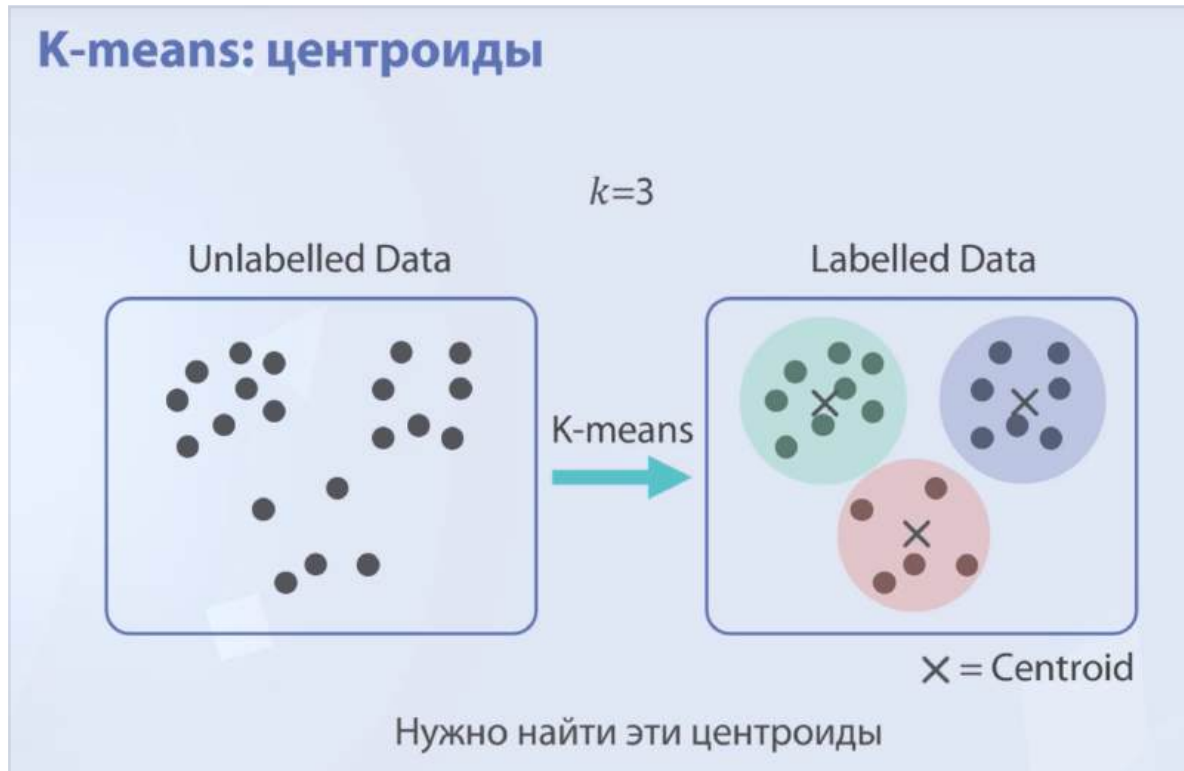


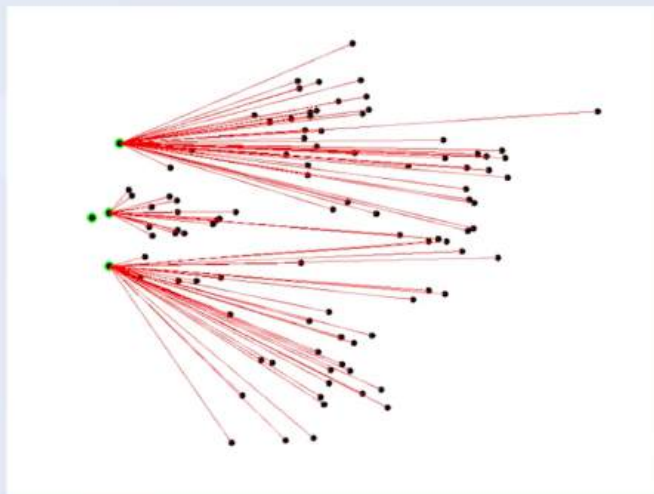
# ML

## K-MEANS

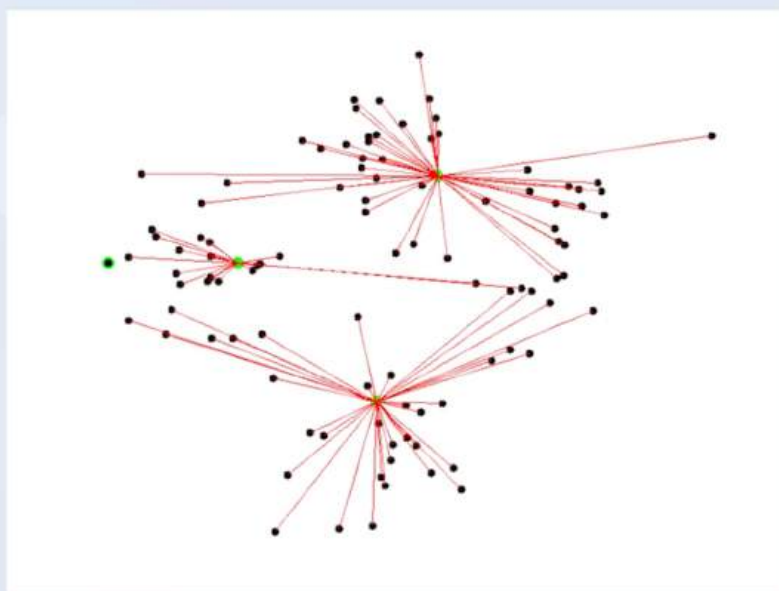
[https://colab.research.google.com/drive/1fg\\_NyOlgcM2ldJuQSBiKLFu1vKpb4qP1?usp=sharing](https://colab.research.google.com/drive/1fg_NyOlgcM2ldJuQSBiKLFu1vKpb4qP1?usp=sharing)



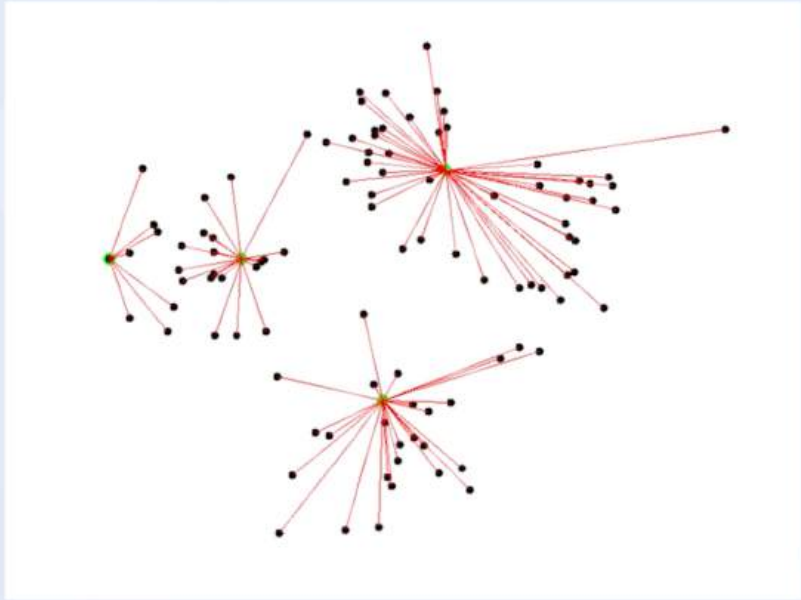
**Для каждой точки находим ближайший центр**



**Пересчитываем центры**



## Опять ищем для каждой точки ближайший центр



- Очень простой алгоритм
- Работает даже на больших данных



- Надо задавать число  $k$  руками
- Не всегда находит кластеры правильно

### DBSCAN

[https://colab.research.google.com/drive/1fg\\_NyOlgcM2ldJuQSBIKLFu1vKpb4qP1?usp=sharing](https://colab.research.google.com/drive/1fg_NyOlgcM2ldJuQSBIKLFu1vKpb4qP1?usp=sharing)

## DBSCAN: идея



Хотим найти скопления точек, не зная заранее количество кластеров.



### Параметры:

$\epsilon$  – размер окрестности точки

**minPts** – минимальное количество точек в  $\epsilon$ -окрестности

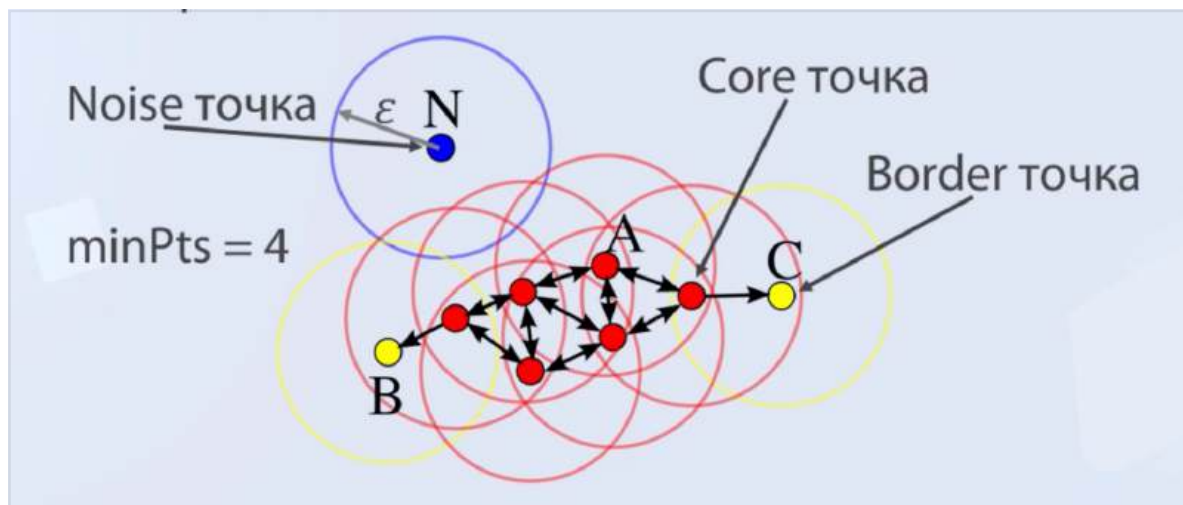


### Все точки делятся на 3 типа:

**Core** точки – от minPts соседей в  $\epsilon$ -окрестности

**Border** точки – не Core точки, но достижимы из Core точек

**Noise** точки – все остальные, меньше minPts соседей в  $\epsilon$ -окрестности



## DBSCAN: алгоритм



Для следующей произвольной точки ищем соседей в  $\varepsilon$ -окрестности

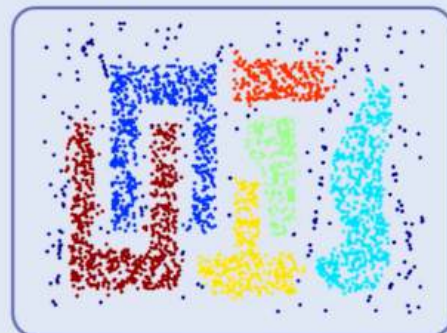
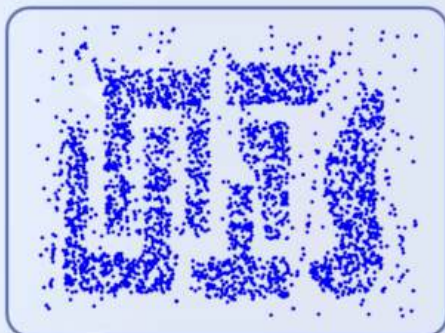


Если их как минимум  $\text{minPts}$ , начинаем поиск связной компоненты из этой Core-точки



Иначе помечаем точку как Noise, она может быть подключена позднее к какой-то Core-точке как Border точка

## DBSCAN: пример



## Резюме: DBSCAN



- Не нужно задавать кол-во кластеров
- Кластеры могут быть любой формы
- Может работать с шумными данными



- Долго работает на больших данных
- Чувствителен к выбору гипер-параметров (лучше использовать HDBSCAN)

## РСА

### РСА: постановка задачи



Пусть у нас есть  $m$  векторов  $x_i \in \mathbb{R}^n$



Мы хотим найти такие вложенные многообразия  $L_k$ , что сумма квадратов расстояний от  $x_i$  до  $L_k$  минимальна:

$L_0 = \{a_0\}$  – *точка*

$L_0 \subset L_1 = \{a_0 + \beta_1 a_1 \mid \beta_1 \in \mathbb{R}\}$  – *линия*

$L_1 \subset L_2 = \{a_0 + \beta_1 a_1 + \beta_2 a_2 \mid \beta_1, \beta_2 \in \mathbb{R}\}$  – *плоскость*

...

где  $a_0 \in \mathbb{R}^n$ , а  $\{a_1, a_2, \dots, a_k\} \subset \mathbb{R}^n$  – ортонормированный набор векторов



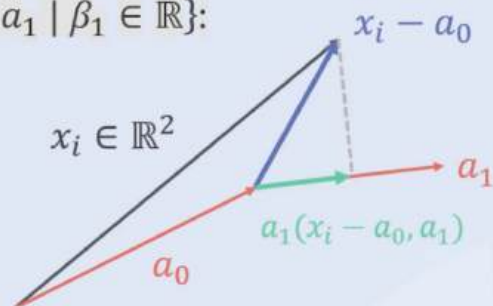
## РСА: расстояние до точки

- Пусть у нас есть  $m$  векторов  $x_i \in \mathbb{R}^n$
- Мы хотим найти такие вложенные многообразия  $L_k$ , что сумма квадратов расстояний от  $x_i$  до  $L_k$  минимальна
- Для точки  $L_0 = \{a_0\}$ :

$$a_0 = \operatorname{argmin}_{a_0 \in \mathbb{R}^n} \left( \sum_{i=1}^m \|x_i - a_0\|^2 \right) = \frac{1}{m} \sum_{i=1}^m x_i$$

## РСА: расстояние до линии

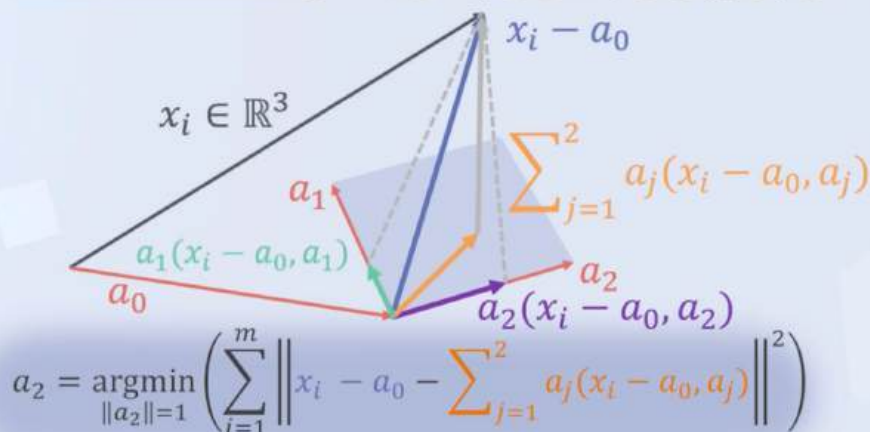
- Пусть у нас есть  $m$  векторов  $x_i \in \mathbb{R}^n$
- Мы хотим найти такие вложенные многообразия  $L_k$ , что сумма квадратов расстояний от  $x_i$  до  $L_k$  минимальна
- Для линии  $L_1 = \{a_0 + \beta_1 a_1 \mid \beta_1 \in \mathbb{R}\}$ :



$$a_1 = \operatorname{argmin}_{\|a_1\|=1} \left( \sum_{i=1}^m \|x_i - a_0 - a_1(x_i - a_0, a_1)\|^2 \right)$$

## РСА: расстояние до плоскости

- Пусть у нас есть  $m$  векторов  $x_i \in \mathbb{R}^n$
- Мы хотим найти такие вложенные многообразия  $L_k$ , что сумма квадратов расстояний от  $x_i$  до  $L_k$  минимальна
- Для плоскости  $L_2 = \{a_0 + \beta_1 a_1 + \beta_2 a_2 \mid \beta_1, \beta_2 \in \mathbb{R}\}$ :



## РСА: расстояние до многообразия

$$a_1 = \operatorname{argmin}_{\|a_1\|=1} \left( \sum_{i=1}^m \|x_i - a_0 - a_1(x_i - a_0, a_1)\|^2 \right)$$

$$a_2 = \operatorname{argmin}_{\|a_2\|=1} \left( \sum_{i=1}^m \left\| x_i - a_0 - \sum_{j=1}^2 a_j (x_i - a_0, a_j) \right\|^2 \right)$$

- В общем виде:

$$a_k = \operatorname{argmin}_{\|a_k\|=1} \left( \sum_{i=1}^m \left\| x_i - a_0 - \sum_{j=1}^k a_j (x_i - a_0, a_j) \right\|^2 \right)$$



## РСА: алгоритм

- 1  $a_0 = \frac{1}{m} \sum_{i=1}^m x_i$
- 2  $x_i := x_i - a_0$
- 3  $a_1 = \operatorname{argmin}_{\|a_1\|=1} (\sum_{i=1}^m \|x_i - a_1(x_i, a_1)\|^2)$
- 4  $x_i := x_i - a_1(x_i, a_1)$
- 5  $a_2 = \operatorname{argmin}_{\|a_2\|=1} (\sum_{i=1}^m \|x_i - a_2(x_i, a_2)\|^2)$  так как  $a_2 \perp a_1$
- 6  $x_i := x_i - a_2(x_i, a_2)$
- 7 Повторяем пока не найдем  $k$  компонент

## РСА: компоненты с максимальной дисперсией

Распишем одно слагаемое в задаче оптимизации:

$$\begin{aligned} \|x_i - a_k(x_i, a_k)\|^2 &= \\ &= \|x_i\|^2 + \|a_k(x_i, a_k)\|^2 - 2(x_i, a_k(x_i, a_k)) = \\ &= \|x_i\|^2 + (x_i, a_k)^2 - 2(x_i, a_k)^2 = \\ &= \|x_i\|^2 - (x_i, a_k)^2 \end{aligned}$$

Перепишем нашу задачу оптимизации:

$$a_k = \operatorname{argmin}_{\|a_k\|=1} \left( \sum_{i=1}^m \|x_i - a_k(x_i, a_k)\|^2 \right) = \operatorname{argmax}_{\|a_k\|=1} \left( \sum_{i=1}^m (x_i, a_k)^2 \right)$$

Так как  $\mathbb{D}\xi = \mathbb{E}\xi^2 - (\mathbb{E}\xi)^2$  и  $x_i$  центрированы в нуле, то мы максимизируем дисперсию длин проекций на  $a_k$

## РСА: компоненты с максимальной дисперсией

Максимизация дисперсии в матричной записи  
( $x_i$  – строки матрицы  $X$ ):

$$\begin{aligned} a_1 &= \operatorname{argmax}_{\|a_1\|=1} \left( \sum_{i=1}^m (x_i, a_1)^2 \right) = \operatorname{argmax}_{\|a_1\|=1} \|Xa_1\|^2 = \\ &= \operatorname{argmax}_{\|a_1\|=1} (Xa_1)^T Xa_1 = \operatorname{argmax}_{\|a_1\|=1} a_1^T X^T X a_1 \end{aligned}$$

Применим метод множителей Лагранжа:

$$\begin{aligned} L &= a_1^T X^T X a_1 - \lambda (a_1^T a_1 - 1) \\ \frac{\partial L}{\partial a_1} &= X^T X a_1 - \lambda a_1 = 0 \end{aligned}$$

То есть  $a_1$  – собственный вектор  $X^T X$ , соответствующий максимальному собственному значению, так как  $a_1^T X^T X a_1 = a_1^T \lambda a_1 = \lambda$

## РСА: компоненты с максимальной дисперсией

Продолжим алгоритм:  $y_i := x_i - a_1(x_i, a_1)$  в матричном виде  $Y := X - Xa_1a_1^T$ , ищем  $a_2 = \operatorname{argmax}_{\|a_2\|=1} a_2^T Y^T Y a_2$ :

$$\begin{aligned} Y^T Y &= (X^T - a_1 a_1^T X^T)(X - Xa_1 a_1^T) = \\ &= X^T X - \textcolor{red}{X^T X a_1 a_1^T} - a_1 a_1^T X^T X + a_1 a_1^T \textcolor{red}{X^T X a_1 a_1^T} = \\ &= X^T X - \textcolor{red}{\lambda_1 a_1 a_1^T} - a_1 a_1^T X^T X + \textcolor{red}{\lambda_1 a_1 a_1^T a_1 a_1^T} = \\ &= X^T X - a_1 a_1^T X^T X \quad \underbrace{\phantom{\lambda_1 a_1 a_1^T a_1 a_1^T}}_1 \end{aligned}$$

$$\begin{aligned} a_2 &= \operatorname{argmax}_{\|a_2\|=1} a_2^T Y^T Y a_2 = \operatorname{argmax}_{\|a_2\|=1} (a_2^T X^T X a_2 - \underbrace{\textcolor{green}{a_2^T a_1 a_1^T X^T X a_2}}_0) = \\ &= \operatorname{argmax}_{\|a_2\|=1, a_2 \perp a_1} a_2^T X^T X a_2, \text{ то есть } a_2 \text{ – следующий с.в. } X^T X \end{aligned}$$

## PCA: расчет через SVD

Любую матрицу  $X$  можно представить в виде SVD разложения:

$$X = U\Sigma V^T$$

$$UU^T = U^T U = I$$

$$VV^T = V^T V = I$$

Распишем

$$X^T X = V\Sigma^T U^T U \Sigma V^T = V\Sigma^2 V^T$$

$$X^T X V = V\Sigma^2$$

То есть в столбцах матрицы  $V$  лежат собственные вектора  $X^T X$ , отсортированные по убыванию собственных значений.

Это и есть наши главные компоненты.

Проекция на главные компоненты:  $XV = U\Sigma$

Вернуться в исходное пространство

(линейная комбинация столбцов  $V$ ):  $U\Sigma V^T = X$



- Быстрый способ уменьшения размерности
- Позволяет визуализировать многомерные данные (при малом  $k$ )



- Новые координаты для точки – *линейные* комбинации исходных координат ( $XV = U\Sigma$ )

t-SNE

## t-SNE: постановка задачи



Пусть у нас есть  $m$  векторов  $x_i \in \mathbb{R}^n$



Хотим снизить размерность каждого вектора  $x_i$ , сохранив информацию о расстоянии до остальных векторов:

1

Зададим распределение пропорциональное расстоянию до остальных точек:

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

2

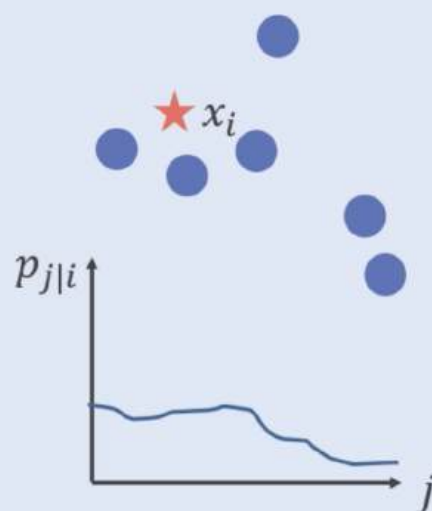
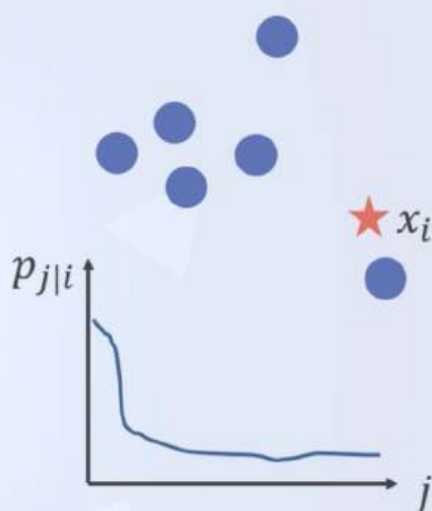
Сделаем его симметричным:

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2}$$

3

Пока что положим  $\sigma_i \equiv \sigma$  (какая-то константа)

## t-SNE: как выглядит $p_{j|i}$ с $\sigma_i \equiv \sigma$



$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$



## t-SNE: уменьшение размерности

- Нужно найти такие  $y_i \in \mathbb{R}^2$  (или  $\mathbb{R}^3$ ), соответствующие  $x_i$ , чтобы распределение  $q_{ij}$  расстояний между  $y_i$  было похоже на распределение  $p_{ij}$ :

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

- Похожесть распределений будем считать при помощи KL-дивергенции:

$$KL(P||Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

- Будем минимизировать  $KL(P||Q)$  по  $y_i$  градиентным спуском:

$$\frac{\partial KL}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j)(1 + \|y_i - y_j\|^2)^{-1}$$

## t-SNE: выбор $\sigma_i$

Для визуализации лучше сохранить информацию о  $k$  ближайших соседях каждой точки  $x_i$

Для этого нужно сделать энтропию распределений

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)} \text{ примерно одинаковой для всех } i$$

Подберем такие  $\sigma_i$ , чтобы  $\text{Perp}(P_i) = 2^{-\sum_j p_{j|i} \log p_{j|i}}$  была примерно одинаковой для всех  $i$

Этого можно добиться бинарным поиском каждой  $\sigma_i$

$\text{Perp}(P_i)$  становится гиперпараметром метода (от 5 до 50)



## t-SNE: визуализация MNIST



## t-SNE: визуализация ImageNet





- Нелинейные преобразования для визуализации данных

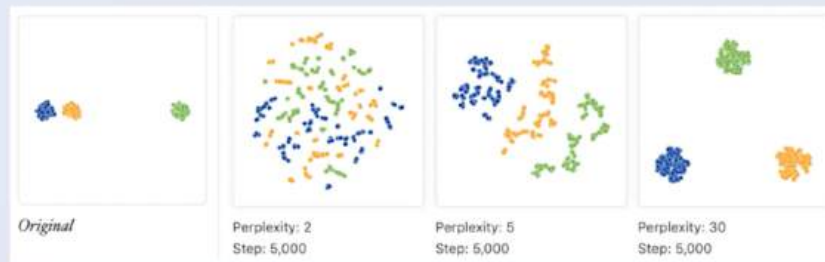
- Долго работает

(<https://github.com/DmitryUlyanov/Multicore-TSNE>)



- Нельзя применить к новой точке (нужно пересчитывать с нуля все точки вместе)

- Неочевидно как выбирать гиперпараметр (перплексия)



для практики <https://colab.research.google.com/drive/18j8sNLcqTscjn6W--LSCHb0DrUofYNQj?usp=sharing>