

# Обзор некоторых статистических методов скорейшего обнаружения

Алексей Артемов

8 мая 2014

# Общая постановка задачи (дискретное время)

- Предполагается, что наблюдаемые данные описываются числовой последовательностью

$$X_1, X_2, \dots, X_{\theta-1}, X_{\theta}, \dots,$$

— это результаты наблюдений над случайными величинами

$$\xi_1, \xi_2, \dots, \xi_{\theta-1}, \xi_{\theta}, \dots$$

- Распределение процесса  $\xi = (\xi_t)_{t \geq 0}$ :
  - до момента  $\theta - P_{\infty}$
  - после момента  $\theta - P_0$
- Момент разладки  $\theta$  заранее неизвестен

# Общая постановка задачи (дискретное время)

- Варианты постановок задач:
  - offline проверка гипотезы о наличии изменения
  - offline обнаружение изменения (сегментация)
  - online обнаружение изменения
    - finite horizon:  $0 \leq t \leq T < \infty$
    - infinite horizon:  $0 \leq t \leq \infty$
- Момент разладки  $\theta$  может быть:
  - случайной величиной с некоторым распределением
  - просто неизвестным параметром
- Существует множество различных критериев оптимальности методов скорейшего обнаружения (= постановок задач)

# Общая постановка задачи (дискретное время)

- Можно использовать свойства наблюдений
  - зависимость и независимость
  - скалярность и векторность
  - распределение
- Бывает необходимо учитывать, что
  - разладка может появляться более одного раза в течение наблюдений
  - параметры процесса после разладки, как правило, неизвестны
- ...и всё это только для стационарных процессов

## Из этого доклада вы не узнаете

Ввиду объема материала, сложности изложения, релевантности поставленным задачам, а также личного опыта докладчика опущены:

- случай finite horizon
- сегментация временного ряда освещается ограниченно
- оценивание значения временного ряда
- случай зависимых наблюдений
- векторный случай

# Понятие разладки

- Свойства данных всюду предполагаются либо постоянными, либо медленно меняющимися
- Многие адаптивные алгоритмы оценивания хорошо работают для медленных изменений
- Разладка — это **быстрое** по сравнению с шагом дискретизации изменение

## Общие обозначения

- $\theta \in \{0, 1, \dots, \infty\}$  — момент разладки
- $P_\theta$  — распределение процесса  $\xi$  в предположении, что разладка происходит в момент  $\theta$
- $P_0$  — распределение процесса  $\xi$  в предположении, что разладка произошла с самого начала ( $\theta = 0$ )
- $P_\infty$  — распределение процесса  $\xi$  в предположении, что разладка не происходит никогда ( $\theta = \infty$ )
- $f_\theta(\cdot), f_0(\cdot), f_\infty(\cdot)$  — плотности мер  $P_\theta, P_0, P_\infty$
- $\mathbf{E}_\theta, \mathbf{E}_0, \mathbf{E}_\infty$  — средние по мерам  $P_\theta, P_0, P_\infty$

# Метод Неймана-Пирсона

- Различение двух гипотез по фиксированному числу наблюдений
- Требуется принять одну из гипотез  $H_0(\theta = 0)$  или  $H_\infty(\theta = \infty)$
- Решающее правило  $d = d(X_1, \dots, X_n)$  характеризуется вероятностями ошибок I и II рода  $\alpha(d) = P_\infty(H_0)$ ,  $\beta(d) = P_0(H_\infty)$ ,
- Возможные постановки задач:
  - $\alpha(d) + \beta(d) \sim \inf_d$
  - $\beta(d) \sim \inf_{d \in D_\alpha}$ ,  $D_\alpha = \{d : \alpha(d) \leq \alpha\}$
- Статистика — отношение правдоподобия

$$L_n = \frac{f_0(X_1, \dots, X_n)}{f_\infty(X_1, \dots, X_n)}$$



## Метод Неймана-Пирсона

- Если  $L_n > h_\alpha$ , то принимается  $H_0$ . Иначе —  $H_\infty$ .
- Удобно также использовать статистику  $Z_n = \log L_n$
- **Пример.** Пусть  $\xi_1, \dots, \xi_n$  — нормальные i.i.d.r.v., причем

$$f_0(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-r_0)^2}{2\sigma^2}}, \quad f_\infty(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-r_\infty)^2}{2\sigma^2}}$$

Тогда

$$Z_n = \frac{r_0 - r_\infty}{\sigma^2} \left[ \bar{X}_n - \frac{r_0 + r_\infty}{2} n \right]$$

Оптимальное решающее правило:

$$d(X_1, \dots, X_n) = \begin{cases} H_0 & \text{если } Z_n \geq h \\ H_\infty & \text{если } Z_n < h. \end{cases}$$

- Величина  $h$  выбирается исходя из заданных величин ошибок I и II рода

# Контрольные карты Шухарта

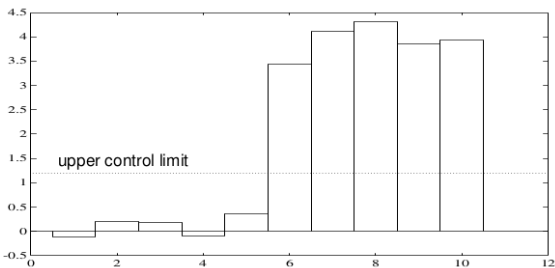
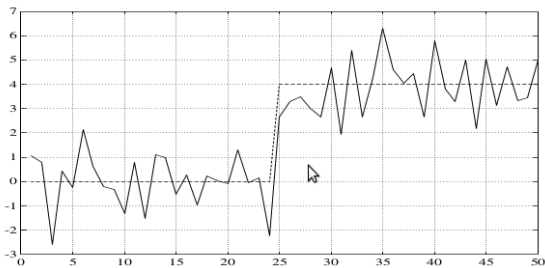
- Наблюдения  $X_1, X_2, \dots$  разбиваются на группы размера  $N$
- Для каждой группы  $X_K = (X_{(K-1)N+1}, \dots, X_{KN})$ ,  $K = 1, 2, \dots$  ставится задача различения двух гипотез:

$$H_0 : \theta \leq (j-1)N + 1 \quad \text{и} \quad H_\infty : \theta > (j-1)N + 1$$

- Решающее правило  $d_K = d(X_K)$  следует из леммы Неймана-Пирсона
- Момент остановки — первый момент принятия гипотезы  $H_0$

$$\tau = N \cdot \min\{K : d_K = H_0\}$$

## Пример [Basseville et al., 1993]



## Последовательное различение гипотез

- Последовательный тест — пара  $(\tau, \varphi)$ :
  - $\tau = \tau(X) \in \{0, 1, \dots, \infty\}$  — момент остановки относительно  $\{\mathcal{F}_n, n \geq 1\}$
  - $\varphi = \varphi(X) \in [0, 1]$  — решающая функция
- Останавливаем наблюдения в момент  $\tau$  и принимаем гипотезу  $H_0$  с вероятностью  $\varphi(X_1, \dots, X_\tau)$
- Характеристики — средние длительности наблюдений  $\mathbf{E}_0\tau$  и  $\mathbf{E}_\infty\tau$ , а также вероятности ошибок I и II рода

$$\alpha(\varphi) = \mathbf{E}_\infty\varphi \quad \text{и} \quad \beta(\varphi) = \mathbf{E}_0(1 - \varphi)$$

- Пусть заданы два числа  $\alpha$  и  $\beta$  и класс тестов

$$\Delta(\alpha, \beta) = \{\delta = (\tau, \varphi) : \alpha(\varphi) \leq \alpha, \beta(\varphi) \leq \beta, \mathbf{E}_0\tau \leq 0, \mathbf{E}_\infty\tau \leq \infty\}$$

Критерий оптимальности теста  $\delta^* = (\tau^*, \varphi^*) \in \Delta(\alpha, \beta)$ :

$$\forall (\tau, \varphi) \in \Delta(\alpha, \beta) \quad \mathbf{E}_0\tau^* \leq \mathbf{E}_0\tau, \quad \mathbf{E}_\infty\tau^* \leq \mathbf{E}_\infty\tau$$

# Метод Вальда

- Выбираются две константы  $A < 0$  и  $B > 0$  и момент

$$\tau_{A,B} = \inf\{n \geq 1 : Z_n \notin (A, B)\}$$

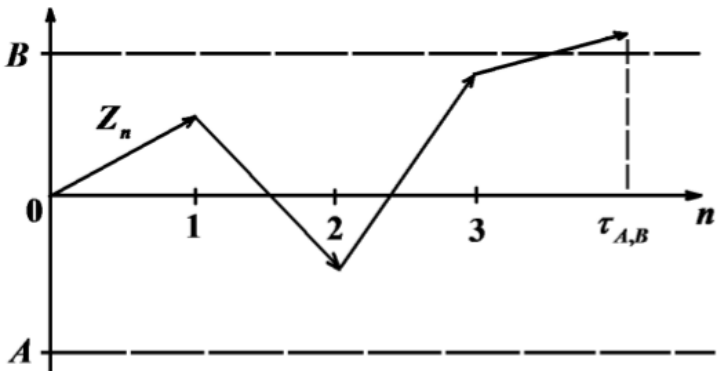
- Решающее правило

$$\varphi_{A,B} = \begin{cases} 1 & \text{если } Z_{\tau_{A,B}} \geq B, \\ 0 & \text{если } Z_{\tau_{A,B}} \leq A. \end{cases}$$

- Пороги  $A, B$  связаны с заданными вероятностями ошибок I и II рода  $\alpha, \beta$  соотношениями

$$A = \log \frac{\beta}{1 - \alpha}, \quad B = \log \frac{1 - \beta}{\alpha}$$

## Пример [Ширяев, 2011]



## Геометрическое среднее (EWMA)

- Подходит для обнаружения изменения среднего значения или дисперсии временного ряда
- Идея в рекурсивной оценке среднего вида

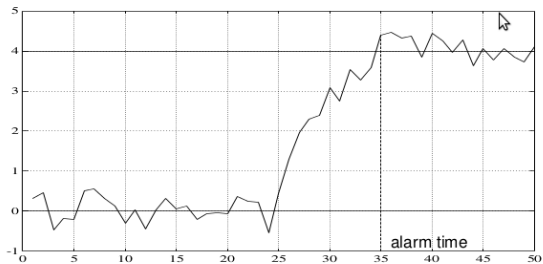
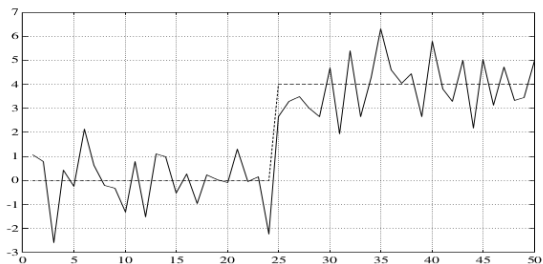
$$\hat{m}_k = (1 - \lambda)\hat{m}_{k-1} + \lambda X_k, \quad k = 1, 2, \dots,$$

где  $\lambda$  («вес» новых данных) — параметр алгоритма.

- Момент остановки — момент первого выхода статистики  $\hat{m}_k$  на заданный уровень  $h$ :

$$\tau = \min\{k \geq 1 : \hat{m}_k \geq h\}$$

## Пример [Basseville et al., 1993]





# Оптимальные алгоритмы

- Рассмотрим некоторые постановки задач скорейшего обнаружения в режиме online
- Обычно ставятся для детектирования появления сноса у броуновского движения (в непрерывном времени):

$$\xi_t = \begin{cases} W_t & \text{если } 0 \leq t < \theta \\ \mu t + W_t & \text{если } t \geq \theta. \end{cases}$$

- Для дискретного времени возможна постановка:

$$\xi_t = \begin{cases} Z_t & \text{если } 0 \leq t < \theta \\ \mu + Z_t & \text{если } t \geq \theta, \end{cases}$$

где  $Z_t, t = 1, 2, \dots$  — нормальные i.i.d.r.v.

# Кумулятивные суммы

- О параметре  $\theta$  не делается никаких предположений
- Фиксируется некоторое число  $T > 0$  и задается класс

$$\mathcal{M}_T = \{\tau : \mathbf{E}_\infty \tau \geq T\}$$

тех моментов остановки, для которых среднее время до ложной тревоги не меньше  $T$ .

- Качество алгоритма задается величиной

$$\mathbf{D}(T) = \sup_{\theta \geq 0} \operatorname{ess\,sup}_\omega \underbrace{\mathbf{E}_\theta((\tau - \theta)^+ | \mathcal{F}_\theta)(\omega)}_{\substack{\text{среднее время} \\ \text{обнаружения разладки}}} \sim \inf_{\tau \in \mathcal{M}_T} \underbrace{\quad}_{\substack{\text{наихудшее среди всех траекторий} \\ \text{и всех моментов } \theta \text{ появления разладки}}}$$

# Кумулятивные суммы

- Вводятся статистики

$$\gamma_n = \sup_{\theta \geq 0} \frac{f_\theta(X_1, \dots, X_n)}{f_\infty(X_1, \dots, X_n)} \quad \text{и} \quad T_n = \log \gamma_n$$

- Если случайные величины  $\xi_1, \dots, \xi_n$  независимы, то

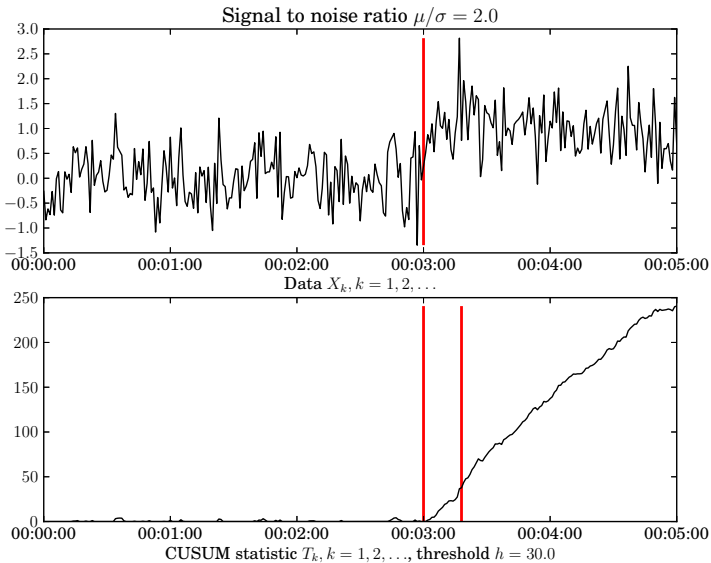
$$\gamma_n = \max \left\{ 1, \max_{1 \leq \theta \leq n} \prod_{k=\theta}^n \frac{f_0(X_k)}{f_\infty(X_k)} \right\},$$

$$T_n = \max \left\{ 0, \max_{1 \leq \theta \leq n} \sum_{k=\theta}^n \log \frac{f_0(X_k)}{f_\infty(X_k)} \right\} = \max \left\{ 0, \max_{1 \leq \theta \leq n} \sum_{k=\theta}^n \zeta_k \right\},$$

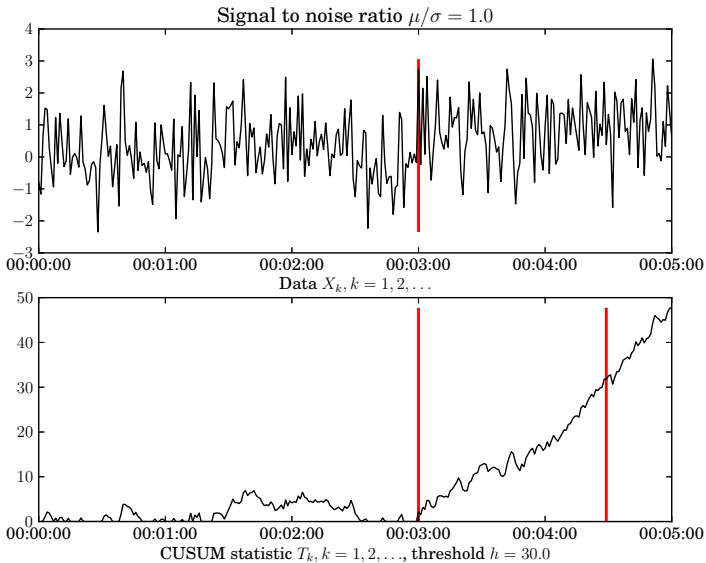
- Статистика  $T_n$  обладает свойством  $T_n = \max(0, T_{n-1} + \zeta_n)$  и называется статистикой кумулятивных сумм (CUMulative SUMs, CUSUM).
- Остановка в момент  $\tau_{\text{CUSUM}}$  минимизирует величину  $\mathbf{D}(T)$

$$\tau_{\text{CUSUM}} = \inf \{ n \geq 0 : T_n \geq h \}$$

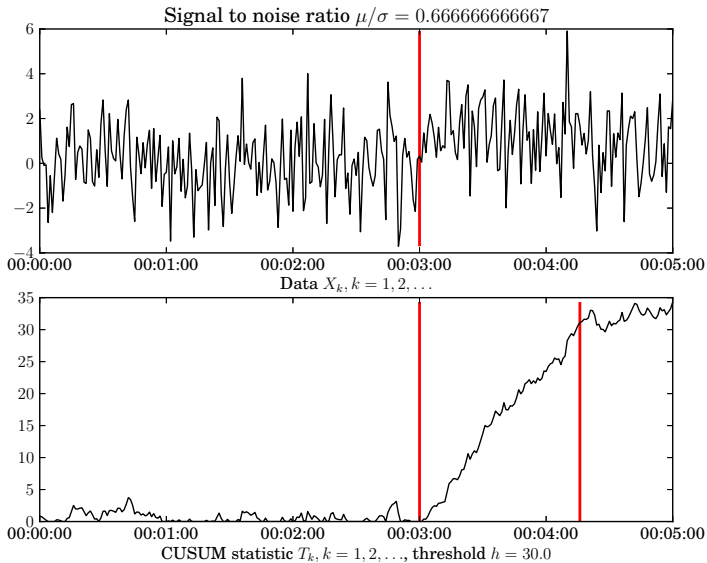
## Пример 1 [Razladki]



## Пример 2 [Razladki]



## Пример 3 [Razladki]



# Статистика Ширяева-Робертса

- Условия для моментов остановки и параметра  $\theta$  в точности совпадают с условиями для процедуры кумулятивных сумм
- Качество алгоритма задается величиной

$$C(T) = \sup_{\theta \geq 0} \underbrace{\mathbf{E}_\theta((\tau - \theta)^+ | \tau \geq \theta)}_{\substack{\text{(условное) среднее время} \\ \text{обнаружения разладки}}} \sim \inf_{\tau \in M_T} \underbrace{\phantom{\mathbf{E}_\theta((\tau - \theta)^+ | \tau \geq \theta)}}_{\substack{\text{наихудшее среди всех} \\ \text{моментов } \theta \text{ появления разладки}}}$$

# Статистика Ширяева-Робертса

- Вводится статистика

$$R_n = \sum_{\theta=1}^n \frac{f_{\theta}(X_1, \dots, X_n)}{f_{\infty}(X_1, \dots, X_n)}$$

- Если случайные величины  $\xi_1, \dots, \xi_n$  независимы, то

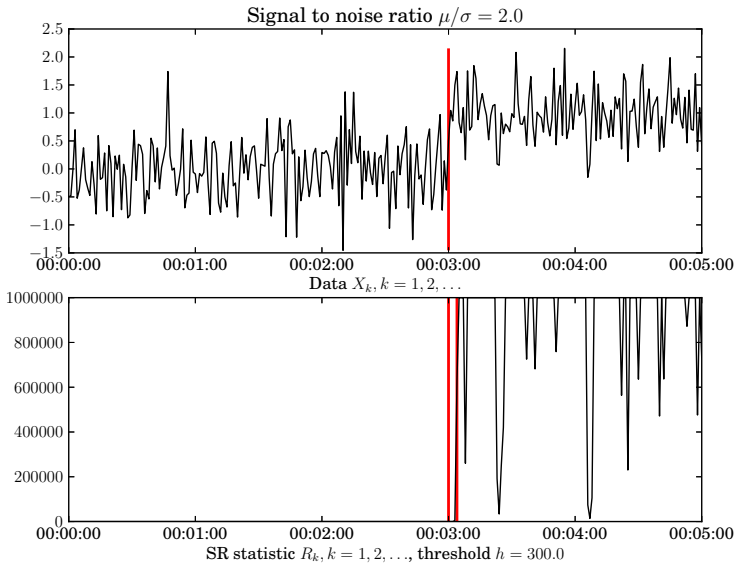
$$R_n = \sum_{\theta=1}^n \prod_{k=\theta}^n \frac{f_0(X_k)}{f_{\infty}(X_k)} = \sum_{\theta=1}^n \prod_{k=\theta}^n l_k.$$

- Статистика  $R_n$  обладает свойством  $R_n = (1 + R_{n-1})l_n$  и называется статистикой Ширяева-Робертса (Shiryaev-Roberts, SR).
- Остановка в момент  $\tau_{\text{SR}}$  минимизирует величину  $\mathbf{C}(T)$

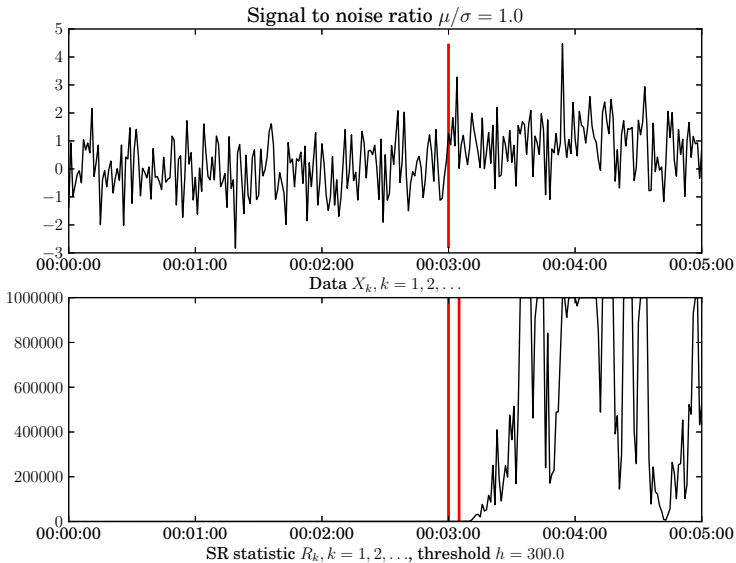
$$\tau_{\text{SR}} = \inf\{n \geq 0 : R_n \geq h\}$$



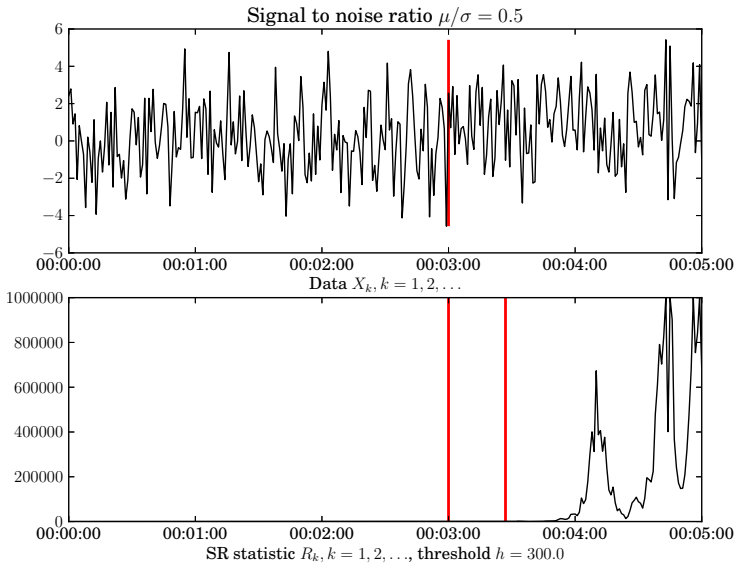
## Пример 1 [Razladki]



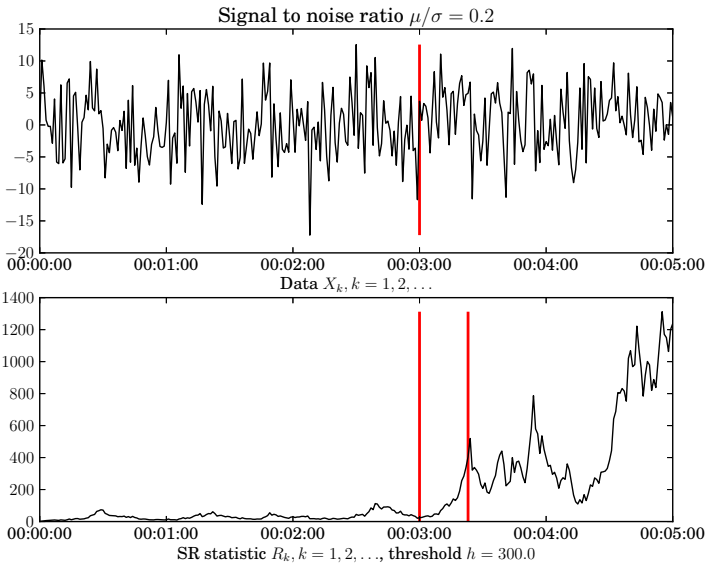
## Пример 2 [Razladki]



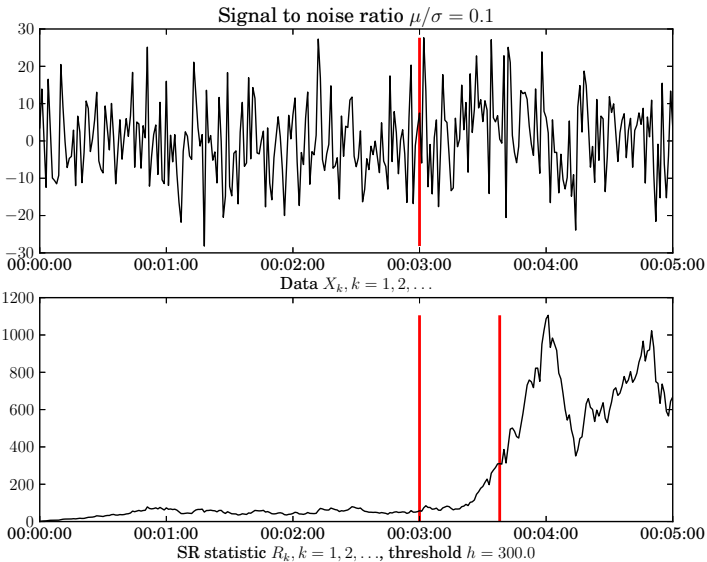
## Пример 3 [Razladki]



## Пример 4 [Razladki]



## Пример 5 [Razladki]



# Статистика $\pi_t$

- Пусть  $\theta = \theta(\omega)$  — случайная величина,  $\theta \perp Z_t$ , имеющая (геометрическое) распределение

$$P(\theta = 0) = \pi, \quad P(\theta = n | \theta > 0) = pq^{n-1},$$

причем  $\pi \in [0, 1)$  и  $p \in (0, 1)$  известны,  $q = 1 - p$ .

- Фиксируется некоторое число  $\alpha \in (0, 1]$  и задается класс

$$\mathcal{M}_\alpha = \{\tau : P(\tau < \theta) \leq \alpha\}$$

тех моментов остановки, для которых вероятность ложной тревоги не выше  $\alpha$ .

# Статистика апостериорной вероятности

- Качество алгоритма задается величиной

$$\mathbf{E}(\tau - \theta | \tau > \theta) \sim \inf_{\tau \in \mathcal{M}_\alpha}$$

- Этот критерий эквивалентен условно-байесовскому критерию

$$\mathbf{A}(c) = \underbrace{P(\tau < \theta)}_{\text{вероятность ложной тревоги}} + \underbrace{c\mathbf{E}(\tau - \theta | \tau > \theta)}_{\text{(условное) среднее время обнаружения разладки}} \sim \inf_{\tau \in \mathcal{M}_\alpha}$$

# Статистика апостериорной вероятности

- Вводится статистика

$$\pi_n = \frac{\varphi_n}{1 + \varphi_n},$$

где для  $\varphi_n$  справедливо рекуррентное соотношение

$$\varphi_{n+1} = (p + \varphi_n) \frac{f_0(X_n)}{qf_\infty(X_n)}.$$

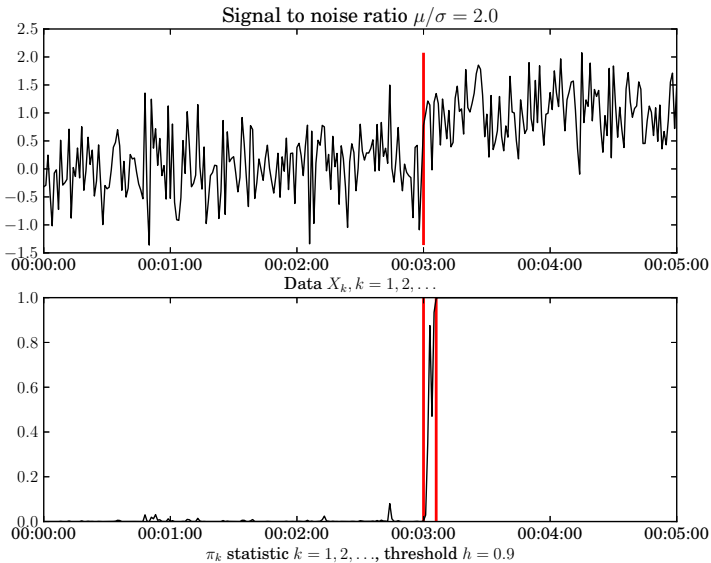
- Остановка в момент  $\tau_\pi$  минимизирует величину  $\mathbf{A}(c)$ :

$$\tau_\pi = \inf\{n \geq 0 : \pi_n \geq h\}$$

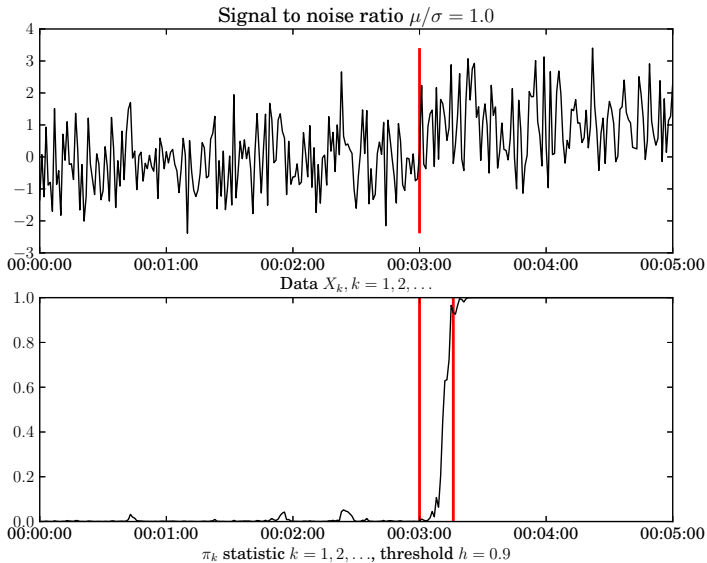
- Значение  $\pi_n$  — апостериорная вероятность появления разладки до момента времени  $n$  в предположении, что получены наблюдения  $\{X_k, k = \overline{1, n}\}$ .



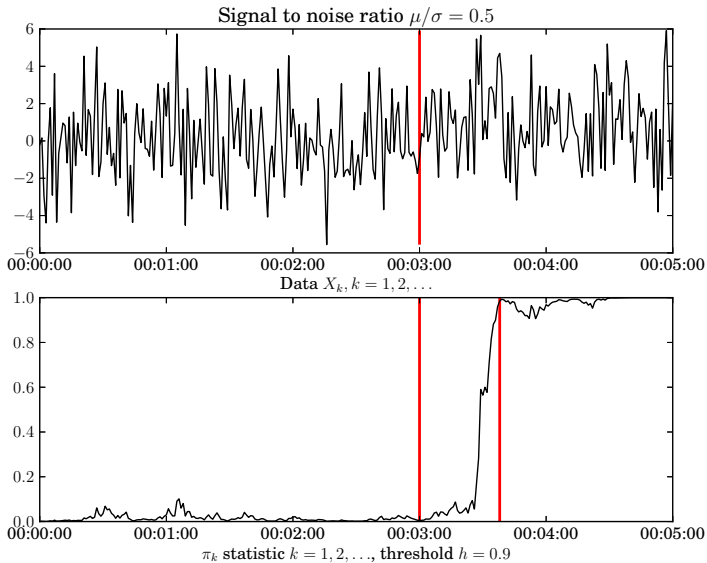
## Пример 1 [Razladki]



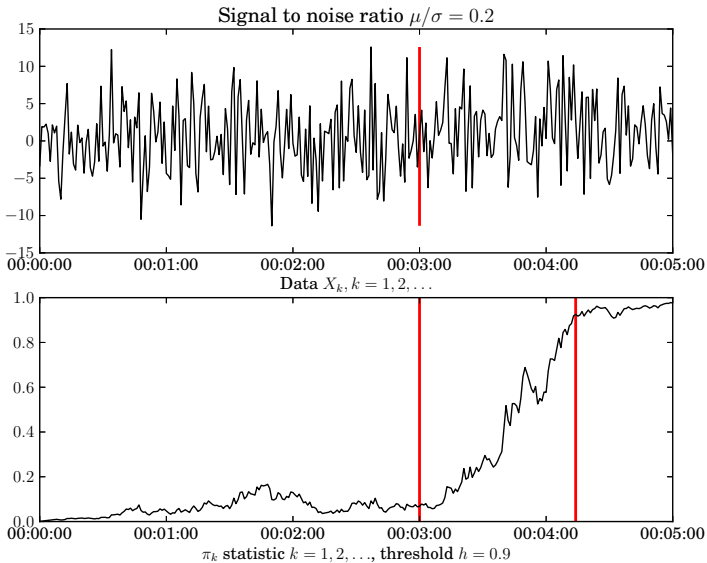
## Пример 2 [Razladki]



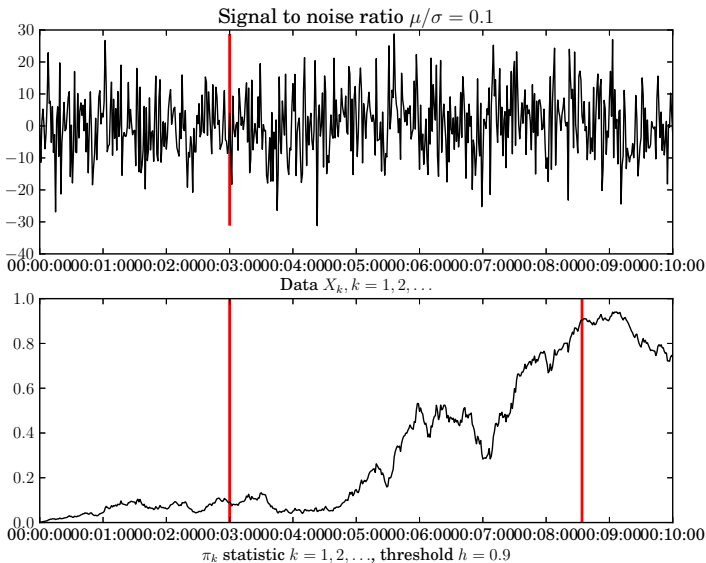
## Пример 3 [Razladki]



## Пример 4 [Razladki]



## Пример 5 [Razladki]



# Общая литература

## Базовая литература по теории оптимальной остановки

- Обзор методов неоптимальной и оптимальной остановки с примерами и обсуждением [Basseville et al., 1993]
- Обзор методов оптимальной остановки с доказательствами и выводом формул [Ширяев, 2011]

# Литература по оптимальным методам

## Кумулятивные суммы

- Введены в [Page, 1954]
- Критерий оптимальности — [Lorden et al., 1971]
- Оптимальность для дискретного случая — [Moustakides et al., 1986], [Ritov, 1990]
- Оптимальность для непрерывного случая — [Ширяев, 1996]

## Статистика Ширяева-Робертса

- Введена независимо в [Ширяев, 1961] и [Roberts, 1966]

## Статистика апостериорных вероятностей

- Введена в [Ширяев, 1969]

# Литература I

Michèle Basseville, Igor V Nikiforov, et al. *Detection of abrupt changes: theory and application*, volume 104. Prentice Hall Englewood Cliffs, 1993.

Gary Lorden et al. Procedures for reacting to a change in distribution. *The Annals of Mathematical Statistics*, 42(6): 1897–1908, 1971.

George V Moustakides et al. Optimal stopping times for detecting changes in distributions. *The Annals of Statistics*, 14(4): 1379–1387, 1986.

ES Page. Continuous inspection schemes. *Biometrika*, pages 100–115, 1954.

Ya'acov Ritov. Decision theoretic optimality of the cusum procedure. *The Annals of Statistics*, pages 1464–1469, 1990.



## Литература II

SW Roberts. A comparison of some control chart procedures.  
*Technometrics*, 8(3):411–430, 1966.

Альберт Николаевич Ширяев. Задача скорейшего обнаружения нарушения стационарного режима. In *Докл. АН СССР*, volume 138, pages 1039–1042, 1961.

Альберт Николаевич Ширяев. *Статистический последовательный анализ: Оптимальные правила остановки*. Наука, 1969.

Альберт Николаевич Ширяев. Минимаксная оптимальность метода кумулятивных сумм (cusum) в случае непрерывного времени. *Успехи математических наук*, 51(4 (310)):173–174, 1996.

Альберт Николаевич Ширяев. *Вероятностно-статистические методы в теории принятия решений*, volume 144. МЦНМО, 2011.