

Занятие 10

Методы снижения размерности и поиск аномалий

Блуменау М.И.

На основе материалов Кантонистовой Е.О.

ВШЭ, 2025

МЕТОД ГЛАВНЫХ КОМПОНЕНТ

Предыдущие методы отбирали из исходных признаков некоторое подмножество признаков.

Теперь мы *хотим придумать новые признаки, каким-то образом выражающиеся через старые, причем новых признаков хочется получить меньше, чем старых*. Сегодня будем рассматривать только случай, когда новые признаки линейно выражаются через старые.

МЕТОД ГЛАВНЫХ КОМПОНЕНТ

Постановка задачи:

x_1, \dots, x_n - исходные числовые признаки, $x_i = f_i(\mathbf{x})$

z_1, \dots, z_d - новые числовые признаки, $d \leq n$, $z_j = g_j(\mathbf{x})$.

Хотим:

1. чтобы новые числовые признаки z_j линейно выражались через исходные признаки x_i

2. чтобы при переходе к новым признакам было потеряно наименьшее количество исходной информации.

Дисперсия выборки, посчитанная в новых признаках, показывает, как много информации нам удалось сохранить после понижения размерности, поэтому дисперсия в новых признаках должна быть максимальной.

МЕТОД ГЛАВНЫХ КОМПОНЕНТ

Новые признаки линейно выражаются через исходные:

$$\begin{cases} z_1 = u_{11}x_1 + \dots + u_{1n}x_n \\ z_2 = u_{21}x_1 + \dots + u_{2n}x_n \\ \dots \\ z_d = u_{d1}x_1 + \dots + u_{dn}x_n \end{cases}$$

Будем искать такие векторы u_1, \dots, u_m , что они:

- ортогональны: $(u_i, u_j) = 0$
- нормированы: $\|u_i\| = 1$

МЕТОД ГЛАВНЫХ КОМПОНЕНТ

Новые признаки линейно выражаются через исходные:

$$\begin{cases} z_1 = u_{11}x_1 + \dots + u_{1n}x_n \\ z_2 = u_{21}x_1 + \dots + u_{2n}x_n \\ \dots \\ z_d = u_{d1}x_1 + \dots + u_{dn}x_n \end{cases}$$

Геометрическая интерпретация: новые признаки z_i — это проекции исходных признаков x_i на некоторые векторы (компоненты) u .

МЕТОД ГЛАВНЫХ КОМПОНЕНТ

Новые признаки линейно выражаются через исходные:

$$\begin{cases} z_1 = u_{11}x_1 + \dots + u_{1n}x_n \\ z_2 = u_{21}x_1 + \dots + u_{2n}x_n \\ \dots \\ z_d = u_{d1}x_1 + \dots + u_{dn}x_n \end{cases}$$

Геометрическая интерпретация: новые признаки z_i — это проекции исходных признаков x_i на некоторые векторы (компоненты) u .

- Проекция объекта x на компоненту u_i : $z_i = (x, u_i) = u_{i1}x_1 + \dots + u_{in}x_n$
- Проекция всей выборки на компоненту u_i : $Z_i = Xu_i$

МЕТОД ГЛАВНЫХ КОМПОНЕНТ

Геометрическая интерпретация: новые признаки z_i — это проекции исходных признаков x_i на некоторые векторы (компоненты) u .

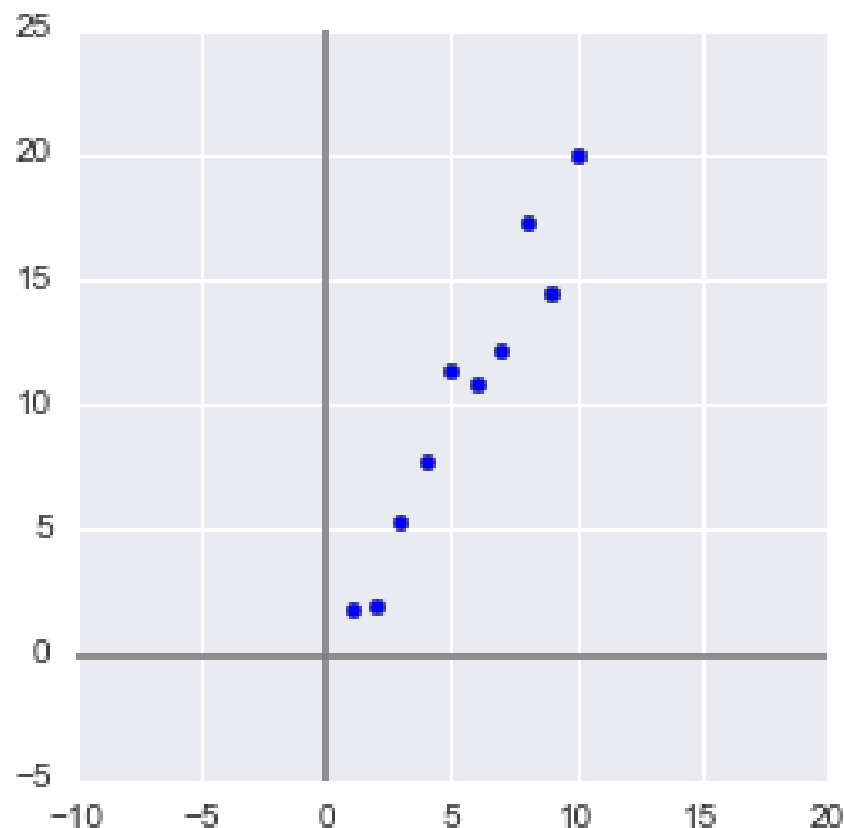
- Проекция объекта x на компоненту u_i : $z_i = (x, u_i) = u_{i1}x_1 + \dots + u_{in}x_n$
- Проекция всей выборки на компоненту u_i : $Z_i = Xu_i$

Наша цель: найти такие компоненты u_i , чтобы дисперсия проекции выборки на них была максимальной:

$$D(Xu_i) \rightarrow \max_{u_i}, \quad i = 1, \dots, d$$

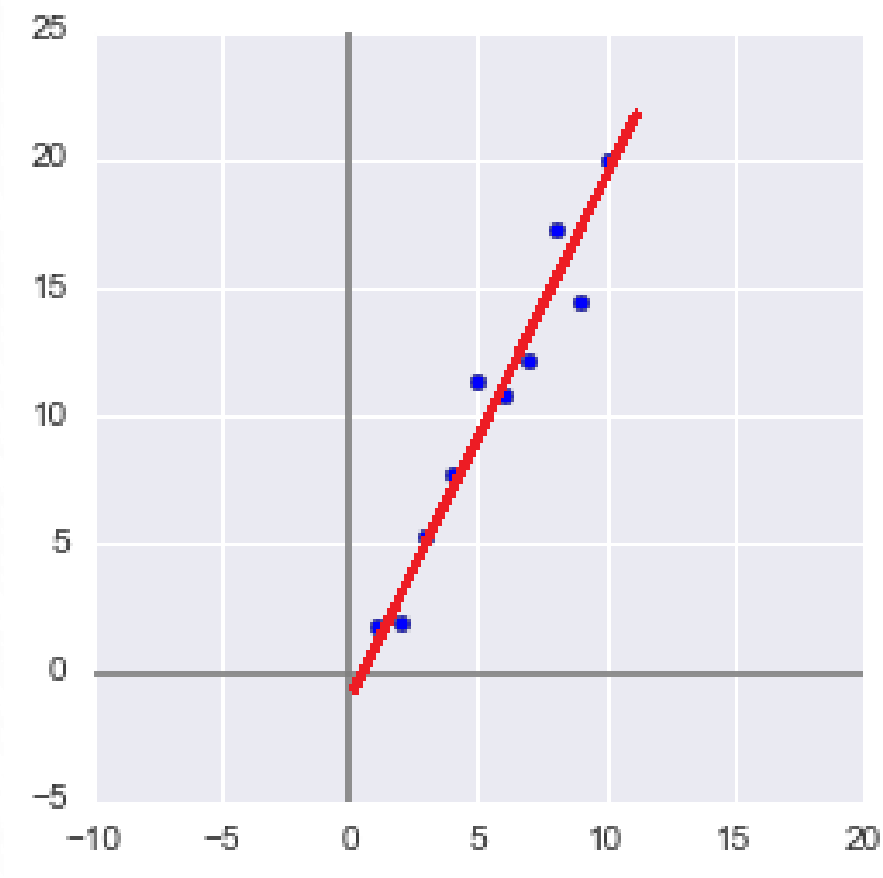
ПРИМЕР

Хотим спроецировать двумерные данные X на одномерный вектор u так, чтобы дисперсия проекции Xu была максимальной:



ПРИМЕР

Хотим спроецировать двумерные данные X на одномерный вектор u так, чтобы дисперсия проекции Xu была максимальной:



ПРОЕКЦИИ МЕТОДА ГЛАВНЫХ КОМПОНЕНТ

- Пусть X – матрица объект-признак для исходных признаков.
- Метод главных компонент делает проекцию исходных объектов на гиперплоскость некоторой размерности d .

Теорема. Базисные векторы этой гиперплоскости – это собственные векторы матрицы $X^T X$ (матрица ковариаций), соответствующие d её наибольшим собственным значениям.

КОНСТРУКТИВНОЕ ПОСТРОЕНИЕ БАЗИСА В РСА

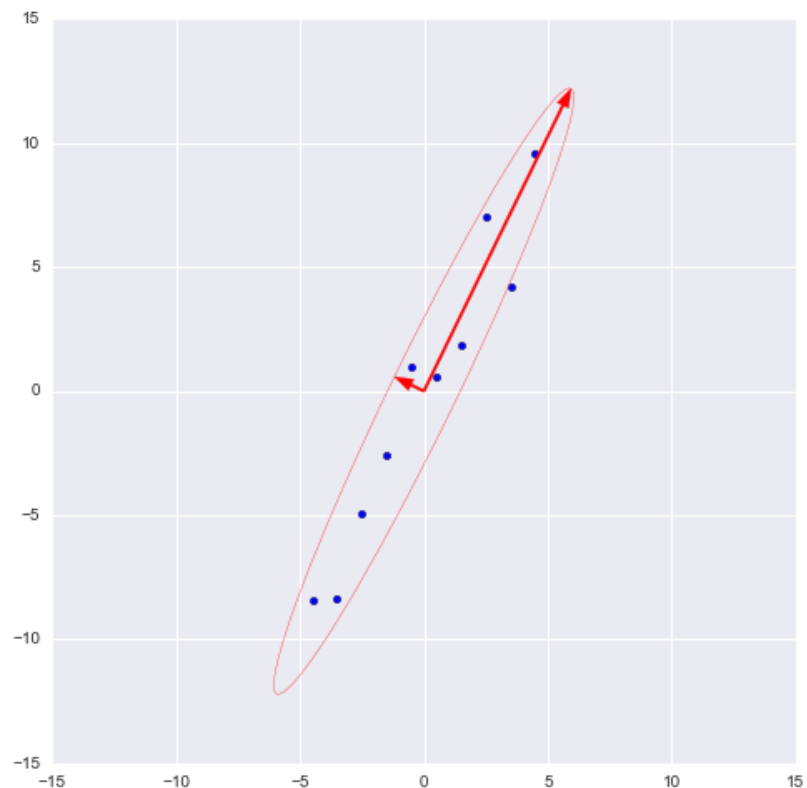
- Находим вектор $u_1 = \operatorname{argmax}_u (D(Xu))$ и нормируем его: $u_1 \rightarrow \frac{u_1}{\|u_1\|}$
- Находим вектор $u_2 = \operatorname{argmax}_u (D(Xu))$ такой, что $(u_1, u_2) = 0$ и нормируем его: $u_2 \rightarrow \frac{u_2}{\|u_2\|}$
- Находим вектор $u_3 = \operatorname{argmax}_u (D(Xu))$ такой, что $(u_1, u_3) = (u_2, u_3) = 0$ и нормируем его: $u_3 \rightarrow \frac{u_3}{\|u_3\|}$.

И т.д.

Получаем ортонормированный базис $\{u_1, u_2, \dots, u_d\}$.

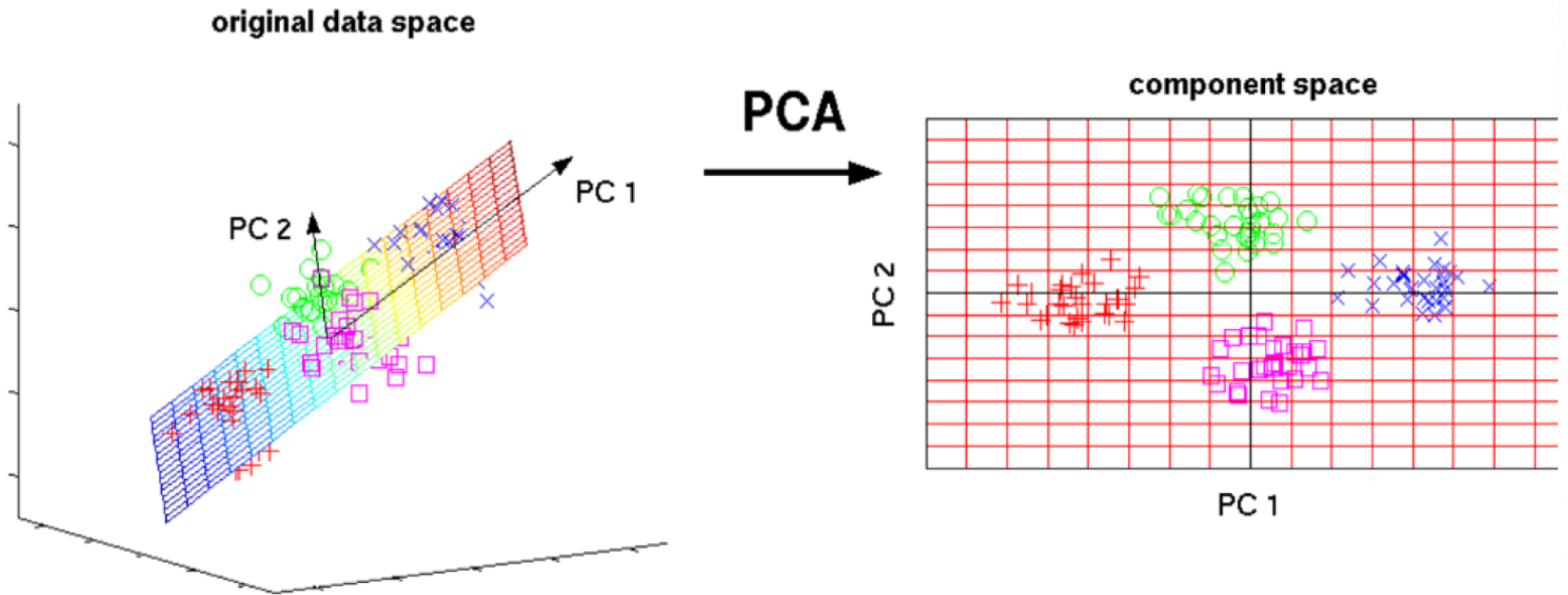
ГЕОМЕТРИЧЕСКИЙ СМЫСЛ PCA

- Нахождение собственных векторов матрицы $X^T X$ позволяет нам аппроксимировать исходные данные эллипсоидом, натянутым на эти векторы



- Затем мы делаем проекцию на подпространство, натянутое на собственные векторы с наибольшими собственными значениями

ПРОЕКЦИЯ НА ГИПЕРПЛОСКОСТЬ

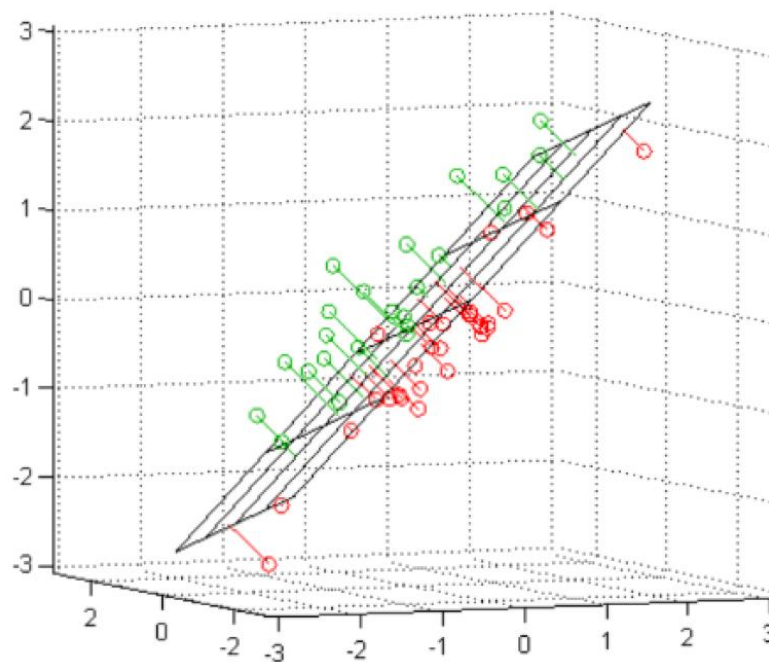


АЛЬТЕРНАТИВНАЯ ПОСТАНОВКА ЗАДАЧИ

Найти новые признаки Z и матрицу проецирования U , наилучшим образом восстанавливающие исходные

признаки: $\|X - ZU^T\|^2 \rightarrow \min_{Z,U}$

Геометрически это эквивалентно нахождению гиперплоскости, сумма квадратов расстояний от которой до точек выборки минимальна:



ДОЛЯ ОБЪЯСНЕННОЙ ДИСПЕРСИИ

- Упорядочим собственные значения матрицы $X^T X$ по убыванию: $\lambda_1 \geq \lambda_2 \geq \dots > \lambda_n \geq 0$.

- Доля дисперсии, объяснённой j -й компонентой (explained variance ratio):

$$\delta_j = \frac{\lambda_j}{\sum_{i=1}^n \lambda_i}$$

- Доля дисперсии, объясняемой первыми k компонентами:

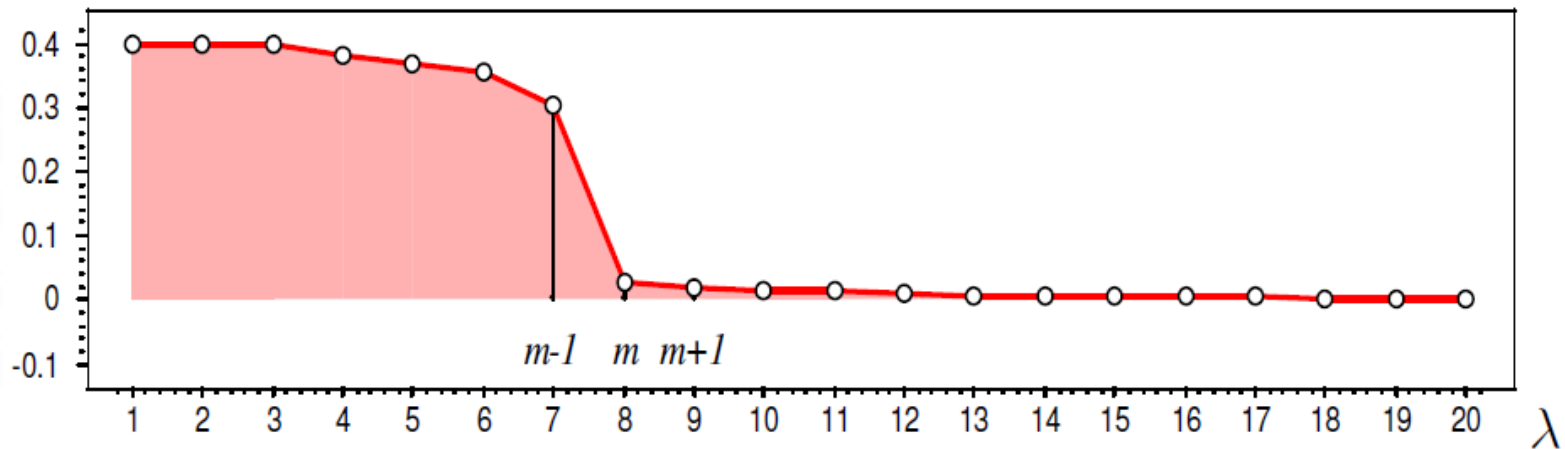
$$\delta = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_n} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n \lambda_i}$$

ВЫБОР ЧИСЛА ГЛАВНЫХ КОМПОНЕНТ

- Эффективная размерность выборки – это наименьшее целое m , при котором *доля необъясненной дисперсии*

$$E_m = \frac{\|ZU^T - X\|^2}{\|X\|^2} = \frac{\lambda_{m+1} + \dots + \lambda_n}{\sum_{i=1}^n \lambda_i} \leq \varepsilon$$

Критерий крутого склона:



ПРИМЕР: FACES DATASET



FACES DATASET (ГЛАВНЫЕ КОМПОНЕНТЫ)



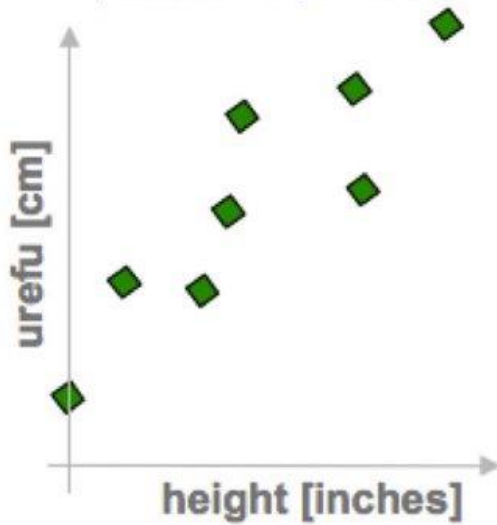
ВОССТАНОВЛЕННОЕ ИЗОБРАЖЕНИЕ



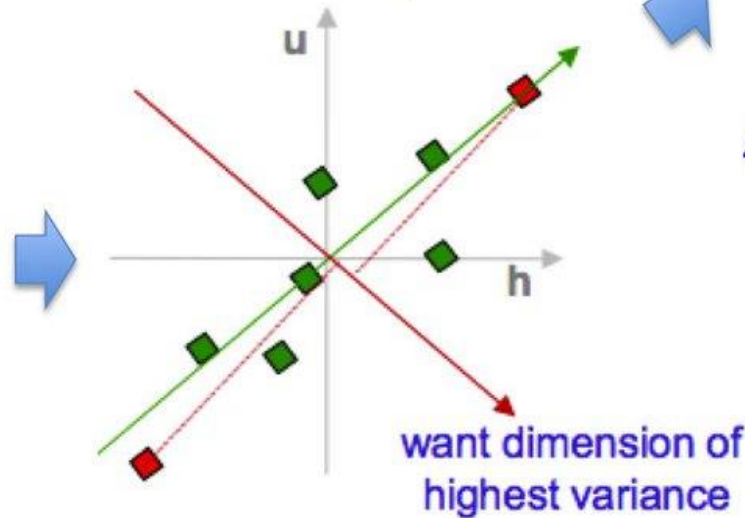
PCA in a nutshell

1. correlated hi-d data

("urefu" means "height" in Swahili)



2. center the points



3. compute covariance matrix

$$\begin{matrix} & h & u \\ h & \begin{pmatrix} 2.0 & 0.8 \end{pmatrix} \\ u & \begin{pmatrix} 0.8 & 0.6 \end{pmatrix} \end{matrix} \rightarrow \text{cov}(h,u) = \frac{1}{n} \sum_{i=1}^n h_i u_i$$

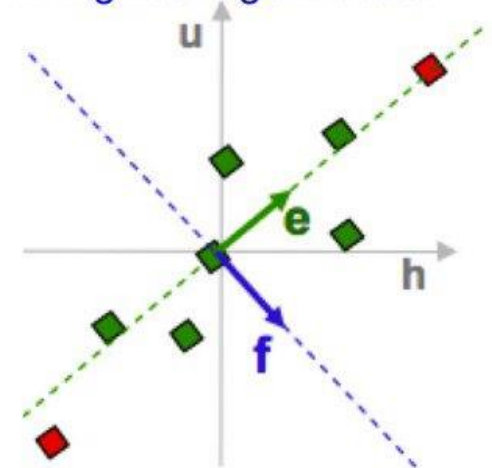
4. eigenvectors + eigenvalues

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix} \begin{pmatrix} e_h \\ e_u \end{pmatrix} = \lambda_e \begin{pmatrix} e_h \\ e_u \end{pmatrix}$$

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix} \begin{pmatrix} f_h \\ f_u \end{pmatrix} = \lambda_f \begin{pmatrix} f_h \\ f_u \end{pmatrix}$$

`eig(cov(data))`

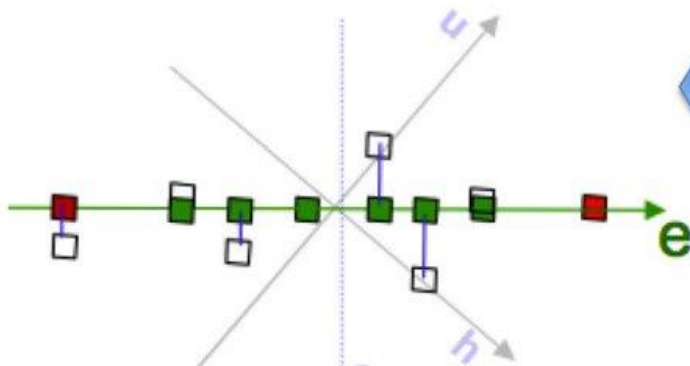
5. pick $m < d$ eigenvectors w. highest eigenvalues



6. project data points to those eigenvectors

$$x'_e = x^T e = \sum_{j=1}^d x_{ij} e_j$$

7. uncorrelated low-d data



СИНГУЛЯРНОЕ РАЗЛОЖЕНИЕ МАТРИЦЫ (SINGULAR VALUE DECOMPOSITION, SVD)

Теорема. Матрицу $A \in \mathbb{R}^{m \times n}$ можно представить в виде

$$A = U \Sigma V^T,$$

- где $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ - ортогональные матрицы,
- $\Sigma \in \mathbb{R}^{m \times n}$ - диагональная матрица с ненулевыми элементами $\sigma_i = \sqrt{\lambda_i}$, где λ_i - собственные значения матрицы $A^T A$.

При этом

- Столбцы матрицы U являются собственными векторами матрицы AA^T
- Столбцы матрицы V являются собственными векторами матрицы $A^T A$.

SINGULAR VALUE DECOMPOSITION

- При $m \leq n$:

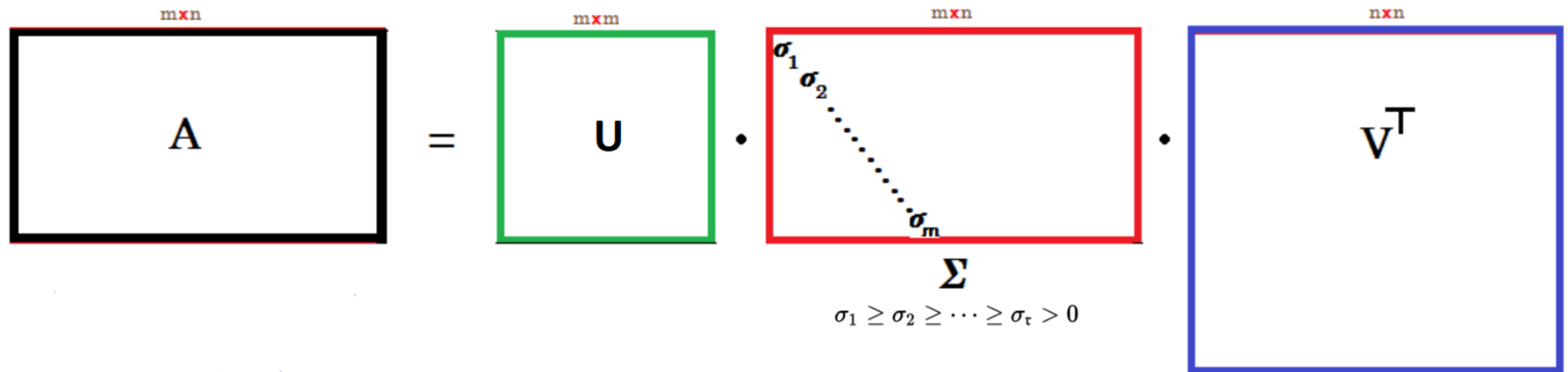


Diagram illustrating the Singular Value Decomposition (SVD) for the case $m \leq n$. The matrix A (size $m \times n$) is decomposed into three matrices: U (size $m \times m$), Σ (size $m \times n$), and V^T (size $n \times n$). The matrix Σ is shown with its diagonal elements $\sigma_1, \sigma_2, \dots, \sigma_m$. Below Σ , the condition $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ is stated.

$$A = U \cdot \Sigma \cdot V^T$$

$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$

- При $m > n$:

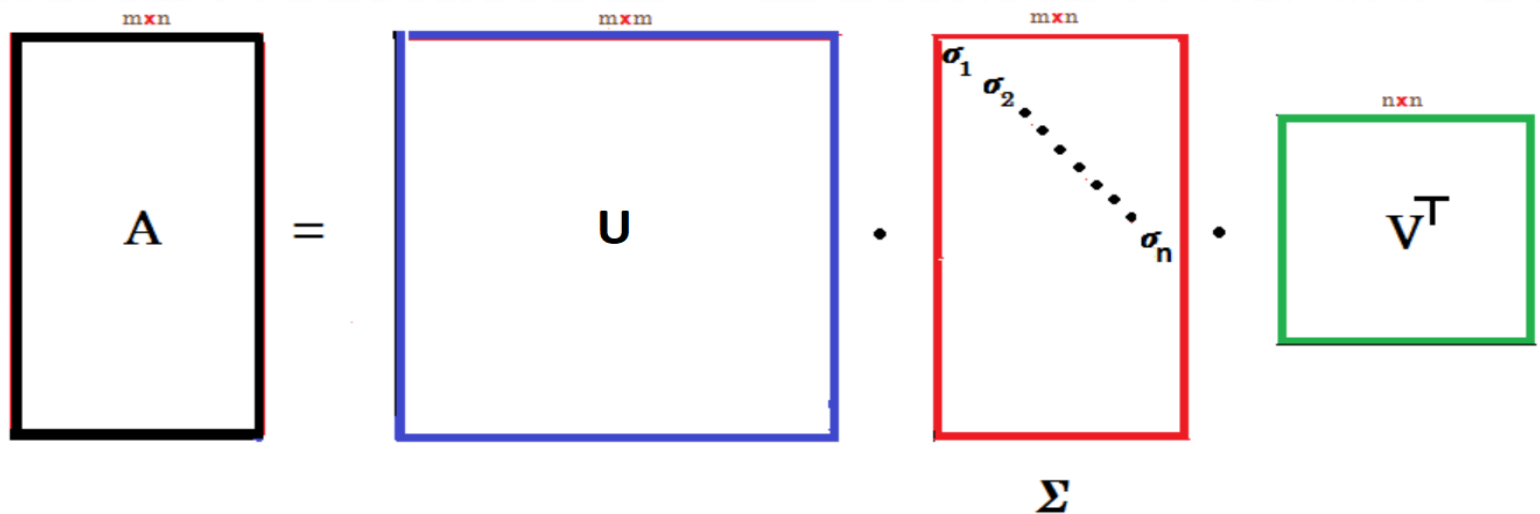


Diagram illustrating the Singular Value Decomposition (SVD) for the case $m > n$. The matrix A (size $m \times n$) is decomposed into three matrices: U (size $m \times m$), Σ (size $m \times n$), and V^T (size $n \times n$). The matrix Σ is shown with its diagonal elements $\sigma_1, \sigma_2, \dots, \sigma_n$.

$$A = U \cdot \Sigma \cdot V^T$$

СВЯЗЬ SVD И PCA

Пусть X – матрица объект-признак, для которой мы хотим снизить размерность и $X = U\Sigma V^T$ её SVD-разложение.

Тогда:

- Столбцы матрицы V – это собственные векторы матрицы $X^T X$, т.е. векторы v_1, \dots, v_n – главные компоненты.
- Столбцы матрицы $U\Sigma$ – это новые признаки, то есть, проекции исходных признаков на главные компоненты $Z = Xv$

$$(X = U\Sigma V^T \Leftrightarrow U\Sigma = XV).$$

- Сингулярные числа матрицы Σ – это корни из собственных чисел матрицы $X^T X$.

СВЯЗЬ SVD И PCA

- Столбцы матрицы V – это собственные векторы матрицы $X^T X$, т.е. векторы v_1, \dots, v_n – главные компоненты.
- Столбцы матрицы $U\Sigma$ – это новые признаки $z = Xv$ ($X = U\Sigma V^T \Leftrightarrow U\Sigma = XV$).
- Сингулярные числа матрицы Σ – это корни из собственных чисел матрицы $X^T X$.

Для снижения размерности берем первые k столбцов матрицы U и верхний $k \times k$ -квадрат матрицы Σ , тогда матрица $U_k \Sigma_k$ содержит k новых признаков, соответствующих первым k главным компонентам.

ПОСТРОЕНИЕ СИНГУЛЯРНОГО РАЗЛОЖЕНИЯ

Ищем сингулярное разложение: $X = U\Sigma V^T$

- Сингулярные числа матрицы Σ – это корни из собственных чисел матрицы $X^T X$

⇒ находим $\lambda_1 \geq \dots \geq \lambda_k$ собственные числа матрицы $X^T X$ и получаем матрицу Σ – у которой на диагонали стоят $\sigma_1 = \sqrt{\lambda_1}, \dots, \sigma_k = \sqrt{\lambda_k}$.

- Столбцы матрицы V – это собственные векторы матрицы $X^T X \Rightarrow$ находим собственные векторы $v_i: (X^T X - \lambda_i I)v_i = 0$.
- Столбцы матрицы $U\Sigma$ – это векторы Xv_1, Xv_2, \dots , т.е.

$$\sigma_i u_i = Xv_i \Rightarrow u_i = \frac{1}{\sigma_i} Xv_i$$

(либо находим $u_i: (XX^T - \lambda_i I)u_i = 0$)

The image features a light gray background with a subtle pattern of concentric circles. In the four corners, there are decorative circuit-like patterns. The top-left and top-right corners have dark blue lines, while the bottom-left and bottom-right corners have light blue lines. These lines form various geometric shapes, including rectangles and circles, resembling a stylized circuit board.

ПОИСК АНОМАЛИЙ

ПОИСК АНОМАЛИЙ С ПОМОЩЬЮ МОДЕЛЕЙ ML

Идея: можно настроить модель машинного обучения так, чтобы на нормальных объектах она принимала значения, близкие к нулю (или, например, положительные значения). Тогда если прогноз на объекте сильно отличается от прогноза на обучающей выборке, то такой объект можно считать аномальным.

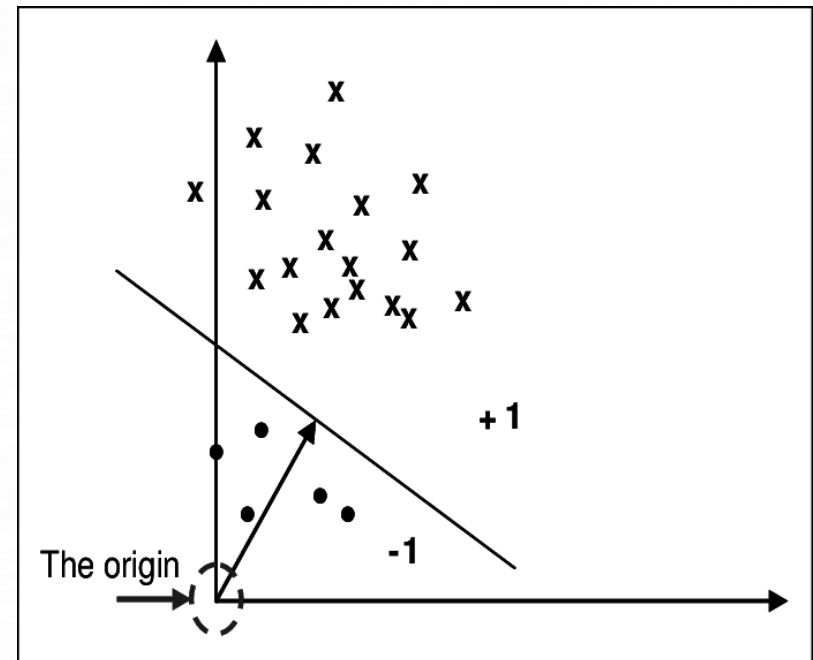
ONE-CLASS SVM

Метод строит линейную функцию $a(x) = \text{sign}(w, x)$ так, чтобы она отделяла выборку от начала координат с максимальным отступом, а именно:

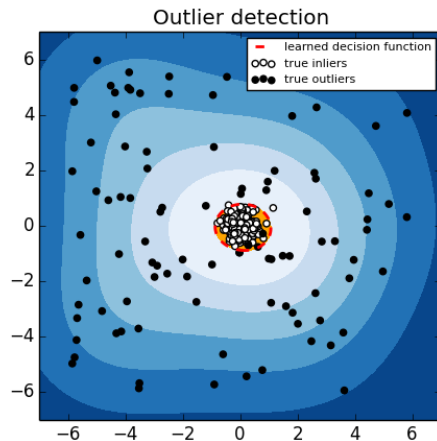
- $a(x)$ отделяет как можно больше объектов выборки от нуля: $a(x) = +1$ на области как можно меньшего объема, содержащей как можно больше объектов выборки
- имеет большой отступ от 0.

Тогда объекты с $a(x) = -1$

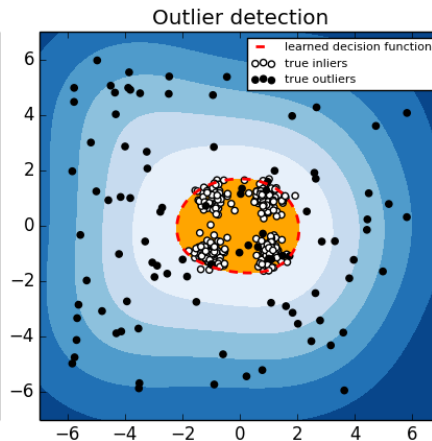
– это аномалии.



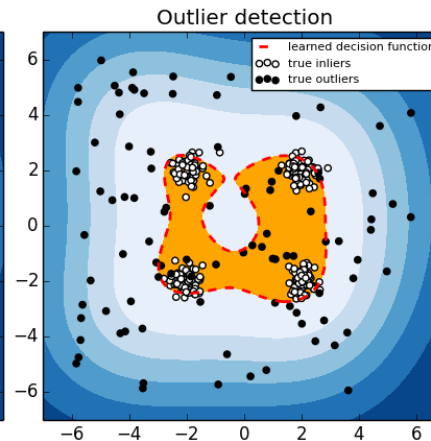
ONE-CLASS SVM С RBF-ЯДРОМ



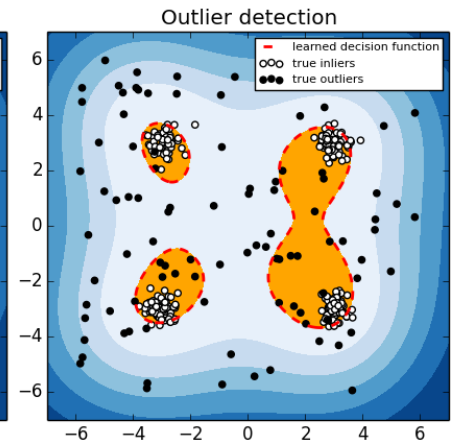
1. one class SVM (errors: 6)



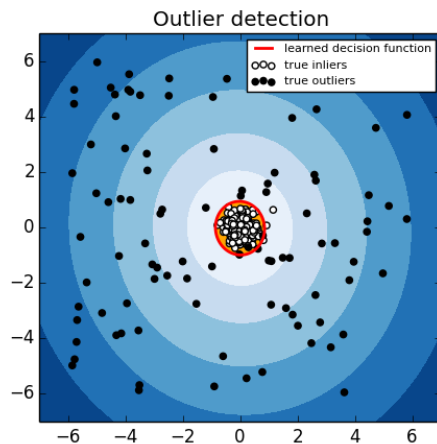
2. one class SVM (errors: 26)



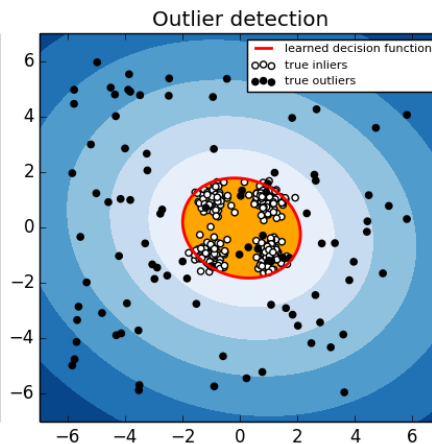
3. one class SVM (errors: 40)



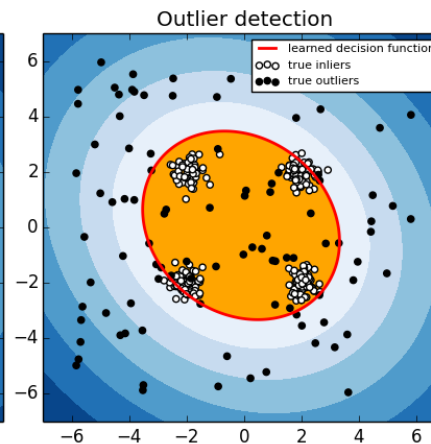
4. one class SVM (errors: 46)



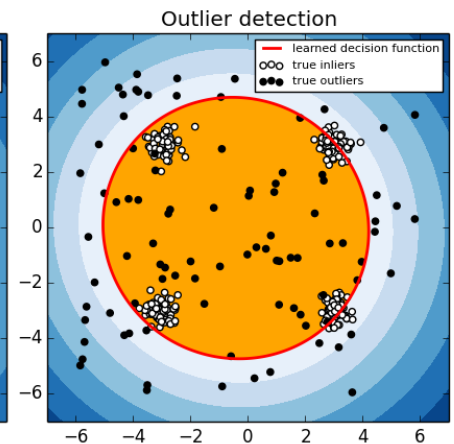
1. covariance estimation (errors: 6)



2. covariance estimation (errors: 26)



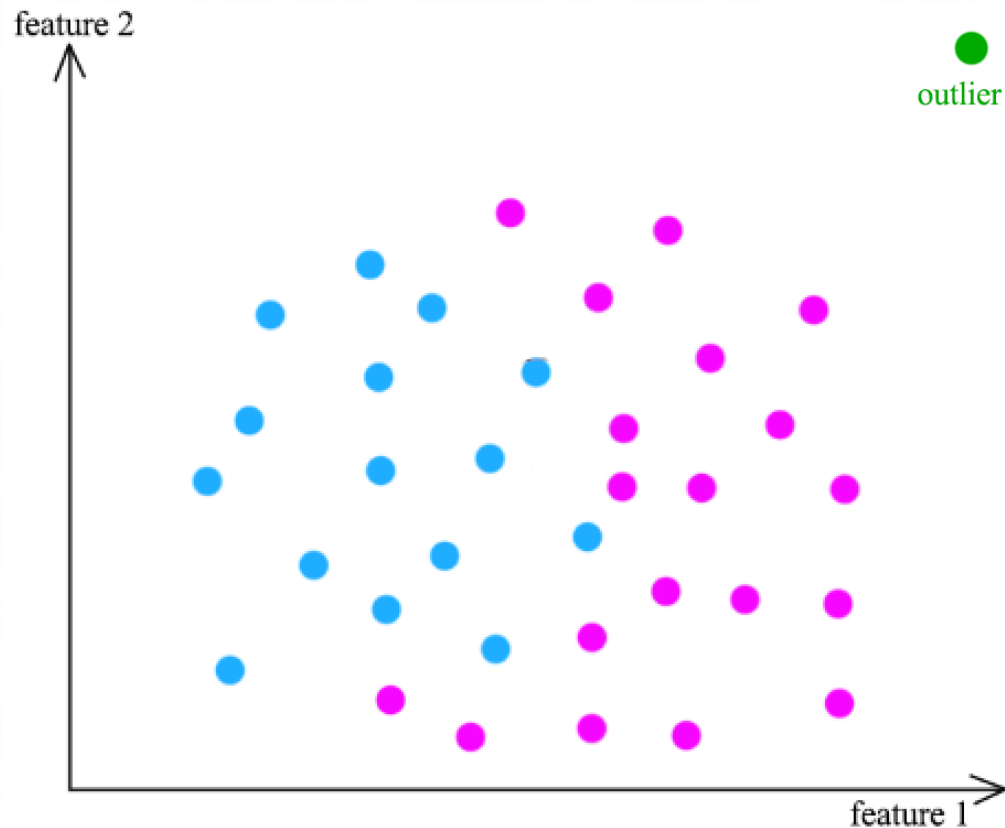
3. covariance estimation (errors: 54)



4. covariance estimation (errors: 98)

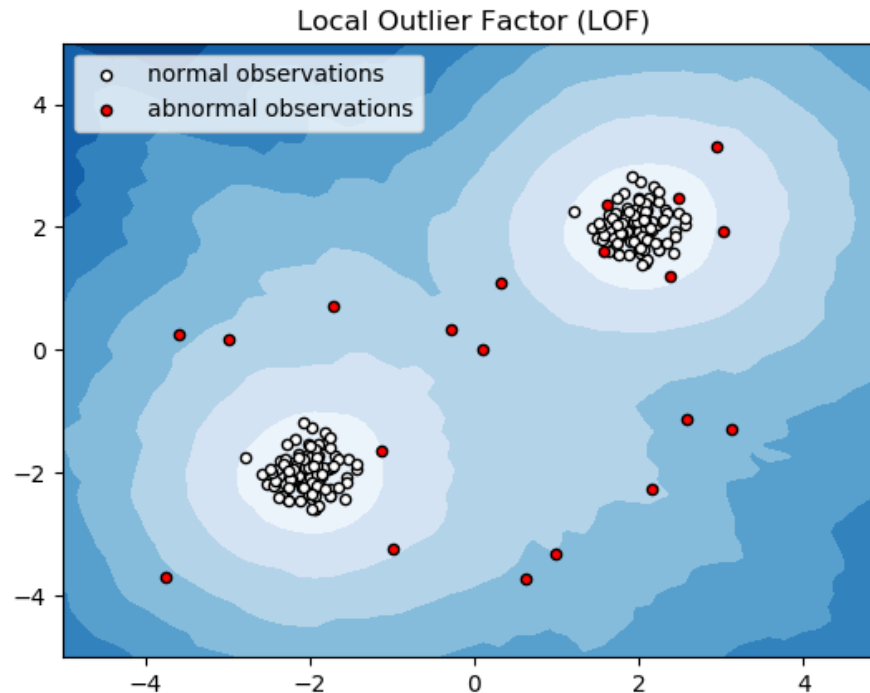
ПОИСК ВЫБРОСОВ С ПОМОЩЬЮ KNN

- Вычисляем среднее расстояние от каждой точки до её ближайших k соседей
- Точки с наибольшим средним расстоянием – выбросы



LOCAL OUTLIER FACTOR

- Задаем плотность распределения в точке, используя k ближайших соседей
- Точки, плотность распределения в которых значительно меньше, чем у соседей – выбросы.



ССЫЛКИ

- <https://dyakonov.org/2017/04/19/поиск-аномалий-anomaly-detection/>
- https://scikit-learn.org/stable/modules/outlier_detection.html
- <https://github.com/yzhao062/pyod>