

Выбросы

Лекция 02



Основные задачи машинного обучения

1. Поиск выбросов (outlier detection). Есть множество объектов M . Найти в нем все аномальные объекты.
2. Поиск новизны (novelty detection). Есть множество объектов M . Определить, является ли объект $A \notin M$ похожим на объекты из M или нет?

Основные задачи машинного обучения

3. Кластеризация (clustering). Дано множество объектов. Их нужно разбить на несколько групп (кластеров), состоящих из похожих друг на друга объектов.

4. Предсказание (prediction). Есть множество объектов M с известными значениями признака Y . Найти значение признака Y для нового объекта $A \notin M$.

Поиск выбросов

Отличие выбросов от новизны и пропусков

выброс VS пропуск в данных

Выброс – это часто реально существующий объект, но обладающий аномальными свойствами, он сильно отличается от других объектов выборки.

выброс VS новизна

Новизна считается по отношению к старой выборке объектов. А выброс является аномальным уже для своих «соседей» по выборке.

Примеры выбросов

1. (из Википедии) если наугад измерять температуру предметов в комнате, получим цифры от 18 до 22 °C, но радиатор отопления будет иметь температуру в 70° - и это выброс, не типичное значение!

Примеры выбросов

2. На матфаке ОмГУ я проводил анкетирование порядка 60 студентов: просил их написать средние баллы за все сессии. Выбросом оказался...

Примеры выбросов

2. На матфаке ОмГУ я проводил анкетирование порядка 60 студентов: просил их написать средние баллы за все сессии. Выбросом оказался КРУГЛЫЙ ОТЛИЧНИК (что было установлено с помощью алгоритмов поиска выбросов)

Зачем нужно искать и уничтожать выбросы?

1. Если данные будут использоваться при решении задачи предсказания, то удаление выбросов, как правило, повышает точность предсказания (ибо правило «мусор на входе – мусор на выходе» никто не отменял).
2. Удаление выбросов позволяет получить нормальные (типичные, эталонные) объекты.

Идеальных методов обнаружения выбросов не бывает потому, что

1. ... не существует формального определения выброса.
2. ... алгоритм, беспощадный к выбросам, будет удалять и часть «нормальных» объектов.
3. ... алгоритм, гуманный к «нормальным» объектам, будет пропускать часть выбросов.

Построить идеальный детектор выбросов – это всё равно что предложить мед. анализ без ложноположительных и ложноотрицательных результатов.

Методы обнаружения выбросов

- I. Поиск аномальных объектов с помощью здравого смысла. Например, если известен нормальный диапазон для значений признака.

Пример: человек с ростом более 200см (такие люди могут в реальности существовать, но их очень мало). Таких людей лучше объявить выбросами.

Методы обнаружения выбросов

- II. Методы, основанные на анализе одного признака (каждый признак берётся отдельно и ищутся объекты аномальными значениями этого признака).
- III. Методы, основанные на одновременном анализе нескольких признаков.

Методы, анализирующие один признак

Вот у вас есть значения признака

$$P = (p_1, p_2, \dots, p_n)$$

Основная идея поиска аномалий:

найти значения p_i , расположенные вдали от среднего значения (медианы).

Далее используем обозначения:

\bar{n} - среднее значение, n - объем выборки,

s_p - отклонение

Простейшие методы

1. Удалить все объекты, у которых величина $|p_i - \bar{p}|$ слишком велика.
2. Удалить все объекты, у которых величина $\frac{|p_i - \bar{p}|}{s_p}$ слишком велика.
3. Метод, использующий медиану, есть в Википедии (см. статью «Выброс»)
4. А еще лучше найти формулу, зависящую от n (это будет критерий Шавене на след. слайдах).

Критерий Шавене (Chauvenet)

Значение p_i является выбросом, если выполнено неравенство

$$\operatorname{erfc}\left(\frac{|p_i - \bar{p}|}{s_p}\right) < \frac{1}{2n}$$

где

$$\operatorname{erfc}(x) =$$

18+ censored

дополнение к функции ошибок.

Критерий Шавене (Chauvenet)

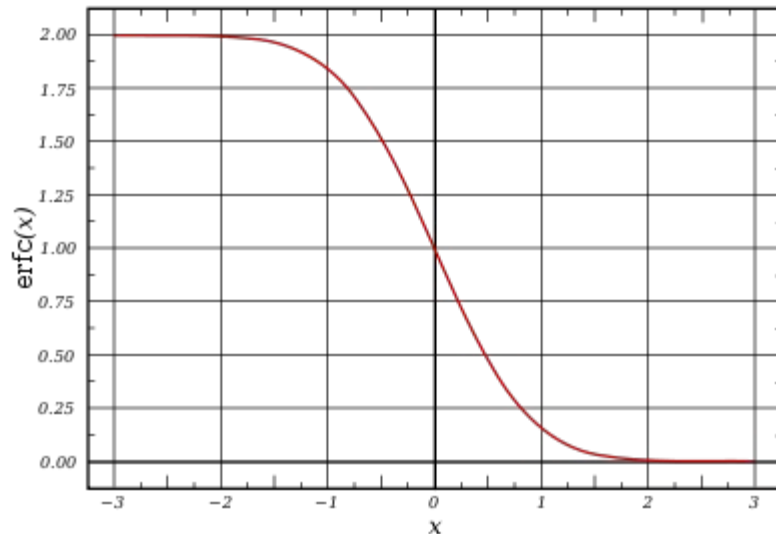
Значение p_i является выбросом, если выполнено неравенство

$$\operatorname{erfc}\left(\frac{|p_i - \bar{p}|}{s_p}\right) < \frac{1}{2n}$$

где

$$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^{\infty} e^{-t} dt$$

дополнение к функции ошибок.



Пример

Есть $n=14$ объектов со следующими значениями признака P

8.02 8.16 3.97 8.64 0.84 4.46 0.81 7.74 8.78 9.26 20.46
29.87 10.38 25.71

Вычисляем: $\bar{p}=10.51$, $s_p=8.77$.

Проверка для 25.71:

$$\text{erfc}\left(\frac{|25.71 - 10.51|}{8.77}\right) = 0.014 < 0.036 = \frac{1}{2 * 14}$$

то есть это выброс!!!

Пример

Есть $n=14$ объектов со следующими значениями признака P

8.02 8.16 3.97 8.64 0.84 4.46 0.81 7.74 8.78 9.26 20.46
29.87 10.38 25.71

Вычисляем: $\bar{p}=10.51$, $s_p=8.77$.

Проверка для 29.87 также говорит: «выброс». А вот 20.46 уже не выброс:

$$\operatorname{erfc}\left(\frac{|20.46 - 10.51|}{8.77}\right) = 0.11 > 0.036 = \frac{1}{2 * 14}$$

Пример (вторая итерация)

Осталось $n=12$ объектов со следующими значениями признака P

8.02 8.16 3.97 8.64 0.84 4.46 0.81 7.74 8.78 9.26 20.46
10.38

Вычисляем: $\bar{p}=7.63$, $s_p=5.17$.

Проверка для 20.46:

$$\text{erfc}\left(\frac{|20.46 - 7.63|}{5.17}\right) = 0.00045 < 0.042 = \frac{1}{2 * 12}$$

то есть теперь 20.46 стало выбросом!!!

Пример (вторая итерация)

Осталось $n=12$ объектов со следующими значениями признака P

8.02 8.16 3.97 8.64 0.84 4.46 0.81 7.74 8.78 9.26 20.46
10.38

Вычисляем: $\bar{p}=7.63$, $s_p=5.17$.

Проверка для 10.38:

$$\text{erfc}\left(\frac{|10.38 - 7.63|}{5.17}\right) = 0.452 > 0.042 = \frac{1}{2 * 12}$$

То есть 10.38 не выорос. Остальные числа также проходят проверку.

Пример (третья итерация)

Осталось $n=11$ объектов со следующими значениями признака P

8.02 8.16 3.97 8.64 0.84 4.46 0.81 7.74 8.78 9.26 10.38

Вычисляем: $\bar{p} = 6.46, s_p = 3.38$.

Проверка для 10.38:

$$\text{erfc}\left(\frac{|10.38 - 6.46|}{3.38}\right) = 0.1 > 0.045 = \frac{1}{2 * 11}$$

то есть 10.38 не выброс. Остальные числа также проходят проверку. КОНЕЦ работы, так как новые выбросы не появляются.

Настройка критерия Шавене

Значение p_i является выбросом, если выполнено неравенство

$$\operatorname{erfc}\left(\frac{|p_i - \bar{p}|}{s_p}\right) < \frac{1}{2n}$$

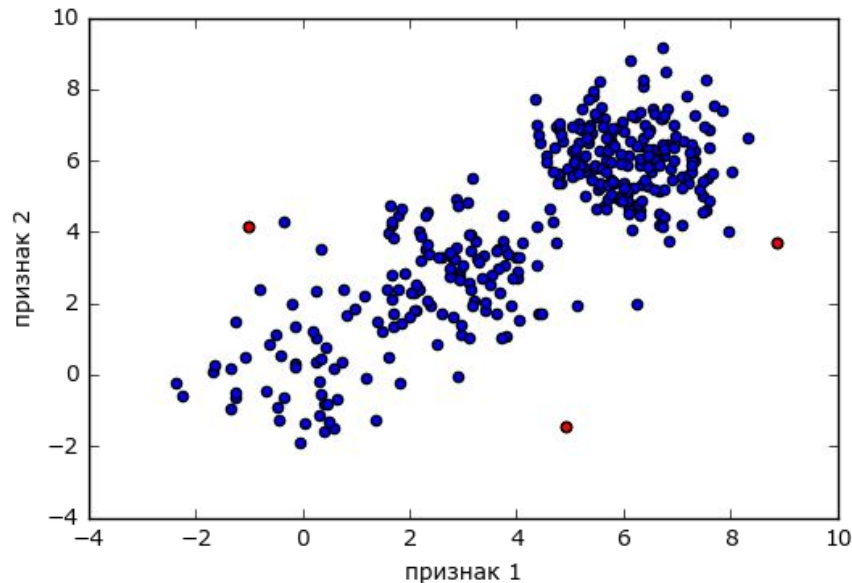
Константу 2 (в формуле) можно заменить на любую другую константу. Это сделает критерий более беспощадным (либо более лояльным) к выбросам.

Недостатки методов, которые анализируют 1 признак (1-й недостаток)

В выборке (1, 50, 2, 1, 101, 2, 1, 100, 1, 3, 101, 1, 3, 100, 101, 100, 100) число 50, очевидно, аномально. Но поскольку 50 очень близко к среднему значению, то все описанные выше методы его не заметят.

Недостатки методов, которые анализируют 1 признак (2-й недостаток)

Аномалия часто характеризуется не только экстремальными значениями отдельных признаков. На картинке каждый выброс имеет «нормальное» значения каждого признака, но их комбинация приводит к аномалии.

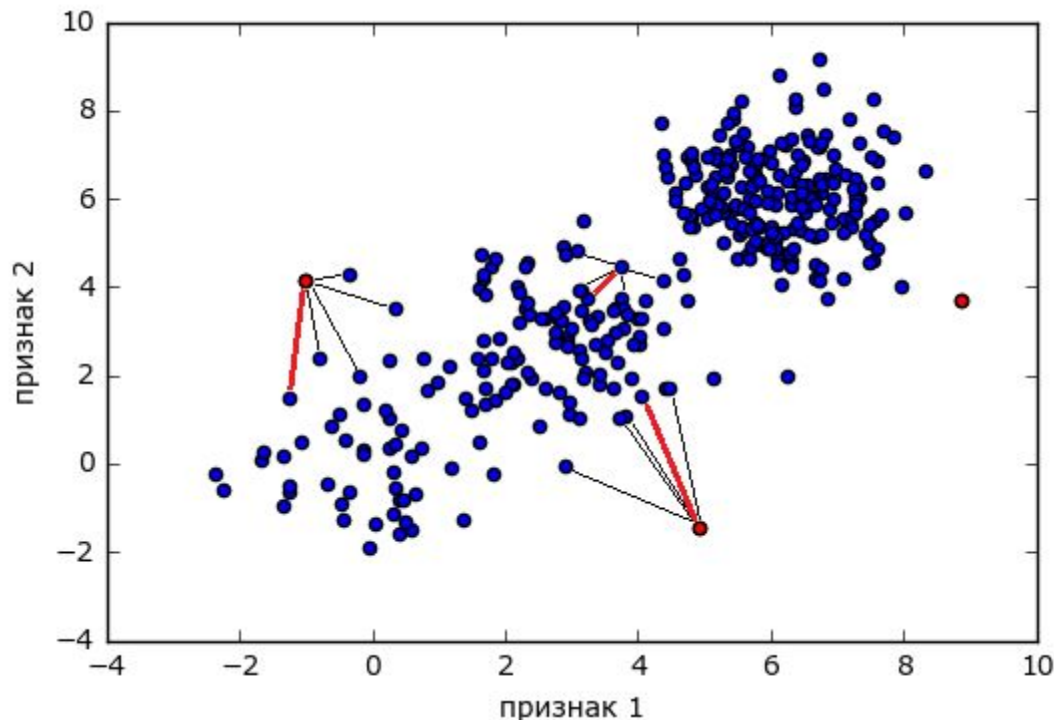


Методы, анализирующие несколько признаков

Метрические методы

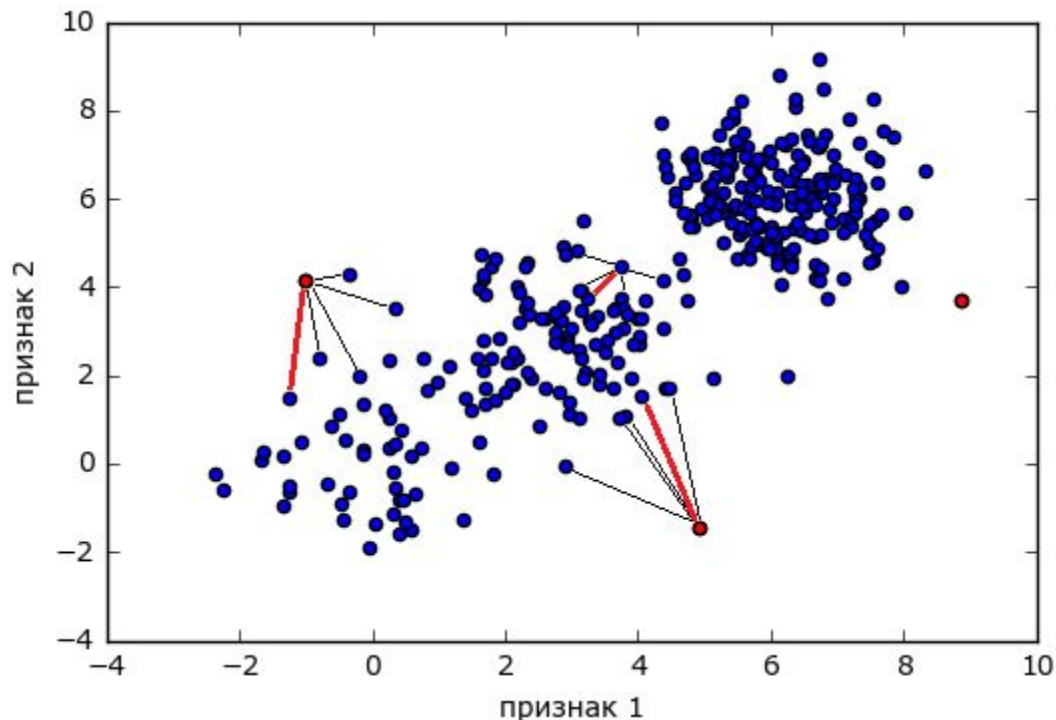
Представим объекты с t признаками с помощью точек в пространстве R^m .

Идея: у выброса мало соседей, а у типичного объекта много.



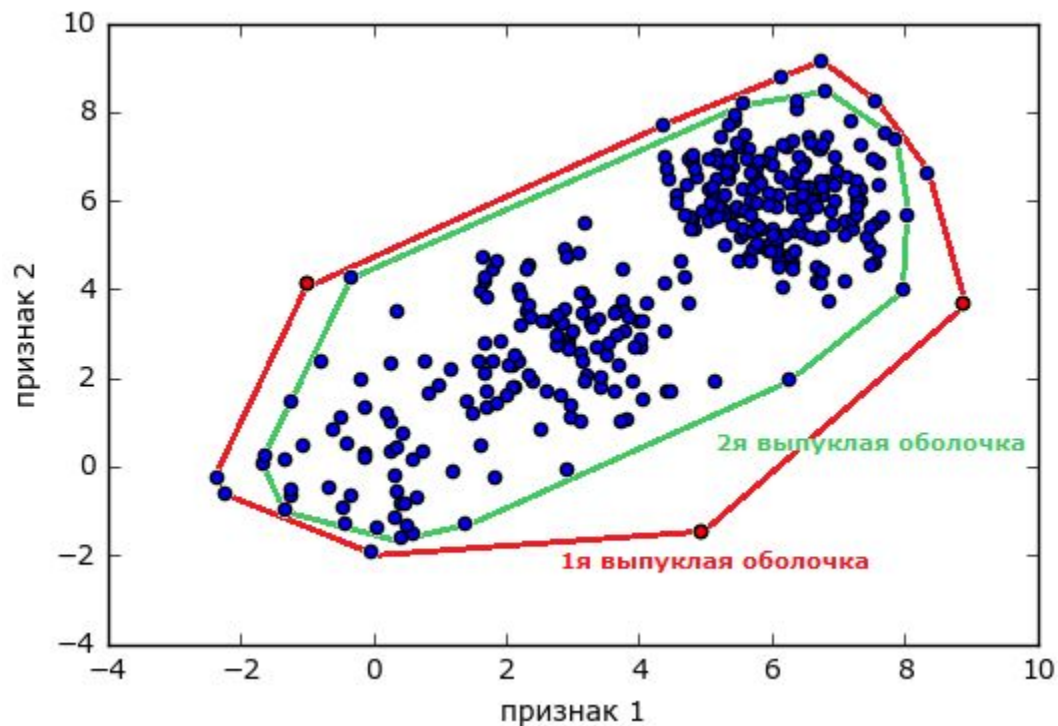
Метрические методы

Можно найти расстояние от каждого объекта до его ближайшего соседа. У выбросов такое расстояние будет большим (можно предложить и другие варианты алгоритма).



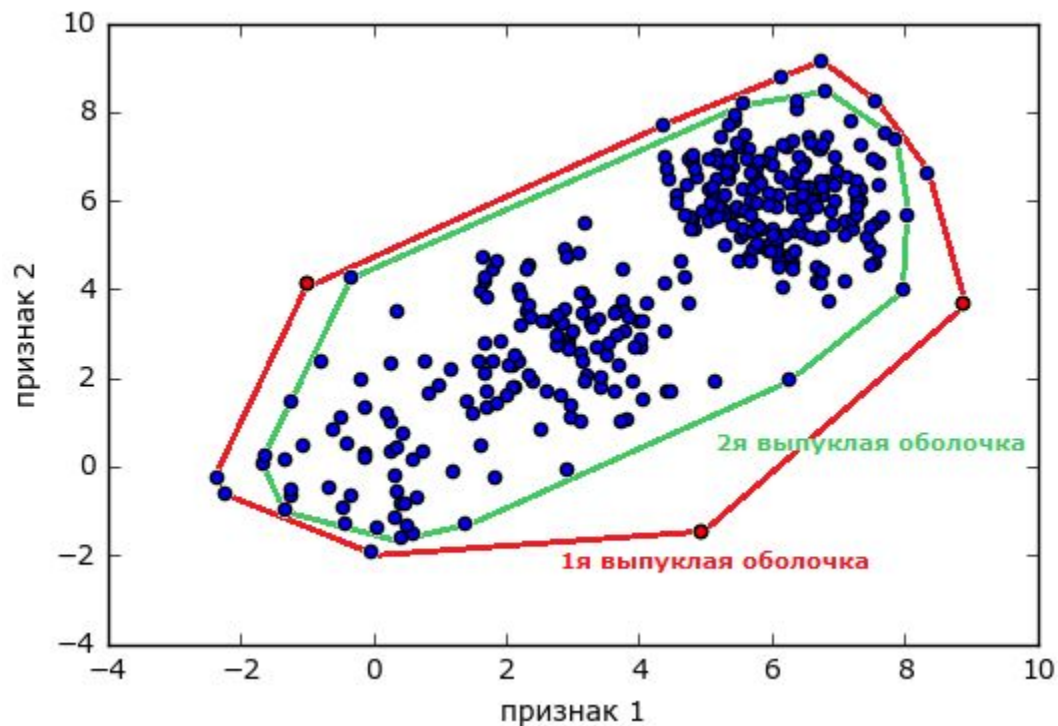
Геометрические методы

Можно вычислить
выпуклую оболочку
объектов (как точек в m -
мерном пространстве).
Выбросами будут ...



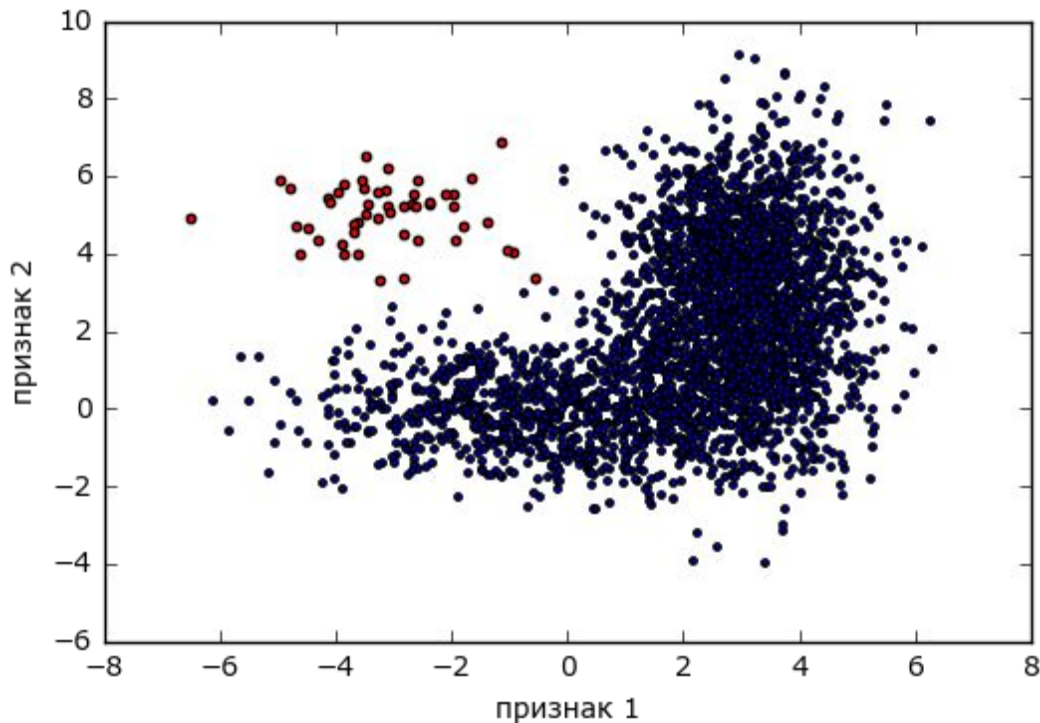
Геометрические методы

Можно вычислить выпуклую оболочку объектов (как точек в m -мерном пространстве). Выбросами будут объекты на границе. (Их можно удалить, а процедуру повторить еще несколько раз).



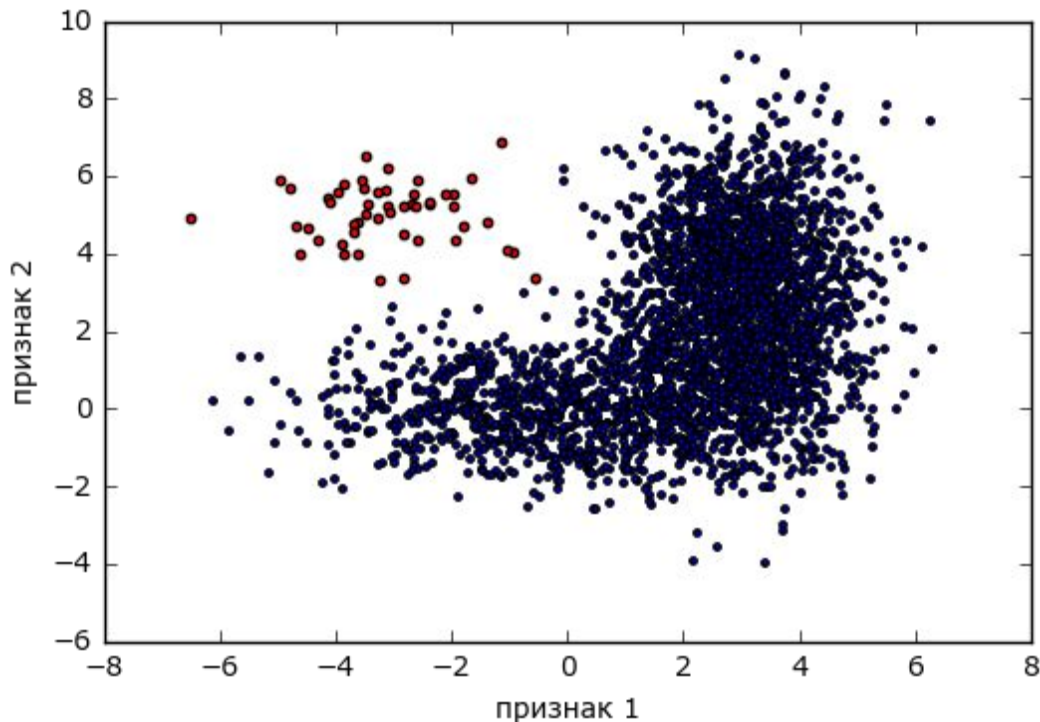
Поиск выбросов с помощью кластеризации

Можно запустить алгоритм кластеризации (подробности в след. лекции). Он разобьет объекты на группы. Выбросы – это ...



Поиск выбросов с помощью кластеризации

Можно запустить алгоритм кластеризации (подробности в след. лекции). Он разобьет объекты на группы. Выбросы – это элементы малых (в том числе и одноэлементных) групп.



Поиск выбросов с помощью моделей предсказания

- 1) Некоторые вариации метода опорных векторов (SVM) позволяют находить выбросы.
- 2) Вариация решающих деревьев (decision trees) под названием «изолирующий лес» может искать выбросы.

Подробнее об этом мы поговорим на соответствующих лекциях.

Поиск новизны

Поиск новизны можно свести к поиску выбросов

(Поиск новизны) Есть множество объектов M .
Определить, является ли объект $A \notin M$ похожим на
объекты из M или нет?

Переход к поиску выбросов: добавляем объект A в
множество M и запускаем алгоритм поиска выброса.
Если A будет детектирован как выброс, то A являлся
новизной.

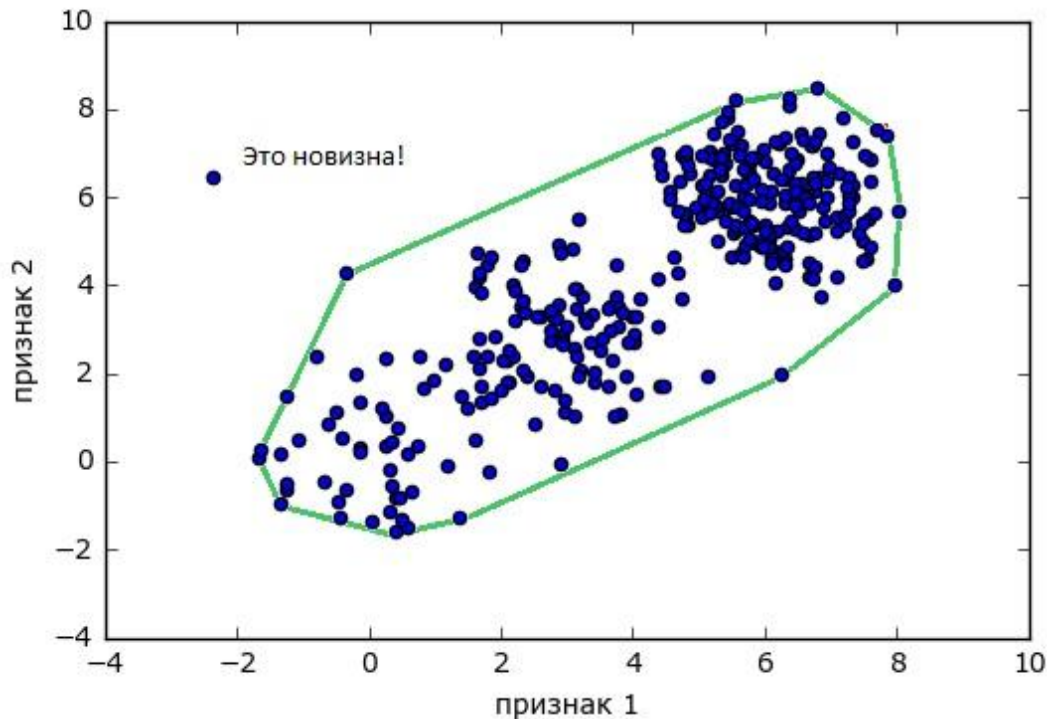
Но есть и алгоритмы, которые ищут новизну, не
используя поиск выбросов.

Поиск новизны без поиска выбросов

Строим выпуклую оболочку всех объектов из выборки. Объект А будет считаться новизной, если...

Поиск новизны без поиска выбросов

Строим выпуклую оболочку всех объектов из выборки. Объект А будет считаться новизной, если он НЕ попадает в эту выпуклую оболочку.



Использованная литература

1. <https://alexanderdyakonov.wordpress.com/2017/04/19/поиск-аномалий-anomaly-detection/> (отсюда же взяты и картинки)
2. Статья «Cleaning Data the Chauvenet Way» by Lily Lin, Paul Sherman