

Фильтрация сигналов с трендом в задачах обнаружения разладки

Алексей Артёмов

5 июня 2016 г.

Deep Machine Intelligence Workshop,
Skolkovo Institute of Science & Technology

Работа выполнена совместно с Е. В. Бурнаевым и др.

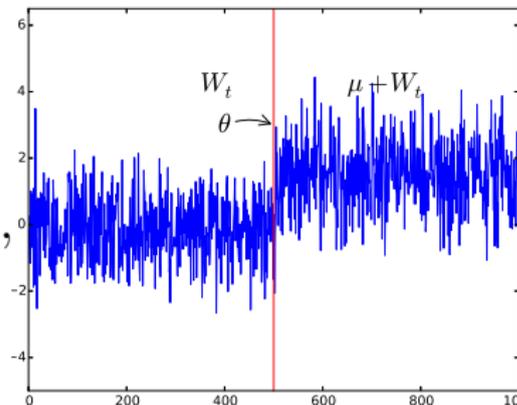
Постановка задачи обнаружения разладки (упрощенная)

- Структура наблюдаемого процесса X_t :

$$\underbrace{X_1, X_2, \dots, X_\theta}_{\sim f_\infty(\cdot)}, \underbrace{X_{\theta+1}, \dots}_{\sim f_0(\cdot)}$$

Пример:

$$X_t = \begin{cases} W_t, & 0 \leq t < \theta, \\ \mu + W_t, & \theta \leq t. \end{cases}$$

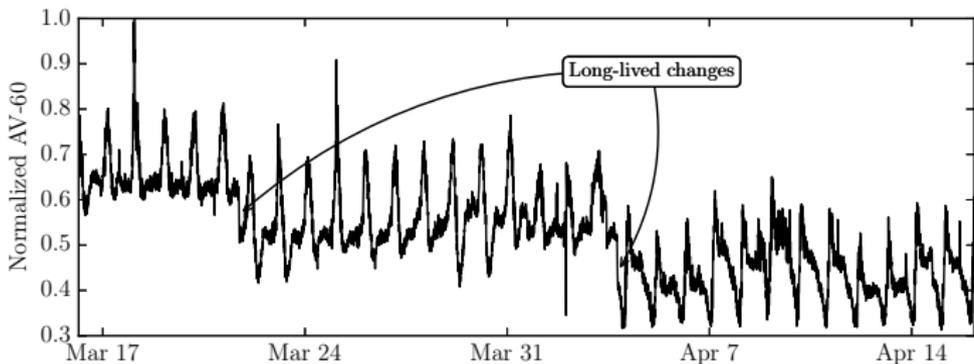
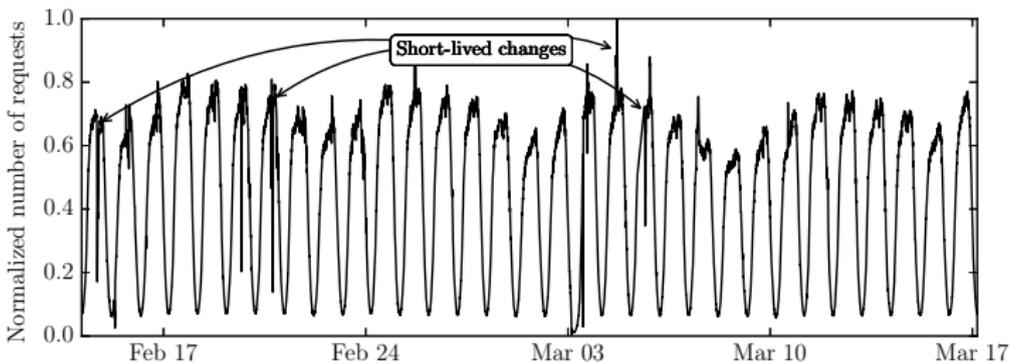


Примеры приложений

- Контроль качества производимого продукта
- Обнаружение целей при радиолокации, эхолокации и видеонаблюдении
- Выявление отказов (сбоев, разладок, утечек) оборудования и программного обеспечения
- Обнаружение злоумышленных действий (внедрений, атак, аномального поведения)
- Финансовый мониторинг (изменение волатильности, возникновение рецессий)
- Мониторинг в медицине (аритмия, ишемия, смертность, эпидемии и др.)
- Ретроспективный анализ временных данных

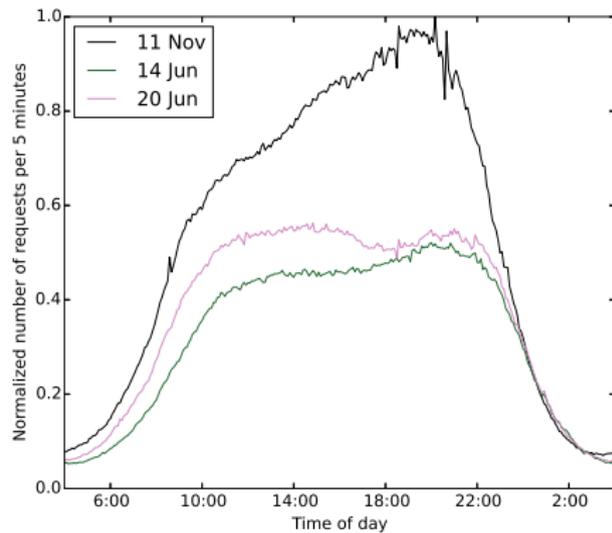
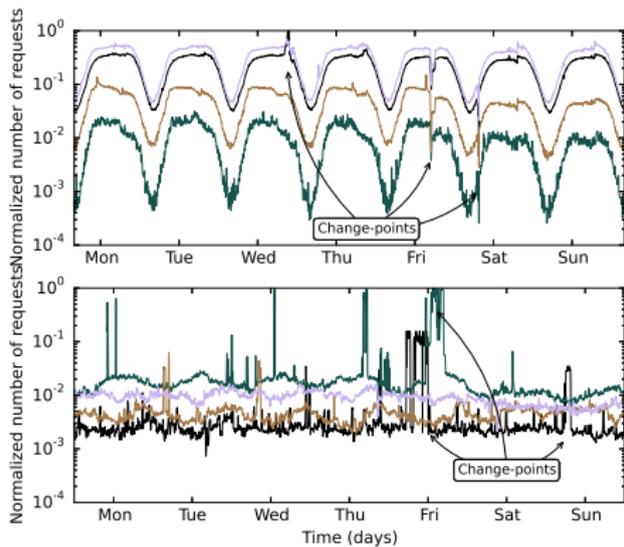
Пример реальных данных

Пример: нагрузка интернет-сервиса.



Пример реальных данных

Пример: нагрузка интернет-сервиса.



Большие информационные системы

- Стохастические циклы нагрузки (дневной, недельный, годовой)
- Всплески нагрузки («длинная память»)
- Априори произвольные типы отказов
- Огромные объемы данных (сотни тысяч или миллионы характеристик)

Большие информационные системы

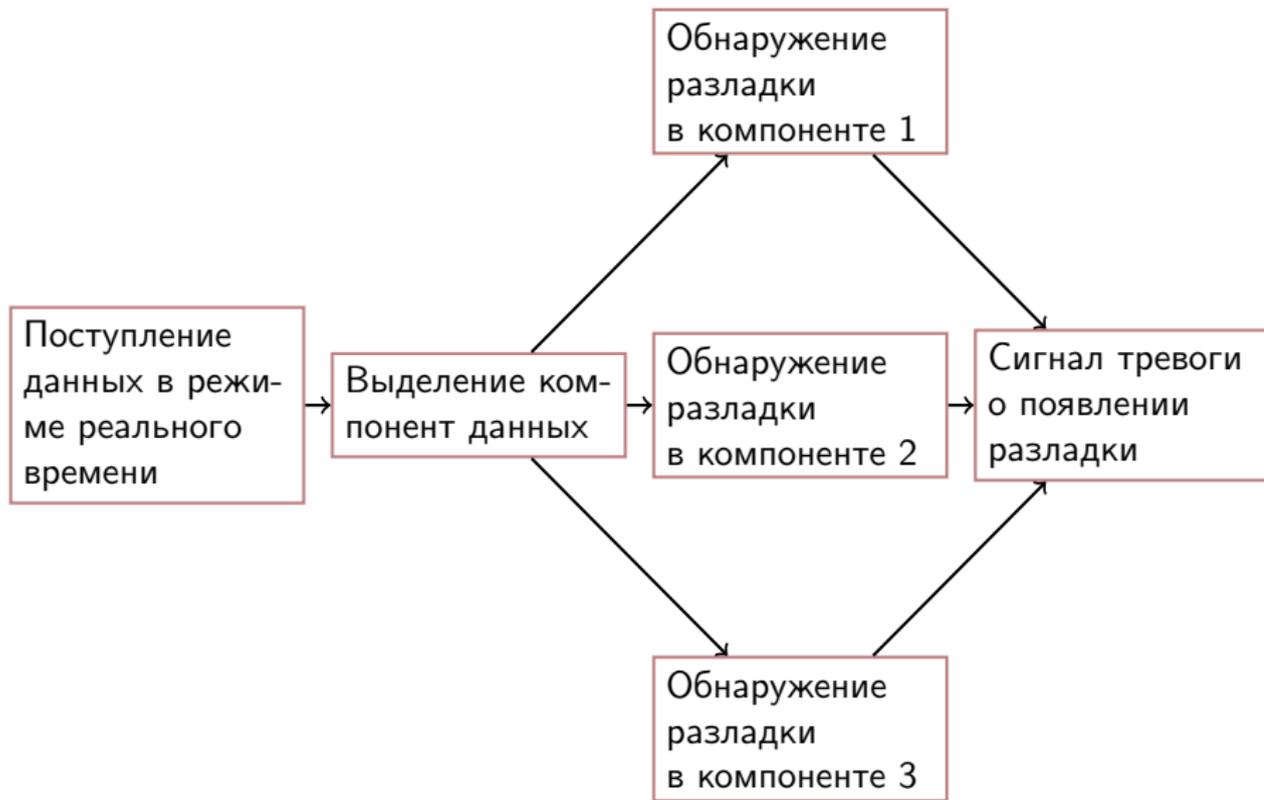
- Стохастические циклы нагрузки (дневной, недельный, годовой)
- Всплески нагрузки («длинная память»)
- Априори произвольные типы отказов
- Огромные объемы данных (сотни тысяч или миллионы характеристик)

Игнорирование этих особенностей →

непрерывный поток ложных тревог!

Цель работы: развитие **эффективной методологии** обнаружения разладок в указанных условиях

Автоматизированная компонента обнаружения разладки



Автоматизированная компонента обнаружения разладки

Ключевая проблема: формализация понятий нормального (целевого) и аномального состояний

Автоматизированная компонента обнаружения разладки

Ключевая проблема: формализация понятий нормального (целевого) и аномального состояний

Наш подход: моделирование тренда нагрузки

$$\xi_t = f(t) + \nu_t,$$

где $f(t)$ — гладкий тренд, ν_t — случайный шум.

Автоматизированная компонента обнаружения разладки

Ключевая проблема: формализация понятий нормального (целевого) и аномального состояний

Наш подход: моделирование тренда нагрузки

$$\xi_t = f(t) + \nu_t,$$

где $f(t)$ — гладкий тренд, ν_t — случайный шум.

Две частные модели:

1. Модель с явным учетом сезонности.
2. Модель с явным учетом длинной памяти.

Модель с явным учетом сезонности

$$\xi_t = q_t s_{\varphi(t)} + \sigma_{\varphi(t)} \varepsilon_t, \quad \text{где:}$$

- q_t : случайная амплитуда
- $\varphi(t) = 2\pi\{t/T\}$: фаза, T : период цикла
- $s_{\varphi(t)}$: неслучайный профиль цикла
- $\nu_t = \sigma_{\varphi(t)} \varepsilon_t$: случайный шум
- $\sigma_{\varphi(t)}$: неслучайная дисперсия шума
- ε_t : стандартный процесс белого шума

Модель с явным учетом сезонности

$$\xi_t = q_t s_{\varphi(t)} + \sigma_{\varphi(t)} \varepsilon_t, \quad \text{где:}$$

- q_t : случайная амплитуда
- $\varphi(t) = 2\pi\{t/T\}$: фаза, T : период цикла
- $s_{\varphi(t)}$: неслучайный профиль цикла
- $\nu_t = \sigma_{\varphi(t)} \varepsilon_t$: случайный шум
- $\sigma_{\varphi(t)}$: неслучайная дисперсия шума
- ε_t : стандартный процесс белого шума

- $X^\ell = (X_k, t_k)_{k=1}^\ell$: данные наблюдений
- $\varphi_k = \varphi(t_k)$: фаза в точке $t_k, k = 1, \dots, \ell$
- **Цель:** оценить $q_t, s_\psi, \sigma_\psi^2$
для каждого значения фазы $\psi \in [0, 2\pi]$

Оценка параметров квазипериодического сигнала

Шаг 0. Положить $\hat{q}_{t_k} = 1, \hat{\sigma}_{\psi_j}^2 = \text{var}(X_1, \dots, X_l)$.

Оценка параметров квазипериодического сигнала

Шаг 0. Положить $\hat{q}_{t_k} = 1, \hat{\sigma}_{\psi_j}^2 = \text{var}(X_1, \dots, X_l)$.

Шаг 1. Обновить $\hat{s}_{\psi_j} = \frac{\sum_{k=1}^l w_k X_k / \hat{q}_k K_h(\varphi_k, \psi_j)}{\sum_{k=1}^l w_k K_h(\varphi_k, \psi_j)}$.

Оценка параметров квазипериодического сигнала

Шаг 0. Положить $\hat{q}_{t_k} = 1, \hat{\sigma}_{\psi_j}^2 = \text{var}(X_1, \dots, X_l)$.

Шаг 1. Обновить $\hat{s}_{\psi_j} = \frac{\sum_{k=1}^l w_k X_k / \hat{q}_k K_h(\varphi_k, \psi_j)}{\sum_{k=1}^l w_k K_h(\varphi_k, \psi_j)}$.

Шаг 2. Обновить $\hat{\sigma}_{\psi_j}^2 = \frac{\sum_{k=1}^l (X_k - \hat{X}_k)^2 K_h(\varphi_k, \psi_j)}{\sum_{k=1}^l K_h(\varphi_k, \psi_j)}$.

Оценка параметров квазипериодического сигнала

Шаг 0. Положить $\hat{q}_{t_k} = 1, \hat{\sigma}_{\psi_j}^2 = \text{var}(X_1, \dots, X_l)$.

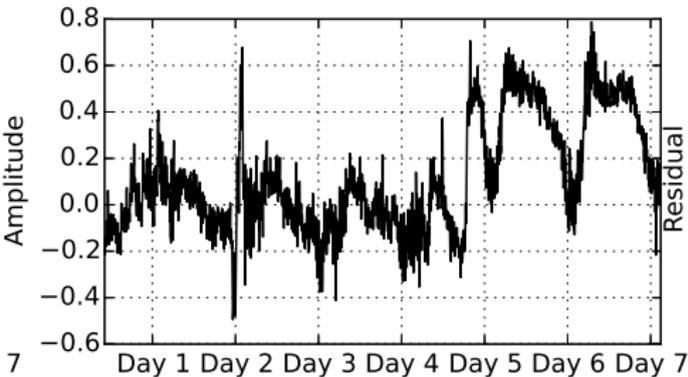
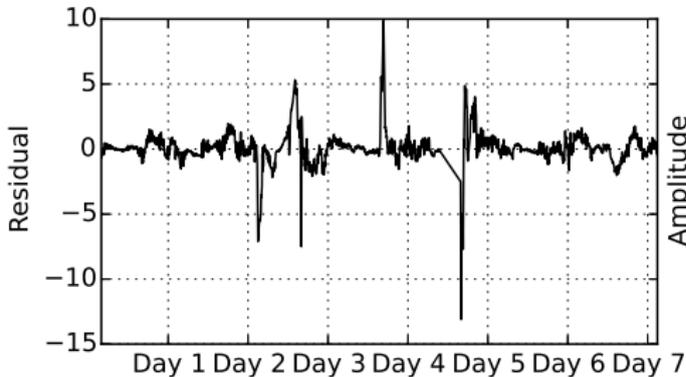
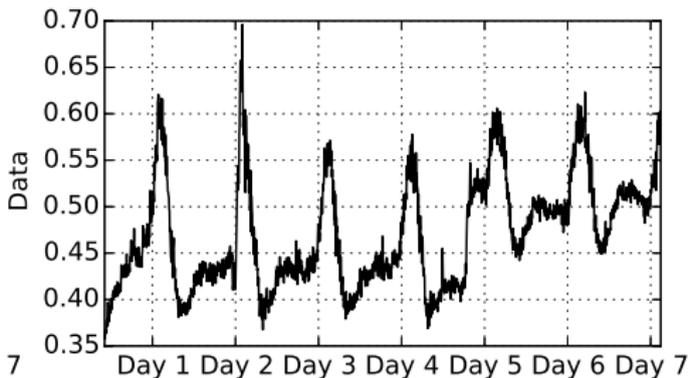
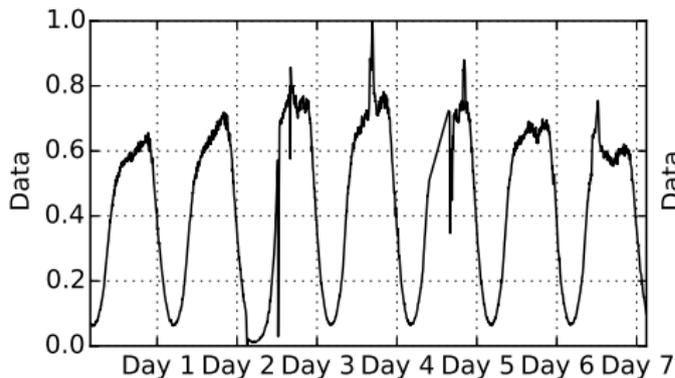
Шаг 1. Обновить $\hat{s}_{\psi_j} = \frac{\sum_{k=1}^l w_k X_k / \hat{q}_k K_h(\varphi_k, \psi_j)}{\sum_{k=1}^l w_k K_h(\varphi_k, \psi_j)}$.

Шаг 2. Обновить $\hat{\sigma}_{\psi_j}^2 = \frac{\sum_{k=1}^l (X_k - \hat{X}_k)^2 K_h(\varphi_k, \psi_j)}{\sum_{k=1}^l K_h(\varphi_k, \psi_j)}$.

Шаг 3. Прогноз $\hat{X}_n = \hat{q}_n \hat{s}_{\varphi(t_n)}$, где \hat{q}_n — оценка методом наименьших квадратов из

$$X_i = q_n \hat{s}_{\varphi(t_i)} + \nu_i, \quad \nu_i \sim \mathcal{N}(0, \hat{\sigma}_{\psi_i}^2)$$

Результаты фильтрации сигнала с циклами



Модель с явным учетом длинной памяти

$$\xi_t = \sum_{i=0}^n \theta_i t^i + \sigma(t) \varepsilon_t^H, \quad t \in [a, b]$$

- $\boldsymbol{\theta} = (\theta_0, \dots, \theta_n)$: неизвестные параметры
- $\sigma(t)$: неслучайная дисперсия шума
- ε_t^H : стандартный процесс фрактального шума

Модель с явным учетом длинной памяти

$$\xi_t = \sum_{i=0}^n \theta_i t^i + \sigma(t) \varepsilon_t^H, \quad t \in [a, b]$$

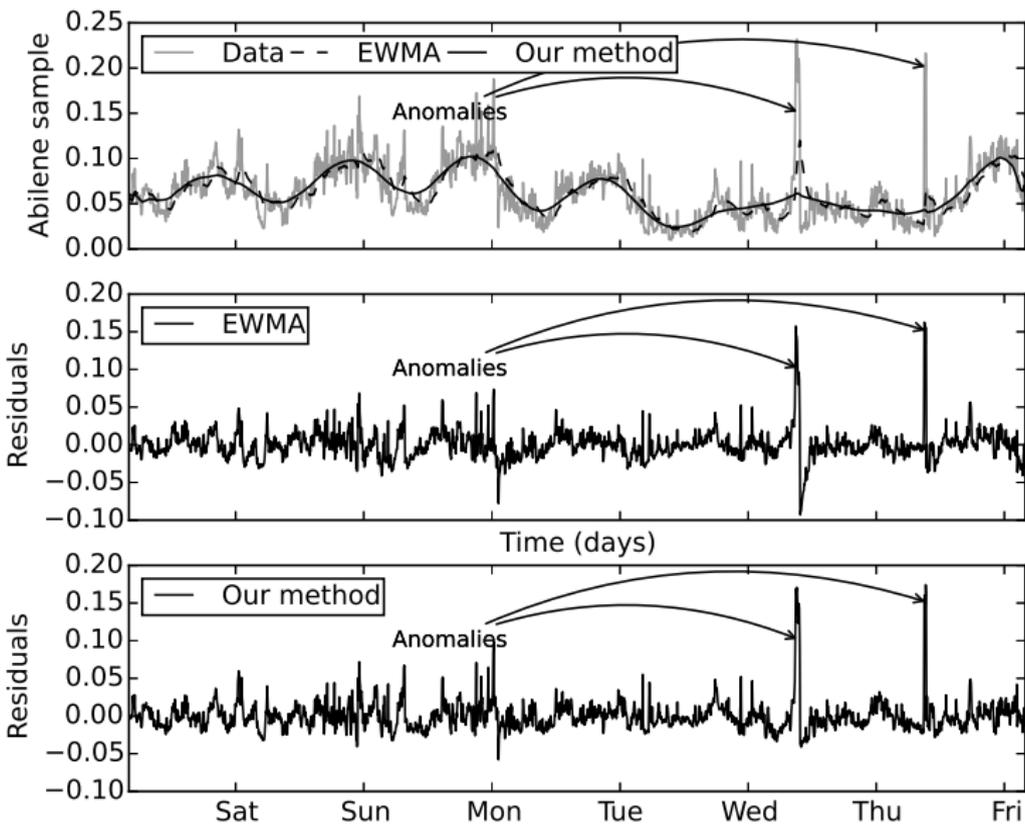
- $\theta = (\theta_0, \dots, \theta_n)$: неизвестные параметры
- $\sigma(t)$: неслучайная дисперсия шума
- ε_t^H : стандартный процесс фрактального шума

Теорема. Оценка максимума правдоподобия:

$$\hat{\theta}_{\text{ML}} = [\mathbf{R}_H(t)]^{-1} \psi_t^H, \quad \text{где}$$

$$(\psi_t^H)_i = \int_0^t \psi_i(s) dM_s^H, \quad (\mathbf{R}_H(t))_{ij} = \int_0^t \psi_i(s) \psi_j(s) dw_s^H$$

Результаты фильтрации сигнала с длинной памятью



Ансамбль процедур обнаружения разрядки

- Π_k : момент тревоги $\tau_k = \inf\{t \geq 0 : s_t^k \geq h_k\}$
- Π_k : «слабый» детектор (ограничительные предположения на модели сигнала и разрядки)

Ансамбль процедур обнаружения разладки

- Π_k : момент тревоги $\tau_k = \inf\{t \geq 0 : s_t^k \geq h_k\}$
- Π_k : «слабый» детектор (ограничительные предположения на модели сигнала и разладки)
- **Ансамбль**: момент $\tau_A = \inf\{t \geq 0 : a_t \geq h_A\}$

$$a_t = \psi(\boldsymbol{\lambda}; \mathbf{S}_t^1, \dots, \mathbf{S}_t^n)$$

- $\mathbf{S}_t^k = \{s_u^k, 0 \leq u \leq t\}$: траектория s^k
- $\boldsymbol{\lambda} \in \mathbb{R}^m$: параметры

Ансамбль процедур обнаружения разладки

- Π_k : момент тревоги $\tau_k = \inf\{t \geq 0 : s_t^k \geq h_k\}$
- Π_k : «слабый» детектор (ограничительные предположения на модели сигнала и разладки)
- **Ансамбль**: момент $\tau_A = \inf\{t \geq 0 : a_t \geq h_A\}$

$$a_t = \psi(\boldsymbol{\lambda}; \mathbf{S}_t^1, \dots, \mathbf{S}_t^n)$$

- $\mathbf{S}_t^k = \{s_u^k, 0 \leq u \leq t\}$: траектория s^k
- $\boldsymbol{\lambda} \in \mathbb{R}^m$: параметры
- **Пример**: логистическая регрессия

$$\psi(\boldsymbol{\lambda}; \mathbf{S}) = \sigma\left(\sum_{j=0}^p \sum_{k=1}^n \lambda_{kj} s_{t-j}^k - \lambda_0\right)$$

Обучение ансамбля «слабых» детекторов

- Средние потери, свойственные процедуре Π

$$F(\Pi) = c_{\infty} \mathbf{E}_{\infty} \left[\frac{\overbrace{\int_{\mathcal{T}_{\infty}} \mathbb{1}_{\{a_t \geq h\}}(t) dt}^{\text{длительность ложного сигнала тревоги}}}{\underbrace{\int_0^T \mathbb{1}_{\mathcal{T}_{\infty}}(t) dt}_{\text{длительность состояния без разрядки}}} \right] + c_0 \mathbf{E}_0 \left[\frac{\overbrace{\int_{\mathcal{T}_0} \mathbb{1}_{\{a_t < h\}}(t) dt}^{\text{длительность ложного молчания при разрядке}}}{\underbrace{\int_0^T \mathbb{1}_{\mathcal{T}_0}(t) dt}_{\text{длительность состояния с разрядкой}}} \right]$$

- Обучение ансамбля: оптимизация потерь

$$F(\Pi) \rightarrow \min_{\lambda \in \mathbb{R}^m}$$

по выборке $X^{\ell} = (X_t^i, Y_t^i)_{i=1}^{\ell}$

Эффективность ансамбля. Процесс fGn

Пример: изменение среднего значения процесса с длинной памятью

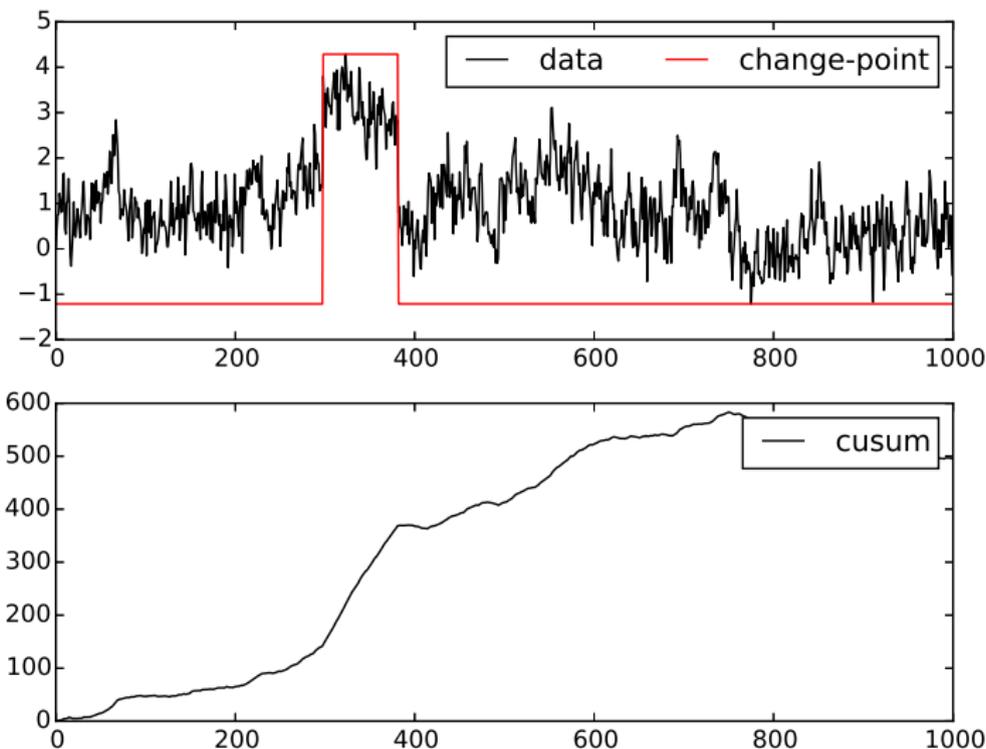
$$X_t = \begin{cases} W_t^H, & \text{если } t \notin [\theta, \theta + \Delta], \\ \mu + W_t^H, & \text{если } t \in [\theta, \theta + \Delta] \end{cases}$$

- $W^H = (W_t^H)_{0 \leq t \leq T}$: фрактальный гауссовский шум с показателем Хёрста $H = 0.95$
- величина разладки $\mu > 0$ неизвестна
- момент разладки $\theta \in [0, T]$ неизвестен

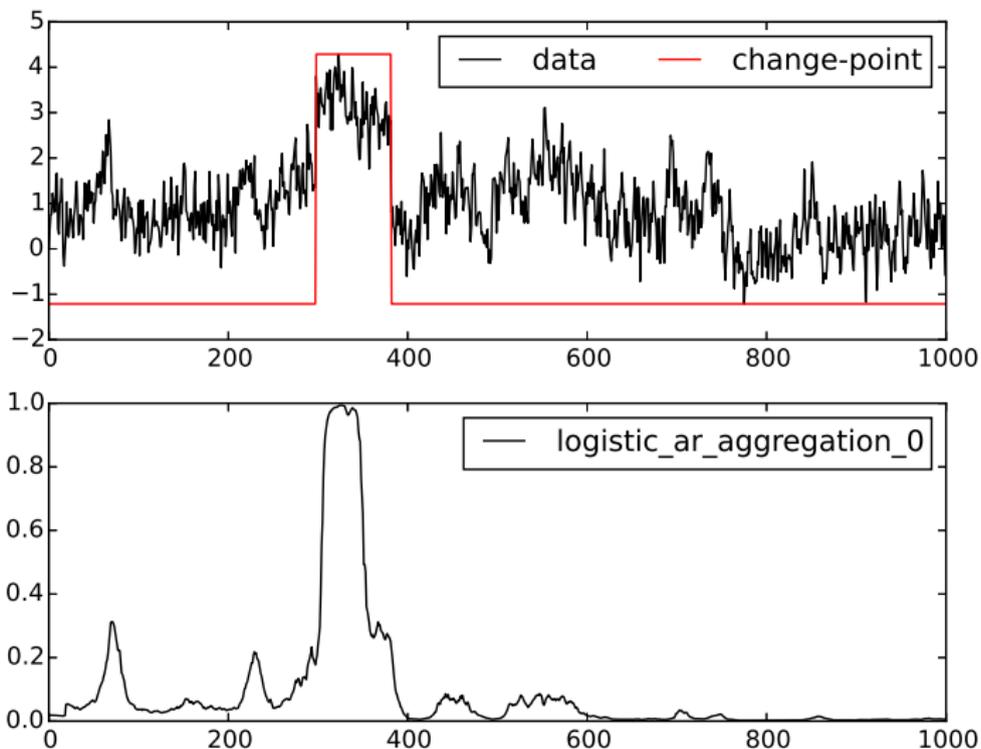
«Слабый» детектор предполагает модель

$$X_t = \mu \mathbb{1}_{\{t \geq \theta\}}(t) + W_t, \quad \mu = 1$$

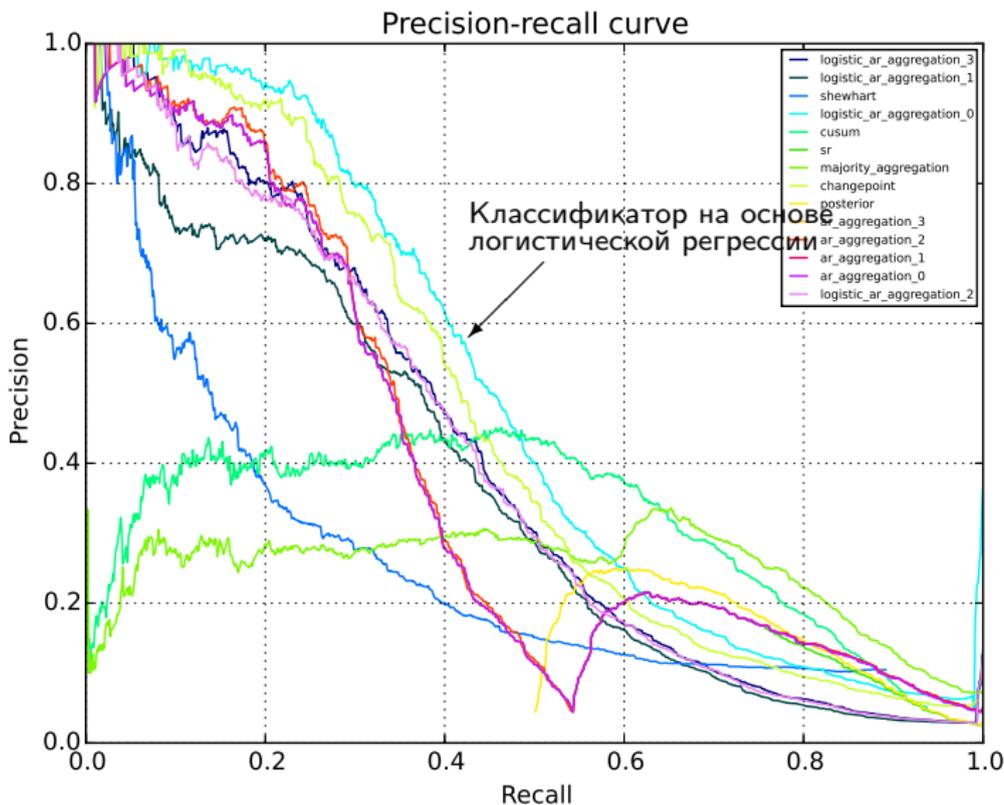
Эффективность ансамбля. Процесс fGn



Эффективность ансамбля. Процесс fGn



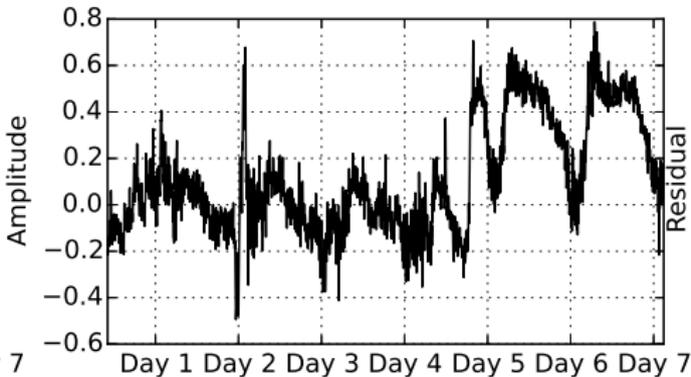
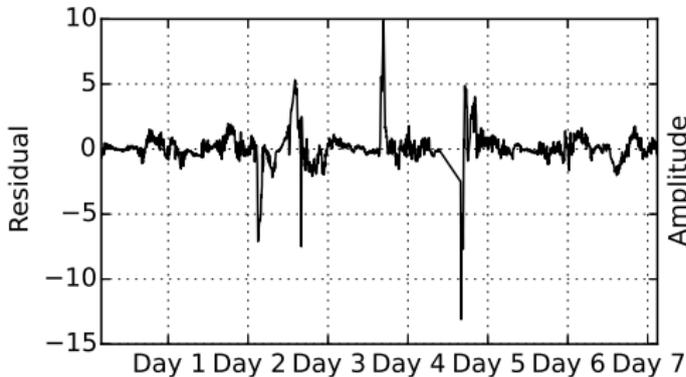
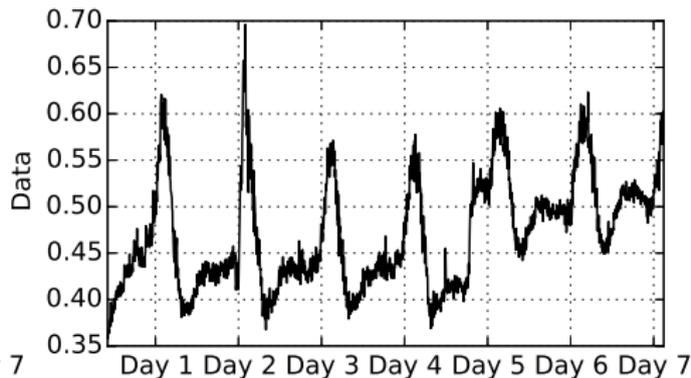
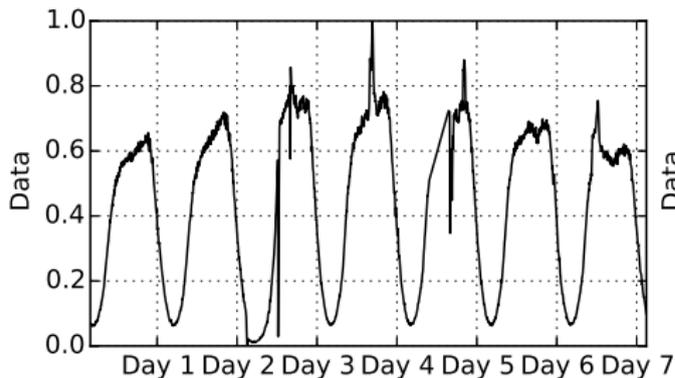
Эффективность ансамбля. Кривые качества



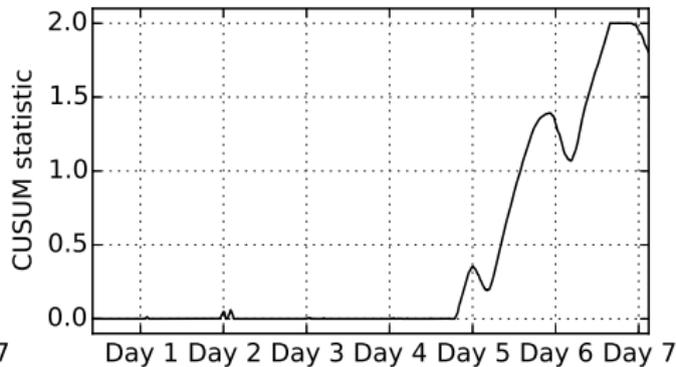
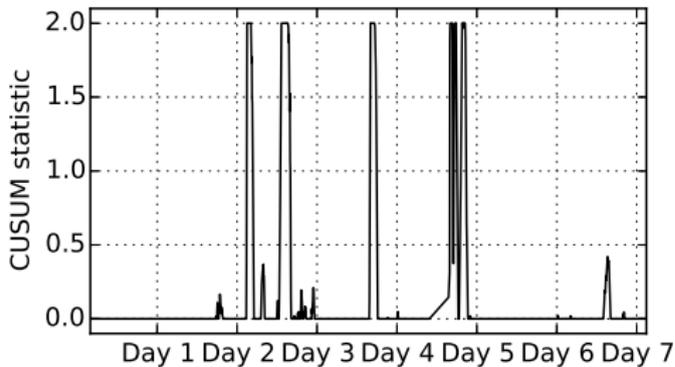
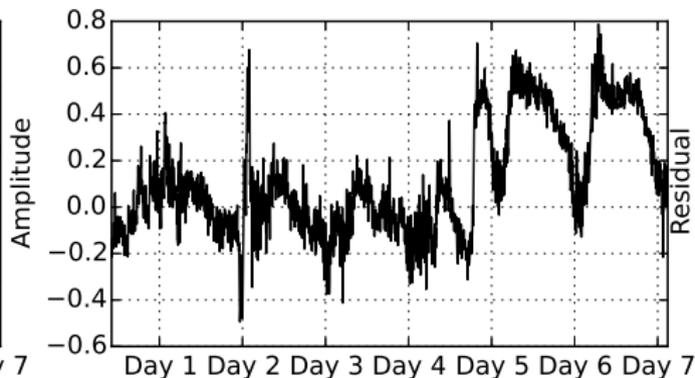
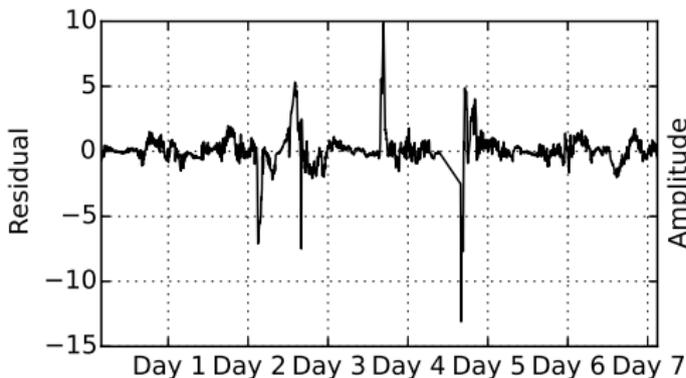
Детектирование отказов в информационных системах

- Обнаружение отказов большой информационной системы
- 100 000 (однотипных) одновременно наблюдаемых характеристик
- Временное разрешение для каждого временного ряда в системе — 5 мин, присутствует суточная цикличность (длина периода $T = 288$ точек/день)
- Обнаружение краткосрочных ($\Delta t \ll T$) и долгосрочных ($\Delta t \sim T$) разладок

Приложение: оценка параметров реальных сигналов



Приложение: обнаружение разладок реальных сигналов



Результаты работы

1. Разработаны новые математические методы оценки параметров сигнала по его измерениям во фрактальном шуме.
2. Разработаны математические модели и алгоритм оценки параметров квазипериодического сигнала.
3. Разработаны алгоритмы обнаружения разладки на основе ансамблей «слабых» детекторов.
4. Разработана внутренняя программная система обнаружения проблемного поведения сервисов компании Яндекс.

Спасибо за внимание!

csmlab.ru
yandexdatafactory.com
ru.linkedin.com/in/artonson
artemov@physics.msu.ru