

Санкт-Петербургский государственный университет  
Прикладная математика и информатика  
Статистическое моделирование

Мерзляков Климент Викторович

ОБНАРУЖЕНИЕ РАЗЛАДКИ ВО ВРЕМЕННЫХ РЯДАХ ПОКАЗОВ МОБИЛЬНОЙ  
РЕКЛАМЫ

Отчет о научно-исследовательской работе

Научный руководитель:  
к. ф.-м. н., доцент Н. Э. Голяндина

Санкт-Петербург

2018

# Оглавление

<b>Введение</b>	3
<b>Глава 1. Обнаружение разладки во временных рядах</b>	4
1.1. Моделирование данных	4
1.2. Методы обнаружения разладки	6
1.2.1. Методы на основе аппроксимации	7
1.2.2. Методы на основе прогнозирования	8
1.3. Оценка качества	8
<b>Заключение</b>	9
<b>Список литературы</b>	10

## Введение

Рекламной сетью называют некоторую площадку или систему, которая является посредником между рекламодателями и собственниками рекламных мест — владельцами сайтов, мобильных приложений и каких-либо других пространств, где можно размещать рекламу.

В интернет-рекламе взаимодействие рекламной сети с пользователем можно описать следующей последовательностью событий. При выполнении некоторых условий (например, пользователь открыл мобильное приложение) с устройства пользователя отправляется запрос на показ рекламы. Если запрос удовлетворяется, то происходит событие „показ“, то есть пользователь непосредственно видит рекламу. После этого может произойти событие „клик“ и далее какое-либо целевое действие. В мобильной интернет-рекламе „показ“ является одним из ключевых событий, поскольку он отражает количество рекламы доставленное до конечного пользователя.

Рекламные интернет-сети являются интересным объектом для исследования с точки зрения обнаружения разладки, поскольку все показатели отслеживаются с точностью до секунды, происходит большое количество событий, а так как рекламные сети, как правило, работают на международном рынке, то существует возможность тестировать гипотезы на большом количестве различных временных рядов.

Одной из текущих проблем, стоящих перед рекламными сетями — это низкая скорость реагирования на любые резкие изменения текущего состояния. Такие изменения безусловно отражаются в данных в виде аномальных значений, резких всплесков и внезапных изменений тренда. Однако проблема заключается в том, что показателей требующих отслеживания могут быть десятки, при этом на каждый показатель может влиять большое количество факторов. Поэтому зачастую, чтобы локализовать и устранить проблему требуется просмотреть сотни графиков. Отсюда следует, что наличие качественного метода обнаружения разладки каждого показателя по каждому измерению позволило бы не только существенно сэкономить ресурсы, но и в целом повысить эффективность бизнеса. Поэтому целью данной работы является разработка методики обнаружения разладки. В работе будут использоваться фактические, данные одной из работающих рекламных сетей.

## Глава 1

## Обнаружение разладки во временных рядах

## 1.1. Моделирование данных

Реальные данные интернет-рекламы имеют стабильную дневную периодичность (на рисунке 1.1 приведен пример типичной динамики в рамках дня). По более длинному ряду изображенному на рисунке 1.2 видно, что в данных время от времени возникают разладки разных видов, при этом сам ряд имеет мультипликативный характер (с изменением среднего уровня ряда пропорционально меняется и амплитуда колебаний). В реальных временных рядах достаточно сложно разметить наличие разладок — зачастую сложно отделить разладку от шума. Поэтому вместо разметки реальных рядов мы будем моделировать искусственные ряды похожие на ряды данных интернет-рекламы с определенным шумом и разладками в известных местах.

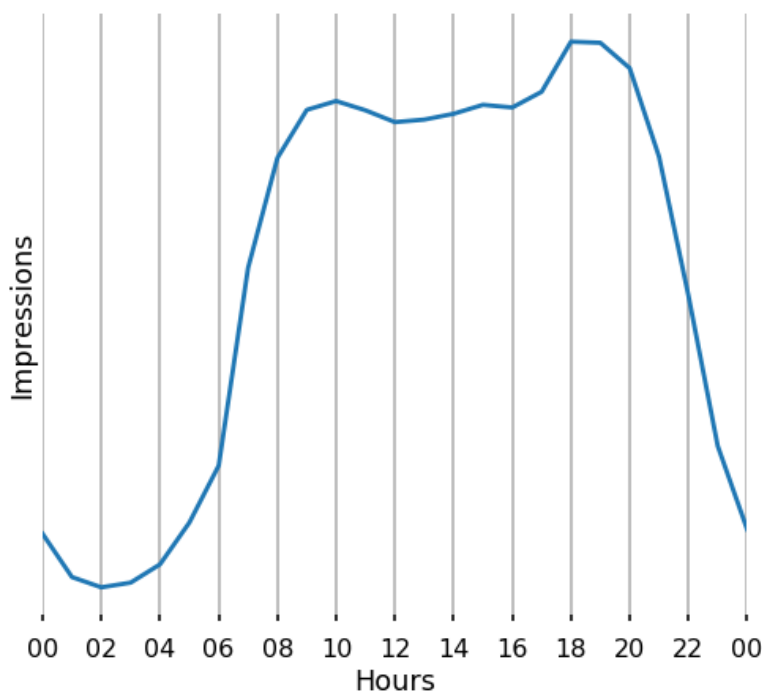


Рис. 1.1. Пример показов рекламы за сутки

Обозначим временной ряд  $Y = (y_1, \dots, y_n)$ . Наблюдаемые значения ряда можно представить в виде суммы компонент:

$$Y = T + S + E,$$

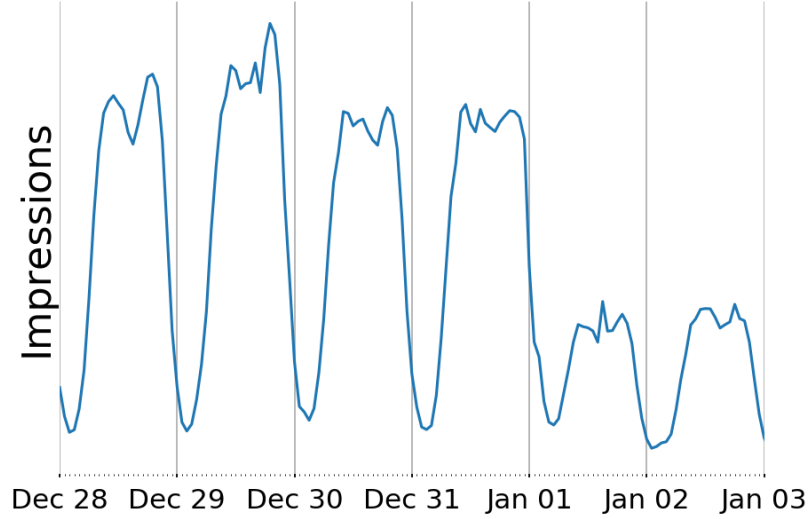


Рис. 1.2. Пример показов рекламы с разладкой за неделю

где  $T = (t_1, \dots, t_n)$  компонента-тренд,  $S = (s_1, \dots, s_n)$  периодическая компонента,  $E = (\epsilon_1, \dots, \epsilon_n)$  остатки или шум. По каждой из этих компонент требуется построить модель. Это можно сделать, например, следующим образом:

$$t_i = c, \quad i = 1, \dots, n,$$

$$s_i = A \cos\left(\frac{2\pi}{a}i + \phi\right), \quad i = 1, \dots, n,$$

$$\epsilon_i \sim N(\mu, \sigma^2), \quad i = 1, \dots, n,$$

где  $i$  индекс элемента ряда;  $c$  — константа,  $A$  — амплитуда;  $a$  — период;  $\phi$  — фаза.

Ряды, модель которых мы хотим построить имеют мультипликативность (амплитуда колебаний меняется пропорционально изменению тренда). Такого эффекта можно достичь, взяв экспоненту от исходной модели ряда:

$$Y^{(\text{mult})} = e^Y.$$

Далее нам нужно построить модель разладки, исходя из следующего:

- Разладка только в одной точке ряда;
- Разладка только в тренде и заключается в сдвиге;
- Разладка может произойти не всегда, а с некоторой вероятностью  $\rho$ .

Формально это можно описать так: пусть  $\tau$  — точка (индекс) разладки, тогда тренд с разладкой разумно будет записать  $\tilde{T} = (\tilde{t}_1, \dots, \tilde{t}_n)$ . И элементы тренда с разладкой будут иметь вид:

$$\tilde{t}_i = \begin{cases} t_i, & i < \tau, \\ t_i + \delta, & i \geq \tau, \end{cases}$$

где  $\delta$  — значение разладки.

Значение разладки является случайной величиной с некоторым распределением. В данной работе значение разладки будет иметь нормальное распределение  $\delta^* \sim N(\mu^{(cp)}, \sigma^{2(cp)})$ , с некоторой вероятностью возникновения  $\rho$ :

$$\delta = \begin{cases} \delta^*, & \text{с вероятностью } \rho, \\ 0, & \text{с вероятностью } 1 - \rho. \end{cases}$$

Таким образом,  $\delta$  является случайной величиной с распределением-смесью. При этом точка разладки  $\tau$  тоже является случайной величиной с равномерным распределением на  $[n_0, \dots, n]$ , где  $n_0$  — самая первая возможная точка разладки, которая задается параметром.  $n_0$  введена намеренно, чтобы разладка при моделировании не возникала в первых точках ряда.

Таким образом, моделируемый ряд с разладкой будет иметь следующий вид:

$$\tilde{Y}^{(mult)} = e^{\tilde{T}+S+E}.$$

В результате модель временного ряда имеет пять параметров:  $A, a, \phi, \mu, \sigma, c$ , а модель разладки имеет еще пять параметров:  $\delta, \mu^{(cp)}, \sigma^{(cp)}, \rho, n_0$ .

## 1.2. Методы обнаружения разладки

Опишем один из подходов к обнаружению разладки. Данный подход не является единственным, хотя включает в себе широкое разнообразие методов. Как правило, у временного ряда есть некоторая структура (сигнал), которая может быть описана той или иной моделью. Идея подхода заключается в том, что около точки разладки модель плохо описывает временной ряд. Используя некоторую меру ошибки мы можем измерять то, насколько хорошо или плохо описывает выбранная модель реальные данные. Как только ошибка (отклонение модели от реальных данных) превышает заданный порог, метод сигнализирует о разладке.

Можно выделить два типа методов в данном подходе:

- Методы на основе предсказания

- Методы на основе аппроксимации

### 1.2.1. Методы на основе аппроксимации

Пусть  $l$  четное вещественное число, называемое шириной окна. При этом  $1 < l < n$ . С помощью ширины окна из исходного ряда образуется последовательность подрядов  $W = \{w_j\}_{j=1}^k$ , где  $k = n - l + 1$  — количество таких подрядов; а  $w_j = (y_j, \dots, y_{j+l+1})$  —  $j$ -ый подряд. Каждый подряд  $w_j$  в свою очередь делится на два подряда одинаковой длины (это возможно поскольку  $l$  четное по условию):  $W^{(left)} = \{w_j^{(left)}\} = (y_j, \dots, y_{\frac{j+l+1}{2}})$  и  $W^{(right)} = \{w_j^{(right)}\} = (y_{\frac{j+l+1}{2}+1}, \dots, y_{j+l+1})$ .

Таким образом, для каждого ряда  $W$  можно сформировать тройки рядов:

$$W^{(all)} = \{w_j^{(all)}\}_{j=1}^k = \{(w_j; w_j^{(left)}; w_j^{(right)})\}_{j=1}^k.$$

Пусть есть функция ошибки  $e(\cdot)$ , такая что:

$$e(X) = \min_{\theta} \sum_{p=1}^m (x_p - f(x_p|\theta))^2,$$

где  $X = (x_1, \dots, x_m)$  некоторый вещественный временной ряд длины  $m$ , а  $f(x|\theta)$  некоторая модель этого временного ряда с параметрами  $\theta$ .

Функция  $f(x|\theta)$  может быть константной ( $\theta = (k)$ ):

$$f(x|k) = k.$$

Либо другой подходящей под наш ряд функцией, например:

$$f(x|P, p, \chi) = P \cos\left(\frac{2\pi}{p}x + \chi\right) + k.$$

Мера ошибки позволяет нам рассчитать, насколько хорошо аппроксимируется отрезок ряда с помощью выбранной модели. Однако, для обнаружения самой разладки необходимо еще ввести функцию разладки:

$$F(W^{(all)}) = \frac{e(W) - e(W^{(left)}) - e(W^{(right)})}{h},$$

где  $h$  — значение нормировки.

Расчет нормирующей константы является открытой проблемой, поскольку имеются разные варианты её расчета со своими плюсами и минусами. Например, можно рассчитывать её как значение функции ошибки на первом отрезке ряда (предполагая, что на этом отрезке не происходило разладок):

$$h = f(w_1^{(all)}) = e(w_1) - e(w_1^{(left)}) - e(w_1^{(right)}).$$

Итого, взяв ряд  $Y$  мы «скользим» по нему окном длины  $l$  и рассчитываем значения функции разладки  $F()$  для каждого из получаемых подрядов  $W^{(all)}$ . Функция разладки начинает расти в окрестности точки разладки  $\tau$ , следовательно можно задать некий порог  $\gamma$ , такой что при превышении функции разладки этого порога в какой-то точке  $\hat{\tau}$  будем считать, что разладка обнаружена в этой точке.

### 1.2.2. Методы на основе прогнозирования

## 1.3. Оценка качества

В рамках данной работы мы разрабатываем систему своевременного оповещения о разладках во временных рядах. При такой постановке задачи важны две характеристики: точность обнаружения разладки и скорость обнаружения разладки. Поскольку мы используем моделированные данные, то мы точно знаем в каких из смоделированных нами рядов произошла разладка, а в каких разладки не было. Более того, мы точно знаем момент разладки. Благодаря этому мы можем строить матрицы сопряжённости и считать метрики качества. При этом надо понимать, что в такой постановке задачи нам важно обнаружить разладку не позднее какого-то срока, иначе оповещение о разладке будет несвоевременным.

Разберем четыре возможных варианта:

- Разладка произошла и метод обнаружил точку разладки **после** фактической точки, но не слишком поздно. Такая ситуация попадает под категорию True positive
- Разладка произошла и метод не обнаружил точку разладки, либо обнаружил её до фактической разладки, либо обнаружил её сильно после фактической разладки. False negative
- Разладки не было, но метод обнаружил разладку. False positive
- Разладки не было и метод не обнаружил разладку. True negative

Договорившись о таком способе оценки качества, можно строить ROC-кривые (изменяя порог  $\alpha$ ) для разных методов обнаружения разладки и таким образом, сравнивая как работают те или иные методы в контролируемой среде эксперимента.



## Заключение

## Список литературы

1. Голяндина Н.Э. Метод „Гусеница“-SSA: прогноз временных рядов. Учебное пособие. СПб., 2003.
2. Armstrong. Principles of Forecasting: A Handbook for Researchers and Practitioners. Kluwer Academic Publishers, 2001.
3. Dagum, Estela, Bianconcini. Seasonal Adjustment Methods and Real Time Trend-Cycle Estimation, 2016.
4. N. Golyandina, V. Nekrutkin, A. Zhigljavsky. Analysis of Time Series Structure - SSA and Related Techniques, 2001.
5. R. Hyndman, G. Athanasopoulos. Forecasting: Principles and Practice. 2013.
6. R. Hyndman. Moving averages. 2009.