

Data Analytics and Management in Data Intensive Domains DAMDID 2025



Тьюториал

Высокопроизводительный поиск аномалий во временных рядах

https://github.com/KraevaYA/Tutorial_DAMDID2025



Яна Александровна Краева, Михаил Леонидович Цымблер

{kraevaya, mzym}@susu.ru

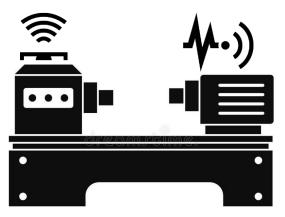
Южно-Уральский государственный университет, Челябинск, Россия

Содержание

- Временные ряды и аномалии в них
- Диссонансы
- Последовательный алгоритм DRAG
- Параллельный алгоритм PD3
- Последовательный алгоритм MERLIN
- Параллельный алгоритм PALMAD

© 2025 М.Л. Цымблер, Я.А. Краева

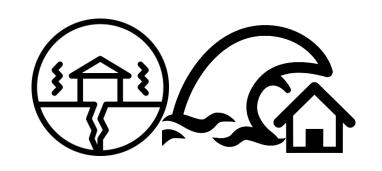
Люди измеряют всё во всех областях с течением времени



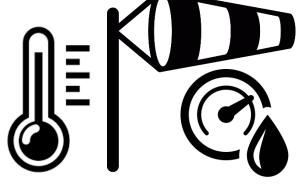
Умное производство, предиктивное ТО



Интернет вещей



Предсказание природных катаклизмов



Прогноз погоды, моделирование климата



Сельское хоз-во, животноводство



Борьба с преступностью



Био- и хемоинформатика



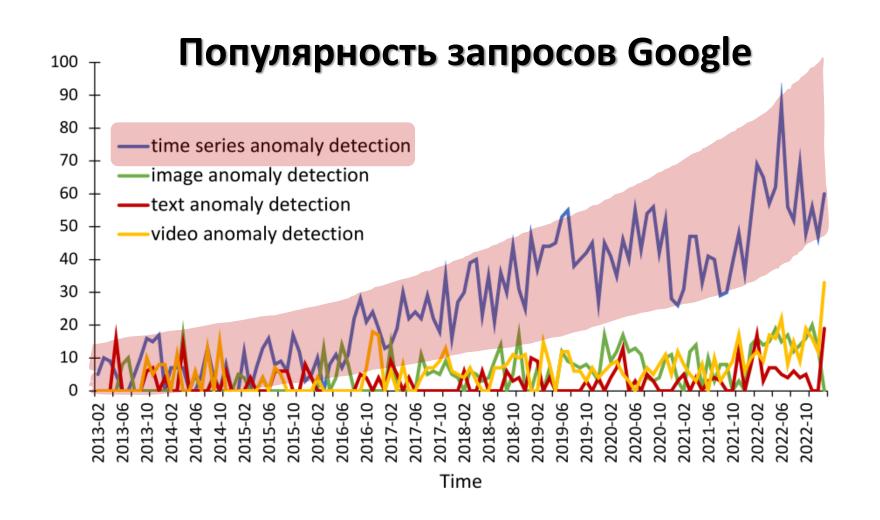
Экономика, бизнес, финансы



Системы электронного обучения

3/52

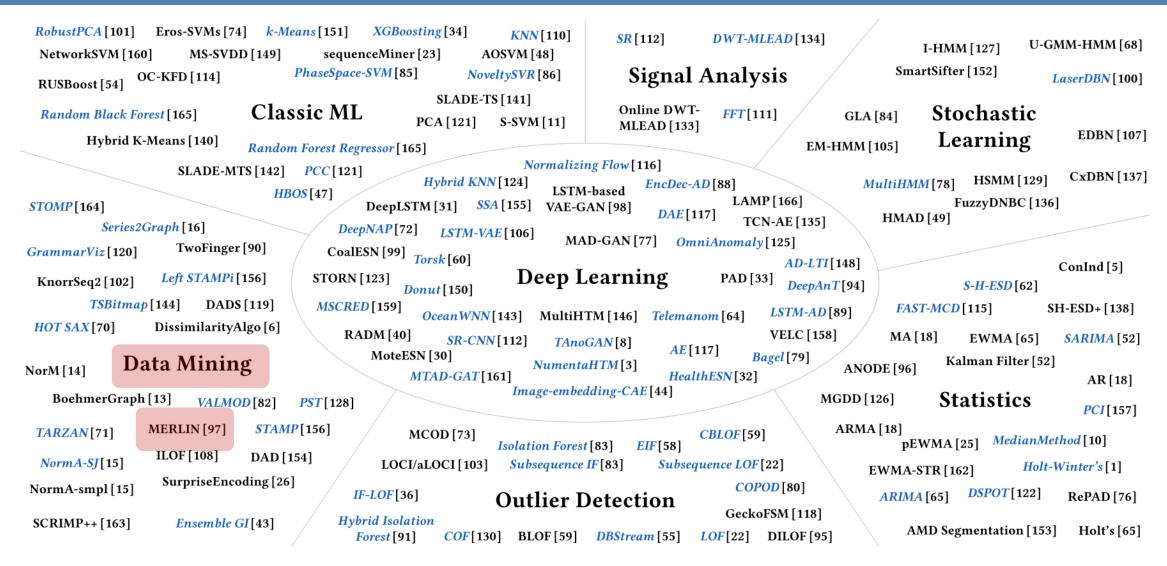
Поиск аномалий во временных рядах – «горячая» тема



^{*} Boniol P. et al. New trends in time-series anomaly detection. EDBT'2023. 847-850 (2023). DOI: 10.48786/edbt.2023.80

29.10.2025

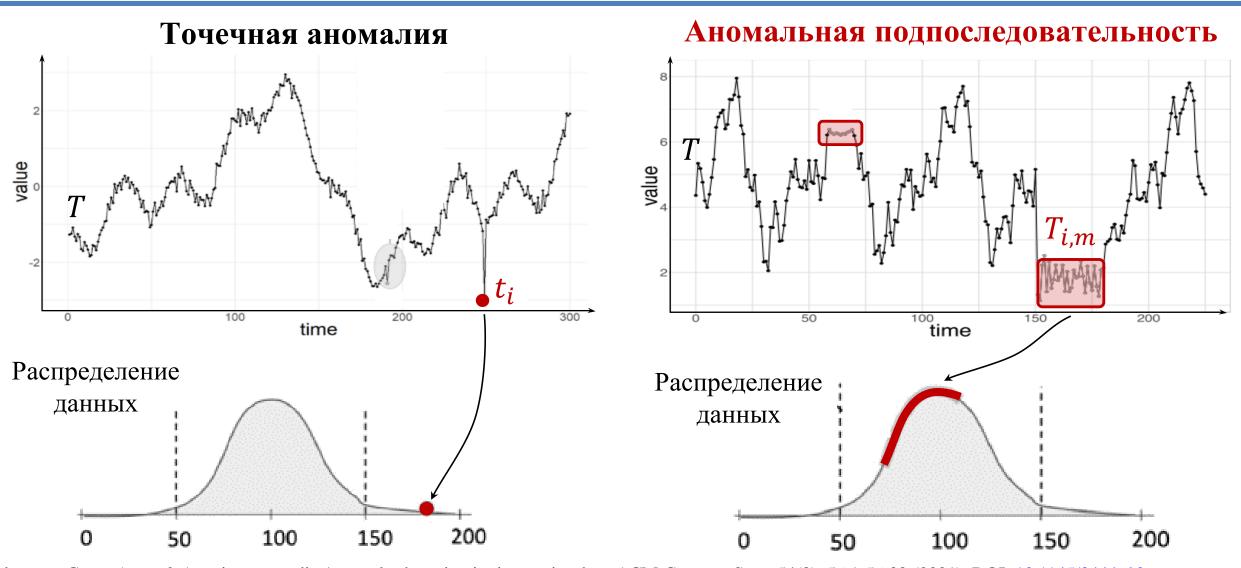
Поиск аномалий во временных рядах – «горячая» тема



^{*} Schmidl S. et al. Anomaly Detection in Time Series: A Comprehensive Evaluation. Proc. VLDB Endow. 15(9), 1779–1797 (2022). DOI: 10.14778/3538598.3538602

29.10.2025

Аномалии временных рядов



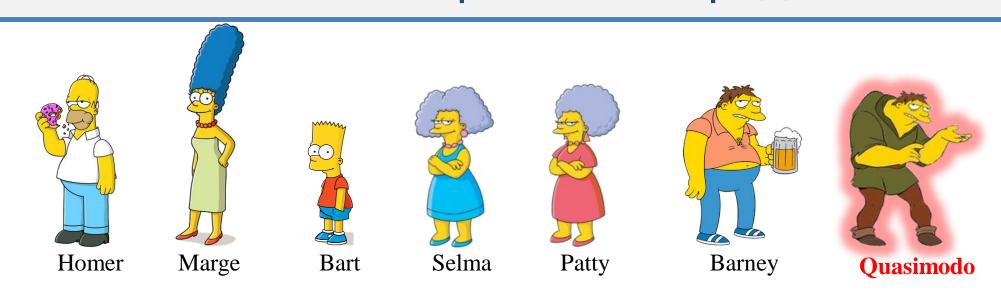
Blazquez-Garca A. et al. A review on outlier/anomaly detection in time series data. ACM Comput. Surv. 54(3), 56:1-56:33 (2021). DOI: 10.1145/3444690

Содержание

- Временные ряды и аномалии в них
- Диссонансы
- Последовательный алгоритм DRAG
- Параллельный алгоритм PD3
- Последовательный алгоритм MERLIN
- Параллельный алгоритм PALMAD

© 2025 М.Л. Цымблер, Я.А. Краева

Аномалия временного ряда



Аномалия – наблюдение, которое настолько сильно отличается от других наблюдений, что вызывает подозрения в том, что оно было создано иным механизмом.

Hawkins D.M. Identification of outliers. Monographs on applied probability and statistics. Springer, 1980. DOI: 10.1007/978-94-015-3994-4.



Диссонанс* – формализация аномалии

Диссонанс (discord) — подпоследовательность заданной длины с максимальным расстоянием до ближайшего соседа

Ближайший сосед — подпоследовательность, которая наиболее похожа на данную и имеет с ней не более половины общих точек



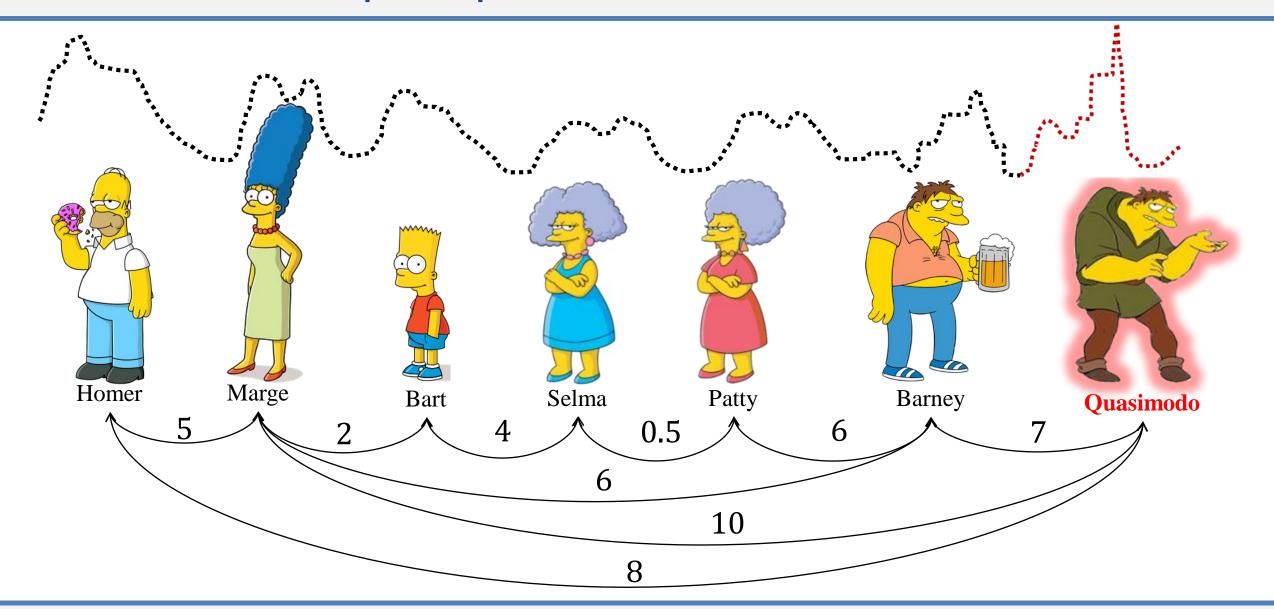
Агностическое и точное понятие:

независимость от предметной области и обучения, результат однозначен и воспроизводим

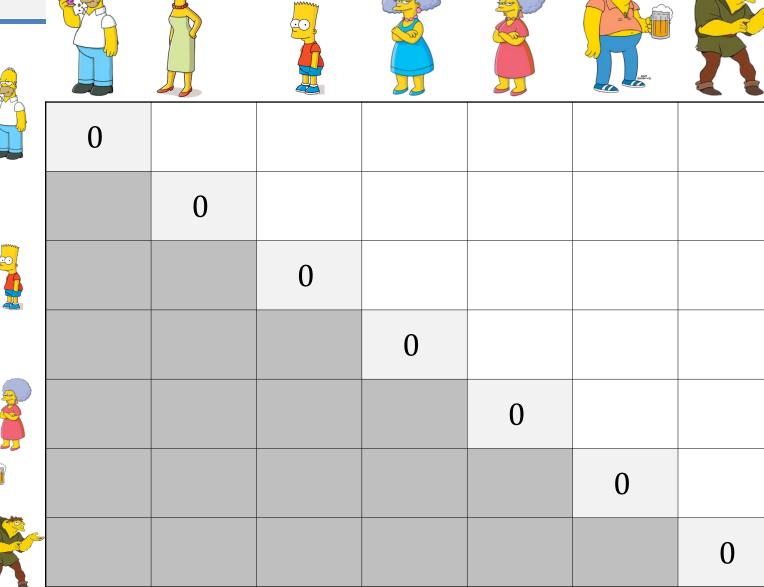
9/52

^{*} Keogh E. et al. HOT SAX: Efficiently finding the most unusual time series subsequence. ICDM 2005. 226-233 (2005). DOI: 10.1109/ICDM.2005.79

Пример поиска диссонанса



Матрица расстояний: чем ближе соседи, тем более они схожи



Selma

Patty

Barney

Marge Bart



Homer

Quasimodo

Матрица расстояний с вычисленными расстояниями



Homer

					2007.0	7
0	5	2	4	4	6	8
5	0	2.5	3	3	6	10
2	2.5	0	4	4	6	9
4	3	4	0	0.5	5	8
4	3	4	0.5	0	5	8
6	6	6	5	5	0	7
8	10	9	8	8	7	0

Selma

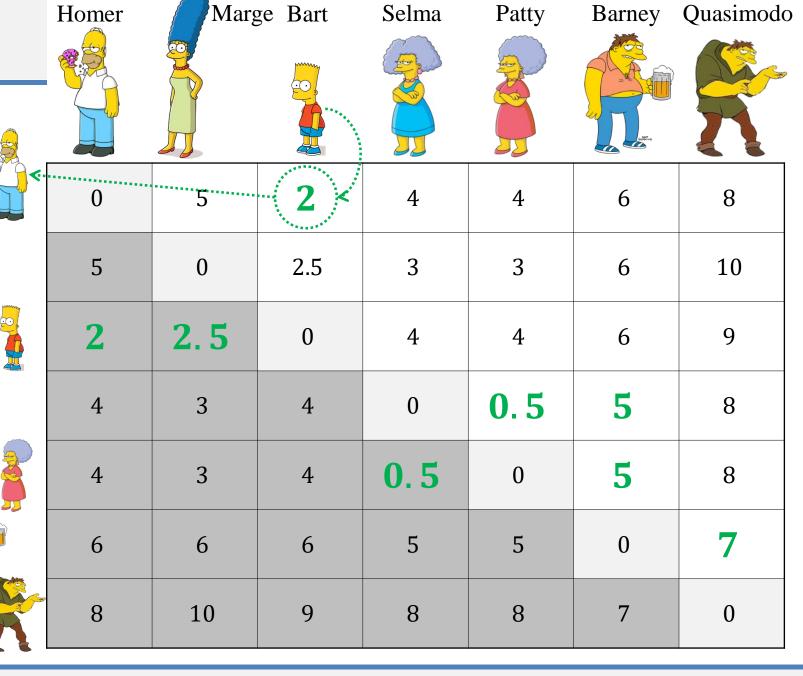
Marge Bart

Quasimodo

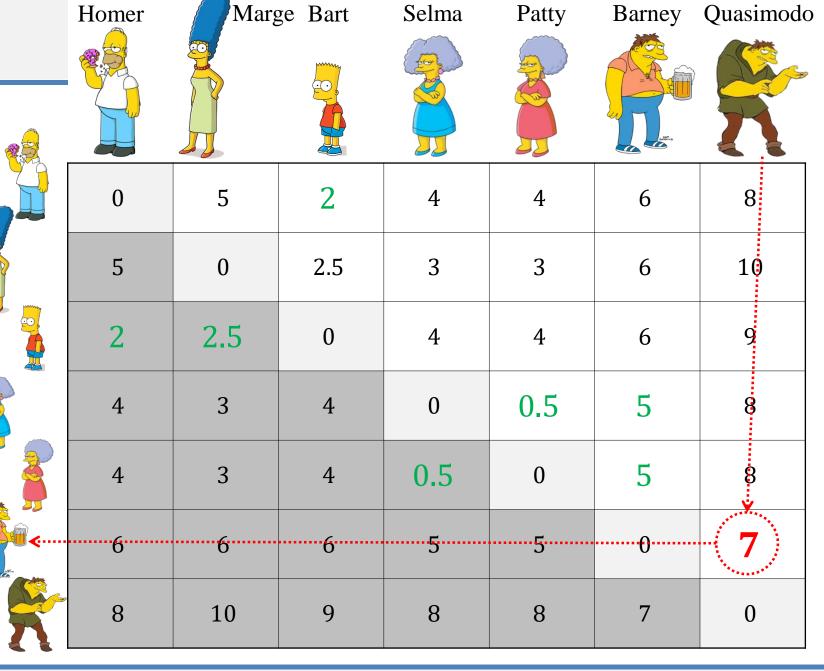
Barney

Patty

Матрица расстояний с расстояниями до ближайших соседей (т.е. минимумы по столбцам)



Матрица расстояний с наибольшим расстоянием до ближайшего соседа (т.е. максимум минимумов по столбцам)



Диссонанс — объект с самым далеким ближайшим соседом (т.е. аргумент максимума минимумов по столбцам)



Homer

					and the second	
0	5	2	4	4	6	8
5	0	2.5	3	3	6	10
2	2.5	0	4	4	6	9
4	3	4	0	0.5	5	8
4	3	4	0.5	0	5	8
6	6	6	5	5	0	7
8	10	9	8	8	7	0

Selma

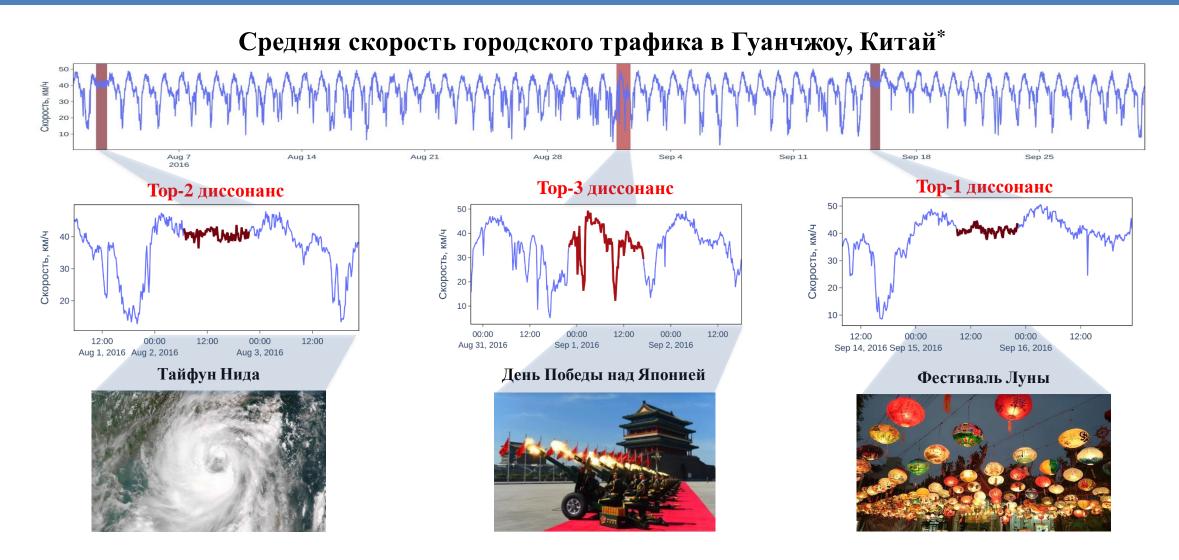
Patty

Barney

Quasimodo

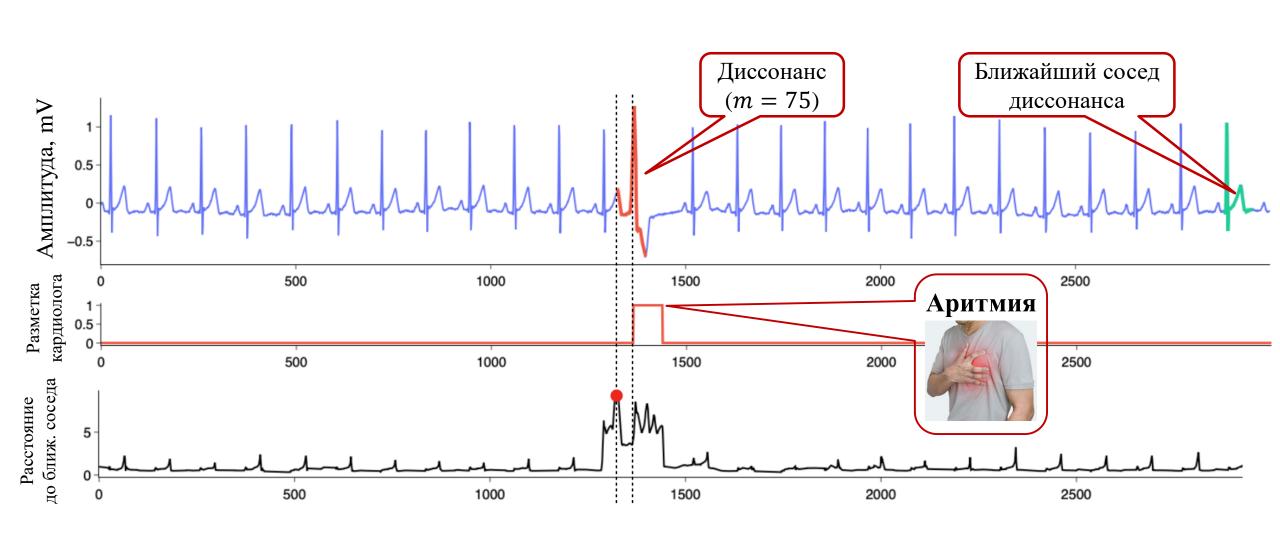
Marge Bart

Диссонанс отражает аномалию в реальной жизни, ...

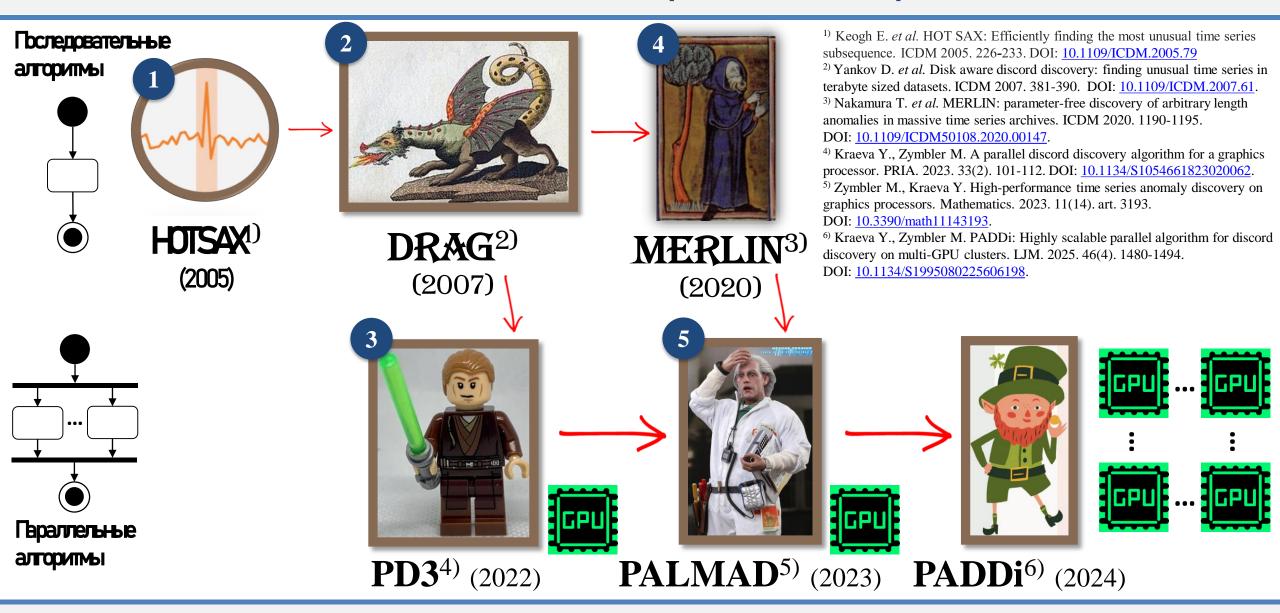


^{*} Chen X, Chen Y, He Z. Urban traffic speed dataset of Guangzhou, China. 2018. DOI: 10.5281/zenodo.1205229.

..., но диссонанс не идентичен аномалии



Диссонансы: дорожная карта



Диссонансы: Практикум

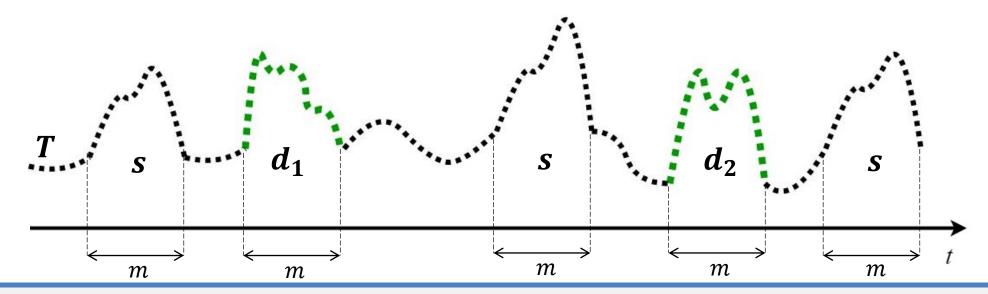
- 1. Проверьте настройки среды
- 2. Отобразите временной ряд
- 3. Идентифицируйте диссонансы

Содержание

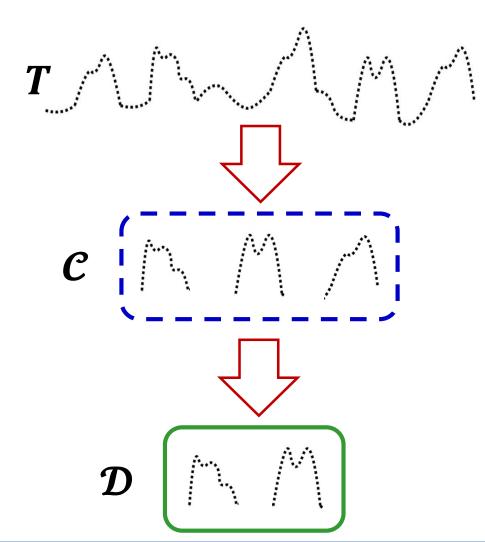
- Временные ряды и аномалии в них
- Диссонансы
- Последовательный алгоритм DRAG
- Параллельный алгоритм PD3
- Последовательный алгоритм MERLIN
- Параллельный алгоритм PALMAD

Диапазонный диссонанс (range discord)

- Диапазонный диссонанс подпоследовательность ряда, расстояние от которой до ее ближайшего соседа не ниже заданного порога
- Дано: ряд T, длина диссонанса m, порог r
- Найти: $\mathcal{D} = \{d_1, d_2, \dots\} \ d_i \in \mathcal{D} \iff \min_{s \in M_{d_i}} \mathrm{Dist}(d_i, s) \geq r$



Алгоритм DRAG (Discord Range Aware Gathering)



1. Отбор

За одно сканирование ряда сформировать множество кандидатов в диссонансы

2. Очистка

За одно сканирование ряда отбросить кандидатов, которые являются ложными диссонансами

DRAG: Отбор кандидатов

```
\mathcal{C} \coloneqq \{T_{1,m}\}
                                                                                                                                                         \mathcal{C} = \{\mathbf{v}\}
while not end of T
   get next subsequence s
   isCandidate := TRUE
   for each c_i \in \mathcal{C} and s \cap c_i = \emptyset
                                                                                                                                                  \operatorname{Dist}(\mathbf{w}, \mathbf{v}) \ge r
         if Dist(s, c_i) < r then
             \mathcal{C} \coloneqq \mathcal{C} \setminus \mathbf{c_i}; is Candidate \coloneqq FALSE
   if isCandidate = TRUE then C := C \cup S
                                                                                                                                                   C = \{v, w\}
                                                                                             X
```

DRAG: Отбор кандидатов

```
C = \{v, w\}
\mathcal{C} \coloneqq \{T_{1,m}\}
while not end of T
    get next subsequence s
    isCandidate := TRUE
                                                                                                                                                  \operatorname{Dist}(\mathbf{x}, \mathbf{v}) < r
    for each c_i \in \mathcal{C} and s \cap c_i = \emptyset
         if Dist(s, c_i) < r then
                                                                                                                                                   \operatorname{Dist}(\mathbf{x}, \mathbf{w}) \geq r
             \mathcal{C} \coloneqq \mathcal{C} \setminus \mathbf{c_i}; is Candidate \coloneqq FALSE
    if isCandidate = TRUE then C := C \cup S
                                                                                                                                                          C = \{w\}
```

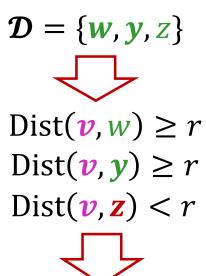
29.10.2025

DRAG: Отбор кандидатов

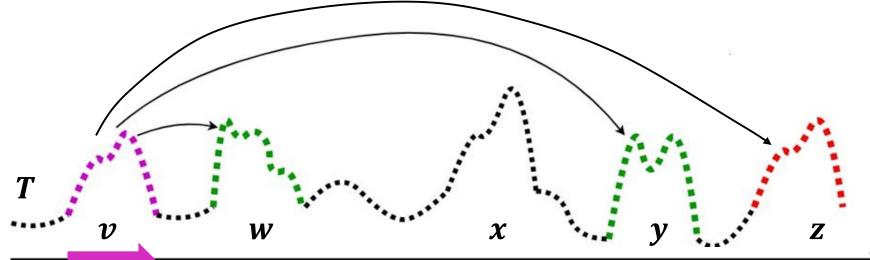
```
C = \{w, y\}
\mathcal{C} \coloneqq \{T_{1,m}\}
                                                                                                       Z – ложный
while not end of T
                                                                                                    диссонанс, т.к.
   get next subsequence s
                                                                                                     Dist(\mathbf{z}, \mathbf{v}) < r
   isCandidate := TRUE
                                                                                                                                             \operatorname{Dist}(\mathbf{z}, \mathbf{w}) \geq r
                                                                                                     Dist(\mathbf{z}, \mathbf{x}) < r.
   for each c_i \in \mathcal{C} and s \cap c_i = \emptyset
                                                                                                                                             \operatorname{Dist}(\mathbf{z},\mathbf{y}) \geq r
         if Dist(s, c_i) < r then
                                                                                                          Ho\boldsymbol{v}и\boldsymbol{x}
             \mathcal{C} \coloneqq \mathcal{C} \setminus \mathbf{c_i}; is Candidate \coloneqq FALSE
                                                                                                    были удалены!
   if isCandidate = TRUE then C := C \cup S
                                                                                                                                              C = \{w, y, z\}
                                                       W
```

DRAG: Очистка кандидатов

```
\mathcal{D} \coloneqq \mathcal{C}
while not end of T
get next subsequence s
for each d_i \in \mathcal{D} and s \cap d_i = \emptyset
if \mathrm{Dist}(s, d_i) < r then
\mathcal{D} \coloneqq \mathcal{D} \setminus d_i
```

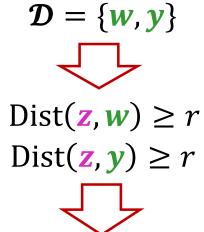


 $\mathcal{D} = \{w, y\}$



DRAG: Очистка кандидатов

```
\mathcal{D} \coloneqq \mathcal{C}
while not end of T
get next subsequence s
for each d_i \in \mathcal{D} and s \cap d_i = \emptyset
if \text{Dist}(s, d_i) < r then
\mathcal{D} \coloneqq \mathcal{D} \setminus d_i
```



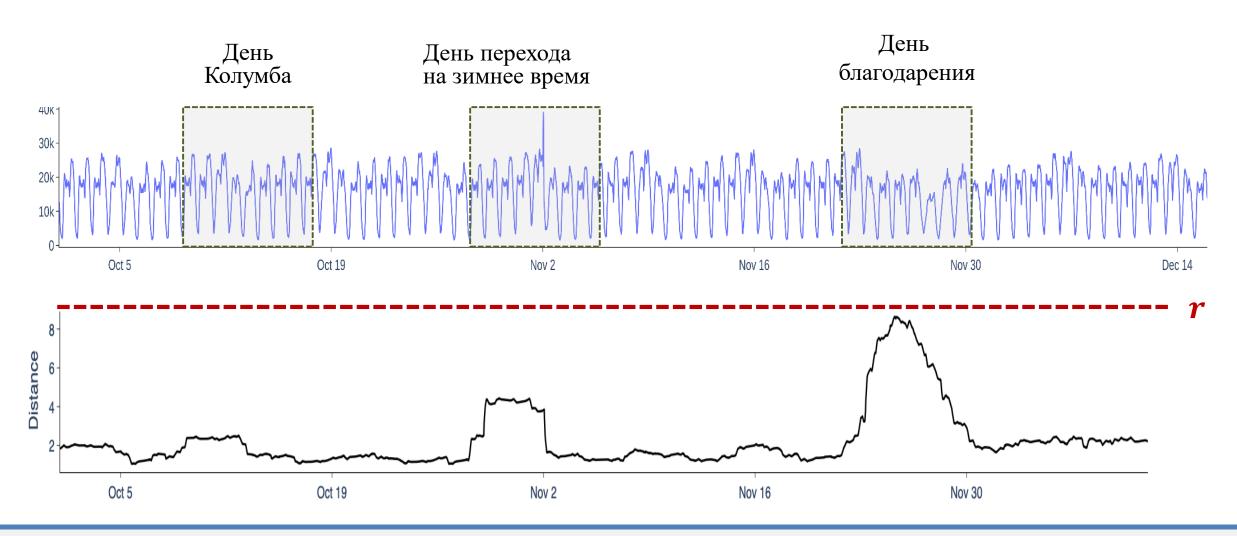
 $\mathcal{D} = \{w, y\}$

Что не так с DRAG: Ручной подбор порога r

Слишком большой порог – нет диссонансов, слишком маленький порог – много ложных диссонансов

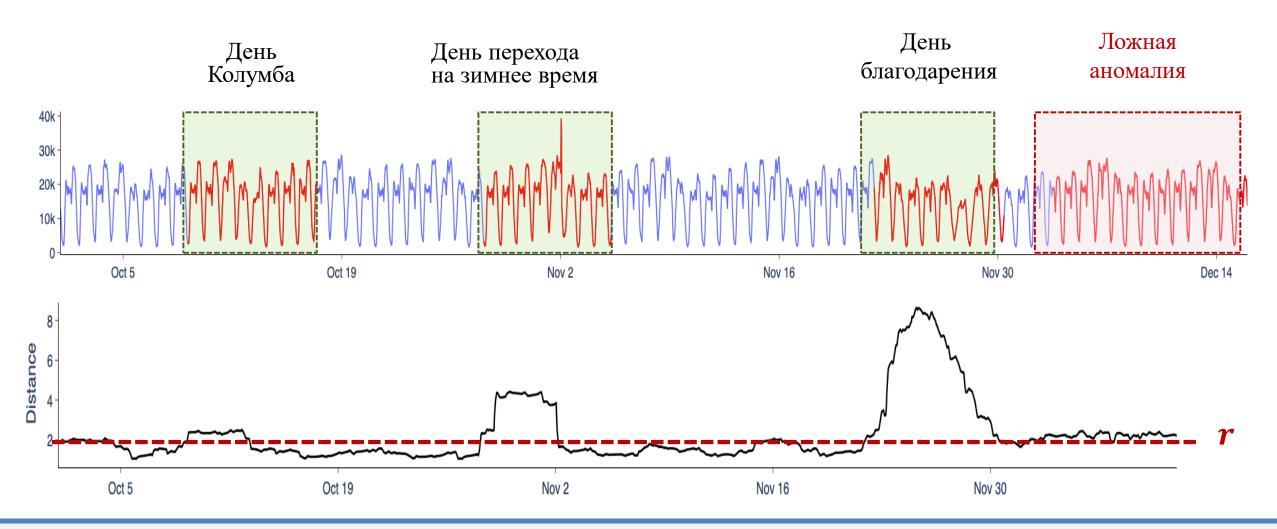
Ручной подбор порога: $r \to +\infty \Rightarrow$ нет диссонансов

Среднее число пассажиров NY такси осенью 2014 г.



Ручной подбор порога: $r \to 0 \implies$ ложные аномалии

Среднее число пассажиров NY такси осенью 2014 г.



Алгоритм DRAG: Практикум

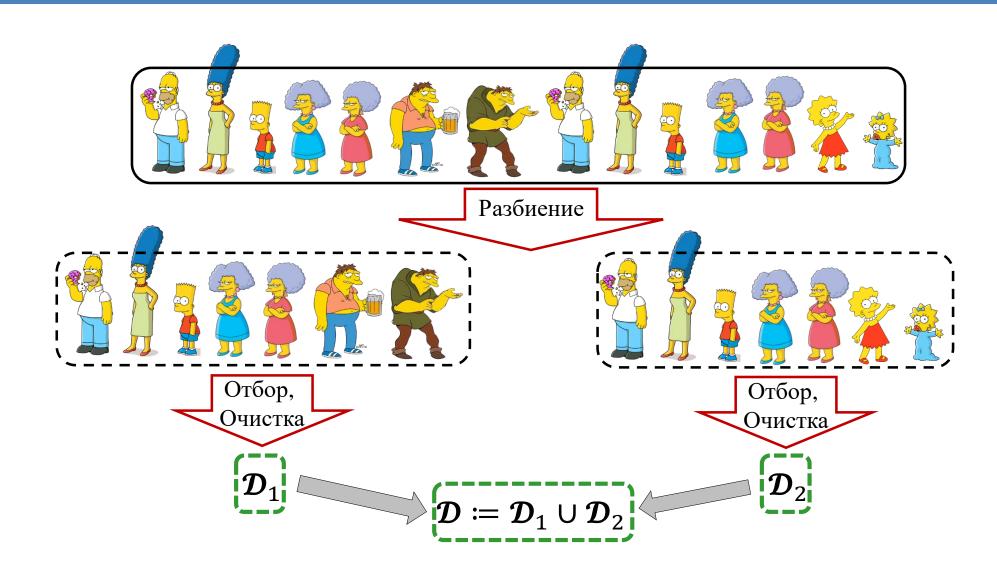
- 1. Найдите диссонансы в заданном временном ряде с помощью алгоритма DRAG
- 2. Отобразите временной ряд и диссонансы

31/52

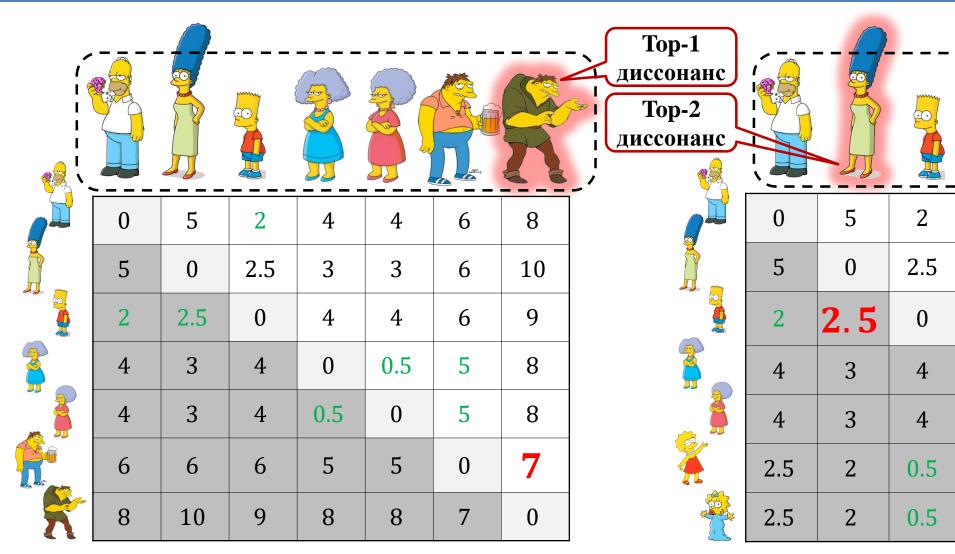
Содержание

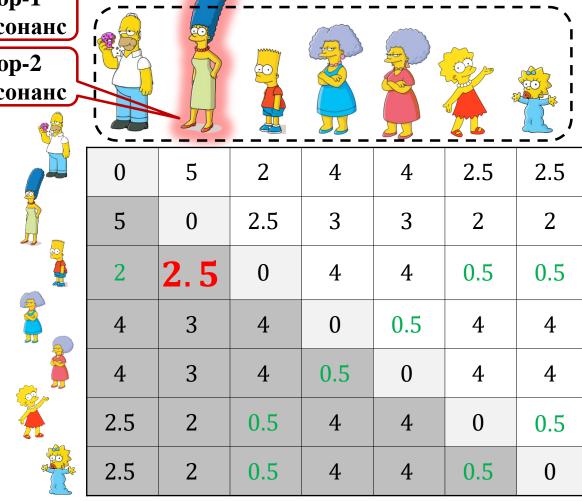
- Временные ряды и аномалии в них
- Диссонансы
- Последовательный алгоритм DRAG
- Параллельный алгоритм PD3
- Последовательный алгоритм MERLIN
- Параллельный алгоритм PALMAD

Наивное распараллеливание ...

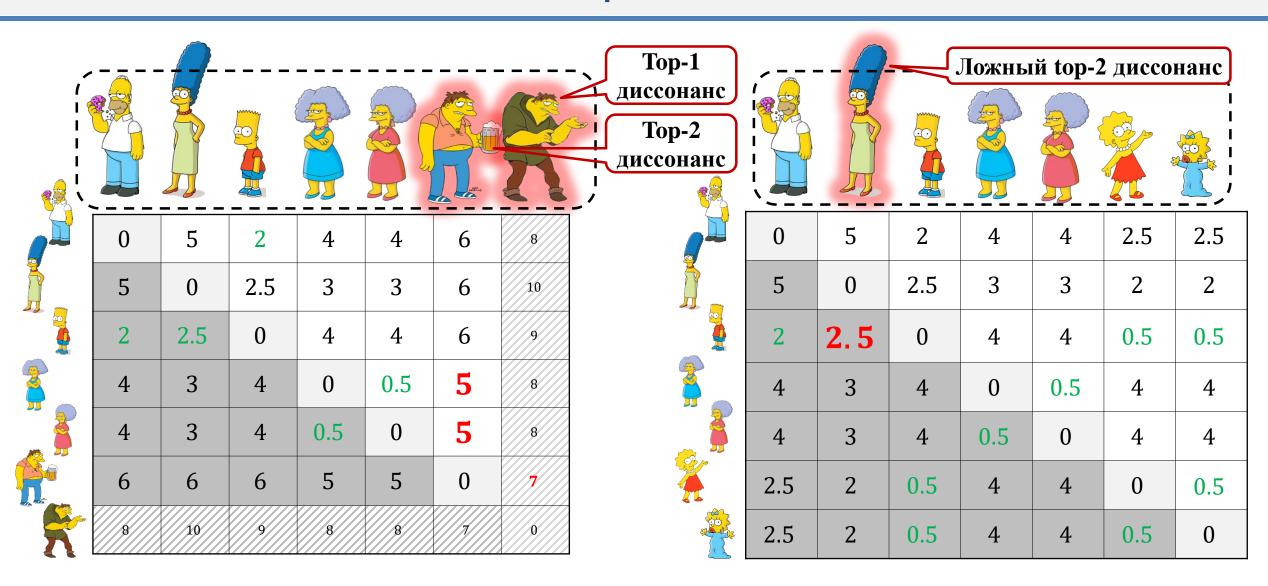


Наивное распараллеливание ...



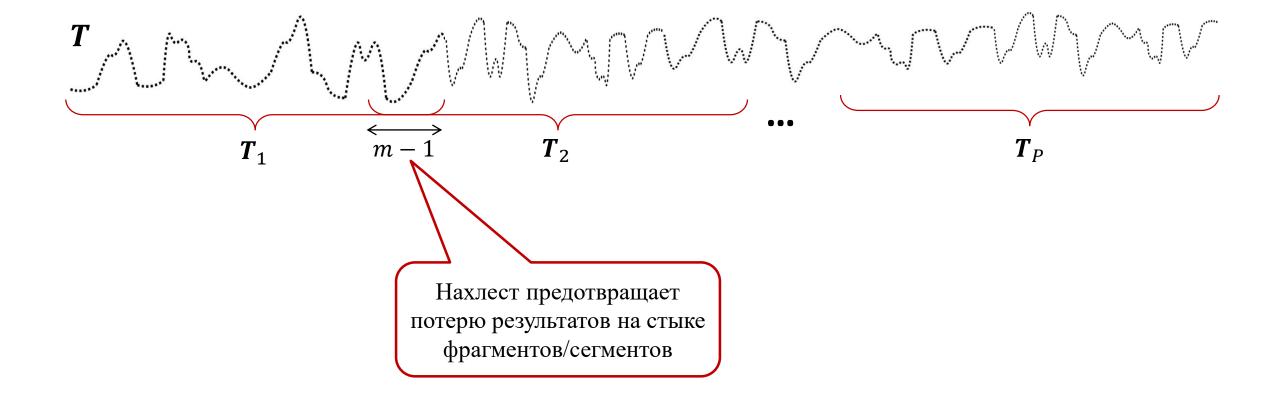


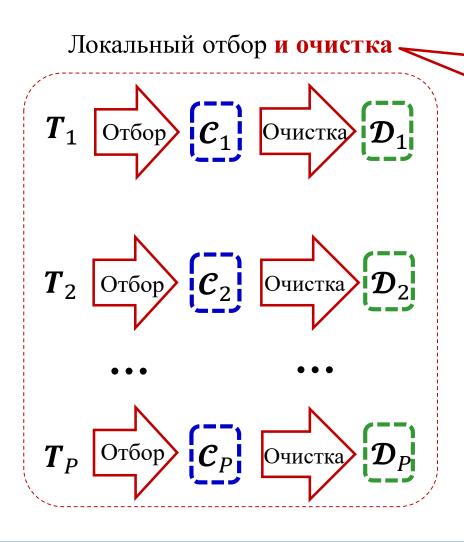
... не работает



29.10.2025

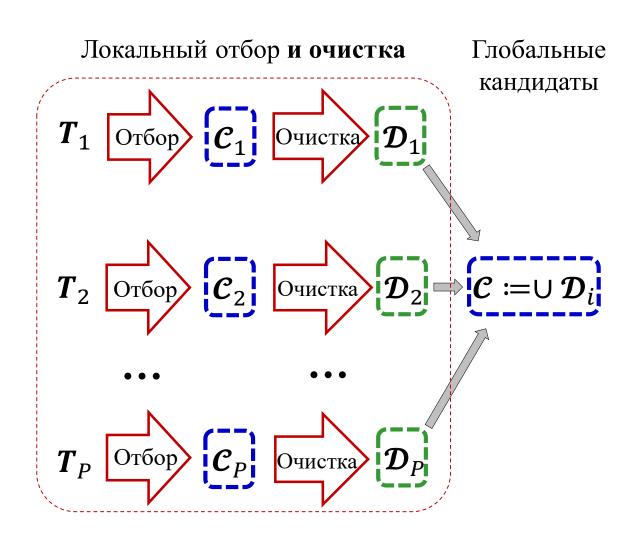
PD3 (Parallel DRAG-based Discord Discovery): Сегментация ряда



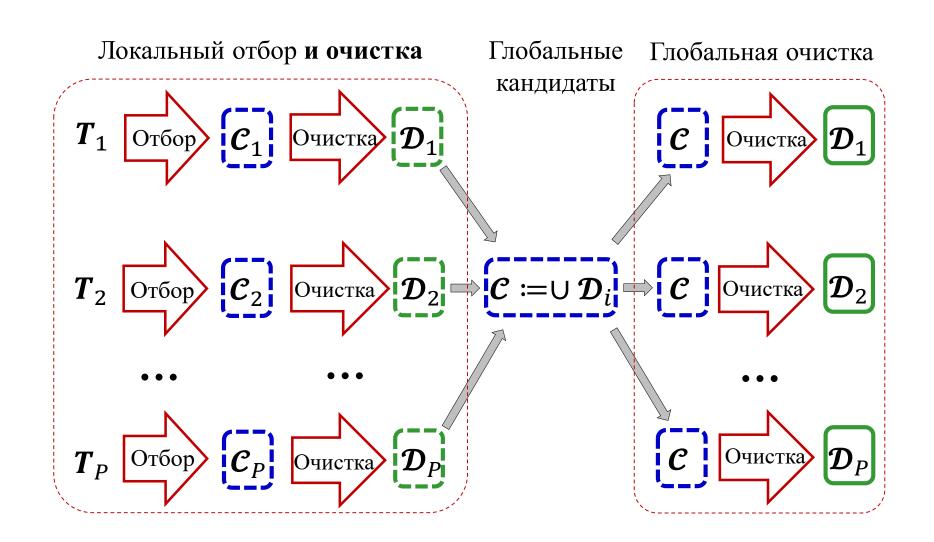


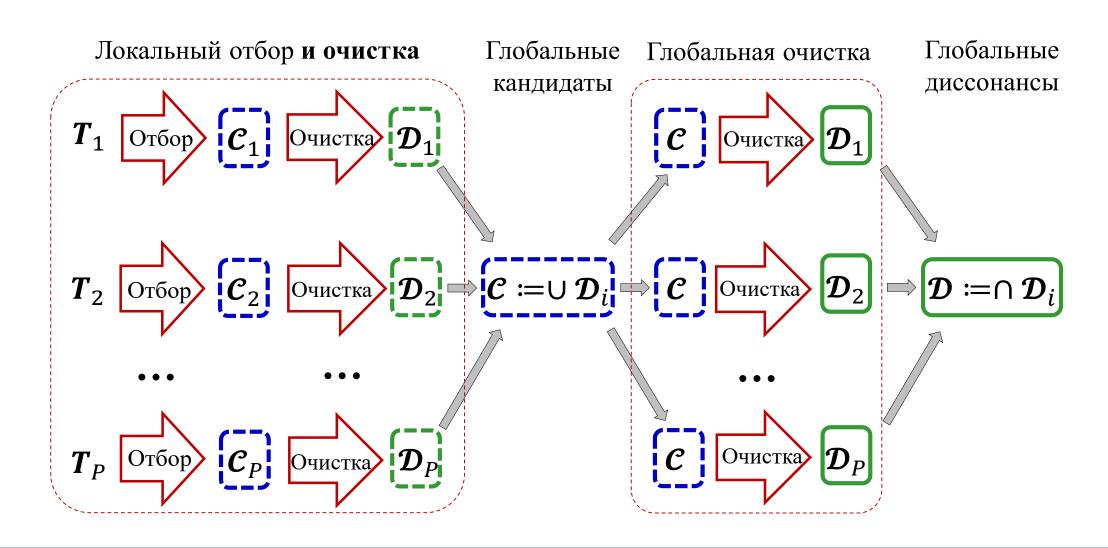
Локальный кандидат не может быть глобальным диссонансом, если он отброшен при очистке хотя бы одного фрагмента/сегмента

37/52



38/52





29.10.2025

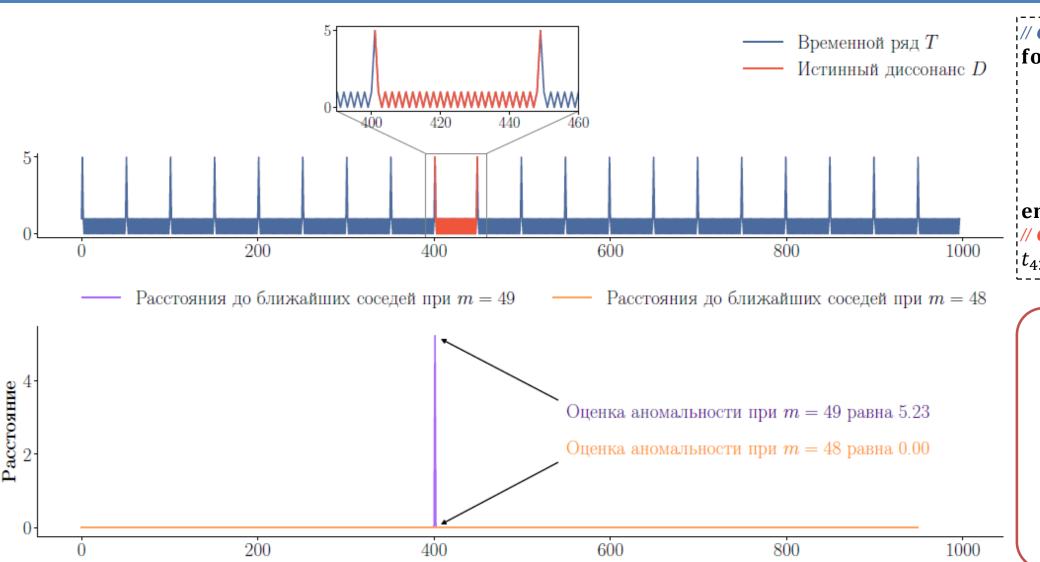
Алгоритм PD3: Практикум

- 1. Найдите диссонансы в заданном временном ряде с помощью алгоритма PD3
- 2. Отобразите временной ряд и найденные диссонансы
- 3. Сравните результаты и время работы алгоритмов DRAG и PD3

Что не так с PD3 (и DRAG): Ручной подбор m

Не всегда заранее известна длина аномалии. Запуск DRAG/PD3 для всех возможных длин вычислительно неосуществим

Чтобы найти все аномалии, нужно проверить все значения m



| Coздание временного ряда | for $i \leftarrow 1$ to 1000 do | if $i \mod 50 \neq 0$ then | $t_i \leftarrow i \mod 2$ else | $t_i \leftarrow 5$ end if | end for | Coздание диссонанса | $t_{430} \leftarrow 0$; $t_{431} \leftarrow 0$

Если $T_{i,m}$ — диссонанс, то не факт, что $T_{i,m\pm 1}$ — диссонанс

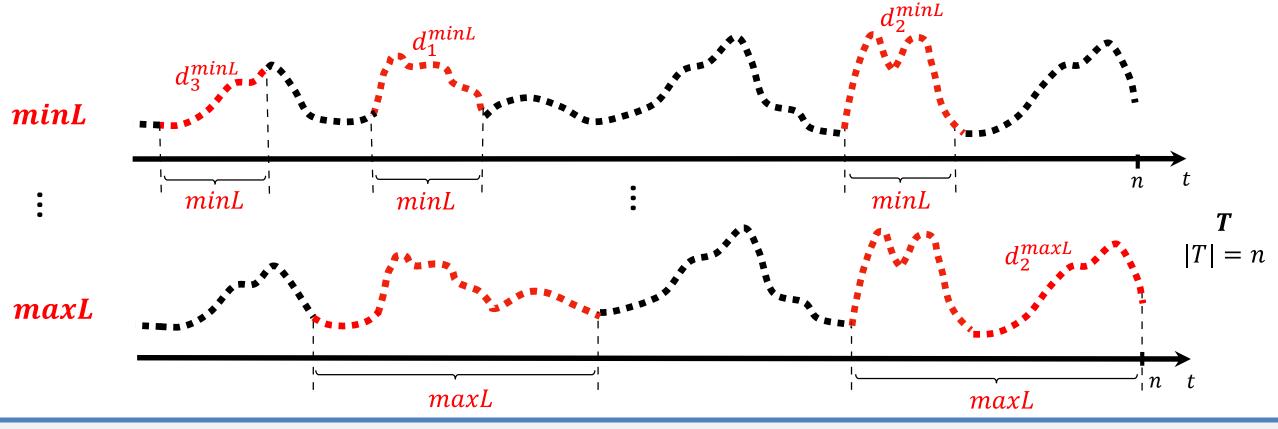
Содержание

- Временные ряды и аномалии в них
- Диссонансы
- Последовательный алгоритм DRAG
- Параллельный алгоритм PD3
- Последовательный алгоритм MERLIN
- Параллельный алгоритм PALMAD

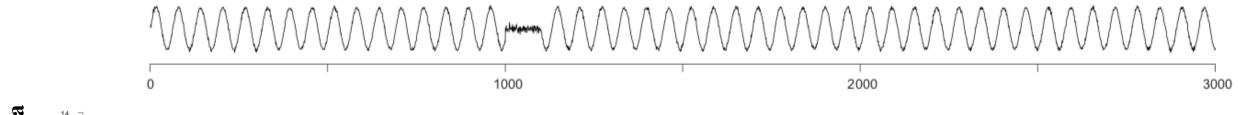
Поиск диссонансов произвольной длины

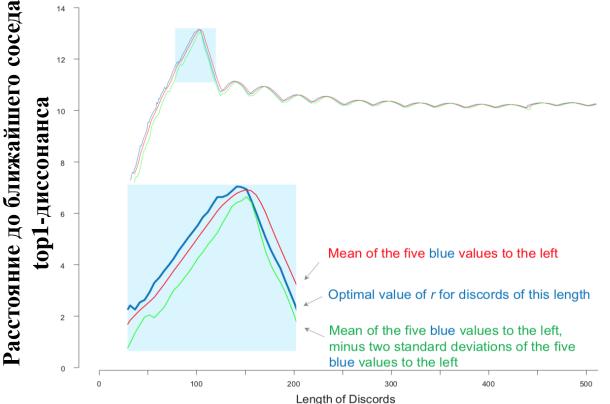
• Дано: временной ряд T, диапазон длины диссонансов $minL \dots maxL$

• Найти:
$$\mathcal{D} = \bigcup_{m=minL}^{maxL} D_m$$
, $D_m = \{d_1^m, d_2^m, ...\}$, где $\exists r_m \min_{s \in M_{d_i^m}} \mathrm{Dist}(d_i^m, s) \geq r_m$



Алгоритм MERLIN: адаптивное вычисление порога r

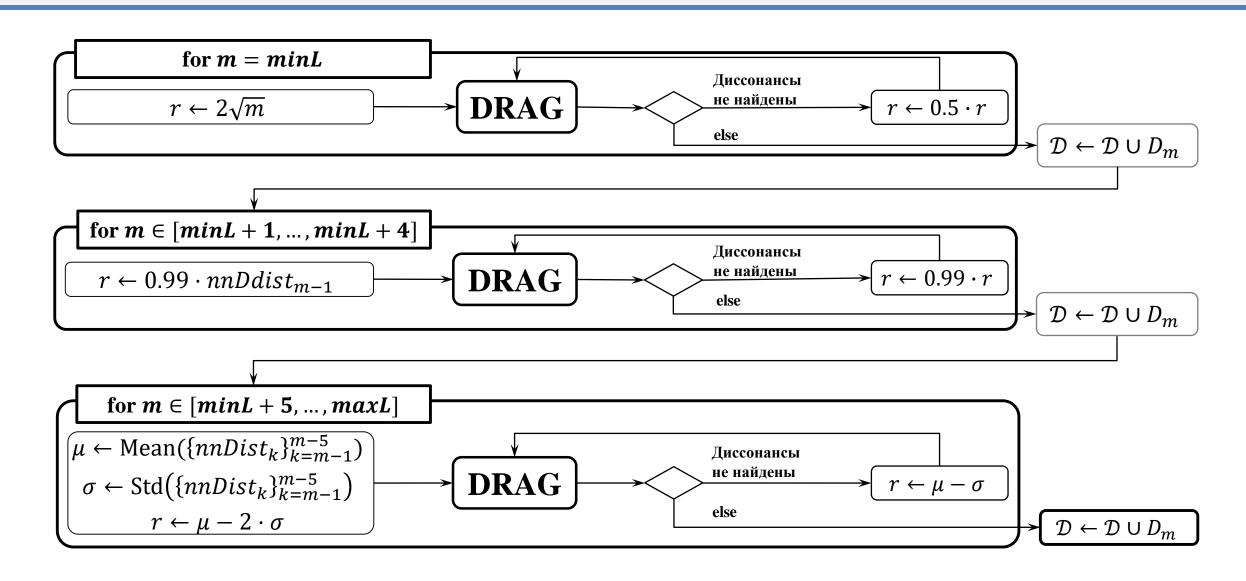




Порог r нужно вычислять по-разному для разных длин m

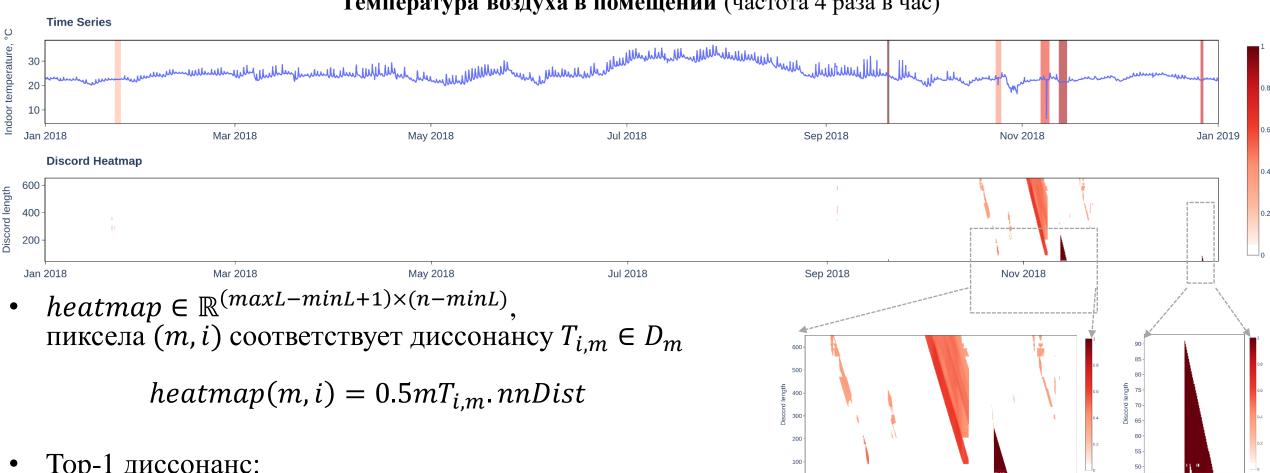
Длина диссонанса, т	Порог, r
minL	$r = 2\sqrt{minL}$
minL + 1,, minL + 4	$r = 0.99 \cdot nnDist_{m-1}$
minL + 5,, maxL	$r = \mu - 2\sigma$

Алгоритм MERLIN



Тепловая карта диссонансов

Температура воздуха в помещении (частота 4 раза в час)



Тор-1 диссонанс:

heatmap(m, i) $\max_{1 \le i \le n-m+1} \max_{minL \le m \le maxL}$

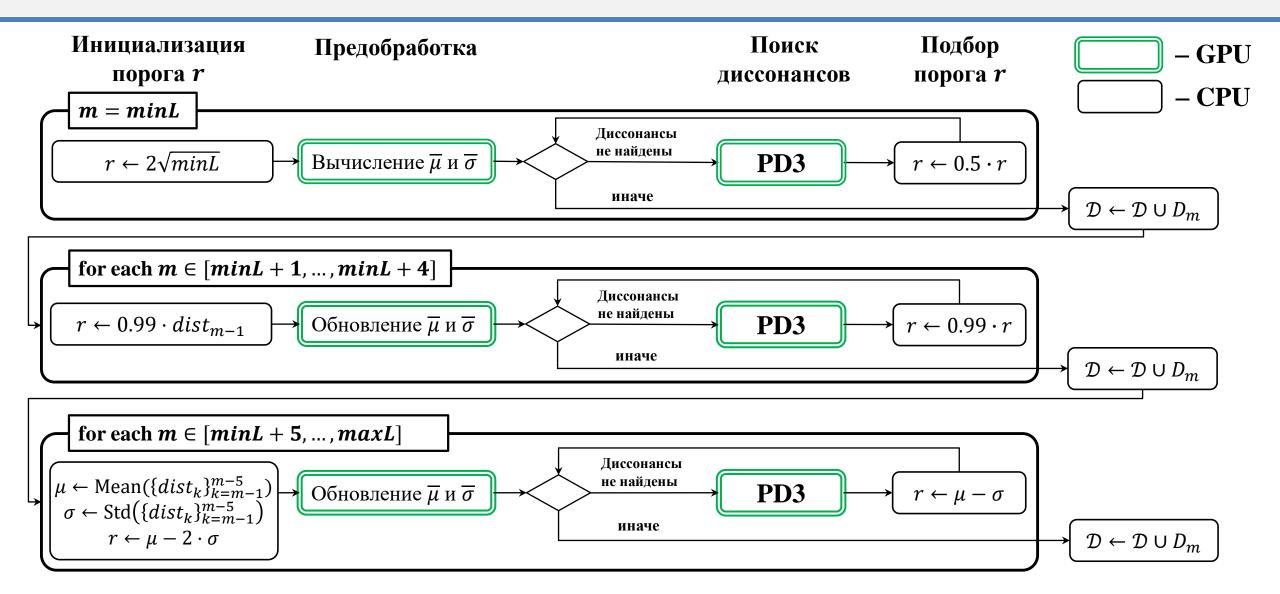
Алгоритм MERLIN: Практикум

- 1. Найдите диссонансы в заданном временном ряде с помощью алгоритма MERLIN
- 2. Отобразите временной ряд и найденные диссонансы
- 3. Сравните результаты и время работы алгоритмов DRAG и MERLIN

Содержание

- Временные ряды и аномалии в них
- Диссонансы
- Последовательный алгоритм DRAG
- Параллельный алгоритм PD3
- Последовательный алгоритм MERLIN
- Параллельный алгоритм PALMAD

PALMAD (Parallel Arbitrary Length MERLIN-based Anomaly Detector)

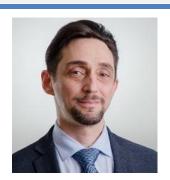


Алгоритм PALMAD: Практикум

- 1. Найдите диссонансы в заданном временном ряде с помощью алгоритма PALMAD
- 2. Отобразите временной ряд и найденные диссонансы
- 3. Сравните результаты и время работы алгоритмов MERLIN и PALMAD

Заключение

- Диссонанс формализует понятие аномалии
- Последовательные алгоритмы
 - **DRAG** поиск диссонансов фиксированной длины
 - MERLIN поиск диссонансов произвольной длины



Михаил Леонидович Цымблер, д.ф.-м.н. mzym@susu.ru https://mzym.susu.ru

- Параллельные алгоритмы для **GPU**
 - PD3 поиск диссонансов фиксированной длины
 - PALMAD поиск диссонансов произвольной длины



Яна Александровна Краева, к.ф.-м.н. kraevaya@susu.ru

53/52