



# > ПОИСК АНОМАЛИЙ (СИСТЕМА АЛЕРТОВ)

## > Что такое аномалия?

**Аномалия** — это необычный характер поведения данных, несвойственные метрике изменения.

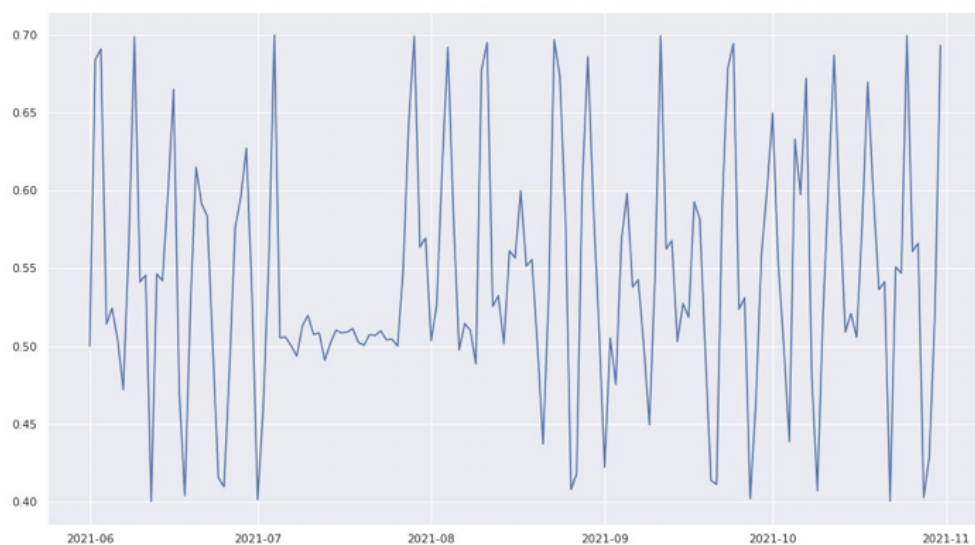
Аномалии в данных могут появляться по разным причинам: технический сбой в работе сервиса, сбои в ETL-процессах, ошибки при обработке данных, внешние факторы и т.д.

Можно условно выделить следующие виды аномалий:

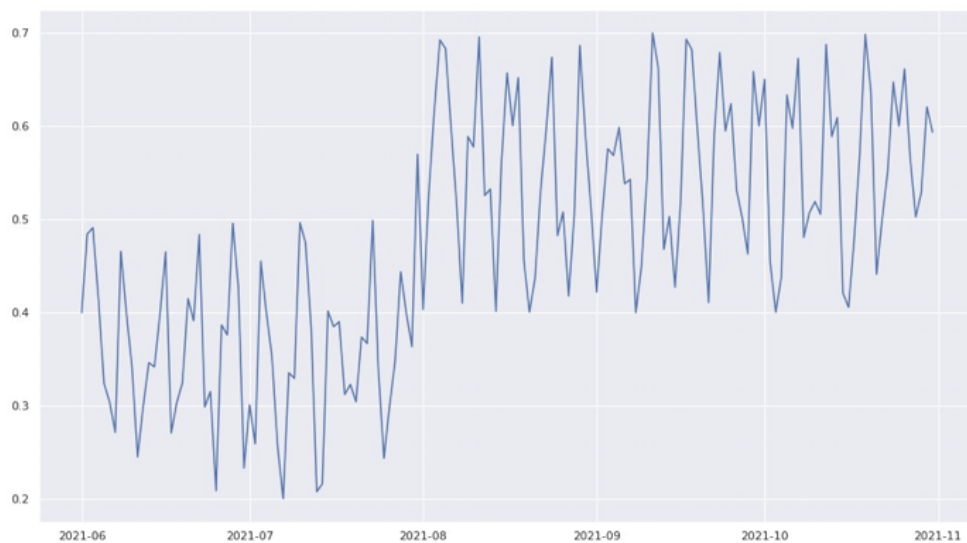
- *Точечная аномалия* — выброс.



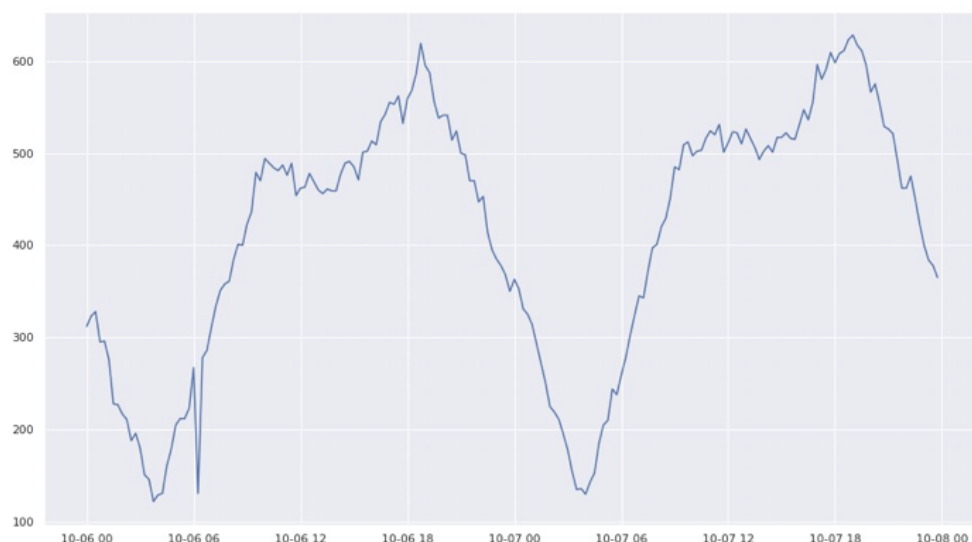
- *Коллективная аномалия* — изменения характера распределения данных. Эта аномалия наблюдается, когда последовательность связанных экземпляров данных (участок временного ряда) распределена иначе. Отдельное наблюдение в такой последовательности может не являться отклонением, однако совместное появление нескольких таких экземпляров является аномальным.



- *Коллективная аномалия — сдвиг.*



- *Контекстуальная аномалия.* Эта аномалия наблюдается, когда метрика в целом принимает приемлемое значение, однако не выполняются другие условия, которые позволили бы считать это значение нормальным. На приведенном графике наблюдается аномальное значение метрики, не характерное для этого момента времени.



## > Метод детектирования аномалий?

В зависимости от характера поведения метрик, наличия исторических данных и условий, в которых необходимо распознавать аномалии, применяются различные методы детектирования. Их можно условно разделить на следующие группы:

- статистические методы
- методы на основе ML

Далее приведены наиболее популярные методы распознавания аномалий из каждой группы.

### Статистические методы

*Правило сигм*

$$x_i < \mu_n - a * \sigma_n \mid x_i > \mu_n + a * \sigma_n$$

$\mu_n$  — скользящее среднее за период  $n$

$\sigma_n$  — стандартное отклонение за период  $n$

$a$  — коэффициент, задающий количество стандартных отклонений

Чаще всего выбирают  $a = 3$ , а метод наиболее известен как “правило трех сигм”

*Межквартильный размах*

$$x_i < q_{25_n} - a * IQR_n \mid x_i > q_{75_n} + a * IQR_n$$

$$\begin{aligned}
 & q_{25_n} - 0.25 \text{ квантиль за период } n \\
 & q_{75_n} - 0.75 \text{ квантиль за период } n \\
 & IQR_n = q_{75_n} - q_{25_n} - \text{межквартильный размах за период } n \\
 & a - \text{коэффициент}
 \end{aligned}$$

## Методы на основе ML

### *DBSCAN (Density-based spatial clustering of applications with noise)*

Метод кластеризует данные, учитывая их плотность и возможный шум. На вход принимает следующие параметры: радиус окрестности и количество соседей. Согласно алгоритму, в окрестности с заданным радиусом вокруг точки должно быть некоторое количество соседей, не меньше заданного числа. Точка будет считаться аномальной, если в обозримой близости у нее нет соседей. Более подробно про механизм работы DBSCAN можно прочитать [тут](#).

### *LOF (Local Outlier Factor)*

Метод неконтролируемого обнаружения аномалий, который основывается на концепции локальной плотности. Локальная плотность оценивается с помощью расстояний до некоторого количества ближайших соседей. Точки с значительно меньшей локальной плотностью считаются выбросами.

## > Работаем с аномалиями в realtime

Давайте представим ситуацию: утром мы получаем ежедневный отчет в Telegram и видим, что количество лайков вчера упало более, чем на 40%! Что же случилось? Мы в срочном порядке начинаем исследовать это аномальное падение, переходим в BI-систему для более детального анализа и видим, что часть пользователей, а точнее те, кто пользуются iOS, не ставили лайки примерно с 13:00. Конечно же, они не могли все сговориться и объявить бойкот. В ходе исследования аномалии мы обнаружили, что у нас на iOS сломалась кнопка лайка. А мы узнали об этом спустя почти сутки!

Конечно, мы могли бы после такого инцидента ввести дежурства в команде аналитиков, чтобы всегда был сотрудник, который внимательно следит за показателями на всех дашбордах, а в случае отклонений незамедлительно начинает расследование. Однако рационально ли так распоряжаться рабочим временем аналитика? Нет. К тому же количество дашбордов и графиков в различных срезах может быть таким большим, что одному сотруднику не уследить за всеми сразу.

Падение количества лайков или просмотров, значительный рост количества новых постов, падение до нуля количества отправленных пользователями сообщений — все это аномалии в данных, которые очень важно как можно раньше заметить, начать исследовать и устранять возможные неполадки. Поэтому нам

необходимо искать подобные аномалии автоматически - то есть нужна **система алертинга**.

**Систему алертинга** можно представить как *набор автоматизированных отчетов*, которые собираются с некоторой периодичностью по всем заданным метрикам и срезам. Отправляется такой *отчет-алерт* только тогда, когда выбранный нами метод зафиксировал аномальное значение метрики.

Как правило, для детектирования аномалий в realtime используют **статистические методы**, так как они, в отличие от методов на основе ML, менее требовательные и быстрее дают ответ.

Так как задача детектирования аномалий в realtime распространенная, существует достаточно большое количество готовых решений, например, [prometheus](#). Эти **системы мониторинга** предлагают инструменты для подключения к базе, отправки сообщений в мессенджер/на электронную почту и позволяют использовать совсем простые методы на основе статистики для мониторинга метрик.

На практике для детектирования аномалий не всегда достаточно самых простых методов, иногда требуется применение более кастомных подходов. Поэтому во многих компаниях распространена практика написание собственных систем мониторинга, позволяющих реализовать **более гибкий подход** к мониторингу метрик их продукта.

Перед началом написания скрипта системы алертов необходимо:

- определить перечень метрик и срезов;
- изучить поведение метрик;
- продумать метод детектирования аномалий и частоту мониторинга значений метрик.

Основные составляющие скрипта кастомной системы алертинга:

- перечень метрик и срезов для мониторинга;
- коннектор к БД для сбора данных;
- алгоритм детектирования аномалии;
- функция для отправки алерта.