}^{(G)} = 1 - \hat{F}_G(t)$ — оценка p .

Замечание 9. Оценка критического значения t_α по квантилям выборки объема G имеет смысл при $\alpha > 1/G$, так как минимальное и максимальное значения выборки соответствуют α - и $(1 - \alpha)$ -квантилям при $\alpha = 1/G$.

Замечание 10. Заметим, что формально оценки квантилей строятся и при $\alpha < 1/G$ некоторой интерполяцией, однако, в этом случае обычно получается отрицательное смещение $\delta(\alpha, G) = E\hat{t}_\alpha - t_\alpha$, причем, $|\delta(\alpha, G)|$ монотонно убывает по α и G . Например, оценим 0.999-квантиль стандартного нормального распределения с помощью выборки размера 100. Тогда истинное значение квантиля равно 3.09, в среднем оценка равна 2.459, а ее медиана — 2.426. Увеличив размер выборки до 1000 получаем уже более точную оценку: в среднем 2.956 и медиана равна 2.932.

Замечание 11. Из замечания 9 следует, что критерий, описанный в алгоритме 1, радикальнее исходного критерия при $\alpha < 1/G$: поскольку $E\hat{t}_\alpha < t_\alpha$, значения t будут чаще попадать в критическую область.

2.1.2. Поправка неточных критериев

Зафиксируем некоторый неточный (консервативный или радикальный) правосторонний критерий и уровень значимости α^* . Пусть дана зависимость ошибки первого рода от уровня значимости $\alpha_I(\alpha) = P_{H_0}(p < \alpha)$. Тогда критерий с формальным уровнем значимости $\tilde{\alpha}^* = \alpha_I^{-1}(\alpha^*)$ является точным: ошибка первого рода $\alpha_I(\tilde{\alpha}^*) = \alpha^*$.

Замечание 12. Если критерий строится по суррогатным данным (алгоритм 1), то по замечаниям 9 и 11 условием точности критерия после поправки является $\tilde{\alpha}^* > 1/G$, что можно рассматривать как ограничение снизу на G . Если $G < 1/\tilde{\alpha}^*$, то после поправки критерий будет радикальным.

Замечание 13. Замечание 12 показывают, что критерий из алгоритма 1, к которому применяется поправка, не должен быть слишком радикальным. Для сильно радикальных критериев с очень маленьким $\tilde{\alpha}^*$ условие $G > 1/\tilde{\alpha}^*$ приводит к очень большим вычислительным затратам и, тем самым, к практической нереализуемости. Везде далее будем рассматривать $G = 1000$ и слова о невозможности построения критерия из-за сильной радикальности будем относить именно к этому выбору G .

Если зависимость $\alpha_I(\alpha)$ неизвестна, она оценивается с помощью моделирования. Приведем алгоритм поправки в этом случае. Помимо критерия и уровня значимости, зафиксируем количество выборок M для оценки $\alpha_I(\alpha)$ и их объем N .

Алгоритм 2. Поправка уровня значимости по зависимости $\alpha_I(\alpha)$ [7]

1. Моделируется M выборок объема N при верной H_0 .
2. По моделированным данным строится оценка зависимость ошибки первого рода от уровня значимости $\alpha_I(\alpha)$.
3. Рассчитывается формальный уровень значимости: $\tilde{\alpha}^* = \alpha_I^{-1}(\alpha^*)$. Критерий с таким уровнем значимости является асимптотически точным при $M \rightarrow \infty$.

2.2. Monte Carlo SSA

2.2.1. Постановка задачи

Рассмотрим задачу поиска сигнала (неслучайной составляющей) во временном ряде. Модель выглядит следующим образом:

$$\mathbf{X} = \mathbf{S} + \boldsymbol{\xi},$$

где \mathbf{S} — сигнал, $\boldsymbol{\xi}$ — стационарный процесс с нулевым средним. Тогда нулевая гипотеза $H_0 : \mathbf{S} = 0$ (отсутствие сигнала, ряд состоит из чистого шума) и альтернатива $H_1 : \mathbf{S} \neq 0$ (ряд содержит сигнал, например, периодическую составляющую).

Определение 14. Случайный процесс $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n, \dots)$ называют красным шумом с параметрами φ и δ , если $\xi_n = \varphi \xi_{n-1} + \delta \varepsilon_n$, где $0 < \varphi < 1$, ε_n — белый гауссовский шум с дисперсией 1 и ξ_1 имеет нормальное распределение с нулевым средним и дисперсией $\delta^2/(1 - \varphi^2)$.

В этой и следующей главах под шумом будем подразумевать именно красный, причем в данной главе с известными параметрами. Также будем рассматривать только односторонние критерии.

2.2.2. Одиночный тест

Пусть $\boldsymbol{\xi}$ — красный шум. Зафиксируем длину окна L и обозначим траекторную матрицу ряда $\boldsymbol{\xi}$ как $\boldsymbol{\Xi}$. Рассмотрим вектор $\mathbf{W} \in \mathbb{R}^L$ такой, что $\|\mathbf{W}\| = 1$. Введем величину

$$p = \|\boldsymbol{\Xi}^T \mathbf{W}\|^2.$$

Статистикой критерия является величина

$$\hat{p} = \|\mathbf{X}^T W\|^2.$$

Если вектор W — синусоида с частотой ω , то \hat{p} отражает вклад частоты ω в исходный ряд.

Рассмотрим алгоритм статистического критерия проверки наличия сигнала в ряде с проекцией на один вектор W , описанный в работе [6].

Алгоритм 3. Одиночный тест [6]

1. Построить статистику критерия \hat{p} .
2. Построить доверительную область случайной величины p : интервал от нуля до $(1 - \alpha)$ -квантиля, где α — уровень значимости.
3. Если \hat{p} не попадает в построенный интервал — H_0 отвергается.

Построенная доверительная область называется *прогнозируемым интервалом* с уровнем доверия $1 - \alpha$.

Замечание 14. В большинстве случаев, распределение p неизвестно. Поэтому оно оценивается методом Монте-Карло: берется G реализаций случайной величины ξ , для каждой вычисляется p и строится эмпирическое распределение. В связи с этим в названии метода присутствуют слова «Monte Carlo».

Замечание 15. Если частота ω сигнала S известна, то в качестве W можно взять синусоиду с частотой ω . Но на практике ω редко бывает известна, что делает данный критерий несостоятельным против H_1 .

2.2.3. Множественный тест

Пусть теперь частоты периодических компонент неизвестны, что не редкость на практике. Тогда подобно одиночному тесту рассмотрим набор W_1, \dots, W_H векторов для проекции, и для каждого $k = 1, \dots, H$ построим статистику критерия \hat{p}_k :

$$\hat{p}_k = \|\mathbf{X}^T W_k\|^2, \quad k = 1, \dots, H. \quad (2.1)$$

В таком случае нужно построить H предсказательных интервалов для каждого W_k по выборкам $P_k = \{p_{ki}\}_{i=1}^G$ с элементами

$$p_{ki} = \|\Xi_i^T W_k\|^2, \quad i = 1, \dots, G; \quad k = 1, \dots, H, \quad (2.2)$$

где G — количество суррогатных реализаций ξ , Ξ_i — траекторная матрица i -й реализации ξ .

В работе [6] подробно описана проблема множественного тестирования, когда вероятность ложного обнаружения периодической составляющей для одной из рассматриваемых частот (групповая ошибка I рода) неизвестна и значительно превышает заданный уровень значимости (частота ошибок одиночного теста), и ее решение. Приведем модифицированный алгоритм построения критерия в случае множественного тестирования, который будем использовать в дальнейшем.

Алгоритм 4. Multiple MC-SSA [6]

1. Для $k = 1, \dots, H$ вычисляется статистика \hat{p}_k , выборка $P_k = \{p_{ki}\}_{i=1}^G$, ее среднее μ_k и стандартное отклонение σ_k .
2. Вычисляется $\eta = (\eta_1, \dots, \eta_G)$, где

$$\eta_i = \max_{1 \leq k \leq H} (p_{ki} - \mu_k) / \sigma_k, \quad i = 1, \dots, G.$$

3. Находится q как выборочный $(1 - \alpha)$ -квантиль η , где α — уровень значимости.
4. Нулевая гипотеза не отвергается, если

$$t = \max_{1 \leq k \leq H} (\hat{p}_k - \mu_k) / \sigma_k < q.$$

5. Если H_0 отвергнута, вклад W_k (и соответствующей частоты) значим, если \hat{p}_k превосходит $\mu_k + q\sigma_k$. Таким образом, $[0, \mu_k + q\sigma_k]$ считаются скорректированными интервалами прогнозирования.

Замечание 16. p -value для этого критерия вычисляется как $1 - \hat{F}_G(t)$, где \hat{F}_G — эмпирическая функция распределения η .

2.2.4. Используемый вариант MC-SSA

В разделе 2.2.3 предполагалось, что векторы W_1, \dots, W_H фиксированные и не зависят от исходного ряда. Такой критерий MC-SSA является точным.

В этой работе будут рассматриваться векторы W_k , порожденные рядом X , при этом по-прежнему при вычислении p_{ki} (2.2) используются те же W_k , что и при вычислении \hat{p}_k (2.1).

Поскольку в этом варианте векторы W_k не заданы заранее, а порождены исходным рядом, критерий MC-SSA становится, вообще говоря, радикальным. Борьба с этой проблемой позволяет метод эмпирической поправки критерия, описанный в разделе 2.1.2.

Замечание 17. Если критерий MC-SSA точный и не требует поправки, то его трудоемкость состоит из трудоемкости одного разложения траекторной матрицы исходного ряда и G вычислений по формуле (2.2). При применении поправки трудоемкость увеличивается в M раз, где M — количество выборок для оценки зависимости $\alpha_I(\alpha)$.

В качестве W_1, \dots, W_H будем брать левые векторы матрицы X (см. определение 8). Такой способ выбора векторов для проекции самый распространенный, поскольку, если есть значимые векторы, можно восстановить сигнал с помощью SSA на их основе. Будем под MC-SSA подразумевать именно этот вариант критерия. Варианты критерия будут определяться конкретным разложением траекторной матрицы. Заметим, что обычно используется сингулярное разложение.

2.2.5. Зависимость мощности от параметров сигнала и шума

Приведем рассуждения относительно зависимости мощности одного из критериев от параметров сигнала и шума. Пусть сигнал является гармоникой с частотой ω , а коэффициент авторегрессии в красном шуме обозначен φ . Можно предположить, что мощность критерия зависит от соотношения амплитуды гармоники и значения спектральной плотности красного шума в точке ω . Отсюда делаем предположение, которое будем использовать далее при выборе параметров примеров.

Предположение 1. При уменьшении φ увеличивается мощность критерия MC-SSA. При увеличении частоты ω мощность также растет.

Обоснование. Рассмотрим одномерный временной ряд $\mathbf{X} = \mathbf{S} + \mathbf{R}$ длины $N = 1000$, где \mathbf{R} — реализация красного шума с параметрами $\varphi = 0.7$ и $\delta = 1$, $\mathbf{S} = \{A \cos(2\pi n\omega)\}_{n=1}^N$ — сигнал с амплитудой $A = 0.6$ и частотой $\omega = 0.1$.

На рис. 2.1 представлены собственные числа теплицева разложения (1.8) и соответствующие 80%-ные доверительные области в предположении, что сигнал отсутствует: по оси Ox отложена частота, которой соответствует вектор, на который производится проекция, по оси Oy отложен вклад этой частоты во временной ряд. Также черной линией проведена спектральная плотность красного шума. Если теперь рассмотреть $\varphi = 0.3$ и посмотреть на рис. 2.2, видно, что, значение спектральной плотности в точке 0.1 уменьшилось, и векторы, соответствующие этой частоте, стали более значимыми, то есть увеличилась мощность критерия.

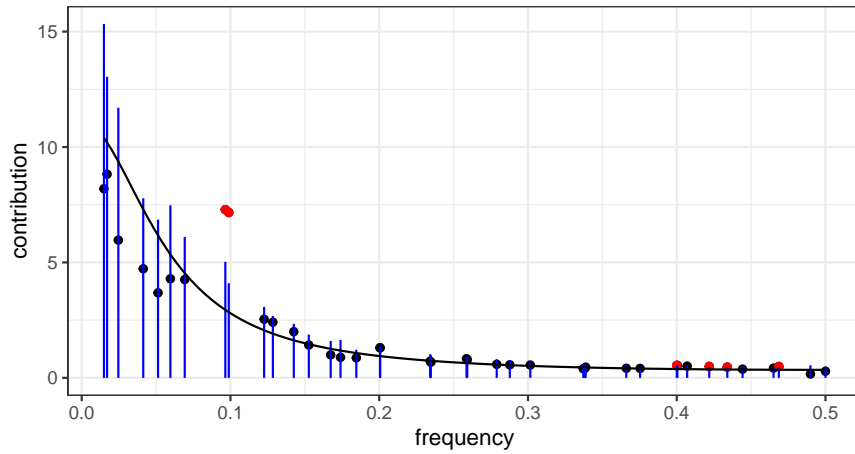


Рис. 2.1. Спектральная плотность при $\varphi = 0.7$ (сигнал на частоте $\omega = 0.1$)

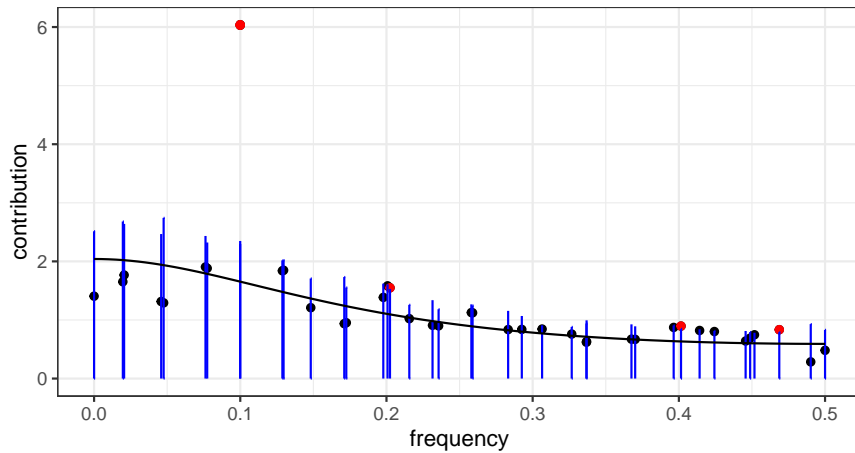


Рис. 2.2. Спектральная плотность при $\varphi = 0.3$ (сигнал на частоте $\omega = 0.1$)

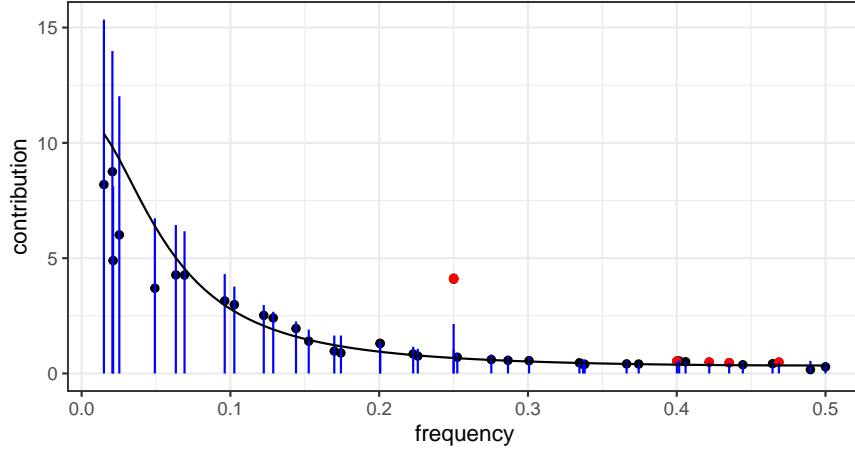


Рис. 2.3. Спектральная плотность при $\varphi = 0.7$ (сигнал на частоте $\omega = 0.25$)

При увеличении частоты ω сигнала значение спектральной плотности в этой точке уменьшается, что приводит к увеличению мощности критерия, это можно увидеть, сравнив рис. 2.1 и 2.3. \square

2.2.6. Численное сравнение MC-SSA с другими критериями

Рассмотрим другие критерии проверки гипотезы о том, что временной ряд состоит из красного шума против альтернативы с наличием сигнала. Пусть дан временной ряд $X = S + \xi$, где S — сигнал, ξ — красный шум с параметрами φ и δ . Большинство критериев проверяет гипотезу о том, что временной ряд — реализация белого шума. Поэтому будем действовать следующим образом:

1. Провести отбеливание: вычислить $Y = \Sigma^{-1/2}X$, где Σ — теоретическая автокорреляционная матрица красного шума с элементами $\varphi^{|i-j|}$, которая при верной H_0 совпадает с автокорреляционной матрицей ряда X .
2. Проверить гипотезу о том, что Y является реализацией белого шума, с применением коррекции при необходимости.

В качестве тестов на белый шум был взят Q-тест Бокса-Пирса [14] и тест с использованием вейвлетов [15] (будем далее называть их box и wavelet). Отметим, что второй тест применим только к рядам длины $N = 2^k$, где $k \in \mathbb{N}$. В связи с этим положим $N = 128$, параметры шума возьмем $\varphi = 0.7$, $\delta = 1$. За альтернативу возьмем $S = \{A \cos(2\pi n\omega)\}_{n=1}^N$ и сравним мощность методов с отбеливанием и MC-SSA при помощи ROC-кривых для разных ω . Для MC-SSA будем рассматривать теплицево разложение (1.8) траекторной

матрицы. В силу предположения 1 для больших частот будем брать меньшую амплитуду сигнала с целью избежать слишком мощных критериев.

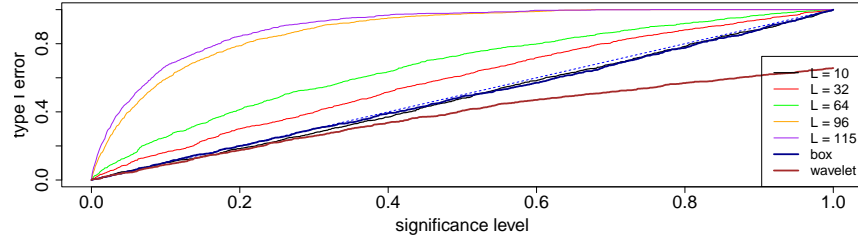
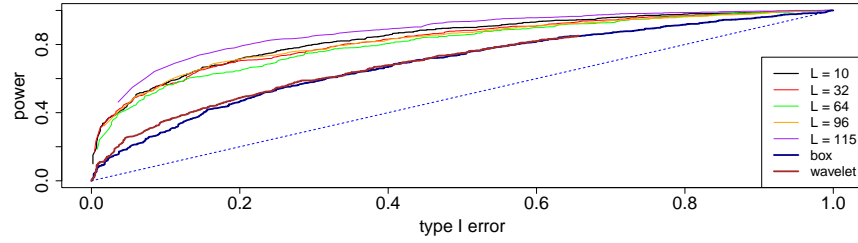
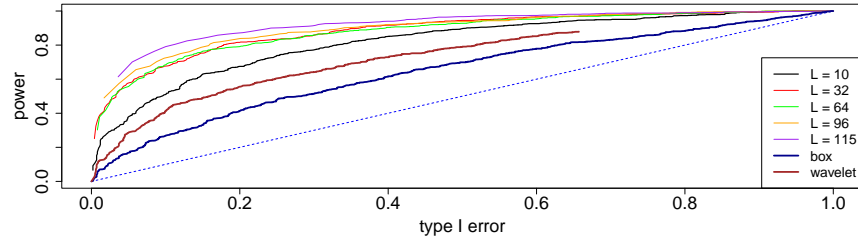


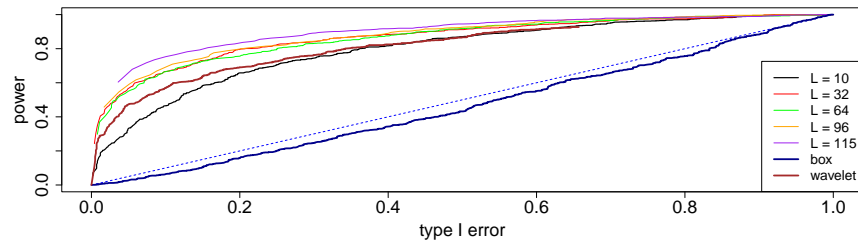
Рис. 2.4. Сравнение ошибки I рода с другими методами



(a) ROC-кривая ($A = 1.5$, $\omega = 0.025$)



(б) ROC-кривая ($A = 0.8$, $\omega = 0.125$)



(в) ROC-кривая ($A = 0.5$, $\omega = 0.225$)

Рис. 2.5. Сравнение ROC-кривых с другими методами

На рис. 2.5, а, 2.5, б, 2.5, в изображены ROC-кривые методов при разных ω . Отметим, что для wavelet построить ROC-кривые для больших ошибок I рода не удалось,

Таблица 2.1. Результаты численного сравнения MC-SSA с другими критериями ($\alpha^* = 0.1$)

Метод	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$	$\beta(\tilde{\alpha}^*)$	$\beta(\tilde{\alpha}^*)$
		$A = 1.5$	$A = 0.8$	$A = 0.5$
		$\omega = 0.025$	$\omega = 0.125$	$\omega = 0.225$
MC-SSA ($L = 10$)	0.101	0.57	0.51	0.465
MC-SSA ($L = 32$)	0.163	0.566	0.678	0.668
MC-SSA ($L = 64$)	0.25	0.556	0.684	0.665
MC-SSA ($L = 96$)	0.593	0.599	0.734	0.709
MC-SSA ($L = 115$)	0.668	0.668	0.791	0.753
box	0.103	0.289	0.269	0.064
wavelet	0.091	0.354	0.414	0.57

поскольку на рис. 2.4 видно, что $\alpha_I(\alpha) < 1 \forall \alpha$. Видно, при разработке критерия wavelet авторы делали его почти точным при $\alpha < 0.2$ и не заботились о том, что будет при бóльших значениях.

Для всех рассмотренных ω MC-SSA при всех длинах окна мощнее, чем box и wavelet, кроме случая с высокой частотой ($\omega = 0.225$), где wavelet оказался немного мощнее MC-SSA с $L = 10$. Отметим также, что на рис. 2.5, в ROC-кривая метода box лежит ниже прямой $y = x$, поэтому этот тест не имеет смысл применять для выявления высоких частот. Для удобства сравнения в таблице 2.1 для каждого критерия указана ошибка первого рода и мощность поправленного критерия для каждой рассмотренной альтернативы при уровне значимости $\alpha^* = 0.1$.

2.3. Monte Carlo MSSA

2.3.1. Отличия от одномерного случая

MC-SSA легко обобщается на многомерный случай: нужно просто заменить SSA на MSSA и генерировать красный шум с тем же количеством каналов, что и у исходного ряда [5].

Стоит отметить, что, в отличие от одномерного случая, левые и правые векторы матрицы отличаются по построению \mathbf{X} (1.3), поэтому в MC-MSSA в качестве векторов для проекции рассмотрены и левые, и правые векторы. Если W_1, \dots, W_H — левые

векторы матрицы \mathbf{X} , метод совпадает с алгоритмом 4. Если рассматривать в качестве векторов для проекции правые векторы, то в формулах (2.1) и (2.2) нужно заменить \mathbf{X} на \mathbf{X}^T и Ξ_i на Ξ_i^T соответственно.

2.3.2. Численное сравнение модификаций MC-MSSA

Поскольку рассматриваемый вариант критерия MC-(M)SSA является радикальным, нужно найти длину окна L , дающую максимально мощный критерий, но по замечанию 13 он не должен быть слишком радикальным, чтобы можно было применить поправку, не увеличивая трудоемкость метода.

Введем понятие условно равномошных критериев. Поскольку оценка мощности критерия $\hat{\beta}$ является долей отвергнутых нулевых гипотез при верной альтернативе, соответствующий доверительный интервал для настоящего значения β выглядит следующим образом:

$$\hat{\beta} \pm z_\alpha \sqrt{\frac{\hat{\beta}(1 - \hat{\beta})}{M}},$$

где z_α — $(1 - \alpha/2)$ -квантиль $N(0, 1)$ и M — количество моделирований. В худшем случае ($\beta = 0.5$) 95%-ный доверительный интервал при $M = 1000$ примерно равен $\hat{\beta} \pm 0.03$. Поэтому будем условно считать критерии равномошными, если их мощности будут отличаться не более чем на 0.03.

Как было показано в [7, Приложение Б.2.4], метод MC-SSA с проекцией на левые (правые) векторы SVD разложения матрицы \mathbf{X} (1.6) дает очень радикальный критерий для больших (малых) значений длины окна L , что делает невозможным построение поправки.

Однако, в одномерном случае было установлено [7], что если вместо SVD разложения матрицы \mathbf{X} использовать теплицево, то радикальность критерия уменьшается, и уже можно применить поправку. Установим, что будет в многомерном случае, если использовать модификации, описанные в разделе 1.3.

Пусть количество каналов $D = 2$, длина каждого канала равна $N = 100$. Рассмотрим модель $\mathbf{X} = \mathbf{S} + \boldsymbol{\xi}$, где $\boldsymbol{\xi}$ — красный шум с параметрами φ и $\delta = 1$, а \mathbf{S} — сигнал с элементами

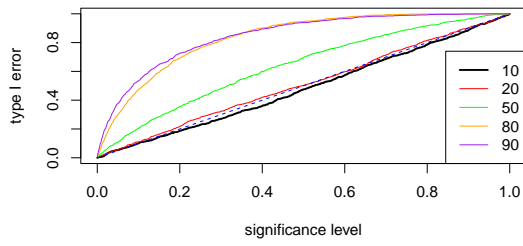
$$s_n^{(1)} = s_n^{(2)} = A \cos(2\pi\omega n), \quad n = 1, \dots, N.$$

Тогда нулевая гипотеза $H_0 : A = 0$ и альтернатива $H_1 : A \neq 0$.

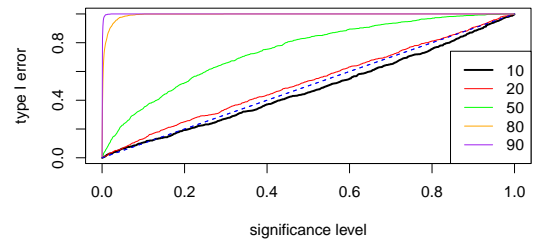
Будем смотреть на графики ошибок первого рода и ROC-кривые критериев для длин окна $L = 10, 20, 50, 80, 90$. Будем воспринимать ROC-кривую как график мощности критерия, к которому был применен алгоритм 2.

Рассмотрим несколько примеров. В данном разделе графики ошибок I рода и ROC-кривые представлены только для первого примера, для остальных графики можно найти в разделе А.1. Таблицы с результатами для каждого примера можно найти в разделе Б.1.

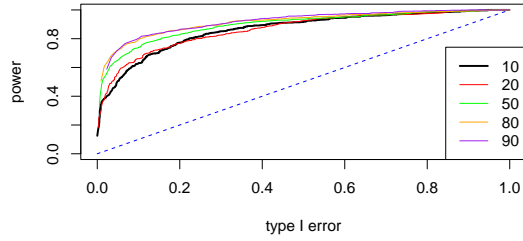
Пример 1. Пусть $\varphi = 0.7$ и параметры сигнала $A = 1, \omega = 0.075$.



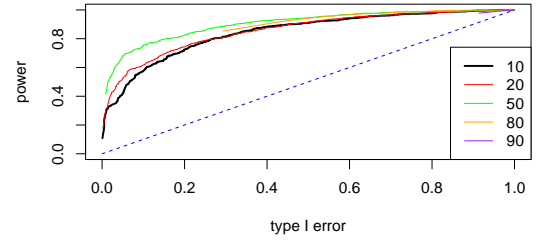
(a) Ошибка первого рода (Sum)



(б) Ошибка первого рода (SVD)

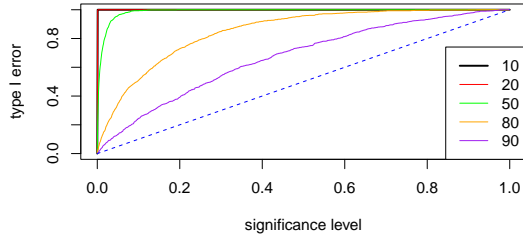


(в) ROC-кривая (Sum)

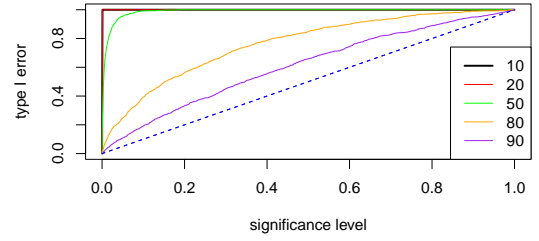


(г) ROC-кривая (SVD)

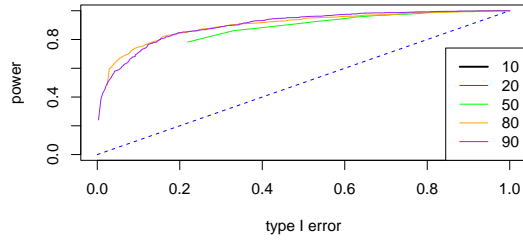
Рис. 2.6. Сравнение методов с проекцией на левые векторы ($\varphi = 0.7, A = 1, \omega = 0.075$)



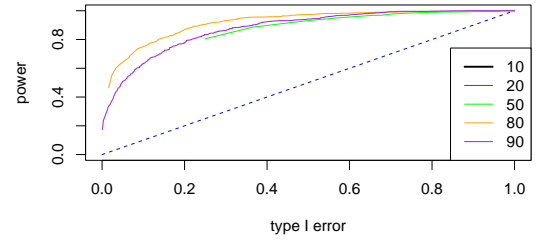
(a) Ошибка первого рода (Sum)



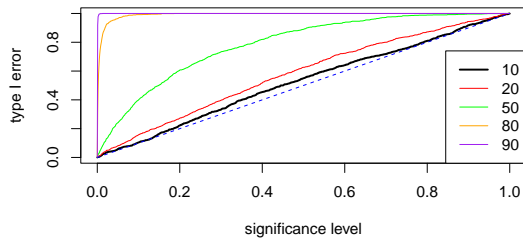
(б) Ошибка первого рода (SVD)



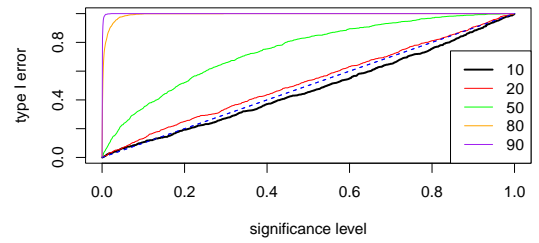
(в) ROC-кривая (Sum)



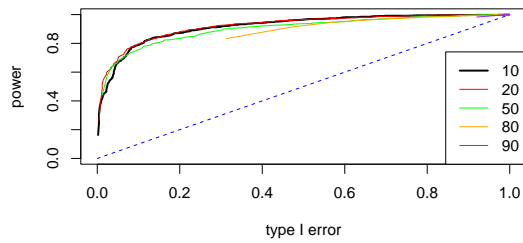
(г) ROC-кривая (SVD)

Рис. 2.7. Сравнение методов с проекцией на правые векторы ($\varphi = 0.7$, $A = 1$, $\omega = 0.075$)

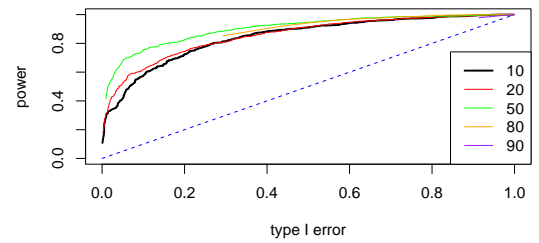
(a) Ошибка первого рода (Block)



(б) Ошибка первого рода (SVD)

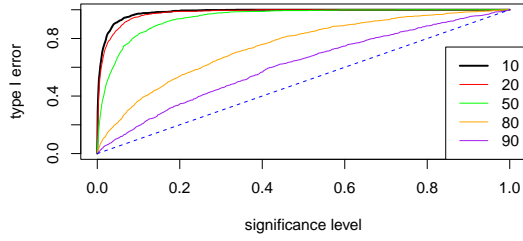


(в) ROC-кривая (Block)

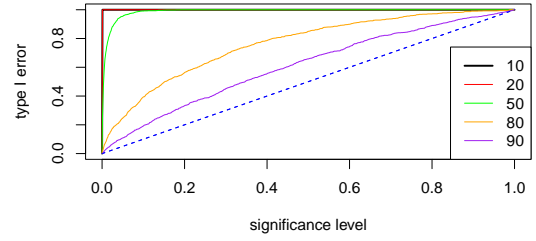


(г) ROC-кривая (SVD)

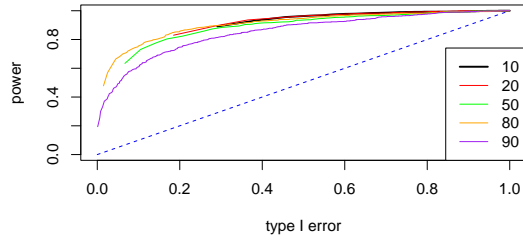
Рис. 2.8. Сравнение методов с проекцией на левые векторы ($\varphi = 0.7$, $A = 1$, $\omega = 0.075$)



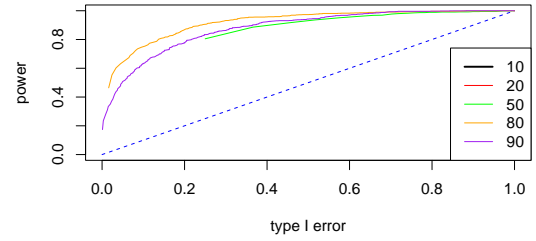
(a) Ошибка первого рода (Block)



(б) Ошибка первого рода (SVD)



(в) ROC-кривая (Block)



(г) ROC-кривая (SVD)

Рис. 2.9. Сравнение методов с проекцией на правые векторы ($\varphi = 0.7$, $A = 1$, $\omega = 0.075$)

На рис. 2.6 и 2.7 векторы для проекции были взяты из разложения (1.8). На рис. 2.6, а видно, что при $L > 20$ метод радикальный, а наибольшая мощность достигается при $L = 90$. На рис. 2.7, а отчетливо заметно, что метод радикальный для всех L . Наибольшая мощность наблюдается при $L = 80$, но отметим, что из-за слишком большой ошибки первого рода построить ROC-кривую на промежутке $[0, 0.22)$ для $L = 50$ и на всем промежутке для $L = 10$ и $L = 20$ не получилось.

На рис. 2.8 и 2.9 векторы для проекции были взяты из разложения (1.7). Если рассматривать проекцию на левые векторы, то на рис. 2.8, а видно, что метод радикальный, а наибольшая мощность достигается при $L = 20$. Проекция на правые векторы также дает радикальный критерий, как видно на рис. 2.9, а. Наибольшая мощность наблюдается при $L = 80$, но из-за слишком большой ошибки первого рода ROC-кривую для $L = 10$ и $L = 20$, для которых метод, предположительно, имеет большую мощность, удалось построить не на всем промежутке.

Пример 2. Пусть $\varphi = 0.3$ и $\omega = 0.075$. В виду предположения 1 уменьшим амплитуду сигнала для этого примера до $A = 0.5$.

Пример 3. В условиях примера 1 увеличим частоту сигнала до $\omega = 0.225$ и в силу предположения 1 уменьшим амплитуду сигнала до $A = 0.4$.

Таблица 2.2. Мощность методов для оптимальных L ($\alpha^* = 0.1$)

Метод	левые/правые векторы	$\beta(\tilde{\alpha}^*)$ (пример 1)	$\beta(\tilde{\alpha}^*)$ (пример 2)	$\beta(\tilde{\alpha}^*)$ (пример 3)
SVD	левые	0.754	0.399	0.573
SVD	правые	0.754	0.382	0.442
Block	левые	0.796	0.398	0.597
Block	правые	0.717	0.389	0.473
Sum	левые	0.806	0.421	0.625
Sum	правые	0.748	0.412	0.613

Таблица 2.3. Размеры матриц методов

Метод	левые/правые векторы	Размер матрицы (пример 1)	Размер матрицы (пример 2)	Размер матрицы (пример 3)
SVD	левые	50	10	20
SVD	правые	80	80	80
Block	левые	102	162	162
Block	правые	42	22	42
Sum	левые	80	20	80
Sum	правые	80	80	80

В таблице 2.2 представлены мощности поправленных критериев при уровне значимости $\alpha^* = 0.1$ для оптимальной длины окна для каждого рассмотренного примера, черным выделены наибольшая мощность и примерно равные ей (отличающиеся не более, чем на 0.03). В таблице 2.3 представлены размеры матриц методов, у которых находятся вектора для проекции для тех длин окна, мощность которых примерно равна мощности при оптимальном L , но критерии с такой длиной окна более эффективные (матрица меньшего размера). Из замечания 17 следует, что трудоемкость метода МС-

MSSA равна трудоемкости разложения матрицы выбранного метода, умноженной на M .

Выводы

Подведем итоги. Для рассмотренных примеров методы Block и Sum с проекцией на левые векторы показывают наибольшую мощность среди всех рассмотренных модификаций MC-MSSA, однако, учитывая таблицу 2.3 и замечания 5, 17, метод Sum во всех трех случаях численно эффективнее метода Block в нахождении векторов для проекции. Также хочется отметить, что метод Sum с проекцией на левые векторы дает наименее радикальный критерий, что важно при построении поправки (см. замечание 13). Поэтому на данный момент, беря в учет рассмотренные примеры, рекомендуется использовать метод Sum с проекцией на левые векторы с $L = 90$ для рядов длины $N = 100$.

Глава 3

Применение метода Monte Carlo SSA на практике

В главе 2 предполагалось, что параметры шума известны и нет мешающего сигнала (например, сезонности или тренда). Мешающий сигнал — это сигнал, который уже не нужно обнаруживать при проверке гипотезы о наличии сигнала. В этой главе рассмотрим случаи, которые более близки к реальным задачам.

В разделе 3.1 исследована зависимость радикальности и мощности MC-SSA от длины окна на различных примерах. В разделе 3.2 рассмотрен один из способов оценки параметров красного шума и численное сравнение MC-SSA с известными и оцененными параметрами в разделе 3.2.1. В разделе 3.3 приведены два алгоритма модификации MC-SSA с мешающим сигналом и проведено численное исследование методов для разных примеров мешающего сигнала. В разделе 3.4 рассмотрены примеры реальных временных рядов и их анализ с помощью SSA и MC-SSA.

В этой главе будем использовать MC-SSA с проекцией на векторы теплицева разложения (1.8) матрицы \mathbf{X} .

3.1. Зависимость радикальности и мощности от параметра L

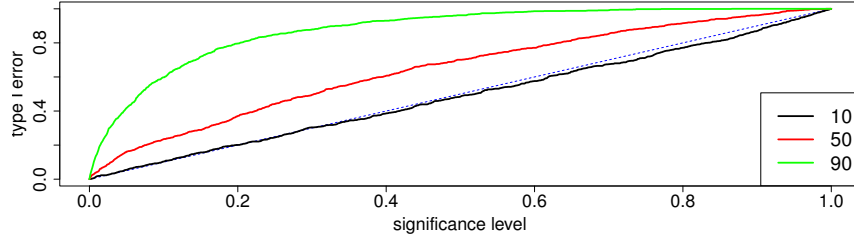
Поскольку рассматриваемый вариант критерия MC-SSA является радикальным, существует проблема выбора такой длины окна L , которая дает максимально мощный критерий, но при этом не слишком радикальный, чтобы можно было применить поправку из алгоритма 2. Однако, в зависимости от длины ряда N и параметров красного шума ξ наблюдаются разные зависимости мощности от L .

Рассмотрим несколько примеров. Пусть дана модель

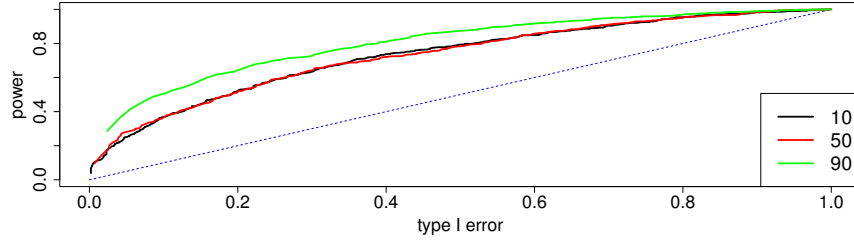
$$\mathbf{X} = \mathbf{S} + \xi,$$

где $\mathbf{S} = \{A \cos(2\pi\omega n)\}_{n=1}^N$, а ξ — красный шум с параметрами φ и $\delta = 1$. Рассмотрим следующие нулевую гипотезу и альтернативу: $H_0 : A = 0$, $H_1 : A \neq 0$. В этом разделе будем предполагать, что параметры красного шума известны. В первых трех примерах рассмотрим частоту сигнала $\omega = 0.075$.

Пример 4. Пусть $\varphi = 0.7$, $N = 100$. По графику ошибок первого рода на рис. 3.1, а



(a) Ошибка I рода

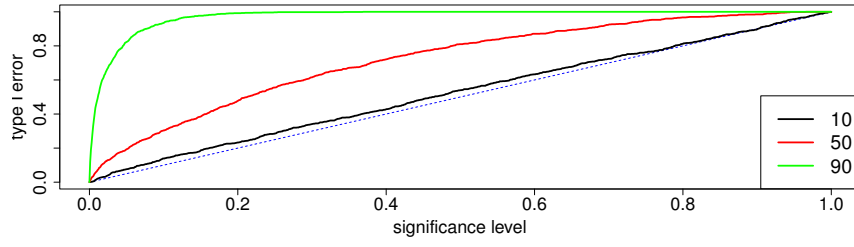
(б) ROC-кривая ($A = 1$, $\omega = 0.075$)Рис. 3.1. Пример 4 ($\varphi = 0.7$, $N = 100$)

видно, что чем больше L , тем более радикальным становится критерий. На рис. 3.1, б изображены ROC-кривые критериев, наибольшую мощность дает критерий с $L = 90$. На этом примере видно, что самым мощным является самый радикальный критерий.

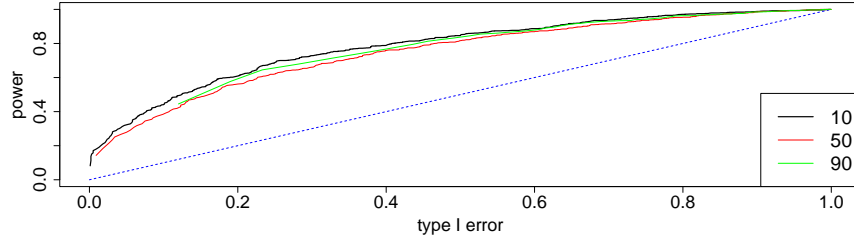
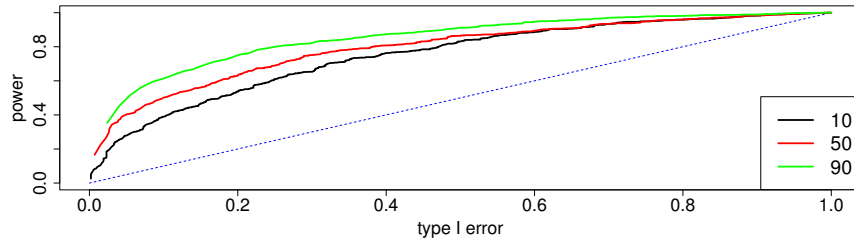
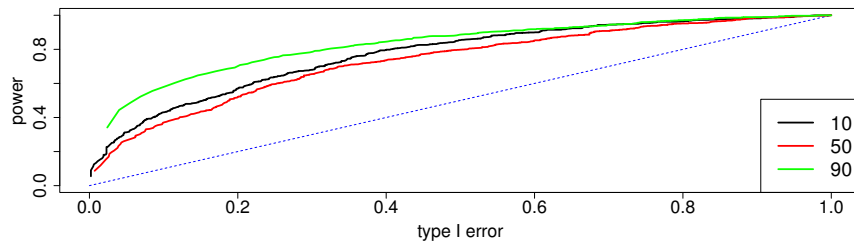
Пример 5. Пусть $\varphi = 0.3$, $N = 100$. На рис. 3.2, а изображен график ошибок первого рода. По нему видно, что, как и в примере 1, чем больше L , тем больше радикальность критерия. Если взглянуть на ROC-кривые на рис. 3.2, б, то видно, что с уменьшением параметра φ уменьшается разброс мощности критериев после поправки в зависимости от длины окна. Лучшей из рассмотренных в этом случае является $L = 10$, хотя разница с $L = 50$ совсем небольшая, а для $L = 90$ для небольших ошибок I рода поправку сделать не удалось из-за радикальности.

Объяснить такое поведение радикальности в зависимости от длины окна L можно теоретически: с увеличением L увеличивается размер автоковариационной матрицы и тем самым увеличивается количество оцениваемых параметров, что влечет за собой большую подгонку собственных векторов к конкретной реализации шума.

Пример 6. В условиях примера 4 рассмотрим разные частоты ω сигнала S и зависимость упорядоченности критериев по мощности от L . На рис. 3.3, б, 3.3, а изображены



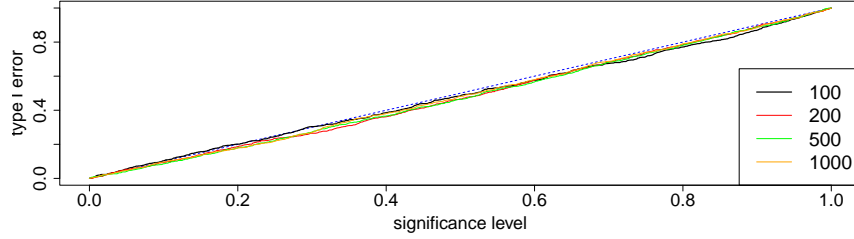
(a) Ошибка I рода

(б) ROC-кривая ($A = 0.7$, $\omega = 0.075$)Рис. 3.2. Пример 5 ($\varphi = 0.3$, $N = 100$)(a) ROC-кривая ($A = 0.6$, $\omega = 0.175$)(б) ROC-кривая ($A = 1.5$, $\omega = 0.025$)Рис. 3.3. Пример 6 ($\varphi = 0.7$, $N = 100$)

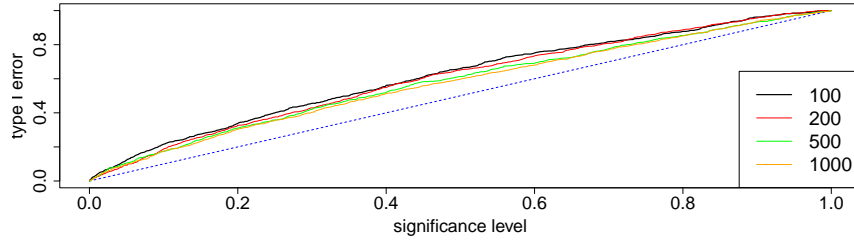
ROC-кривые критериев при разных альтернативах. Видно, что упорядоченность L нарушается при маленьких частотах сигнала. Если упорядочить рис. 3.3, б, рис. 3.1, б и рис. 3.3, а по частоте ω , то видна динамика по соотношению ROC-кривых для $L = 10$

и $L = 50$.

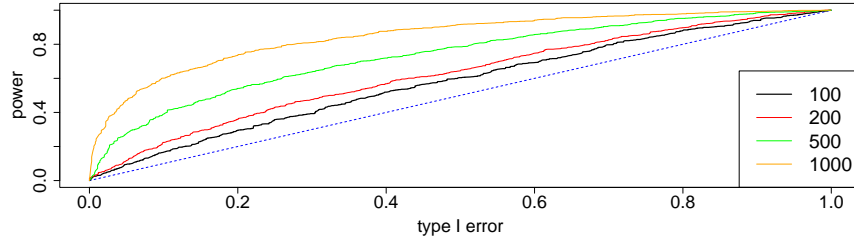
Пример 7. В условиях примера 4 рассмотрим $N = 100, 200, 500, 1000$ и посмотрим на графики ошибок I рода и ROC-кривые при $L = 10$ и $L = 40$.



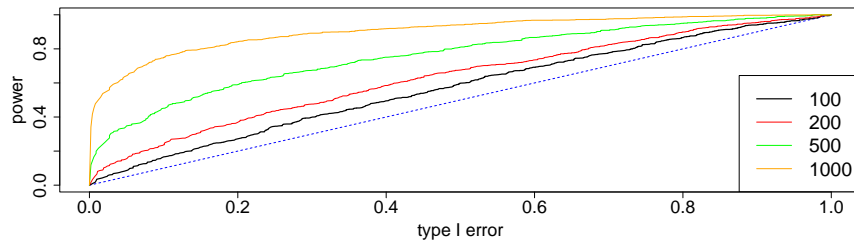
(a) Ошибка I рода ($L = 10$)



(б) Ошибка I рода ($L = 40$)



(в) ROC-кривая ($L = 10$)



(г) ROC-кривая ($L = 40$)

Рис. 3.4. Пример 7 ($\varphi = 0.7$, $A = 0.8$, $\omega = 0.075$, разные N)

Из рис 3.4, а видно, что критерий с $L = 10$ примерно точный для всех N , а на рис. 3.4, б наблюдается очень медленное уменьшение радикальности критерия с ростом

N . По ROC-кривым на рис. 3.4, *в* и 3.4, *г* видна состоятельность критерия против рассмотренной альтернативы.

Численные эксперименты показали, что длина окна L , дающая максимальную мощность критерия после поправки, зависит от параметров шума, длины ряда и, главное, от частоты сигнала в альтернативной гипотезе. Поэтому при выборе длины окна возможны следующие варианты:

1. При больших N применение поправки, описанной в разделе 2.1.2, является трудоемкой задачей даже для небольших L . Поэтому из рассмотренных примеров получено, что без поправки можно использовать MC-SSA только с $L = 10$. Это нетрудозатратно, но возможна некоторая потеря в мощности.
2. В поведении оптимальной по мощности длины окна L в зависимости от параметров ряда наблюдается некоторая регулярность. Поэтому можно построить зависимость оптимальной длины окна от параметров ряда с помощью численного моделирования, оценив параметры красного шума. Однако было показано, что упорядоченность критериев по мощности зависит от частоты сигнала в альтернативе, поэтому эта рекомендация имеет практический смысл, только если есть дополнительная информация о диапазоне возможных частот в альтернативе.

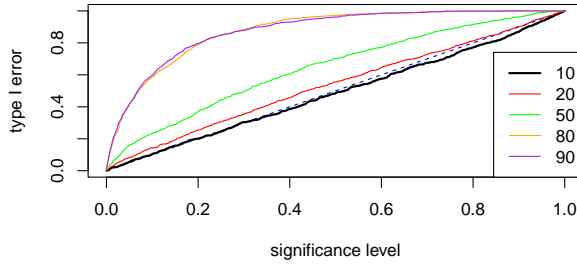
3.2. Оценка параметров красного шума

До сих пор предполагалось, что параметры красного шума φ и δ известны, но в реальных задачах редко возникает такая ситуация. В этой ситуации можно воспользоваться методом bootstrapping, который позволяет использовать оцененные параметры шума для построения критерия [6]. Тогда параметры оцениваются на основе исходного ряда методом максимального правдоподобия, где для нахождения начальных значений используется метод CSS [16].

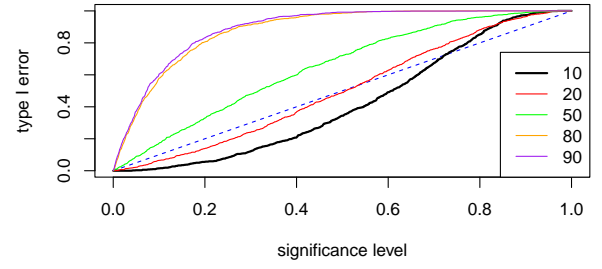
3.2.1. Искажение критерия при использовании оценок

В условиях примеров 1, 2 и 3 раздела 2.3.2 и посмотрим, насколько сильно искажается критерий MC-SSA, если оценивать параметры красного шума.

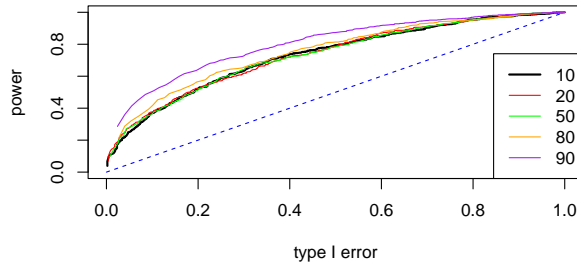
На рис. 3.5 представлено сравнение критерия с известными и критерия с оцененными параметрами красного шума. На рис. 3.5, *а* и 3.5, *б* изображены графики ошибок



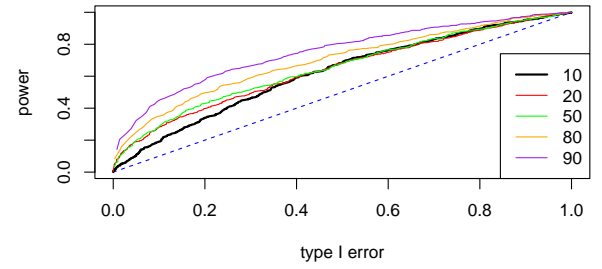
(a) Ошибка I рода (известные параметры)



(б) Ошибка I рода (оцененные параметры)



(в) ROC-кривая (известные параметры)



(г) ROC-кривая (оцененные параметры)

Рис. 3.5. Сравнение критериев ($\varphi = 0.7$, $A = 1$, $\omega = 0.075$)

первого рода для разных длин окна, а на рис. 3.5, в и 3.5, г — ROC-кривые. По ним видно, что оценка параметров снижает радикальность критерия, делая его консервативным при $L = 10$ и $L = 20$ для небольших α . Мощность поправленного критерия при этом тоже снижается. Графики для остальных примеров можно найти в разделе А.2.

В таблице 3.1 представлены результаты численного сравнения методов: длина окна, дающая наибольшую мощность, ошибка первого рода и мощность поправленного критерия при уровне значимости $\alpha^* = 0.1$. Как видно из таблицы, при маленьких φ оценка параметров довольно сильно искажает критерий: мощность для примера 2 сильно уменьшилась. Однако при больших φ уменьшение мощности не такое значительное.

3.3. Наличие мешающего сигнала

Пусть известно, что во временном ряде присутствует некоторый сигнал, но, возможно, еще есть какой-то другой. Тогда модель выглядит следующим образом:

$$X = F + S + \xi,$$

Таблица 3.1. Мощность методов для оптимальных L ($\alpha^* = 0.1$)

Пример 1 ($\varphi = 0.7, A = 1, \omega = 0.075$)	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Известные параметры	90	0.598	0.509
Оцененные параметры	90	0.605	0.454
Пример 2 ($\varphi = 0.3, A = 0.7, \omega = 0.075$)	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Известные параметры	80	0.814	0.487
Оцененные параметры	80	0.839	0.229
Пример 3 ($\varphi = 0.7, A = 0.4, \omega = 0.225$)	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Известные параметры	90	0.598	0.373
Оцененные параметры	90	0.605	0.352

где \mathbf{F} — мешающий сигнал, \mathbf{S} — неизвестный сигнал и $\boldsymbol{\xi}$ — красный шум. Тогда проверяется следующая нулевая гипотеза с альтернативой:

$$H_0 : \mathbf{S} = 0,$$

$$H_1 : \mathbf{S} \neq 0.$$

Алгоритм 5. MC-SSA с мешающим сигналом

1. Находится приближенное значение мешающего сигнала $\hat{\mathbf{F}}$ и оцениваются параметры $\boldsymbol{\xi}$ на основе остатка $\tilde{\mathbf{X}} = \mathbf{X} - \hat{\mathbf{F}}$.
2. Находятся левые векторы P_1, \dots, P_L траекторной матрицы временного ряда $\tilde{\mathbf{X}}$, полученные из разложения (1.8).
3. Применяется MC-SSA к исходному ряду \mathbf{X} с проекцией на векторы P_1, \dots, P_L , при этом суррогатными рядами являются реализации случайной величины $\boldsymbol{\eta}$:

$$\boldsymbol{\eta} = \boldsymbol{\xi} + \hat{\mathbf{F}}.$$

Цель данной модификации алгоритма — устранить влияние мешающего сигнала на вектора, на которые делается проекция.

Рассмотрим вместе с алгоритмом 5 другой вариант MC-SSA с мешающим сигналом.

Алгоритм 6. MC-SSA с мешающим периодическим сигналом

1. Находится приближенное значение мешающего сигнала \hat{F} и оцениваются параметры ξ на основе остатка $\tilde{X} = X - \hat{F}$.
2. Находятся левые векторы P_1, \dots, P_L траекторной матрицы временного ряда X , полученные из разложения (1.8).
3. Применяется MC-SSA к исходному ряду X с проекцией на векторы P_1, \dots, P_L , при этом суррогатными рядами являются реализации случайной величины η :

$$\eta = \xi + \hat{F}.$$

Замечание 18. Алгоритм 6 отличается от алгоритма 5 тем, что рассматривается проекция на левые векторы траекторной матрицы исходного временного ряда, а не его остатка. Поэтому, поскольку для нахождения собственных векторов используется теплицево разложение, использовать в качестве мешающего сигнала для алгоритма 6 можно только стационарный сигнал.

Рассмотрим два примера мешающего сигнала в условиях примеров 1 и 2 раздела 2.3.2.

3.3.1. Периодическая компонента

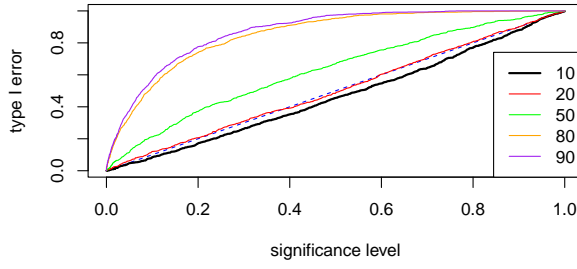
Рассмотрим в качестве мешающего сигнала синусоиду

$$f_n = A \cos(2\pi\omega n), \quad n = 1, \dots, N,$$

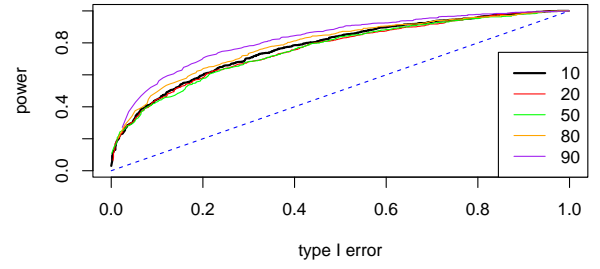
с амплитудой $A = 3$ и частотой $\omega = 0.25$.

Будем выделять периодическую компоненту при помощи SSA: будем оценивать доминирующую частоту левых векторов с помощью метода ESPRIT [17, Раздел 3.1] и на шаге группировки (раздел 1.2.3) будем брать две компоненты с наиболее близкими к ω частотами.

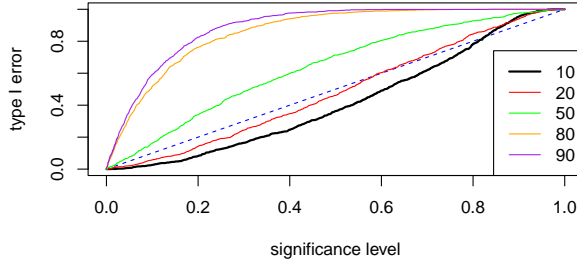
На рис. 3.6 представлены графики ошибок первого рода и ROC-кривые следующих случаев алгоритма 5: когда мешающий сигнал и параметры шума известны точно, когда F известен точно, но параметры шума оцениваются, и когда и мешающий сигнал, и параметры шума оцениваются. Графики ошибок первого рода на рис. 3.6, а, 3.6, в и 3.6, д похожи друг на друга, а отклонение от случая, когда все известно, можно объяснить погрешностью при оценке неизвестных параметров. После применения поправок из раздела 2.1.2 критерии становятся примерно точными для любой длины окна и



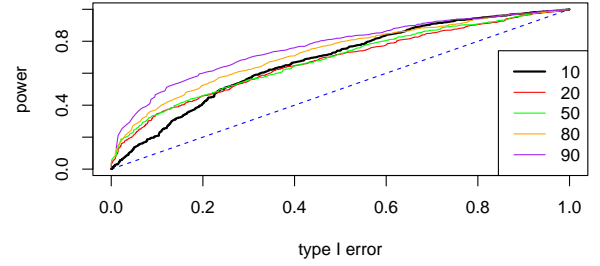
(a) Ошибка I рода



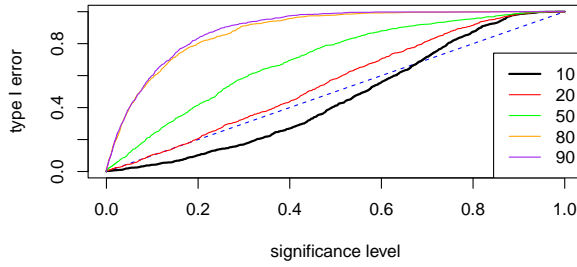
(б) ROC-кривая



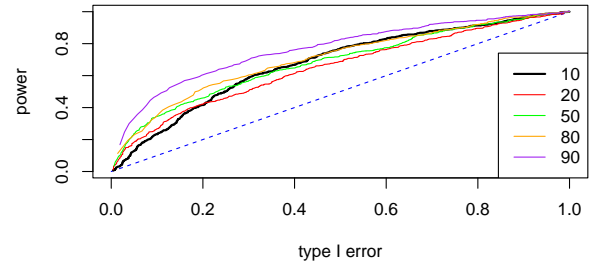
(в) Ошибка I рода (оцененные параметры шума)



(г) ROC-кривая (оцененные параметры шума)



(д) Ошибка I рода (оцененный мешающий сигнал и параметры шума)



(е) ROC-кривая (оцененный мешающий сигнал и параметры шума)

Рис. 3.6. Анализ алгоритма 5 (мешающий сигнал — периодика) ($\varphi = 0.7$, $A = 1$, $\omega = 0.075$)

ROC-кривые на рис. 3.6, б, 3.6, г и 3.6, е представляют собой графики мощности этих критериев. Таким образом, наибольшая мощность во всех трех случаях достигается при $L = 90$, однако заметно снижение мощности при оценивании неизвестных параметров. Графики для остальных примеров можно найти в разделе А.3.1.

В таблице 3.2 представлены результаты сравнения двух алгоритмов, а именно оптимальная длина окна, ошибка первого рода и мощность поправленного критерия при уровне значимости $\alpha^* = 0.1$. Таблицы сравнения для остальных примеров можно найти

Таблица 3.2. Сравнение алгоритма 5 и алгоритма 6 при $\alpha^* = 0.1$ ($\varphi = 0.7$, $A = 1$, $\omega = 0.075$)

Алгоритм 5	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	90	0.57	0.542
Оцененные параметры шума	90	0.593	0.48
Оцененные параметры шума и мешающий сигнал	90	0.6	0.475
Алгоритм 6	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	90	0.594	0.532
Оцененные параметры шума	90	0.588	0.468
Оцененные параметры шума и мешающий сигнал	90	0.624	0.521

в разделе Б.2.

3.3.2. Тренд

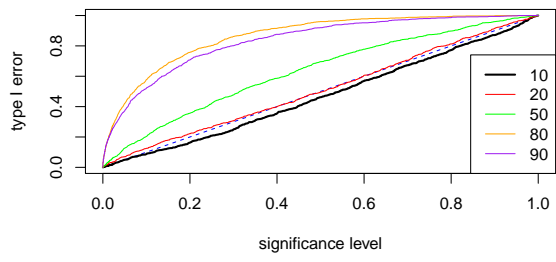
Отдельно рассмотрим вариант, когда мешающий сигнал — тренд, т.е. медленно меняющаяся компонента. Рассмотрим следующий экспоненциальный ряд:

$$f_n = Ae^{\alpha n}, \quad n = 1, \dots, N,$$

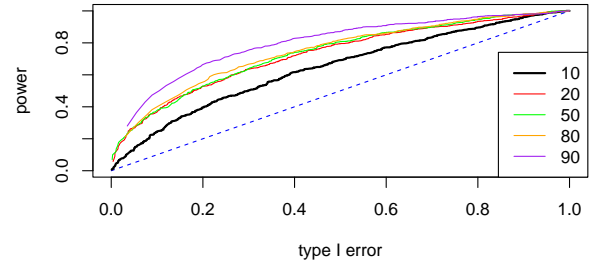
где $A = 0.2$, $\alpha = 0.05$.

Выделять тренд будем с помощью SSA: поскольку в сингулярном разложении (1.6) сингулярные числа, соответствующие тренду, будут самыми большими среди всех сингулярных чисел, на шаге группировки (раздел 1.2.3) будем брать первые r элементарных компонент, где r — ранг тренда. В данном случае $r = 1$.

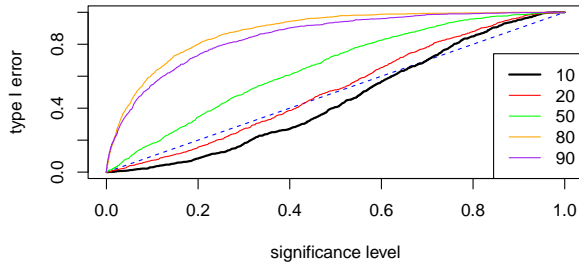
На рис. 3.7 представлены графики ошибок первого рода и ROC-кривые следующих критериев: когда тренд и параметры шума φ и δ известны точно, когда тренд известен точно, но параметры шума оцениваются, и когда и тренд, и параметры шума оцениваются. Как и в разделе 3.3.1, графики ошибок первого рода на рис. 3.7, а, 3.7, в и 3.7, д сохраняют общую тенденцию при оценке неизвестных параметров. По ROC-кривым на рис. 3.7, б, 3.7, г и 3.6, е видно, что оценка неизвестных параметров снижает мощность, но оптимальной длиной окна для этого примера является в любом случае $L = 90$. Графики для остальных примеров можно найти в разделе А.3.2.



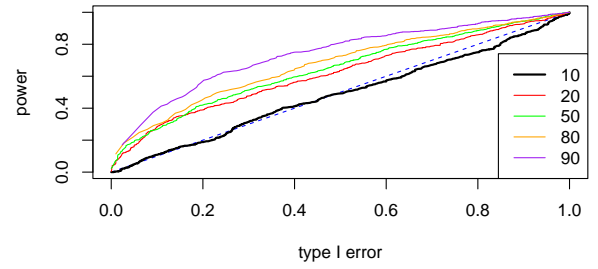
(a) Ошибка I рода



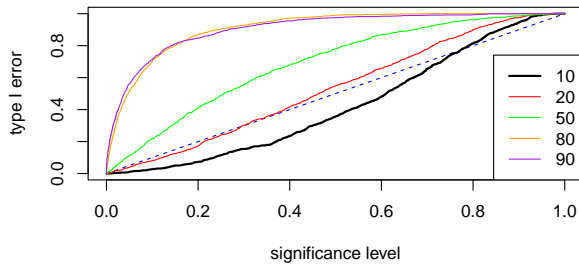
(б) ROC-кривая



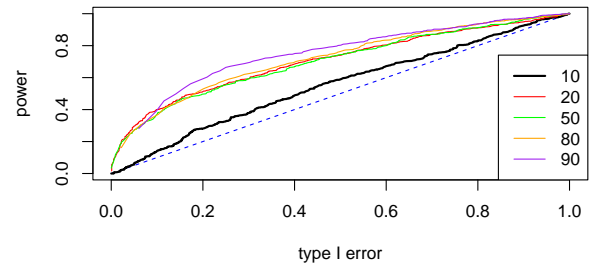
(в) Ошибка I рода (оцененные параметры шума)



(г) ROC-кривая (оцененные параметры шума)



(д) Ошибка I рода (оцененный мешающий сигнал и параметры шума)



(е) ROC-кривая (оцененный мешающий сигнал и параметры шума)

Рис. 3.7. Анализ алгоритма 5 (мешающий сигнал — тренд) ($\varphi = 0.7$, $A = 1$, $\omega = 0.075$)

В таблице 3.3 представлены оптимальная длина окна, ошибка первого рода и мощность поправленного критерия при уровне значимости $\alpha^* = 0.1$ для всех рассмотренных примеров. Как видно из таблицы, оценка тренда значительно увеличивает радикальность критерия, мощность при этом в случае $\varphi = 0.7$ падает примерно на 10%, а в случае $\varphi = 0.3$ уже довольно значительно, как и в разделе 3.2. Также отметим, что в случае $\varphi = 0.3$ критерий слишком радикальный и построить ROC-кривую для больших длин окна не удалось полностью.

Таблица 3.3. Результаты алгоритма 5 (мешающий сигнал — тренд) при $\alpha^* = 0.1$

Пример 1 ($\varphi = 0.7, A = 1, \omega = 0.075$)	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	90	0.521	0.501
Оцененные параметры шума	90	0.546	0.413
Оцененные параметры шума и мешающий сигнал	90	0.714	0.389
Пример 2 ($\varphi = 0.3, A = 0.7, \omega = 0.075$)	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	50	0.304	0.416
Оцененные параметры шума	50	0.255	0.223
Оцененные параметры шума и мешающий сигнал	50	0.358	0.243
Пример 3 ($\varphi = 0.7, A = 0.4, \omega = 0.225$)	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	90	0.521	0.393
Оцененные параметры шума	90	0.546	0.351
Оцененные параметры шума и мешающий сигнал	80	0.613	0.327

3.4. Применение к реальным данным

На рис. 3.8 представлена ежемесячная температура поверхности моря в центральной тропической части Тихого океана в период с 1950 по 2024 год (888 месяцев). В данном регионе происходит явление под названием Эль-Ниньо, характеризующееся аномальным потеплением поверхностных вод. Эти колебания температуры оказывают заметное влияние на погодные условия во всем мире, поэтому важно изучить их поведение.

Сразу заметим, что в этом временном ряде присутствует небольшой тренд, поэтому перед применением MC-SSA удалим его. Для начала воспользуемся базовым SSA с длиной окна $L = N/2 = 444$, как рекомендуется в [2]. На рис. 3.9 изображены первые 6 собственных векторов сингулярного разложения. Видно, что первый вектор соответствует тренду. Посмотрев на двумерные графики собственных векторов на рис. 3.10, видно, что вторая и третья компоненты образуют двенадцатиугольник. Это означает, что они соответствуют периодике с периодом 12 [2]. С учетом всей полученной инфор-

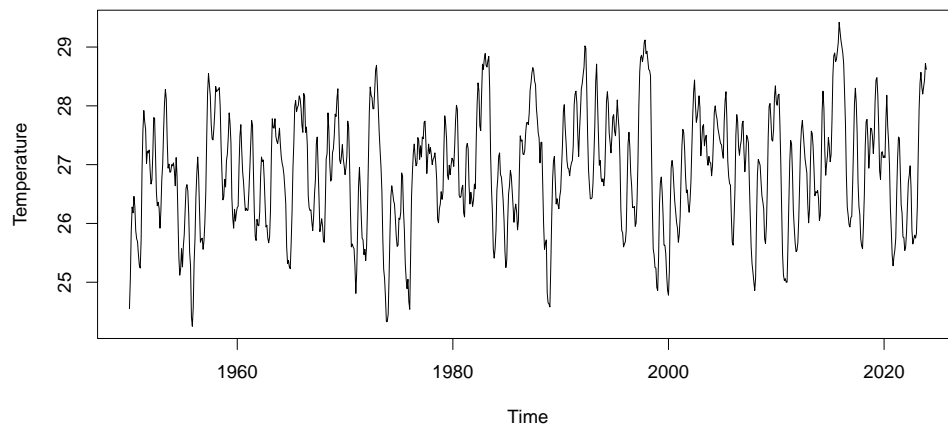


Рис. 3.8. Температура поверхности моря в центральной тропической части Тихого океана

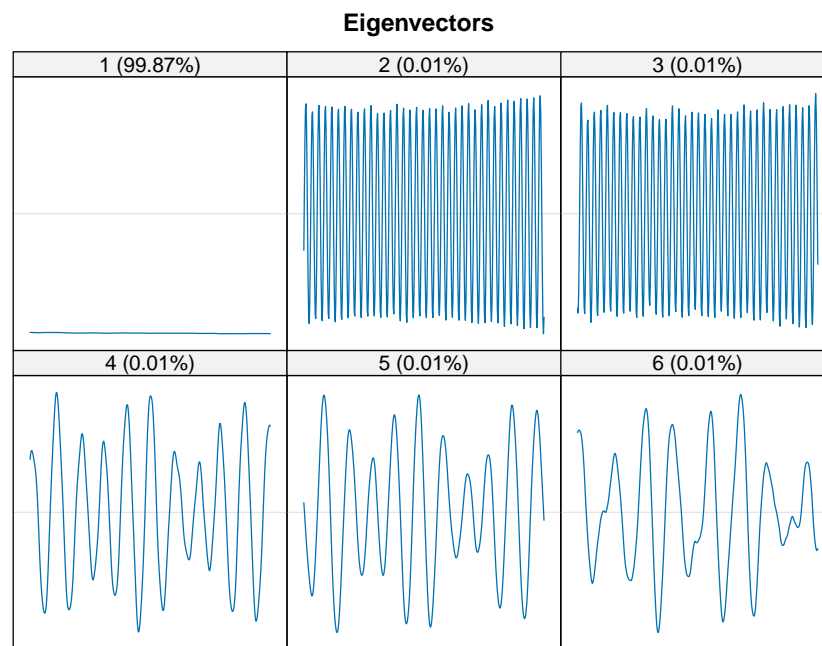


Рис. 3.9. Графики собственных векторов

мации, возьмем длину окна для выделения тренда небольшой, но делящейся на период периодической компоненты для обеспечения разделимости [2]. На рис. 3.11 изображен выделенный тренд при $L = 120$. Для наиболее точного выделения периодической компоненты длина окна должна быть близкой к половине длины ряда и должна делиться

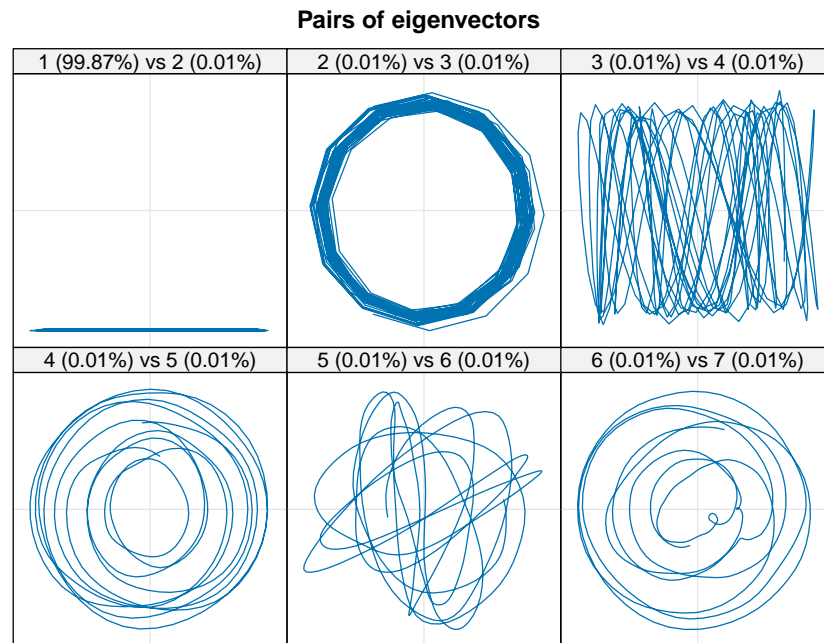


Рис. 3.10. Двумерные графики собственных векторов

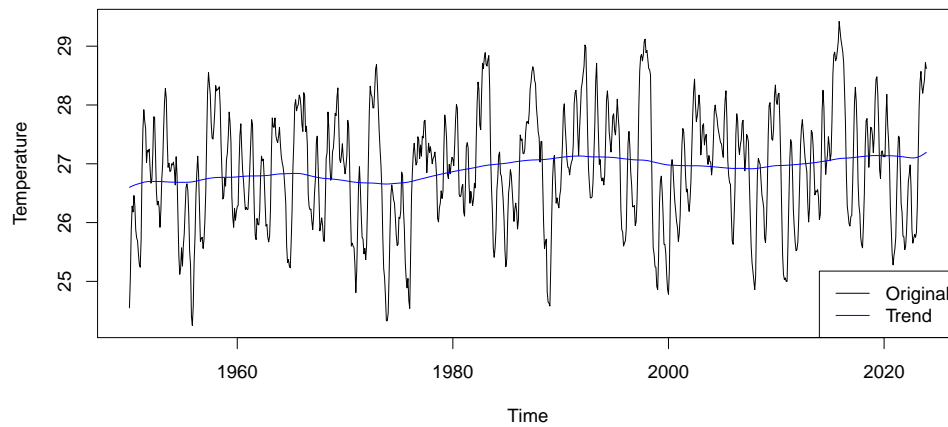


Рис. 3.11. Выделенный тренд

на ее период [2], поэтому применим Toeplitz SSA с длиной окна $L = 444$. На рис. 3.12 изображена выделенная годовая сезонность.

Применим поправленный MC-SSA с $L = 40$ к ряду без тренда с годовой периодичностью в качестве мешающего сигнала (алгоритм 6). Оцененные параметры красного

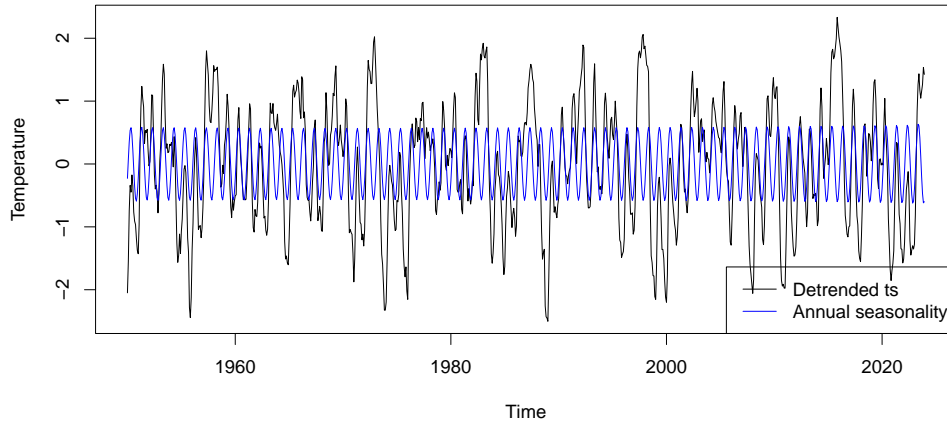


Рис. 3.12. Выделенная годовая сезонность

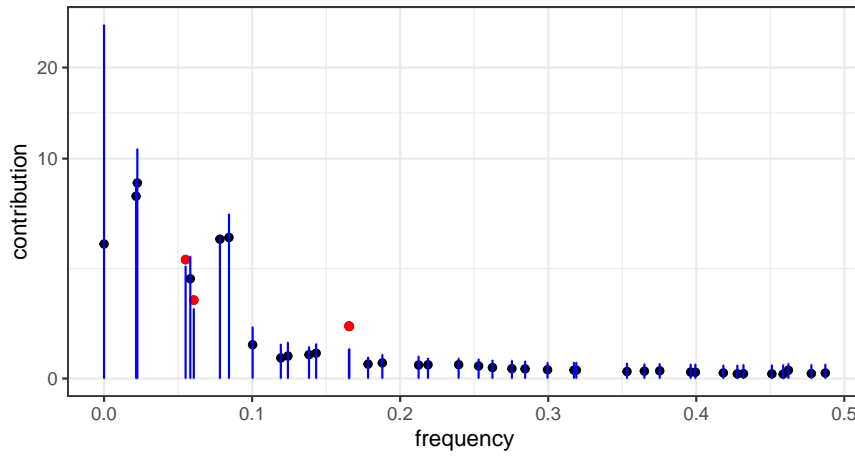


Рис. 3.13. Результат работы MC-SSA

шума следующие: $\varphi = 0.94$ и $\delta = 0.305$. Получено p-value, равное 0. На рис. 3.13 изображены 95%-ные доверительные интервалы статистики \hat{p}_k , $k = 1, \dots, L$ (2.1). Значимыми являются четыре компоненты, две компоненты, имеющие период приблизительно 6, легко интерпретируются — это замеченная полугодовая сезонность. С помощью Toeplitz SSA с той же длиной окна эта сезонность была выделена, ее вид изображен на рис 3.14. Оставшиеся значимые компоненты имеют периоды 18.16 и 16.5, которые довольно сложно интерпретировать.

Значимые векторы, интерпретация которых не представляется возможной, нельзя относить к сигналу, поскольку MC-SSA проверяет гипотезу о том, что временной ряд представляет собой реализацию красного шума, то есть возможно модель является неверной. Предположим, что рассматриваемый ряд без тренда и годовой периодичности

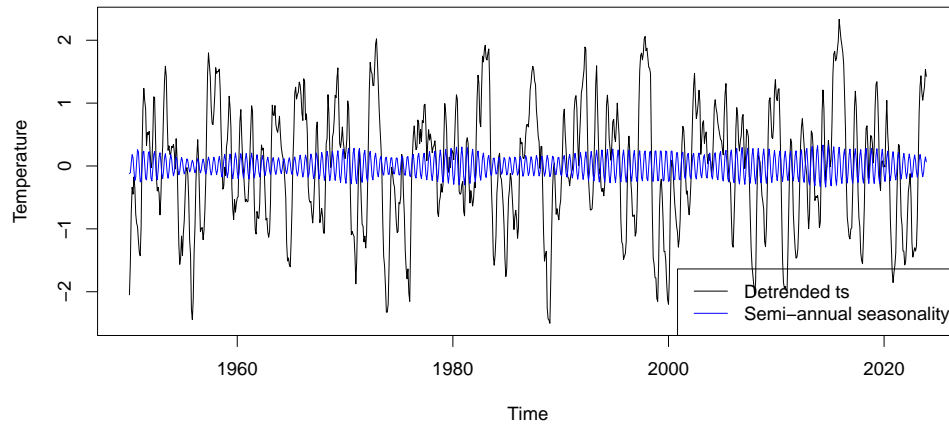


Рис. 3.14. Выделенная полугодовая сезонность

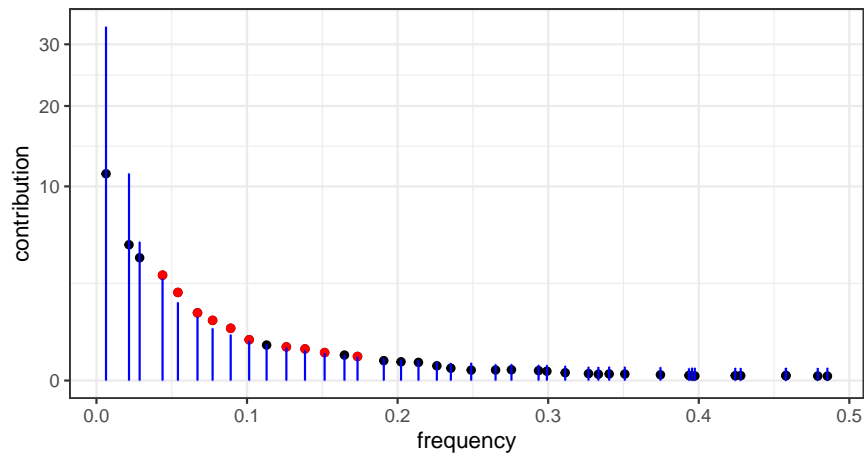


Рис. 3.15. Результат работы MC-SSA для модели ARMA(1, 2)

является моделью $ARMA(p, q)$, где p — порядок авторегрессии, q — порядок скользящего среднего. Тогда наиболее подходящей моделью является $ARMA(1, 2)$ [18]. Отметим, что красному шуму соответствует модель $ARMA(1, 0)$. Промоделировав ряд в соответствии с полученной моделью и посмотрев на доверительные интервалы статистик \hat{p}_k на рис. 3.15, получаем много значимых компонент. Таким образом, неправильно выбранная модель может исказить выводы, полученные в результате применения MC-SSA, поэтому важно внимательно относиться к выбору модели при проверке гипотезы.

Заключение

В ходе данной работы были реализованы два метода Toeplitz MSSA. На основе численных экспериментов были сделаны следующие выводы: Sum и Block версии Toeplitz MSSA для стационарного ряда точнее выделяют гармонический сигнал, чем Basic MSSA, причем метод Toeplitz Sum MSSA для рассмотренных примеров оказался наиболее эффективным в построении автоковариационной матрицы и нахождении ее спектрального разложения.

На основе методов Sum и Block Toeplitz MSSA были реализованы модификации критерия MC-MSSA. Было получено, что во всех рассмотренных примерах метод Toeplitz Sum MC-MSSA с проекцией на левые векторы дает самый мощный критерий, причем не слишком радикальный. Помимо этого, он численно эффективнее по сравнению с другими модификациями MC-MSSA. Поэтому этот метод более предпочтителен, чем метод Block, что важно ввиду его простоты в реализации и структуры, подходящей под пакет Rssa [10].

Проведено численное сравнение MC-SSA с другими критериями проверяющими гипотезу об отсутствии сигнала в красном шуме. Получено, что MC-SSA намного мощнее других критериев, особенно при малых частотах сигнала. Учитывая вычислительную трудоемкость метода, этот результат показывает, что она оправдана.

Была исследована зависимость радикальности и мощности MC-SSA от длины окна, однако оказалось, что длина окна, дающая максимальную мощность критерия после поправки, зависит от частоты сигнала в альтернативной гипотезе. Поэтому рекомендации по выбору параметра были даны только в общем виде.

Также было исследовано поведение MC-SSA в случае оценивании неизвестных параметров красного шума на основе исходного временного ряда. Получено, что степень искажения критерия зависит от параметра авторегрессии, чем он меньше, тем больше уменьшение мощности метода. Помимо оценки параметров было реализовано два алгоритма MC-SSA с мешающим сигналом и были разобраны два варианта мешающего сигнала. Для каждого примера были рассмотрены три случая: когда мешающий сигнал и параметры красного шума известны точно, когда параметры шума оцениваются, и когда вместе с параметрами оценивается мешающий сигнал.

При анализе алгоритмов была выявлена и сформулирована проблема применения

поправки, делающей радикальный критерий точным, для критериев, использующих суррогатные данные. На основе этого было показано, что радикальность критерия приводит к существенному увеличению численных затрат, в связи с чем при сравнении критериев нужно учитывать не только мощность, но и радикальность, выбирая менее радикальные критерии.

Исходный код предложенных в работе алгоритмов опубликован в [19].

Список литературы

1. Broomhead D. S., King G. P. Extracting qualitative dynamics from experimental data // *Physica D: Nonlinear Phenomena*. — 1986. — Vol. 20, no. 2–3. — P. 217–236.
2. Golyandina N., Nekrutkin V., Zhigljavsky A. *Analysis of Time Series Structure*. — Chapman and Hall/CRC, 2001. — ISBN: 9780367801687.
3. Allen M. R., Smith L. A. Monte Carlo SSA: Detecting irregular oscillations in the Presence of Colored Noise // *Journal of Climate*. — 1996. — Vol. 9, no. 12. — P. 3373–3404.
4. Бояров А. А. Исследование статистических свойств метода Монте-Карло SSA : дипломная работа ; СПбГУ. — 2012.
5. Groth A., Ghil M. Monte Carlo Singular Spectrum Analysis (SSA) Revisited: Detecting Oscillator Clusters in Multivariate Datasets // *Journal of Climate*. — 2015. — Vol. 28, no. 19. — P. 7873–7893.
6. Golyandina N. Detection of signals by Monte Carlo singular spectrum analysis: multiple testing // *Statistics and Its Interface*. — 2023. — Vol. 16, no. 1. — P. 147–157.
7. Ларин Е. С. Метод SSA для проверки гипотезы о существовании сигнала во временном ряде : квалификационная работа магистра ; СПбГУ. — 2022.
8. Spectral characteristics and predictability of the NAO assessed through Singular Spectral Analysis / Gámiz-Fortis S. R., Pozo-Vázquez D., Esteban-Parra M. J., and Castro-Díez Y. // *Journal of Geophysical Research: Atmospheres*. — 2002. — Vol. 107, no. D23.
9. Paluš M., Novotná D. Enhanced Monte Carlo Singular System Analysis and detection of period 7.8 years oscillatory modes in the monthly NAO index and temperature records // *Nonlinear Processes in Geophysics*. — 2004. — Vol. 11, no. 5/6. — P. 721–729.
10. Rssa: A Collection of Methods for Singular Spectrum Analysis. — R package version 1.0.5. Access mode: <https://CRAN.R-project.org/package=Rssa>.
11. Multivariate and 2D Extensions of Singular Spectrum Analysis with theRssaPackage / Golyandina N., Korobeynikov A., Shlemov A., and Usevich K. // *Journal of Statistical Software*. — 2015. — Vol. 67, no. 2.
12. Plaut G., Vautard R. Spells of Low-Frequency Oscillations and Weather Regimes in the Northern Hemisphere // *Journal of the Atmospheric Sciences*. — 1994. — Vol. 51, no. 2. — P. 210–236.

13. Korobeynikov A. Computation- and Space-Efficient Implementation of SSA // Statistics and Its Interface. — 2010. — Vol. 3, no. 3. — P. 257–268.
14. Box G. E. P., Pierce D. A. Distribution of Residual Autocorrelations in Autoregressive-Integrated Moving Average Time Series Models // Journal of the American Statistical Association. — 1970. — Vol. 65, no. 332. — P. 1509–1526.
15. Nason G. P., Savchev D. White noise testing using wavelets // Stat. — 2014. — Vol. 3, no. 1. — P. 351–362.
16. Gardner G., Harvey A. C., Phillips G. D. A. Algorithm AS 154: An Algorithm for Exact Maximum Likelihood Estimation of Autoregressive-Moving Average Models by Means of Kalman Filtering // Applied Statistics. — 1980. — Vol. 29, no. 3. — P. 311.
17. Golyandina N., Korobeynikov A., Zhigljavsky A. Singular Spectrum Analysis with R. Use R! — Springer Berlin Heidelberg, 2018. — ISBN: 9783662573808.
18. Hyndman R. J., Khandakar Y. Automatic Time Series Forecasting: TheforecastPackage forR // Journal of Statistical Software. — 2008. — Vol. 27, no. 3.
19. Poteshkin E. R-scripts for Toeplitz MSSA and Monte Carlo MSSA. — 2024. — Access mode: <https://doi.org/10.5281/zenodo.11372757>.

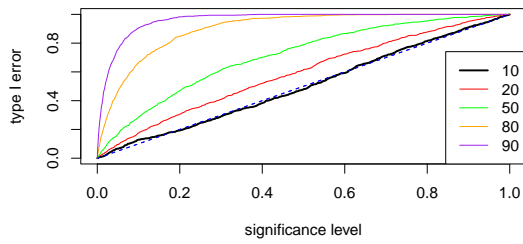
Приложение А

Графики

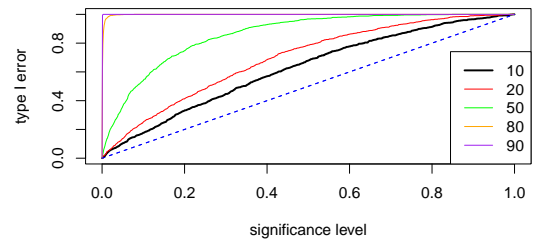
А.1. Численное сравнение модификаций MC-MSSA

Здесь приведены графики к разделу 2.3.2 основного текста.

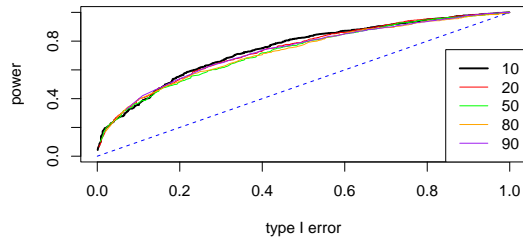
$$\varphi = 0.3, \omega = 0.075$$



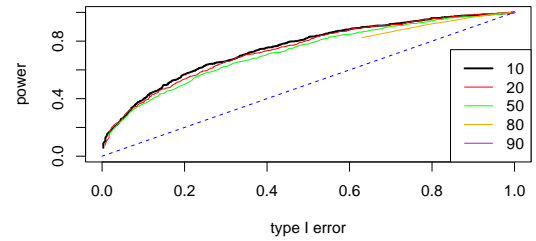
(а) Ошибка I рода (Sum)



(б) Ошибка I рода (Block)



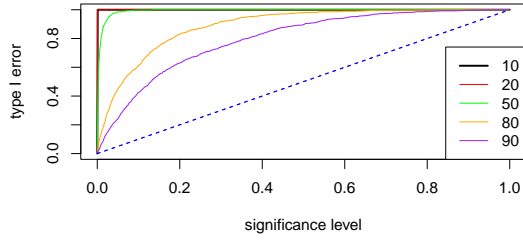
(в) ROC-кривая (Sum)



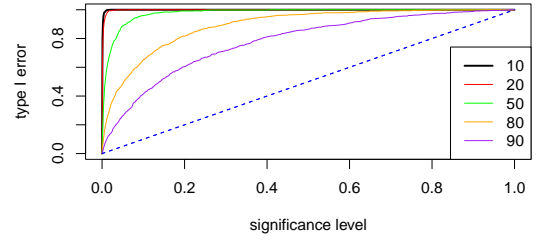
(г) ROC-кривая (Block)

Рис. А.1. Сравнение методов с проекцией на левые векторы ($\varphi = 0.3$, $A = 0.5$, $\omega = 0.075$)

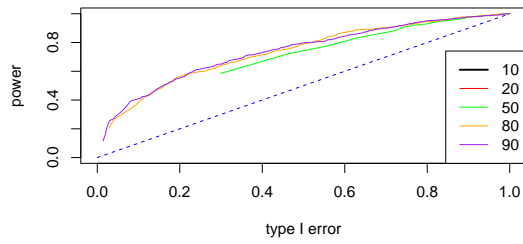
На рис. А.1 изображены графики ошибок первого рода и ROC-кривые методов Sum и Block с проекцией на левые векторы, а на рис. А.2 — на правые векторы. Как видно по рисункам, при уменьшении параметра φ , во-первых, увеличивается радикальность критериев для любой длины окна, а во-вторых уменьшается разброс мощности после поправки. Отметим, что наименее радикальным критерием среди рассмотренных является Sum с проекцией на левые векторы (рис. А.1, а).



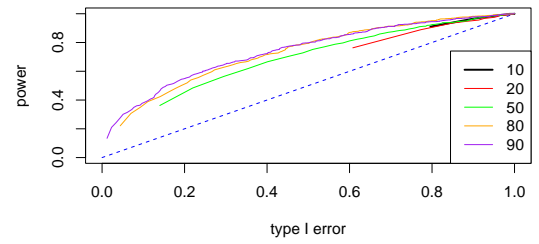
(a) Ошибка I рода (Sum)



(б) Ошибка I рода (Block)



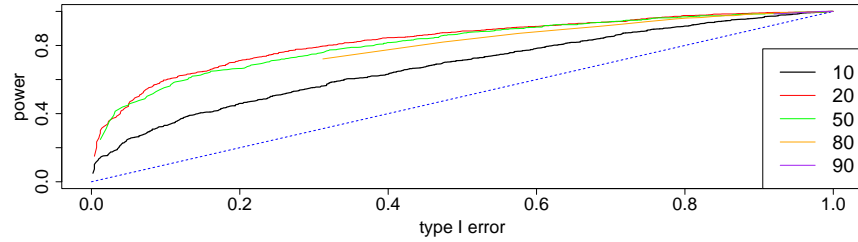
(в) ROC-кривая (Sum)



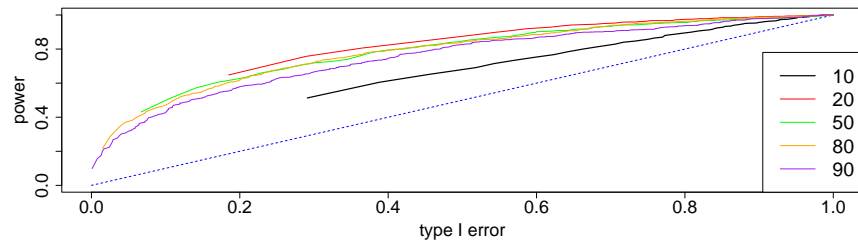
(г) ROC-кривая (Block)

Рис. А.2. Сравнение методов с проекцией на правые векторы ($\varphi = 0.3$, $A = 0.5$, $\omega = 0.075$)

$\varphi = 0.7$, $\omega = 0.225$

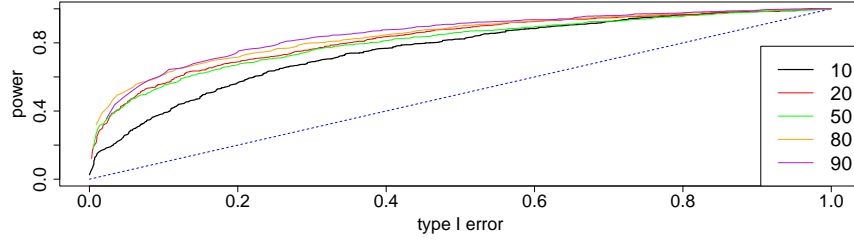


(a) Проекция на левые векторы

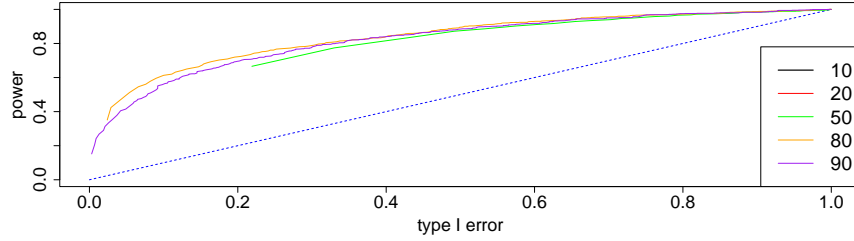


(б) Проекция на правые векторы

Рис. А.3. ROC-кривые метода Block ($\varphi = 0.7$, $A = 0.4$, $\omega = 0.225$)



(a) Проекция на левые векторы



(б) Проекция на правые векторы

Рис. А.4. ROC-кривые метода Sum ($\varphi = 0.7$, $A = 0.4$, $\omega = 0.225$)

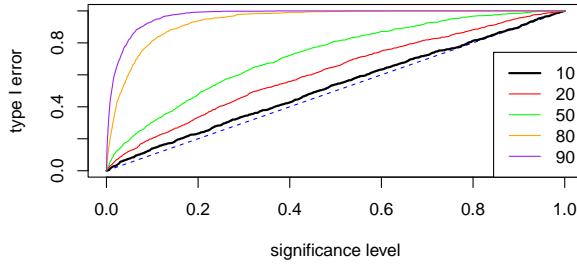
На рис. А.4, А.3 изображены ROC-кривые критериев, графики ошибок первого рода можно найти в разделе 2.3.2. Если сравнить получившиеся графики с графиками на рис 2.7, в, 2.7, в, 2.8, в и 2.9, в из раздела 2.3.2, видно, что с ростом частоты сигнала в альтернативе увеличивается различие в мощностях критериев в зависимости от параметра L .

А.2. Искажение критерия при использовании оценок

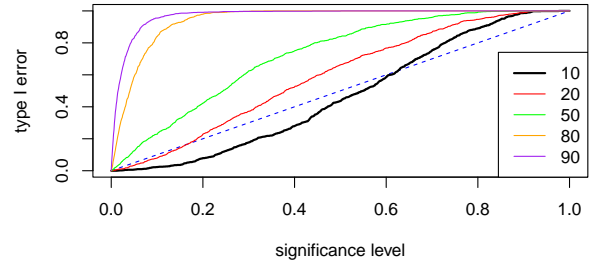
Здесь приведены графики к разделу 3.2.1 основного текста.

$$\varphi = 0.3, \omega = 0.075$$

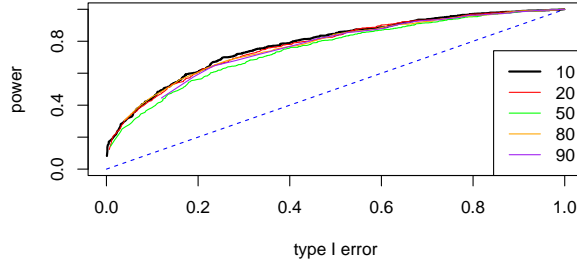
На рис. А.5 приведено сравнение критерия с известными и критерия с оцененными параметрами красного шума. По графикам ошибок первого рода на рис. А.5, а и А.5, б видно, что для длины окна $L = 10$ критерий при оценивании параметров шума становится консервативным при $\alpha < 0.6$, а для $L = 20$ критерий становится примерно точным при небольших α . При остальных L оценка параметров сильно на радикальность критерия не повлияла. Если взглянуть на ROC-кривые методов на рис. А.5, в и А.5, г, можно заметить значительное снижение мощности поправленного критерия



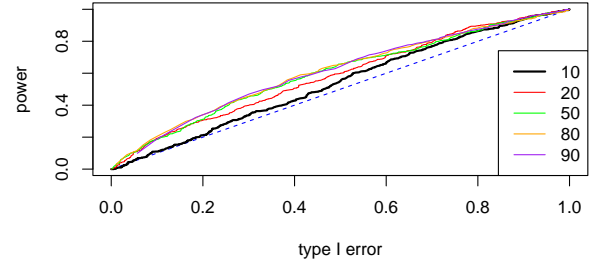
(a) Ошибка I рода (известные параметры)



(б) Ошибка I рода (оцененные параметры)



(в) ROC-кривая (известные параметры)



(г) ROC-кривая (оцененные параметры)

Рис. А.5. Влияние оценки параметров на критерий ($\varphi = 0.3$, $A = 0.7$, $\omega = 0.075$)

для всех длин окна.

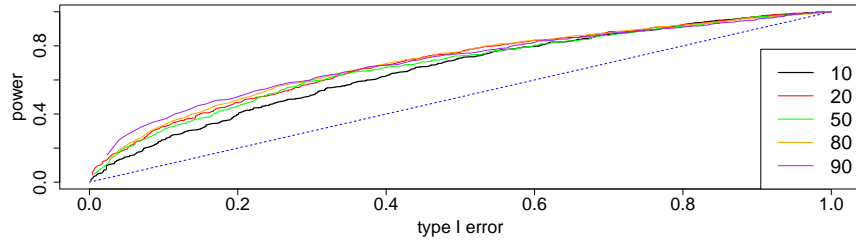
$$\varphi = 0.7, \omega = 0.225$$

На рис. А.6 изображены ROC-кривые метода при известных параметрах и при оцененных параметрах красного шума, графики ошибок первого рода можно найти в разделе 3.2.1. Также, как и в случае $\varphi = 0.7$ и $\omega = 0.075$ (см. раздел 3.2.1), оценка параметров не слишком сильно искажает критерий.

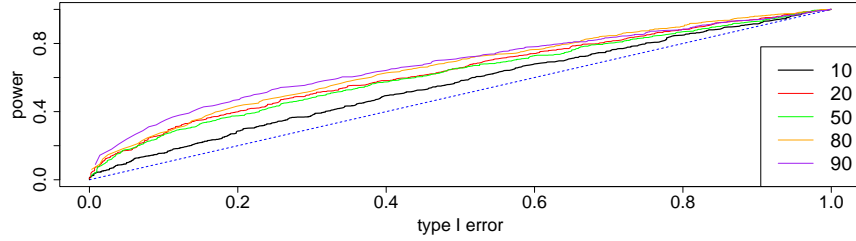
А.3. Наличие мешающего сигнала

А.3.1. Периодическая компонента

Здесь приведены графики к разделу 3.3.1 основного текста.



(a) ROC-кривая (известные параметры)



(б) ROC-кривая (оцененные параметры)

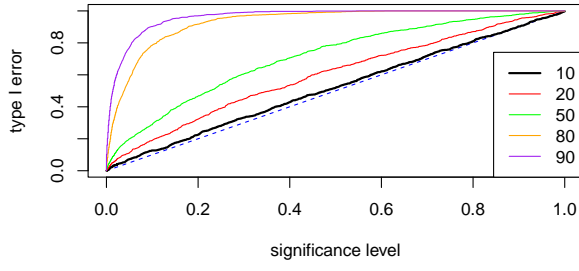
Рис. А.6. Влияние оценки параметров на критерий ($\varphi = 0.7$, $A = 0.4$, $\omega = 0.225$)

$$\varphi = 0.3, \omega = 0.075$$

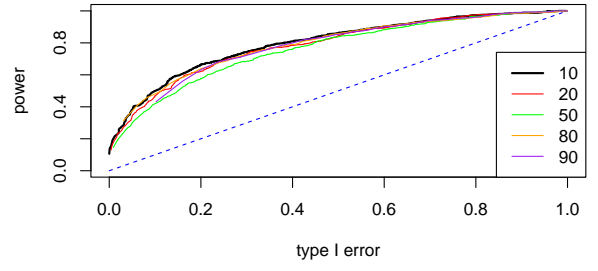
Уменьшим параметр φ до 0.3 и амплитуду сигнала в альтернативе до $A = 0.7$. На рис. А.7 приведено сравнение следующих критериев: когда мешающий сигнал и параметры красного шума известны точно, когда параметры шума оцениваются, и когда вместе с параметрами оценивается мешающий сигнал. Если сравнить графики ошибок первого рода на рис. А.7, а, А.7, в и А.7, д с рис. 3.6, а, 3.6, в и 3.6, д, то видно, что уменьшение параметра φ приводит к увеличению радикальности критерия. Если взглянуть на ROC-кривые критериев на рис. А.7, б, А.7, г и А.7, е, заметно сильное снижение мощности при оценивании параметров красного шума.

$$\varphi = 0.7, \omega = 0.225$$

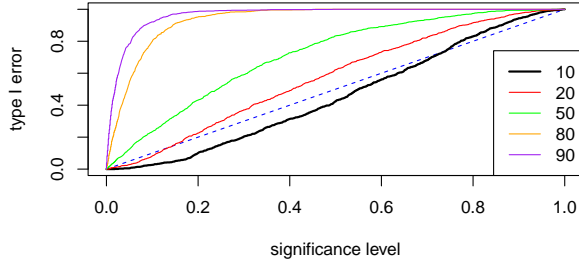
В условиях примера, рассмотренного в разделе 3.3.1, увеличим частоту в альтернативе до $\omega = 0.225$, уменьшив амплитуду до $A = 0.4$. На рис. А.8 приведено сравнение ROC-кривых. Графики ошибок первого рода для этого примера можно найти в разделе 3.3.1. Напомним, что частота мешающего сигнала $\omega = 0.25$, поэтому довольно низкую мощность метода можно объяснить близостью частоты сигнала в альтернативе с частотой мешающего сигнала.



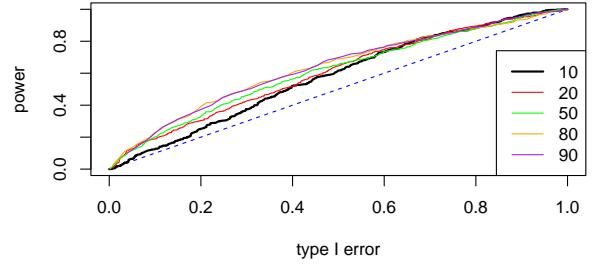
(a) Ошибка I рода



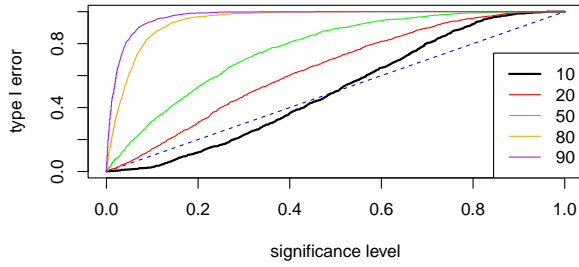
(б) ROC-кривая



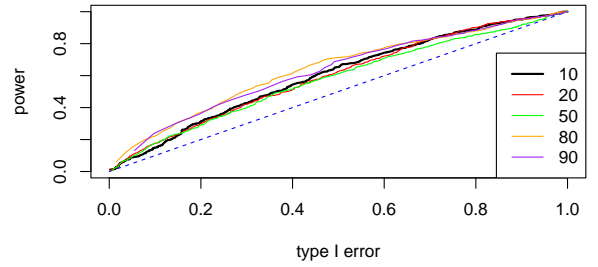
(в) Ошибка I рода (оцененные параметры шума)



(г) ROC-кривая (оцененные параметры шума)



(д) Ошибка I рода (оцененный мешающий сигнал и параметры шума)



(е) ROC-кривая (оцененный мешающий сигнал и параметры шума)

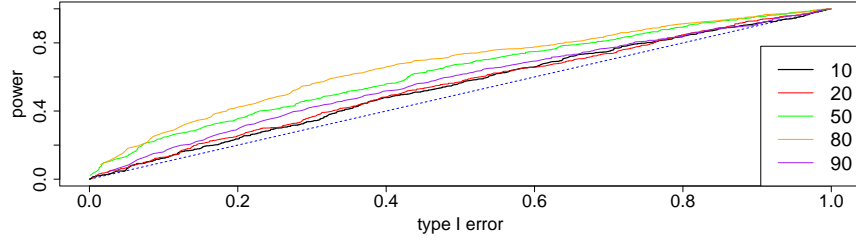
Рис. А.7. Анализ алгоритма 5 (мешающий сигнал — периодика) ($\varphi = 0.3$, $A = 0.7$, $\omega = 0.075$)

А.3.2. Тренд

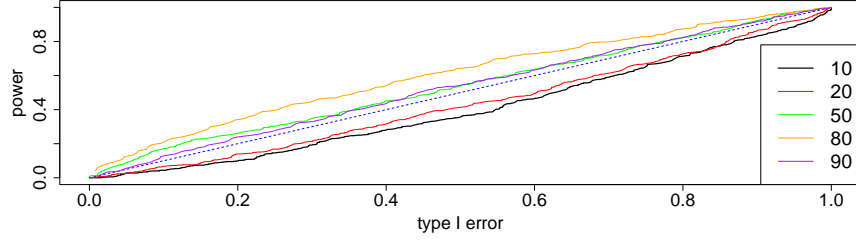
Здесь приведены графики к разделу 3.3.2.

$$\varphi = 0.3, \omega = 0.075$$

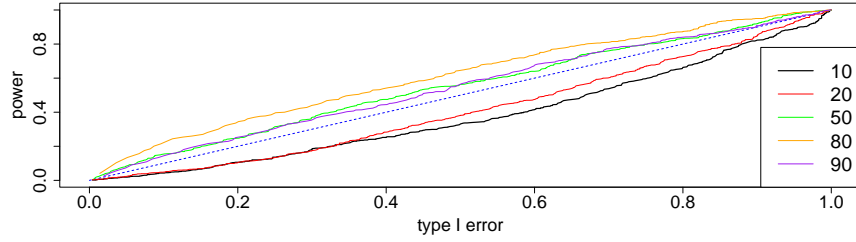
Уменьшим параметр φ до 0.3, амплитуду сигнала в альтернативе до $A = 0.7$ и рассмотрим 3 случая: когда мешающий сигнал и параметры красного шума известны



(a) ROC-кривая



(б) ROC-кривая (оцененные параметры шума)



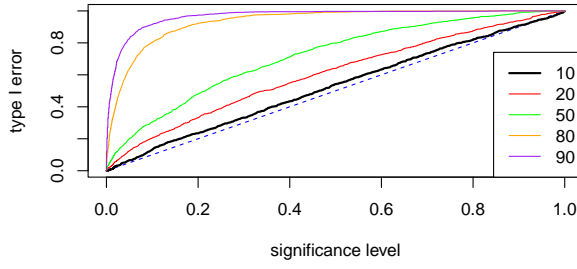
(в) ROC-кривая (оцененные параметры шума и мешающий сигнал)

Рис. А.8. Анализ алгоритма 5 (мешающий сигнал — периодика) ($\varphi = 0.7$, $A = 0.4$, $\omega = 0.225$)

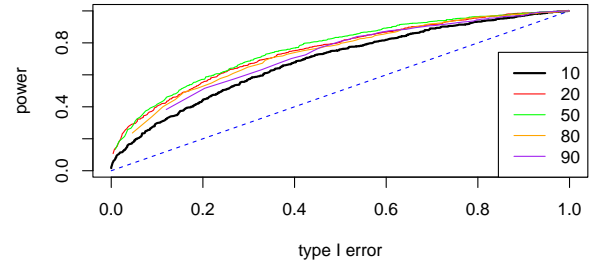
точно, когда параметры шума оцениваются, и когда вместе с параметрами оценивается мешающий сигнал. Заметим, что как и в разделе А.3.1, уменьшение параметра φ приводит к увеличению радикальности критерия, а также приводит к уменьшению разницы в мощностях поправленных критериев. Также, как в случае с периодикой в качестве мешающего сигнала, при небольших φ оценка неизвестных параметров сильно снижает мощность критериев. Это можно объяснить погрешностью при оценке параметров и погрешностью в выделении мешающего сигнала.

$$\varphi = 0.7, \omega = 0.225$$

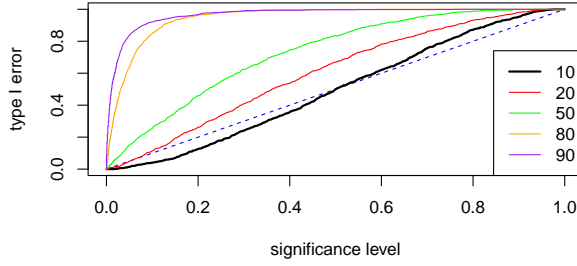
Теперь в условиях примера, рассмотренного в разделе 3.3.2, и увеличим частоту в альтернативе до $\omega = 0.225$, уменьшив амплитуду до $A = 0.4$. Также рассмотрим 3



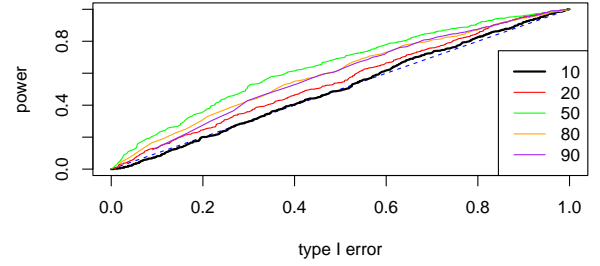
(a) Ошибка I рода



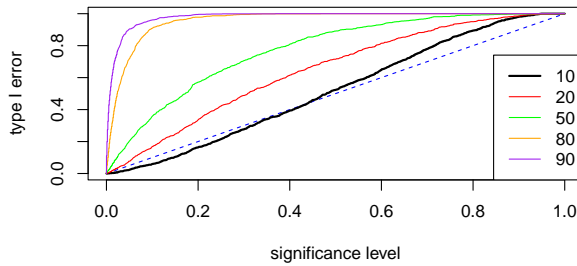
(б) ROC-кривая



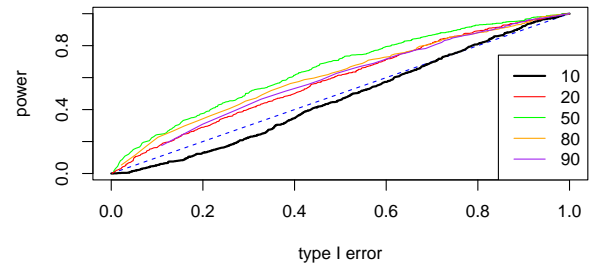
(в) Ошибка I рода (оцененные параметры шума)



(г) ROC-кривая (оцененные параметры шума)



(д) Ошибка I рода (оцененный мешающий сигнал и параметры шума)

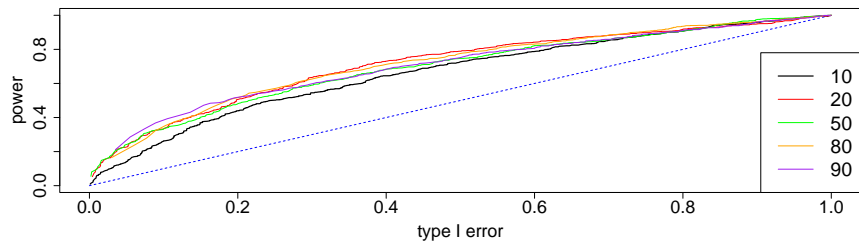


(е) ROC-кривая (оцененный мешающий сигнал и параметры шума)

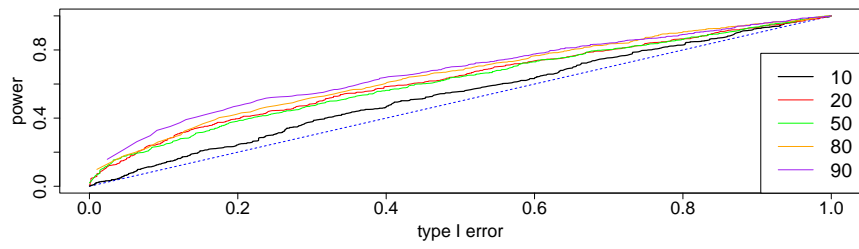
Рис. А.9. Анализ алгоритма 5 (мешающий сигнал — тренд) ($\varphi = 0.3$, $A = 0.7$, $\omega = 0.075$)

случая.

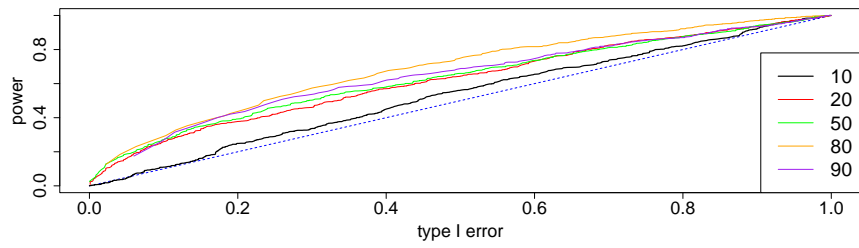
На рис. А.10 приведено сравнение ROC-кривых. Графики ошибок первого рода для этого примера можно найти в разделе 3.3.1. Графики ошибок первого рода для этого примера можно найти в разделе 3.3.2. При оценивании параметров наблюдается потеря в мощности, но она не такая существенная, как в случае $\varphi = 0.3$.



(a) ROC-кривая



(б) ROC-кривая (оцененные параметры шума)



(в) ROC-кривая (оцененные параметры шума и мешающий сигнал)

Рис. А.10. Анализ алгоритма 5 (мешающий сигнал — тренд) ($\varphi = 0.7$, $A = 0.4$, $\omega = 0.225$)

Приложение Б

Таблицы

Б.1. Численное сравнение модификаций MC-MSSA

Здесь приведены таблицы к разделу 3.2.1 основного текста.

$$\varphi = 0.7, \omega = 0.075$$

Таблица Б.1. Численное сравнение методов для оптимальных L ($\varphi = 0.7$, $A = 1$, $\omega = 0.075$)

Метод	левые/правые векторы	L	KD	длина векторов	кол-во векторов	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
SVD*	левые	50	102	50	50	0.316	0.754
SVD*	правые	80	42	42	42	0.394	0.754
Block*	левые	20, 50	162, 102	162	20	0.157	0.796
Block*	правые	80	42	80	42	0.367	0.717
Sum	левые	80, 90	42, 22	90	22	0.535	0.806
Sum*	правые	80	42	42	42	0.51	0.748

В таблице Б.1 представлены результаты для примера 1. Для каждого метода указана оптимальная по мощности длина окна (для которой удалось построить ROC-кривую, звездочкой помечены те методы, для которых эта L может не являться оптимальной), значение KD , длина проекционных векторов, их количество, а также ошибка первого рода и мощность поправленного критерия при уровне значимости $\alpha^* = 0.1$. Черным выделены максимальная мощность и примерно равные ей (отличающиеся не более, чем на 0.03). Также для каждого критерия дополнительно добавлена длина окна, дающая примерно ту же мощность, что и оптимальная, но метод с таким L более эффективный в нахождении левых или правых векторов автоковариационной матрицы (такая длина окна и соответствующее значение KD выделены синим). Видно, что MC-MSSA с проекцией на левые или правые векторы обеих модификаций мощнее, чем с проекцией на векторы базового MSSA. Самыми мощными являются методы Block и Sum с проекцией на левые векторы, причем метод Sum более эффективнее (матрица

размера 80 против матрицы размера 162).

$$\varphi = 0.3, \omega = 0.075$$

Таблица Б.2. Численное сравнение методов для оптимальных L ($\varphi = 0.3$, $A = 0.5$, $\omega = 0.075$)

Метод	левые/правые векторы	L	KD	длина векторов	кол-во векторов	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
SVD*	левые	10, 20	182, 162	20	20	0.199	0.399
SVD*	правые	80, 90	42, 22	22	22	0.449	0.382
Block*	левые	10, 20	182, 162	182	10	0.177	0.398
Block*	правые	80, 90	42, 22	90	22	0.414	0.389
Sum	левые	20, 90	162, 22	20	20	0.905	0.421
Sum*	правые	80, 90	42, 22	22	22	0.425	0.412

В таблице Б.2 представлены результаты для примера 2. По таблице видно, что различие в мощностях у методов совсем небольшое. Если сравнивать методы с наибольшей мощностью по трудоемкости, самым эффективным является метод Sum с проекцией на левые векторы с $L = 20$.

$$\varphi = 0.7, \omega = 0.225$$

Таблица Б.3. Численное сравнение методов для оптимальных L ($\varphi = 0.7$, $A = 0.4$, $\omega = 0.225$)

Метод	левые/правые векторы	L	KD	длина векторов	кол-во векторов	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
SVD*	левые	20	162	20	20	0.122	0.573
SVD*	правые	80	42	80	42	0.394	0.442
Block*	левые	20	162	162	20	0.157	0.597
Block*	правые	50, 80	102, 42	80	42	0.83	0.509
Sum	левые	80, 90	42, 22	10	90	0.535	0.625
Sum*	правые	80	42	42	42	0.51	0.613

В таблице Б.3 представлены результаты для примера 3. Метод Block с проекцией на левые и метод Sum с проекцией на левые и правые векторы в этом случае дают

наибольшую мощность, причем метод Sum с проекцией на левые или правые векторы при $L = 80$ наиболее эффективный среди наиболее мощных критериев.

Б.2. Сравнение двух алгоритмов MC-SSA с мешающим сигналом

Здесь приведены таблицы к разделу 3.3.1 основного текста.

$$\varphi = 0.3, \omega = 0.075$$

Уменьшим параметр φ до 0.3, амплитуду сигнала в альтернативе до $A = 0.7$ и сравним алгоритм 5 и 6 в 3 случаях: когда мешающий сигнал и параметры красного шума известны точно, когда параметры шума оцениваются, и когда вместе с параметрами оценивается мешающий сигнал.

Таблица Б.4. Сравнение алгоритма 5 и алгоритма 6 при $\alpha^* = 0.1$ ($\varphi = 0.3, A = 0.7, \omega = 0.075$)

Алгоритм 5	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	10	0.127	0.497
Оцененные параметры шума	90	0.921	0.261
Оцененные параметры шума и мешающий сигнал	90	0.94	0.239
Алгоритм 6	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	90	0.842	0.489
Оцененные параметры шума	90	0.867	0.292
Оцененные параметры шума и мешающий сигнал	90	0.887	0.27

В таблице Б.4 представлены результаты сравнения двух алгоритмов, а именно оптимальная длина окна, ошибка первого рода и мощность поправленного критерия при уровне значимости $\alpha^* = 0.1$. По таблице видно, что в этом примере алгоритм 6 менее радикальное алгоритма 5, что дает использовать поправку для больших L , тем самым давая бóльшую мощность при оценке неизвестных параметров.

$$\varphi = 0.7, \omega = 0.225$$

В условиях примера, рассмотренного в разделе 3.3.1, увеличим частоту в альтернативе до $\omega = 0.225$, уменьшив амплитуду до $A = 0.4$.

Таблица Б.5. Сравнение алгоритма 5 и алгоритма 6 при $\alpha^* = 0.1$ ($\varphi = 0.7, A = 0.4, \omega = 0.225$)

Алгоритм 5	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	80	0.525	0.273
Оцененные параметры шума	80	0.518	0.214
Оцененные параметры шума и мешающий сигнал	80	0.586	0.214
Алгоритм 6	L	$\alpha_I(\alpha^*)$	$\beta(\tilde{\alpha}^*)$
Точная модель	50	0.215	0.231
Оцененные параметры шума	50	0.149	0.174
Оцененные параметры шума и мешающий сигнал	50	0.195	0.166

По таблице Б.5 видно, что оба алгоритма дают довольно малую мощность. Это можно объяснить близостью частоты мешающего сигнала (0.25) и частоты в альтернативе. Поэтому применять эти алгоритмы для выявления сигнала с частотой, близкой к частоте мешающего сигнала, не рекомендуется.