

Московский Физико-Технический Институт
(Государственный Университет)

Факультет Управления и Прикладной Математики
Кафедра «Интеллектуальные Системы»

ДИПЛОМНАЯ РАБОТА СТУДЕНТА 774 ГРУППЫ

«Выбор иерархических моделей в авторегрессионном прогнозировании»

Выполнил:

студент 6 курса 774 группы

Фадеев Илья Владимирович

Научный руководитель:

к.ф-м.н., н.с. ВЦ РАН

Стрижов Вадим Викторович

Аннотация

В работе предложены методы анализа и прогнозирования периодических временных рядов с использованием полупараметрических моделей и инвариантных преобразований. Сегменты временного ряда кластеризуются на группы со схожей формой, для каждой группы оценивается форма и параметры полупараметрической модели. Для поиска моментов изменения формы сегментов с течением времени адаптируется алгоритм обнаружения разладок с помощью дискретной производной. Для прогнозирования используется двухуровневая иерархическая модель. В качестве примера использования предложенных методов рассматривается прогнозирование почасового потребления электроэнергии.

1 Введение

В работе рассматриваются методы построения прогностических моделей, описывающих периодические временные ряды и включающие инвариантные преобразования. Временной ряд разбивается на сегменты, длина которых равна периоду. Сегменты рассматриваются как выборка в полупараметрической регрессионной модели.

Полупараметрические модели были предложены в работе [1] для анализа выборок, состоящих из регрессионных кривых. Примером такой выборки является набор зависимостей целевой переменной от времени: каждому объекту соответствует своя зависимость, заданная в некоторые моменты времени. В полупараметрической модели в пространстве кривых определяется параметрическое семейство преобразований. Предполагается, что существует кривая, называемая формой и общая для всех элементов выборки, такая, что каждый элемент выборки является результатом преобразования этой кривой с некоторыми параметрами. Таким образом, построение модели включает задачу непараметрической регрессии (вычисление формы) и параметрической регрессии (вычисление параметров преобразования).

Частным случаем полупараметрических моделей являются модели, инвариантные относительно формы. В таких моделях параметрическое семейство преобразований содержит четыре параметра: растяжение и сдвиг кривой вдоль оси времени и оси целевой переменной. В работах [2, 3, 4] предложен итеративный алгоритм вычисления формы и параметров такой модели: на каждой итерации при фиксированной форме параметры вычисляются методом наименьших квадратов, далее оценивается форма с помощью усреднения по формам всех кривых выборки. Модели, инвариантные относительно формы, исследуются также в работах [5, 6, 7]. Подклассом таких моделей являются модели сдвига вдоль оси времени [8, 9, 10].

Для оценки формы и параметров полупараметрических моделей используется минимизация функционала ошибки [11, 3, 6] или максимизация правдоподобия для заданной статистической модели [12, 5]. Кривые, составляющие выборку, являются аппроксимациями зависимости целевой переменной от времени с помощью ядерного сглаживания [3] или сплайнов [13, 6, 7], или рассматриваются как марковский случайный процесс [12]. Форма в полупараметрической модели оценивается с помощью метода Надарая-Ватсона [9], усреднением по кривым выборки [4] или рассматри-

вается как матожидание марковского случайного процесса [12]. Для периодических кривых форма оценивается рядом Фурье [8, 5].

В данной работе предлагаются методы оценки формы и параметров полупараметрических моделей для широкого класса семейств преобразований. Предложенные оценки не состоятельны, однако доказывається, что для достаточно больших выборок ошибка оценок не превосходит некоторой величины. Для более узкого класса преобразований схожие оценки использовались в [4].

Предполагается, что сегменты временного ряда могут быть разбиты на группы так, что каждой группе соответствует своя форма в полупараметрической модели. Для нахождения разбиения предлагается использовать кластеризацию или априорные предположения о последовательности сегментов временного ряда. При этом для принятия решения о разбиении группы сегментов на две подгруппы предлагается статистический критерий различимости форм для двух подгрупп.

В работе также решается задача поиска моментов времени, в которых происходит изменение формы сегментов временного ряда. Для этого адаптируется алгоритм обнаружения разладок с помощью дискретной производной и вычислением достижимого уровня значимости, предложенный в [14]. Алгоритм заключается в вычислении дискретной производной формы сегментов по времени, отборе потенциальных точек разладки среди локальных максимумов модуля дискретной производной, исключении из потенциальных точек разладки "ложных тревог" с помощью статистического критерия.

Работа построена следующим образом. В главе 2 определяется полупараметрическая модель для сегментов временного ряда. Как и в работах [5, 4, 7], рассматривается задача однозначного определения формы и параметров модели, но для более широкого класса преобразований. В главе 3 предлагаются методы оценки формы полупараметрической модели, доказывається их устойчивость. Глава 4 посвящена задаче разбиения сегментов временного ряда на группы с общей формой. В главе 5 адаптируется алгоритм обнаружения разладок для поиска моментов изменения формы сегментов временного ряда. В главе 6 приведён критерий качества параметрического семейства преобразований, который используется при наличии разбиения сегментов на группы с общей формой, заданной экспертом. Этот критерий позволяет выбрать

семейство преобразований, наилучшим образом отражающее экспертное представление о близости форм временных рядов. В главе 7 предлагается иерархическая двухуровневая модель для прогнозирования значений временного ряда. В главе 8 приводится пример использования предложенных методов для анализа и прогнозирования почасового потребления электроэнергии.

2 Модель порождения данных

Рассматривается квазипериодический временной ряд с периодом n

$$y_t, \quad t = 1, \dots, T,$$

где $T = Nn$ — длина временного ряда, кратная периоду n , $N = \frac{T}{n}$ — число периодов.

Временной ряд y_t разбивается на сегменты, длина которых равна периоду n :

$$\mathbf{x}_i = (y_{(i-1)n+1}, y_{(i-1)n+2}, \dots, y_{in}), \quad i = 1, \dots, N.$$

Пусть $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ — параметрическое семейство преобразований. Рассмотрим следующую модель порождения данных:

$$\mathbf{x}_i = \mathbf{f}(\mathbf{z}_0, \boldsymbol{\alpha}_i) + \boldsymbol{\varepsilon}_i, \quad \boldsymbol{\alpha}_i \in \mathbb{R}^m, \quad \boldsymbol{\varepsilon}_i \sim \mathcal{N}(0, E_n \sigma^2), \quad (1)$$

где \mathbf{z}_0 — некоторый временной ряд, $\boldsymbol{\alpha}_i$ — параметры преобразования, соответствующие каждому временному ряду, E_n — единичная матрица, σ^2 — дисперсия ошибки. Таким образом, временные ряды \mathbf{x}_i являются результатом преобразования временного ряда \mathbf{z}_0 с аддитивным нормальным шумом.

Предположим также, что преобразования \mathbf{f} разбивают пространство \mathbb{R}^n на классы эквивалентности следующим образом: $\mathbf{x}_i \sim \mathbf{x}_j$ если и только если существует вектор $\boldsymbol{\alpha}$ такой, что $\mathbf{x}_j = \mathbf{f}(\mathbf{x}_i, \boldsymbol{\alpha})$.

Форма \mathbf{z}_0 и параметры $\boldsymbol{\alpha}_i$ модели (1) не заданы однозначно: любая форма из класса эквивалентности исходной $\tilde{\mathbf{z}}_0 \sim \mathbf{z}_0$ с соответствующим набором параметров $\tilde{\boldsymbol{\alpha}}_i$ определяет ту же модель. Действительно, из $\tilde{\mathbf{z}}_0 \sim \mathbf{z}_0 \sim \mathbf{f}(\mathbf{z}_0, \boldsymbol{\alpha}_i)$ следует, что существует вектор параметров $\tilde{\boldsymbol{\alpha}}_i$ такой, что $\mathbf{f}(\tilde{\mathbf{z}}_0, \tilde{\boldsymbol{\alpha}}_i) = \mathbf{f}(\mathbf{z}_0, \boldsymbol{\alpha}_i)$, и модель

$$\mathbf{x}_i = \mathbf{f}(\tilde{\mathbf{z}}_0, \tilde{\boldsymbol{\alpha}}_i) + \boldsymbol{\varepsilon}_i, \quad \boldsymbol{\varepsilon}_i \sim \mathcal{N}(0, E_n \sigma^2)$$

эквивалентна исходной модели (1) порождения данных.

Пусть множество $\mathbb{Z} \subset \mathbb{R}^n$ содержит ровно по одному представителю от каждого класса эквивалентности преобразования \mathbf{f} . Тогда любой вектор $\mathbf{x}_i \in \mathbb{R}^n$ однозначно представим в виде

$$\mathbf{x}_i = \mathbf{f}(\mathbf{z}_i, \boldsymbol{\alpha}_i), \quad \mathbf{z}_i \in \mathbb{Z},$$

что позволяет ввести преобразование $\mathbf{u} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ и отображение $\mathbf{v} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, такие, что

$$\mathbf{x}_i = \mathbf{f}(\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x})), \quad \mathbf{u}(\mathbf{x}) \in \mathbb{Z},$$

Введение ограничения $\mathbf{z}_0 \in \mathbb{Z}$ однозначно определяет \mathbf{z}_0 и $\boldsymbol{\alpha}_i$ модели (1):

$$\mathbf{x}_i = \mathbf{f}(\mathbf{z}_0, \boldsymbol{\alpha}_i) + \boldsymbol{\varepsilon}_i, \quad \boldsymbol{\varepsilon}_i \sim \mathcal{N}(0, E_n \sigma^2), \quad \mathbf{z}_0 \in \mathbb{Z}. \quad (2)$$

Вектор \mathbf{z}_0 будем называть формой, соответствующей выборке $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, а вектор $\mathbf{z}_i = \mathbf{u}(\mathbf{x}_i)$ формой вектора \mathbf{x}_i . Заметим, что преобразование \mathbf{u} инвариантно относительно \mathbf{f} , а форма $\mathbf{u}(\mathbf{x}_i)$ является инвариантом вектора \mathbf{x}_i .

2.1 Примеры преобразований f

Наиболее часто встречающиеся в литературе и хорошо интерпретируемые семейства инвариантных преобразований — растяжение-сдвиг вдоль оси значений временного ряда $\mathbf{f}(\mathbf{x}, \boldsymbol{\alpha}) = \alpha_1 \mathbf{x} + \alpha_2$, растяжение-сдвиг вдоль оси времени, а также их суперпозиции.

Обобщением сдвига вдоль оси значений является прибавление полинома k -го порядка:

$$f_j(\mathbf{x}, \boldsymbol{\alpha}) = x_j + \sum_{m=0}^k \alpha_m j^m, \quad j = 1, \dots, n,$$

Множество \mathbb{Z} можно задать преобразованиями

$$\mathbf{v}(\mathbf{x}) = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \sum_{j=1}^n \left(\sum_{m=0}^k \alpha_m j^m - x_j \right)^2,$$

$$u_j(\mathbf{x}) = x_j - \sum_{m=0}^k v_m(\mathbf{x}) j^m,$$

при этом $\mathbf{v}(\mathbf{x})$, $\mathbf{u}(\mathbf{x})$ вычисляются методом наименьших квадратов.

Обобщением растяжения являются преобразования вида

$$f_j(\mathbf{x}, \boldsymbol{\alpha}) = g(x_j, \boldsymbol{\alpha}), \quad j = 1, \dots, n,$$

где g — семейство монотонных функций; например,

$$f_j(\mathbf{x}, \boldsymbol{\alpha}) = \alpha_0 x_j^{\alpha_1}, \quad j = 1, \dots, n.$$

3 Оценка формы \mathbf{z}_0

Предположим, что для любого набора векторов $\mathbf{z}_l \in \mathbb{Z}$ выполнено условие

$$\frac{1}{L} \sum_{l=1}^L \mathbf{z}_l \in \mathbb{Z}, \quad \mathbf{z}_l \in \mathbb{Z}, \quad l = 1, \dots, L. \quad (3)$$

Тогда в качестве оценки формы \mathbf{z}_0 из модели (2) предлагается использовать среднее значение форм векторов \mathbf{x}_i

$$\mathbf{z}^* = \frac{1}{N} \sum_{i=1}^N \mathbf{u}(\mathbf{x}_i). \quad (4)$$

Теорема 1. Пусть преобразование \mathbf{u} удовлетворяет условию Липшица с константой L , т. е. для любых $\mathbf{x}_i, \mathbf{x}_j$

$$\|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\| \leq L \|\mathbf{x}_i - \mathbf{x}_j\|.$$

Тогда почти наверно существует минимальный размер выборки N_0 такой, что

$$\|\mathbf{z}^* - \mathbf{z}_0\| < L\sigma n + \varepsilon_0, \quad \forall N : N > N_0,$$

где \mathbf{z}^* — оценка (4), ε_0 — любое положительное число.

Доказательство. Заметим, что

$$\begin{aligned} |u_j(\mathbf{x}_i) - z_{0j}| &\leq \|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\| \leq \\ &\leq L \|\mathbf{x}_i - f(\mathbf{z}_0, \boldsymbol{\alpha}_i)\| = L \|\boldsymbol{\varepsilon}_i\|. \end{aligned}$$

$$|E[u_j(\mathbf{x}_i) - z_{0j}]| \leq E|u_j(\mathbf{x}_i) - z_{0j}| \leq E[L \|\boldsymbol{\varepsilon}_i\|] = LE \|\boldsymbol{\varepsilon}_i\|. \quad (5)$$

$$V[u_j(\mathbf{x}_i) - z_{0j}] = E(u_j(\mathbf{x}_i) - z_{0j})^2 - (E[u_j(\mathbf{x}_i) - z_{0j}])^2 \leq \quad (6)$$

$$\leq \mathbf{E}(u_j(\mathbf{x}_i) - z_{0j})^2 \leq \mathbf{E}[L^2 \|\boldsymbol{\varepsilon}_i\|^2] = L^2 \mathbf{E}\|\boldsymbol{\varepsilon}_i\|^2 = L^2 \sigma^2 n.$$

Из модели порождения данных (2) следует, что

$$\frac{\boldsymbol{\varepsilon}_i}{\sigma} \sim \mathcal{N}(0, E_n \sigma^2),$$

— вектор с независимыми компонентами, имеющими стандартное нормальное распределение. Следовательно,

$$\left\| \frac{\boldsymbol{\varepsilon}_i}{\sigma} \right\|^2 \sim \chi^2(n).$$

Как известно, $E\chi^2(n) = n$, следовательно,

$$\mathbf{E}\|\boldsymbol{\varepsilon}_i\|^2 = \sigma^2 n,$$

$$\mathbf{E}\|\boldsymbol{\varepsilon}_i\| = \sqrt{\mathbf{E}\|\boldsymbol{\varepsilon}_i\|^2 - \mathbf{V}\|\boldsymbol{\varepsilon}_i\|} \leq \sqrt{\mathbf{E}\|\boldsymbol{\varepsilon}_i\|^2} = \sigma\sqrt{n}.$$

Из (5) и последнего неравенства следует, что

$$|\mathbf{E}[u_j(\mathbf{x}_i) - z_{0j}]| \leq L\sigma\sqrt{n}.$$

Поскольку матожидание среднего значения нескольких случайных величин не превосходит максимального значения матожидания, то

$$\left| \mathbf{E} \left[\frac{1}{N} \sum_{i=1}^N u_j(\mathbf{x}_i) - z_{0j} \right] \right| \leq L\sigma\sqrt{n}.$$

Воспользуемся усиленным законом больших чисел в формулировке Колмогорова, согласно которому

$$\frac{1}{N} \sum_{i=1}^N u_j(\mathbf{x}_i) - z_{0j} \xrightarrow{\text{almost sure}} \mathbf{E} \left[\frac{1}{N} \sum_{i=1}^N u_j(\mathbf{x}_i) - z_{0j} \right]$$

при выполнении условия

$$\sum_{i=1}^{\infty} \frac{1}{i^2} \mathbf{V}[u_j(\mathbf{x}_i) - z_{0j}] < \infty.$$

Последнее неравенство выполняется в силу (6) и того, что

$$\begin{aligned} & \sum_{i=1}^{\infty} \frac{1}{i^2} \mathbf{V}[u_j(\mathbf{x}_i) - z_{0j}] \leq \\ & \leq L^2 \sigma^2 n \sum_{i=1}^{\infty} \frac{1}{i^2} = \frac{L^2 \sigma^2 n \pi^2}{6} < \infty. \end{aligned}$$

Из сходимости почти на вероятности единица существует число N_0 такое, что для любого $N > N_0$

$$\frac{1}{N} \sum_{i=1}^N u_j(\mathbf{x}_i) - z_{0j} < L\sigma\sqrt{n} + \varepsilon_1,$$

где ε_1 — любое положительное число. В итоге получаем

$$\begin{aligned} \|\mathbf{z}^* - \mathbf{z}_0\| &= \left\| \frac{1}{N} \sum_{i=1}^N u(\mathbf{x}_i) - \mathbf{z}_0 \right\| = \sqrt{\sum_{j=1}^n \left[\frac{1}{N} \sum_{i=1}^N u_j(\mathbf{x}_i) - z_{0j} \right]^2} \leq \\ &\leq \sqrt{n}(L\sigma\sqrt{n} + \varepsilon_1) = L\sigma n + \varepsilon_0 \quad \bullet \end{aligned}$$

Не для любого множества \mathbb{Z} выполнено условие (3). В общем случае в качестве оценки формы \mathbf{z}_0 предлагается использовать форму одного из векторов выборки:

$$\mathbf{z}^* = \operatorname{argmin}_{\mathbf{u}(\mathbf{x}_j)} \sum_{i=1}^N \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|^2. \quad (7)$$

Теорема 2. Пусть преобразование \mathbf{u} удовлетворяет условию Липшица с константой L , т. е. для любых $\mathbf{x}_i, \mathbf{x}_j$

$$\|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\| \leq L\|\mathbf{x}_i - \mathbf{x}_j\|.$$

Тогда почти на вероятности единица существует минимальный размер выборки N_0 такой, что

$$\|\mathbf{z}^* - \mathbf{z}_0\| < \frac{31}{2}L\sigma\sqrt{n} + \varepsilon_0, \quad \forall N : N > N_0,$$

где \mathbf{z}^* — оценка (7), ε_0 — любое положительное число.

Доказательство. Из модели порождения данных (2) следует, что

$$\frac{\boldsymbol{\varepsilon}_i}{\sigma} \sim \mathcal{N}(0, E_n\sigma^2),$$

— вектор с независимыми компонентами, имеющими стандартное нормальное распределение. Следовательно,

$$\frac{\|\mathbf{x}_i - f(\mathbf{z}_0, \boldsymbol{\alpha}_i)\|^2}{\sigma^2} = \left\| \frac{\boldsymbol{\varepsilon}_i}{\sigma} \right\|^2 \sim \chi^2(n).$$

По условию теоремы,

$$\|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\|^2 \leq L^2\|\mathbf{x}_i - f(\mathbf{z}_0, \boldsymbol{\alpha}_i)\|^2,$$

следовательно,

$$E\|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\|^2 \leq L^2\sigma^2 E\chi^2(n) = L^2\sigma^2 n.$$

Согласно неравенству Маркова, для любого a

$$P(\|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\|^2 \geq a^2) \leq \frac{E\|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\|^2}{a^2} \leq \frac{L^2\sigma^2 n}{a^2}.$$

Выберем $a > 0$ так, чтобы

$$P(\|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\|^2 \geq a^2) \leq \frac{L^2\sigma^2 n}{a^2} = \frac{1}{4}, \quad (8)$$

$$a^2 = 4L^2\sigma^2 n.$$

Определим разбиение индексов объектов выборки $I = \{1, \dots, N\} = I_1 \cup I_2$ следующим образом:

$$I_1 = \{i : \|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\|^2 \leq a^2\},$$

$$I_2 = \{i : \|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\|^2 > a^2\}.$$

Из условия (8) следует, что существует число N_1 такое, что для любого размера выборки $N > N_1$

$$|I_2| < \frac{N}{3}. \quad (9)$$

В последующих рассуждениях предполагаем, что $N > N_1$ и, соответственно, выполнено условие (9).

Пусть \mathbf{x}_β — некоторый объект выборки, $b = \|\mathbf{u}(\mathbf{x}_\beta) - \mathbf{z}_0\|$. Покажем, что при выполнении некоторого условия для b вектор $\mathbf{u}(\mathbf{x}_\beta)$ не может быть равен оценке \mathbf{z}^* из (7). Для этого необходимо найти такой элемент выборки \mathbf{x}_α , для которого

$$\sum_{i=1}^N \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\alpha)\|^2 < \sum_{i=1}^N \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\beta)\|^2. \quad (10)$$

Пусть α — произвольный индекс из I_1 . Тогда для любого $i \in I_1$ из неравенства треугольника

$$\|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\beta)\| \geq b - a,$$

$$\|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\alpha)\| \leq 2a.$$

Учитывая данные неравенства и условие (9) получаем

$$\sum_{i \in I_1} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\beta)\|^2 - \sum_{i \in I_1} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\alpha)\|^2 \geq$$

$$\geq \sum_{i \in I_1} [(b-a)^2 - (2a)^2] > \frac{2N}{3} [(b-a)^2 - (2a)^2]. \quad (11)$$

Для сокращения записи введём следующие обозначения:

$$\rho_{ij} = \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|,$$

$$\rho_{0i} = \|\mathbf{u}(\mathbf{x}_i) - \mathbf{z}_0\|.$$

В новых обозначениях

$$\begin{aligned} & \sum_{i \in I_2} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\beta)\|^2 - \sum_{i \in I_2} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\alpha)\|^2 = \\ & = \sum_{i \in I_2} (\rho_{\beta i}^2 - \rho_{\alpha i}^2) = \sum_{i \in I_2} (\rho_{\beta i} - \rho_{\alpha i})(\rho_{\beta i} + \rho_{\alpha i}). \end{aligned} \quad (12)$$

Из неравенства треугольника

$$|\rho_{\beta i} - \rho_{\alpha i}| \leq \rho_{\alpha\beta},$$

$$\rho_{\beta i} - \rho_{\alpha i} \geq -\rho_{\alpha\beta} \geq -(a+b). \quad (13)$$

С учётом условия (9)

$$\begin{aligned} & \sum_{i \in I_2} (\rho_{\beta i} + \rho_{\alpha i}) \leq \sum_{i \in I_2} (\rho_{0i} + b + \rho_{0i} + a) \leq \\ & \leq \frac{N}{3}(a+b) + 2 \sum_{i \in I_2} \rho_{0i}. \end{aligned} \quad (14)$$

Т. к. $\rho_{0i} > a$ для $i \in I_2$, то

$$\sum_{i \in I_2} \rho_{0i} \leq \sum_{i \in I_2} \frac{\rho_{0i}^2}{a} = \frac{1}{a} \sum_{i \in I_2} \rho_{0i}^2 \leq \frac{1}{a} \sum_{i=1}^N \rho_{0i}^2.$$

Далее из условия Липшица получаем

$$\begin{aligned} \sum_{i \in I_2} \rho_{0i} & \leq \frac{1}{a} \sum_{i=1}^N \rho_{0i}^2 \leq \frac{1}{a} \sum_{i=1}^N (L^2 \|\mathbf{x}_i - f(\mathbf{z}_0, \boldsymbol{\alpha}_i)\|^2) = \\ & = \frac{L^2 \sigma^2}{a} \sum_{i=1}^N \left\| \frac{\boldsymbol{\varepsilon}_i}{\sigma} \right\|^2. \end{aligned}$$

Согласно усиленному закону больших чисел,

$$\frac{1}{N} \sum_{i=1}^N \left\| \frac{\boldsymbol{\varepsilon}_i}{\sigma} \right\|^2 \xrightarrow{\text{almost sure}} \mathbb{E} \left\| \frac{\boldsymbol{\varepsilon}_i}{\sigma} \right\|^2 = \mathbb{E} \chi^2(n) = n,$$

следовательно, для любого $\varepsilon_1 > 0$ существует такое N_2 , что

$$\frac{1}{N} \sum_{i=1}^N \left\| \frac{\varepsilon_i}{\sigma} \right\|^2 < n + \varepsilon_1, \quad N > N_2.$$

В итоге получаем

$$\sum_{i \in I_2} \rho_{0i} < \frac{L^2 \sigma^2 (n + \varepsilon_0) N}{a}. \quad (15)$$

Подставим в (12) неравенства (13), (14) и (15):

$$\begin{aligned} & \sum_{i \in I_2} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\beta)\|^2 - \sum_{i \in I_2} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\alpha)\|^2 \geq \\ & \geq -(a+b) \left[\frac{N}{3}(a+b) + \kappa \right], \quad (16) \\ & \kappa = \frac{L^2 \sigma^2 (n + \varepsilon_0) N}{a}. \end{aligned}$$

Объединяя (11) и (16), получаем

$$\begin{aligned} & \sum_{i=1}^N \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\beta)\|^2 - \sum_{i=1}^N \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\alpha)\|^2 = \sum_{i \in I_1} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\beta)\|^2 - \sum_{i \in I_1} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\alpha)\|^2 + \\ & + \sum_{i \in I_2} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\beta)\|^2 - \sum_{i \in I_2} \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_\alpha)\|^2 \geq \frac{2N}{3} [(b-a)^2 - (2a)^2] - (a+b) \left[\frac{N}{3}(a+b) + \kappa \right]. \end{aligned}$$

Необходимо определить, при каких ограничениях на b выполнено

$$\frac{2N}{3} [(b-a)^2 - (2a)^2] - (a+b) \left[\frac{N}{3}(a+b) + \kappa \right] > 0.$$

Обозначив

$$\tilde{\kappa} = \frac{\kappa}{aN/3}, \quad v = \frac{b}{a},$$

перепишем последнее неравенство в виде

$$2[(v-1)^2 - 4] - (1+s)[(1+s) + \tilde{\kappa}] > 0,$$

$$v^2 - (6 - \tilde{\kappa})v - 7 - \tilde{\kappa} > 0,$$

$$v > 7 + \tilde{\kappa},$$

$$b > a(7 + \tilde{\kappa}) = a \left(7 + \frac{3(n + \varepsilon_1)}{4n} \right) = \frac{31a}{4} + \varepsilon_0 = \frac{31}{2} L \sigma \sqrt{n} + \varepsilon_0. \quad (17)$$

Таким образом, при выполнении условия (17), $N > \max(N_1, N_2)$, выполнено неравенство (10) и, следовательно,

$$\mathbf{u}(\mathbf{x}_\beta) \neq \operatorname{argmin}_{\mathbf{u}(\mathbf{x}_j)} \sum_{i=1}^N \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|^2. \quad \bullet$$

Оценки (4) и (7) не являются состоятельными оценками формы \mathbf{z}_0 , т. к. при увеличении размеров выборки ошибка оценок не стремится к нулю. Однако теоремы 1 и 2 утверждают, что для достаточной большой выборки ошибка $\|\mathbf{z}^* - \mathbf{z}_0\|$ не превосходит некоторого порога, размер которого пропорционален константе Липшица инвариантного преобразования \mathbf{u} .

4 Разбиение сегментов на группы со схожей формой

В модели (2) рассматривалась единая для всей выборки \mathbf{x}_i форма \mathbf{z}_0 . Теперь будем предполагать, что существует разбиение индексов $\mathcal{I} = \{1, \dots, N\} = \mathcal{I}_1 \cup \dots \cup \mathcal{I}_s$ такое, что

$$\mathbf{x}_i = \mathbf{f}(\mathbf{z}_{0k}, \boldsymbol{\alpha}_i) + \boldsymbol{\varepsilon}_i, \quad i \in I_k, \quad \mathbf{z}_{0k} \in Z, \quad k = 1, \dots, s,$$

т. е. каждому набору индексов \mathcal{I}_k соответствует своя форма \mathbf{z}_{0k} .

Для нахождения разбиения $\mathcal{I}_1 \cup \dots \cup \mathcal{I}_s$ предлагается использовать один из алгоритмов кластеризации, определив функцию расстояния между индексами $\rho_{ij} = \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|$, $i, j \in I$ как расстояние между формами i -го и j -го временного ряда. Далее по каждой подвыборке $\mathbf{x}_i, i \in I_k$ оценивается форма \mathbf{z}_{0k} по формуле (4) или (7). С увеличением количества кластеров s уменьшается количество временных рядов в каждом кластере и, следовательно, увеличивается ошибка при оценке формы \mathbf{z}_{0k} .

Для целей прогнозирования важна интерпретируемость разбиения множества индексов \mathcal{I} , т. к. она позволяет предсказать форму, соответствующую некоторому периоду в будущем. Например, разбиение $\mathcal{I} = \mathcal{I}_1 \cup \mathcal{I}_2$ временного ряда с суточной периодикой на будни \mathcal{I}_1 и выходные \mathcal{I}_2 позволяет при прогнозировании предсказать форму \mathbf{z}_{01} для буднего дня и \mathbf{z}_{02} для выходного дня. В качестве примеров таких разбиений можно рассмотреть следующие:

- фазы календарных и суточных периодик: время суток, будни/выходные, время года, праздники/рабочие дни;

- разбиение на последовательности подряд идущих индексов: $\mathcal{I} = \{1, \dots, N_1\} \cup \{N_1 + 1, \dots, N\}$.

Возникает вопрос о том, существует ли значимое различие между формами векторов с индексами из \mathcal{I}_1 и \mathcal{I}_2 для некоторого априорного разбиения $\mathcal{I}_1 \cup \mathcal{I}_2$. Если различие выявлено, то предлагается оценивать формы \mathbf{z}_{01} и \mathbf{z}_{02} для временных рядов с индексами из \mathcal{I}_1 и \mathcal{I}_2 соответственно вместо общей оценки формы для всех временных рядов.

Чтобы формализовать задачу, рассмотрим формы векторов $\mathbf{u}(\mathbf{x}_i), i \in \mathcal{I}_1$ и $\mathbf{u}(\mathbf{x}_i), i \in \mathcal{I}_2$ как выборки некоторых случайных векторов J_1 и J_2 . Необходимо определить, существует ли статистически значимое различие между распределением J_1 и J_2 . В качестве нулевой гипотезы рассмотрим предположение о равенстве матожидания межклассового расстояния $E_{12} = E[\rho_{ij}|i \in \tilde{\mathcal{I}}_1, j \in \tilde{\mathcal{I}}_2]$ среднему значению матожиданий внутриклассовых расстояний $\frac{1}{2}(E[\rho_{ij}|i, j \in \mathcal{I}_1] + E[\rho_{ij}|i, j \in \mathcal{I}_2])$:

$$M = E_{12} - \frac{E_{11} + E_{22}}{2} = 0.$$

Оценки максимального правдоподобия матожиданий

$$\hat{E}_{12} = \frac{\sum_{i \in \mathcal{I}_1, j \in \mathcal{I}_2} \rho_{ij}}{|\mathcal{I}_1||\mathcal{I}_2|},$$

$$\hat{E}_{11} = \frac{\sum_{i, j \in \mathcal{I}_1, i < j} \rho_{ij}}{|\mathcal{I}_1|(|\mathcal{I}_1| - 1)/2}, \quad \hat{E}_{22} = \frac{\sum_{i, j \in \mathcal{I}_2, i < j} \rho_{ij}}{|\mathcal{I}_2|(|\mathcal{I}_2| - 1)/2}.$$

Используя оценки максимального правдоподобия для дисперсий $V\rho_{ij}$, найдём дисперсию оценок матожиданий

$$V\hat{E}_{12} = \frac{V[\rho_{ij}|\mathcal{I}_1, j \in \mathcal{I}_2]}{|\mathcal{I}_1||\mathcal{I}_2|} \approx \frac{\sum_{i \in \mathcal{I}_1, j \in \mathcal{I}_2} (\rho_{ij} - \hat{E}_{12})^2}{|\mathcal{I}_1|^2|\mathcal{I}_2|^2},$$

$$V\hat{E}_{11} \approx \frac{\sum_{i, j \in \mathcal{I}_1, i < j} (\rho_{ij} - \hat{E}_{11})^2}{|\mathcal{I}_1|^2(|\mathcal{I}_1| - 1)^2/4},$$

$$V\hat{E}_{22} \approx \frac{\sum_{i, j \in \mathcal{I}_2, i < j} (\rho_{ij} - \hat{E}_{22})^2}{|\mathcal{I}_2|^2(|\mathcal{I}_2| - 1)^2/4}.$$

Среднеквадратичное отклонение \hat{M}

$$\hat{s}e_M = \sqrt{V\hat{E}_{12} + \frac{1}{4}(V\hat{E}_{11} + V\hat{E}_{22})}.$$

Считая, что $\hat{M} = \hat{E}_{12} - (\hat{E}_{11} + \hat{E}_{22})/2 \sim \mathcal{N}(0, \widehat{se}_M^2)$, находим

$$\text{p-value} = 1 - \Phi\left(\frac{\hat{M}}{\widehat{se}_M}\right) \quad (18)$$

для нулевой гипотезы $M = 0$ против альтернативы $M > 0$, где Φ — функция стандартного нормального распределения. При уровне значимости 0,05 значение $\text{p-value} < 0,05$ свидетельствует о значимом различии в распределении форм векторов с индексами из \mathcal{I}_1 и \mathcal{I}_2 . В этом случае предлагается оценивать формы \mathbf{z}_{01} для $\mathbf{u}(\mathbf{x}_i), i \in \mathcal{I}_1$ и \mathbf{z}_{02} для $\mathbf{u}(\mathbf{x}_i), i \in \mathcal{I}_2$ вместо общей оценки формы для всей выборки.

5 Поиск моментов изменения формы

Предположим, что существует набор целых чисел $0 = \tau_0 < \tau_1 < \dots < \tau_K < \tau_{K+1} = N$ такой, что

$$\mathbf{x}_i = \mathbf{f}(\mathbf{z}_{0k}, \boldsymbol{\alpha}_i) + \boldsymbol{\varepsilon}_i, \quad \tau_k < i \leq \tau_{k+1}, \quad \mathbf{z}_{0k} \in \mathbb{Z}, \quad k = 0, \dots, K,$$

Числа τ_1, \dots, τ_K будем называть точками разладки; они разбивают последовательность временных рядов $\mathbf{x}_i, i \in \mathcal{I} = \{1, \dots, N\}$ на подпоследовательности с индексами $\mathcal{I}_0 = \{1, \dots, \tau_1\}, \dots, \mathcal{I}_{K+1} = \{\tau_K + 1, \dots, N\}$, и каждой подпоследовательности соответствует некоторая форма \mathbf{z}_{0k} .

Необходимо найти точки разладки τ_1, \dots, τ_K при условии, что их количество K неизвестно.

Для решения задачи предлагается адаптировать алгоритм поиска разладки с помощью дискретной производной и вычислением достижимого уровня значимости (Filtered Derivative with p-Values, FDP-V), предложенный в работе [14].

Определим дискретную производную в точке t как разницу между оценками форм \mathbf{z}^* , вычисленных по выборкам из двух скользящих окон ширины A слева и справа от точки t :

$$D(t, A) = \mathbf{z}^*(t, A) - \mathbf{z}^*(t - A, A),$$

где $\mathbf{z}^*(t, A)$ — оценка формы по выборке $\{\mathbf{x}_i : \tau_t < i \leq \tau_t + A\}$, вычисленная по формуле (4) или (7).

В соответствии с алгоритмом FDr-V локальные максимумы функции $\|D(t, A)\|$ по t рассматриваются как возможные точки разладки. Однако некоторые локальные максимумы могут не соответствовать точкам разладки и возникать вследствие неточности оценок $\mathbf{z}^*(t, A)$.

Алгоритм FDr-V состоит из двух шагов:

1. Поиск потенциальных точек разладки $\tilde{\tau}_k$ как локальных максимумов дискретной производной.
2. Отделение истинных точек разладки от "ложных тревог", применяя для каждой найденной точки $\tilde{\tau}_k$ статистический критерий.

5.1 Поиск потенциальных точек разладки

В качестве потенциальных точек разладки $\tilde{\tau}_k$ выбираются точки локальных максимумов функции $\|D(t, A)\|$ по t , в которых значение $\|D(t, A)\|$ превышает некоторый порог λ . Порог необходимо выбрать таким образом, чтобы вероятность появления в выборке "ложной тревоги" не превышала заданной величины p_1 . В частности, при отсутствии истинных точек разладки необходимо выполнение условия

$$P(\max_t \|D(t, A)\| > \lambda) = p_1. \quad (19)$$

Чтобы оценить значение λ , удовлетворяющее условию (19), введём следующую статистическую модель. Рассмотрим $\mathbf{u}(\mathbf{x}_i)$ как независимые случайные вектора с соответствующими распределениями ξ_i . Нулевая гипотеза

$$H_0 : \quad \xi_1 = \dots = \xi_N$$

отвергается в пользу альтернативы

$$H_1 : \quad \exists K, 0 = \tau_0 < \dots < \tau_{K+1} = N :$$

$$\xi_1 = \dots = \xi_{\tau_1} \neq \xi_{\tau_1+1} = \dots = \xi_{\tau_2} \neq \dots \neq \xi_{\tau_K+1} = \dots = \xi_N$$

при условии

$$\max_t \|D(t, A)\| > \lambda.$$

Порог λ , задающий необходимую вероятность ошибки I рода

$$P(\max_t \|D(t, A)\| > \lambda) = p_1,$$

предлагается оценить с помощью бутстрепа. Для всех $i = 1, \dots, M$, $M \sim 10^4 - 10^5$, выполним следующие шаги:

- Сгенерируем последовательность $\tilde{\mathbf{x}}_i$ с помощью случайной перестановки индексов исходной последовательности \mathbf{x}_i .
- Вычислим $S_i = \max_t \|D(t, A)\|$, где $D(t, A)$ подсчитана по последовательности $\tilde{\mathbf{x}}_i$.

В качестве оценки порога λ возьмём

$$\lambda = S_{(N(1-p_1))}, \tag{20}$$

где $S_{(1)}, \dots, S_{(N)}$ — отсортированные по возрастанию значения S_i .

Таким образом, поиск потенциальных точек разладки включает следующие этапы:

1. Выбор из априорных соображений уровня значимости p_1 и размера окна A , оценка порога λ из (20).
2. Инициализация чисел $d_t = \|D(t, A)\|$, $t \in [A, N - A]$, счётчика $k = 0$.
3. До тех пор, пока $\max_t d_t > \lambda$, выполнять:
 - $k = k + 1$;
 - $\tilde{\tau}_k = \operatorname{argmax}_t d_t$;
 - $d_t = 0$ для всех $t \in (\tilde{\tau}_k - A; \tilde{\tau}_k + A)$.
4. Сортировка чисел $\tilde{\tau}_k$.

5.2 Исключение ложных тревог

Пусть $\tilde{\tau}_1, \dots, \tilde{\tau}_K$ — отсортированная последовательность потенциальных точек разладки. Необходимо отфильтровать из последовательности "ложные тревоги"

— точки, включенные в эту последовательность вследствие флукциаций функции $\|D(t, A)\|$. Поскольку точки $\tilde{\tau}_1, \dots, \tilde{\tau}_{\tilde{K}}$ разбивают последовательность временных рядов \mathbf{x}_i на подпоследовательности с индексами $\mathcal{I}_0 = \{1, \dots, \tilde{\tau}_1\}, \dots, \mathcal{I}_{\tilde{K}+1} = \{\tilde{\tau}_{\tilde{K}} + 1, \dots, N\}$, то исключение точки $\tilde{\tau}_k$ из списка потенциальных точек разладки эквивалентно объединению подпоследовательностей временных рядов $\mathbf{x}_i, i \in \mathcal{I}_{k-1}$ и $\mathbf{x}_i, i \in \mathcal{I}_k$. В качестве критерия объединения используется критерий различимости форм на подвыборках (18).

Выполним следующие шаги:

1. Инициализируем $k = 1$, выбираем уровень значимости p_2 .
2. До тех пор, пока $k \leq \tilde{K}$:
 - Вычисляем p-value по формуле (18), выбрав в качестве множества индексов $\mathcal{I}_1 = \{\tilde{\tau}_{k-1} + 1, \dots, \tilde{\tau}_k\}$, $\mathcal{I}_2 = \{\tilde{\tau}_k + 1, \dots, \tilde{\tau}_{k+1}\}$. (Считаем, что $\tilde{\tau}_0 = 0$, $\tilde{\tau}_{\tilde{K}+1} = N$.)
 - Если p-value $\geq p_2$, то исключаем $\tilde{\tau}_k$ из последовательности, уменьшая \tilde{K} на единицу. В противном случае $k = k + 1$.

6 Критерий качества семейства преобразований \mathbf{f}

Выбор семейства преобразований \mathbf{f} и множества \mathbb{Z} должен отражать экспертные представления о близости форм временных рядов.

Предположим, что существует разбиение индексов

$$\mathcal{I} = \{1, \dots, N\} = \mathcal{I}_1 \cup \dots \cup \mathcal{I}_s,$$

заданное экспертом, такое, что вектора $\mathbf{x}_i, i \in \mathcal{I}_j$ из одного класса имеют схожую форму, в то время как вектора из разных классов имеют разную форму. Необходимо выбрать преобразование \mathbf{f} и множество \mathbb{Z} так, чтобы расстояния между формами временных рядов $\|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|$ наиболее точно отражали представления эксперта, т. е. минимизировать расстояния между формами векторов из одного класса

$$\|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|, \quad J(i) = J(j),$$

и максимизировать расстояния между формами векторов из разных классов

$$\|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|, \quad J(i) \neq J(j),$$

где $J(i) = k \iff i \in \mathcal{I}_k$.

Предлагается вычислить среднее внутриклассовое расстояние

$$F_1 = \frac{\sum_{i < j} [J(i) = J(j)] \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|}{\sum_{i < j} [J(i) = J(j)]}$$

и среднее межклассовое расстояние

$$F_2 = \frac{\sum_{i < j} [J(i) \neq J(j)] \|\mathbf{u}(\mathbf{x}_i) - \mathbf{u}(\mathbf{x}_j)\|}{\sum_{i < j} [J(i) \neq J(j)]}.$$

В качестве критерия качества семейства используется отношение:

$$S(f) = \frac{F_2}{F_1}.$$

7 Иерархическая прогностическая модель

Предположим, что значения

$$y_t, \quad t = 1, \dots, T_0, \quad T_0 = N_0 n$$

временного ряда y_t известны, значения

$$y_t, \quad t = T_0 + 1, \dots, T = N n$$

неизвестны и их необходимо спрогнозировать. Предлагается следующая двухуровневая прогностическая модель.

Временной ряд разбивается на сегменты, как в п. 2:

$$\mathbf{x}_i = (y_{(i-1)n+1}, y_{(i-1)n+2}, \dots, y_{in}), \quad i = 1, \dots, N,$$

при этом значения $\mathbf{x}_i, i = 1, \dots, N_0$ известны, значения $\mathbf{x}_i, i = N_0 + 1, \dots, N$ неизвестны.

Предположим, что существует разбиение индексов $\mathcal{I} = \{1, \dots, N\} = \mathcal{I}_1 \cup \dots \cup \mathcal{I}_s$ такое, что

$$\mathbf{x}_i = \mathbf{f}(\mathbf{z}_{0k}, \boldsymbol{\alpha}_i) + \boldsymbol{\varepsilon}_i, \quad i \in \mathcal{I}_k, \quad \mathbf{z}_{0k} \in Z, \quad k = 1, \dots, s.$$

Построение разбиения $\mathcal{I}_1 \cup \dots \cup \mathcal{I}_s$ выполняется в два шага:

1. Построение априорного разбиения $\tilde{\mathcal{I}}_1 \cup \dots \cup \tilde{\mathcal{I}}_s$ в соответствии с особенностями периодов исходного временного ряда (примеры таких разбиений представлены в п. 4). При этом множество индексов $\tilde{\mathcal{I}}_i$ разбивается на подмножества $\tilde{\mathcal{I}}_j$ и $\tilde{\mathcal{I}}_k$ только тогда, когда между формами сегментов с соответствующими индексами существует значимое различие. Достижимый уровень значимости вычисляется по формуле (18).
2. Каждая группа $\tilde{\mathcal{I}}_k$ разбивается на подгруппы с помощью алгоритма обнаружения разладок, предложенного в п. 5. Индексы, соответствующие сегментам \mathbf{x}_i с неизвестными значениями временного ряда, относятся к последней по времени подгруппе.

Для каждой группы индексов \mathcal{I}_k рассматривается многомерный временной ряд

$$\mathbf{v}(\mathbf{x}_i), i \in \mathcal{I}_k,$$

значения которого известны для $\mathbf{v}(\mathbf{x}_i), i \leq N_0$. Значения $\mathbf{v}(\mathbf{x}_i), N_0 < i \leq N$ оцениваются с помощью одного из алгоритмов прогнозирования многомерных временных рядов.

Таким образом, первый уровень иерархической модели образует полупараметрическая модель для сегментов временного ряда, связывающая целевые значения временного ряда с формой и параметрами полупараметрической модели. На втором уровне рассматриваются последовательности форм и параметров, соответствующие исходной последовательности сегментов.

8 Вычислительный эксперимент

Рассмотрим пример использования алгоритмов, предложенных в работе, для анализа и прогнозирования потребления электроэнергии.

Пусть временной ряд y_t содержит данные о почасовом потреблении электроэнергии в Новосибирске. На рис. 1 показана зависимость потребления электроэнергии от времени в течение недели, с понедельника по воскресенье.

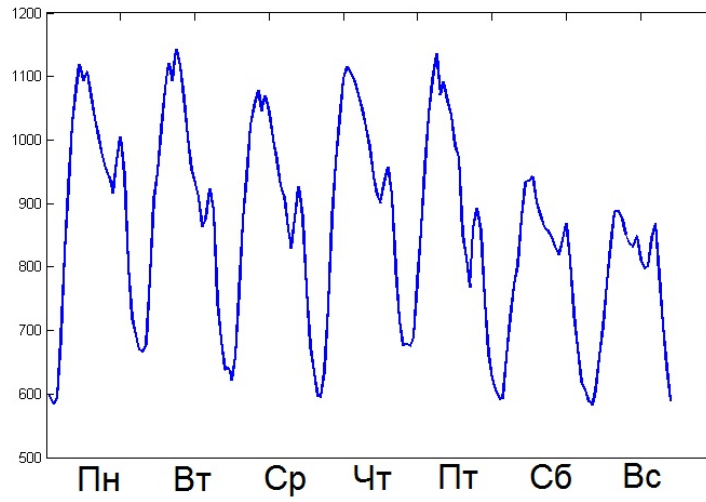


Рис. 1: Потребление электроэнергии в течение недели

Временной ряд y_t разбивается на сегменты \mathbf{x}_i так, что компоненты \mathbf{x}_i содержат данные о потреблении электроэнергии в течении каждого часа i -х суток в году. На рис. 2 показаны сегменты \mathbf{x}_i для будних дней восьми последовательных недель.

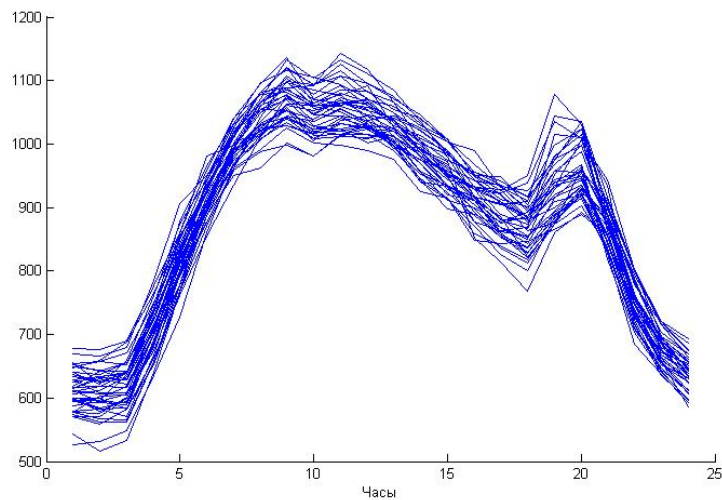


Рис. 2: Потребление электроэнергии в будние дни

Для построения полупараметрической модели выберем в качестве семейства преобразований сдвиг вдоль оси значений временного ряда:

$$\mathbf{f}(\mathbf{x}, \alpha) = \mathbf{x} + (\alpha, \dots, \alpha)^\top.$$

Таким образом, согласно модели (2), все сегменты рассматриваются как один временной ряд, смещённый вдоль оси значений на некоторые величины α_i . Определим множество Z как

$$Z = \{\mathbf{z} : \sum_{j=1}^{24} z_j = 0\}.$$

Тогда для сегмента \mathbf{x}_i смещение

$$\alpha_i = v(\mathbf{x}_i) = \frac{1}{24} \sum_{j=1}^{24} x_j,$$

форма

$$u(\mathbf{x}_i) = \mathbf{x}_i - (v(\mathbf{x}_i), \dots, v(\mathbf{x}_i))^T.$$

Очевидно, для множества Z выполнено условие (3), поэтому в качестве оценки формы \mathbf{z}_0 выборки временных рядов \mathbf{x}_i из модели (2) используем оценку (4). На рис. 3 изображены сегменты \mathbf{x}_i для будних дней восьми последовательных недель (синие линии), и вычисленная по ним форма \mathbf{z}_0 (красная линия), смещённая вверх для компактности графика.

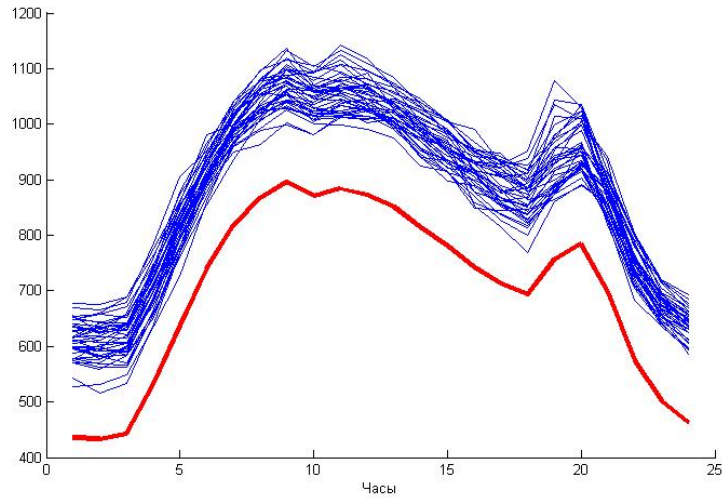


Рис. 3: Форма сегментов временного ряда

Смещения α_i изображены на рис. 4 (выходные на графике пропущены).

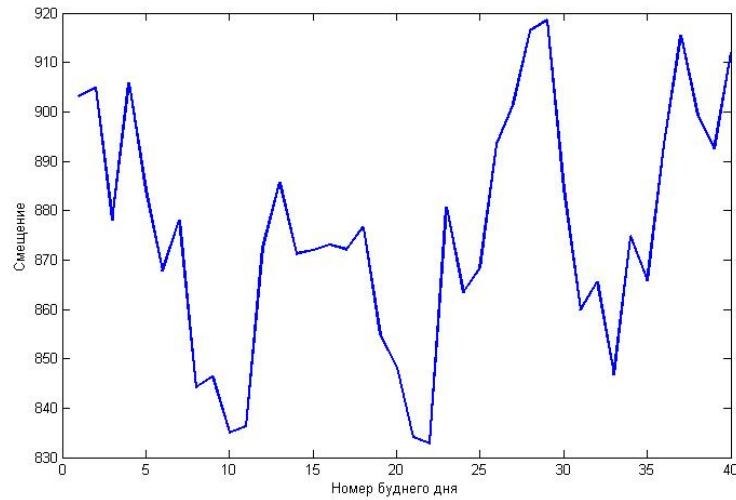


Рис. 4: Смещения сегментов временного ряда

Список литературы

- [1] M. Maggio W. Lawton, E. Sylvestre. Self modeling nonlinear regression. *Technometrics*, 1972.
- [2] Theo Gasser Alois Kneip. Convergence and consistency results for self-modeling nonlinear regression. *The Annals of Statistics*, 1988.
- [3] J. S. Marron W. Hardle. Semiparametric comparison of regression curves. *The Annals of Statistics*, 1990.
- [4] Joachim Engel Alois Kneip. Model estimation in nonlinear regression under shape invariance. *The Annals of Statistics*, 1995.
- [5] Myriam Vimond. Efficient estimation for a subclass of shape invariant models. *The Annals of Statistics*, 2010.
- [6] Daniel Gervini Holger Hurtgen. Semiparametric shape-invariant models for periodic data. 2008.
- [7] Mary J. Lindstrom. Self modeling with random shift and scale parameters and a free-knot spline shape function. *Department of Biostatistics Technical Report*, 1992.
- [8] Jean-Michel Loubes Fabrice Gamboa. Semi-parametric estimation of shifts. *Electronic Journal of Statistics*, 2007.

- [9] Philippe Fraysse Bernard Bercu. A robbins–monro procedure for estimation in semiparametric regression models. *The Annals of Statistics*, 2012.
- [10] Morton B. Brown Yuedong Wang, Chunlei Ke. Shape invariant modelling of circadian rhythms with random effects and smoothing spline anova decompositions. 2002.
- [11] Theo Gasser Kongming Wang. Synchronizing sample curves nonparametrically. *The Annals of Statistics*, 1999.
- [12] B. W. Silverman John A. Rice. Estimating the mean and covariance structure nonparametrically when the data are curve. *Journal of the Royal Statistical Society*, 1991.
- [13] John Staudenmayer Brent A. Coull. Self-modeling regression for multivariate curve data. *Statistica Sinica*, 2004.
- [14] Arnaud Guillin Pierre Bertrand, Mehdi Fhima. Off-line detection of multiple change points with the filtered derivative with p-value method. *Sequential Analysis*, 2010.