

**«Машинное обучение»**

# **Обучение без учителя**

**Александр Дьяконов**

**22 марта 2022 года**



## План

### Задачи обучения без учителя

#### Понижение (сокращение) размерности / Вложение в поверхности (Manifold Learning)

~~SVD, PCA, kernel PCA~~

~~LLE (Locally Linear Embedding)~~

~~SNE (Stochastic Neighbor Embedding)~~

~~t-SNE (t-distributed Stochastic Neighbor Embedding)~~

~~IsoMap (Isometric Mapping)~~

~~MDS (MultiDimensional Scaling)~~

~~Maximum Variance Unfolding~~

~~Spectral Embedding / Laplacian Eigenmap~~

~~ICA~~

## Обучение без учителя (Unsupervised Learning)

**Главный вопрос исследователя – «как всё устроено?»**

**вход: неразмеченные данные**

**выход: описание структуры данных / упрощение данных / объяснение данных**

**Неформально: понимание, как данные устроены**

Задачи обучения без учителя (Unsupervised Learning)

<div>Кластеризация (Clustering) <b>отдельная лекция</b></div> <div>нахождение таргетированных групп для акций, планирование эксперимента, визуализация, определение мутаций вируса</div>	
<div>Поиск аномалий <b>отдельная лекция</b></div> <div>обнаружение фрода, поломок, инсайдеров и т.п.</div>	
<div>Тематические модели (Topic discovery)</div> <div>аналог кластеризации в текстах</div>	<p>Жил <u>старик</u> со своею <u>старухой</u>. У самого <u>синего моря</u>. Они жили в <u>ветхой землянке</u>. Ровно тридцать лет и три года. <u>Старик</u> ловил <u>неводом</u> <u>рыбу</u>. <u>Старуха</u> <u>пряла</u> свою <u>пряжу</u>. Раз он в <u>море</u> <u>закинул невод</u>. Пришел <u>невод</u> с одною <u>тиной</u>.</p>
<div>Оценка плотностей (Density Estimation)</div> <div>понимание вероятностной природы данных</div>	
<div>Генерация данных (Data Generation)</div> <div>создание объектов генеральной совокупности</div>	

Задачи обучения без учителя (Unsupervised Learning)

<p><b>Снижение размерности (Dimensionality Reduction)</b> предобработка данных, сжатие информации, устранение избыточности</p>	
<p><b>Вложение в поверхности (Manifold Learning)</b> понимание структурной природы данных, визуализация</p>	
<p><b>Получение представлений (Representation Learning)</b> представление сложных объектов в простом пространстве</p>	
<p><b>Устранение шума (noise reduction)</b> повышение качества данных</p>	
<p><b>Заполнение пропусков / дополнение матриц (matrix completion)</b> матричная факторизация</p>	
<p><b>Поиск ассоциативных правил (association rule learning)</b> <b>отдельная лекция</b></p>	<p><b>A → B</b></p>



## Обучение без учителя – причины

- **неразмеченные данные проще получить**
  - **методы USL можно применять до SL**  
(в том числе, для получения новых признаков)  
**при этом нет риска переобучения, т.к. не видим метки,**  
**но можем подглядывать в будущее**
- **⇒ повышение качества / экономия памяти / интерпретация**  
(в том числе, для последующей визуализации)

## Дальше в этой лекции

<p><b>Снижение размерности (Dimensionality Reduction)</b> предобработка данных, сжатие информации, устранение избыточности</p>	
<p><b>Вложение в поверхности (Manifold Learning)</b> понимание структурной природы данных, визуализация</p>	

**всё это полезно для визуализации и генерации новых признаков**

## Понижение (сокращение) размерности

$$X \in \mathbb{R}^{m \times n} \rightarrow Z \in \mathbb{R}^{m \times k}, k < n$$

### подходы:

- **выразимость  $X$  через  $Z$  (м.б. и наоборот)**  
**ех: возможность восстановления (в DL автокодировщики)**
- **сохранение расстояний (или порядка расстояний)**

### **меньше признаков пространство =>**

- **борьба с переобучением**
- **интерпретация**
- **визуализация**
- **скорость работы алгоритмов**
- **автоматическое удаление шума**
- **ниже стоимость признаков пространства**



## Понижение (сокращение) размерности

$$X \in \mathbb{R}^{m \times n} \rightarrow Z \in \mathbb{R}^{m \times k}, k < n$$

**Но отличается от отбора признаков!**  
получаем вообще говоря новую матрицу...

**если признаков слишком много,  
то найдётся случайно коррелирующий с целевым...**  
ex: случайная матрица размера  $n \times n$  п.н. невырождена

## Понижение (сокращение) размерности с помощью SVD

у нас было сингулярное разложение

$$X_{m \times n} = U \Lambda V^T$$

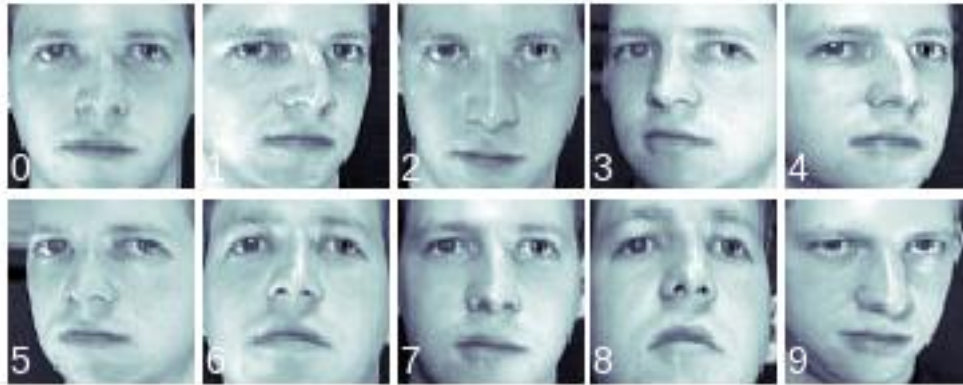
и усечённое сингулярное разложение

$$X_{m \times n} \approx X'_{m \times n} = U[:, 1:k] \cdot \underbrace{\text{diag}(\lambda_1, \dots, \lambda_k)}_{\Lambda[1:k, 1:k]} \cdot V[1:k, :]^T$$

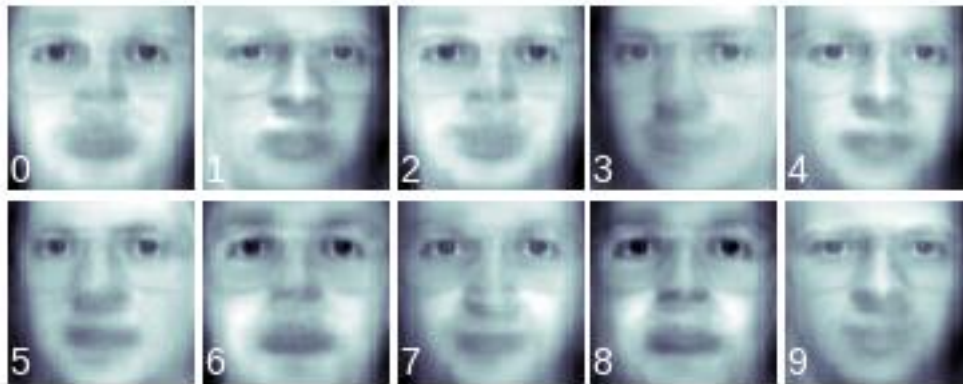
$X'$  лучшее (в каком смысле?) приближение матрицы  $X$   
логично переходить к признаковому пространству  $U[:, 1:k]$

## Эксперименты с лицами «Olivetti faces dataset»

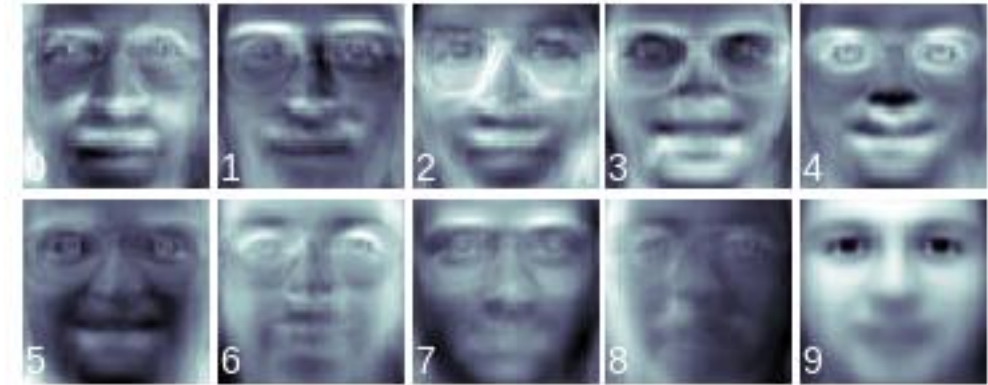
**датасет – 400 картинок 64×64**



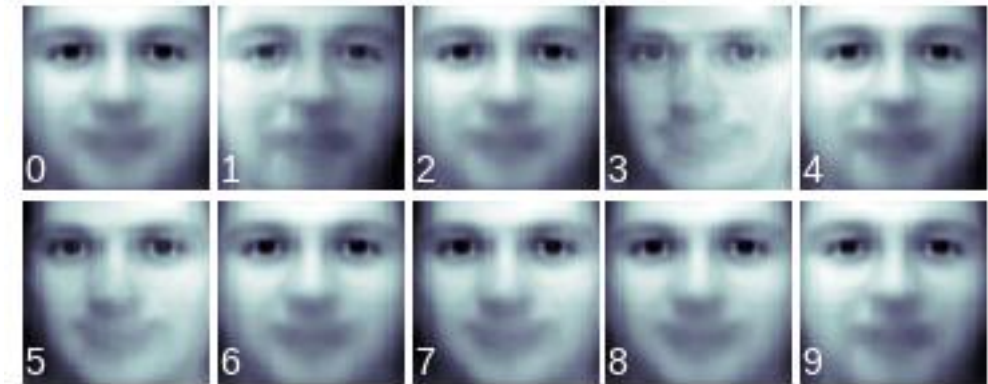
**изображения, восстановленные по  
10 компонентам**



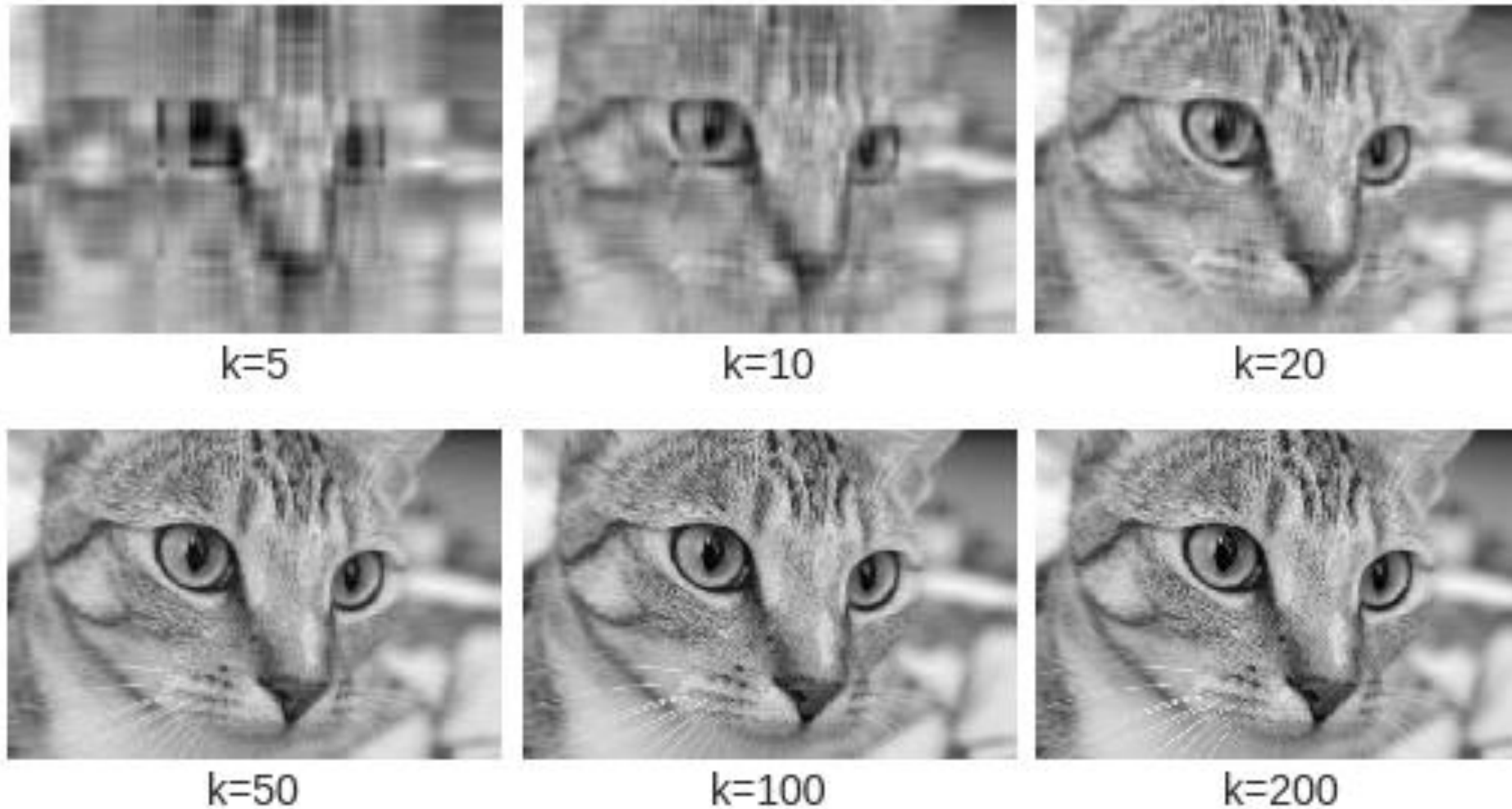
**главные направления**



**изображения, восстановленные по  
2 компонентам**

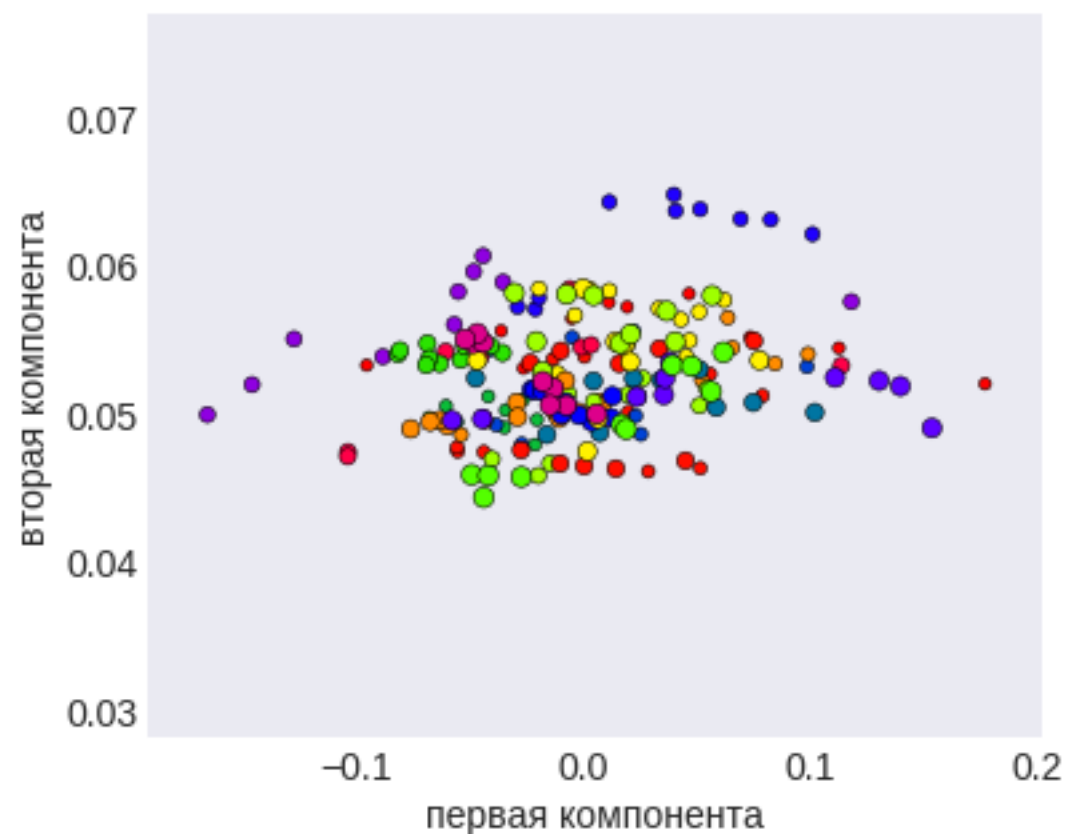


**Вспомним – реконструкция изображений с помощью SVD**  
**это отличается от применения SVD к изображениям, которое было ранее**



**Изначальный размер изображения  $300 \times 451 = 135\,300$**   
 **$300 \times 50 + 50 \times 451 + 50 = 37\,600$**

## Эксперименты с лицами «Olivetti faces dataset»



```
k = 2 # сколько компонент  
from scipy.sparse.linalg import svds  
U, L, V = svds(faces.data, k=k)
```

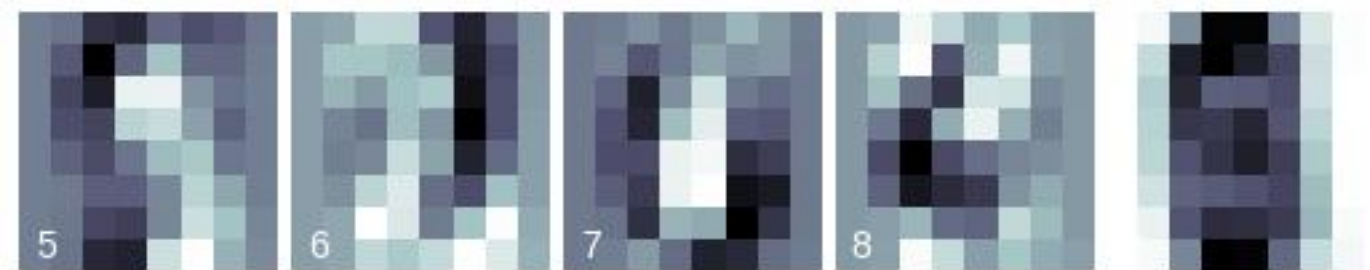
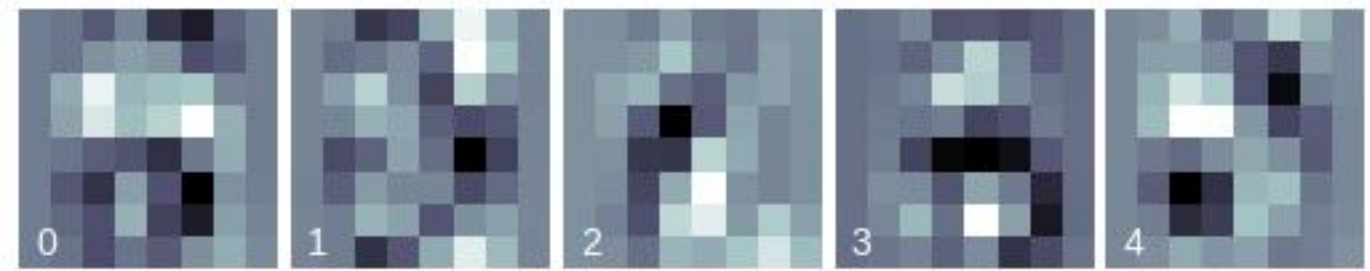
```
x2 = U @ np.diag(L) @ V
```



# Эксперименты с лицами «digits»

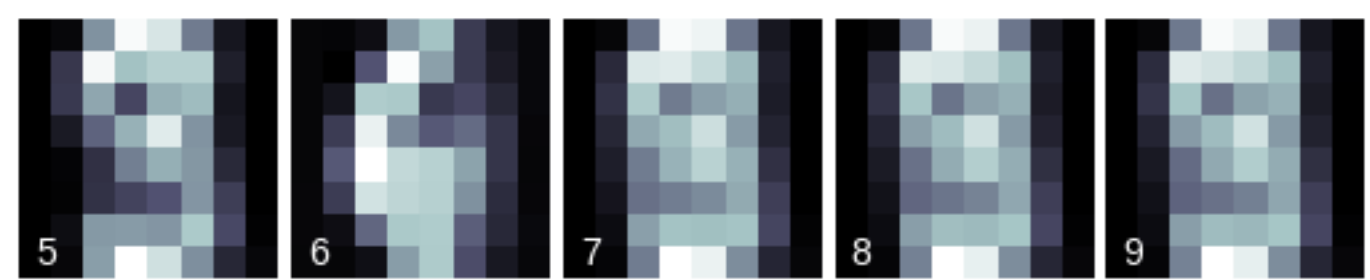
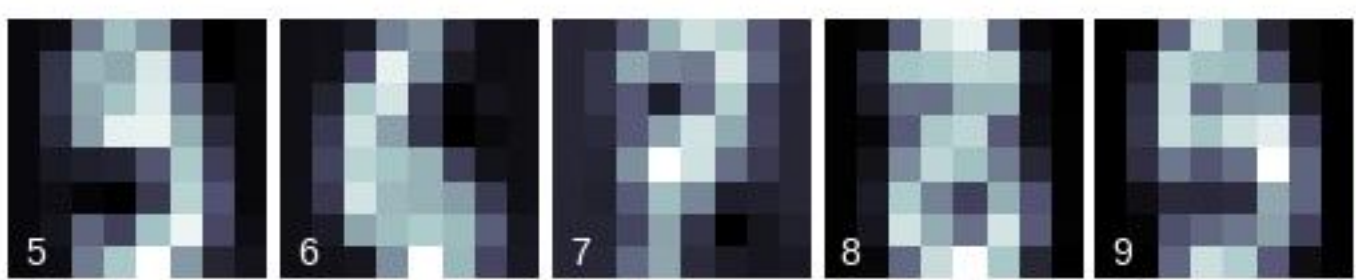
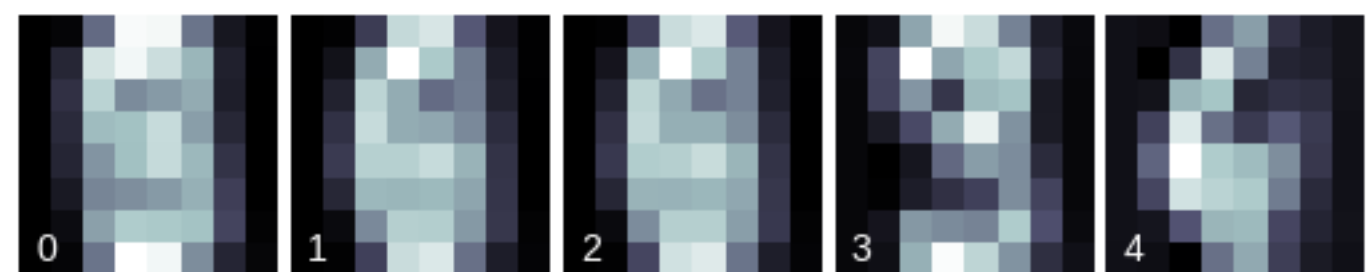
Датасет: 1797 картинок 8×8

главные направления



восстановленные по 10 компонентам

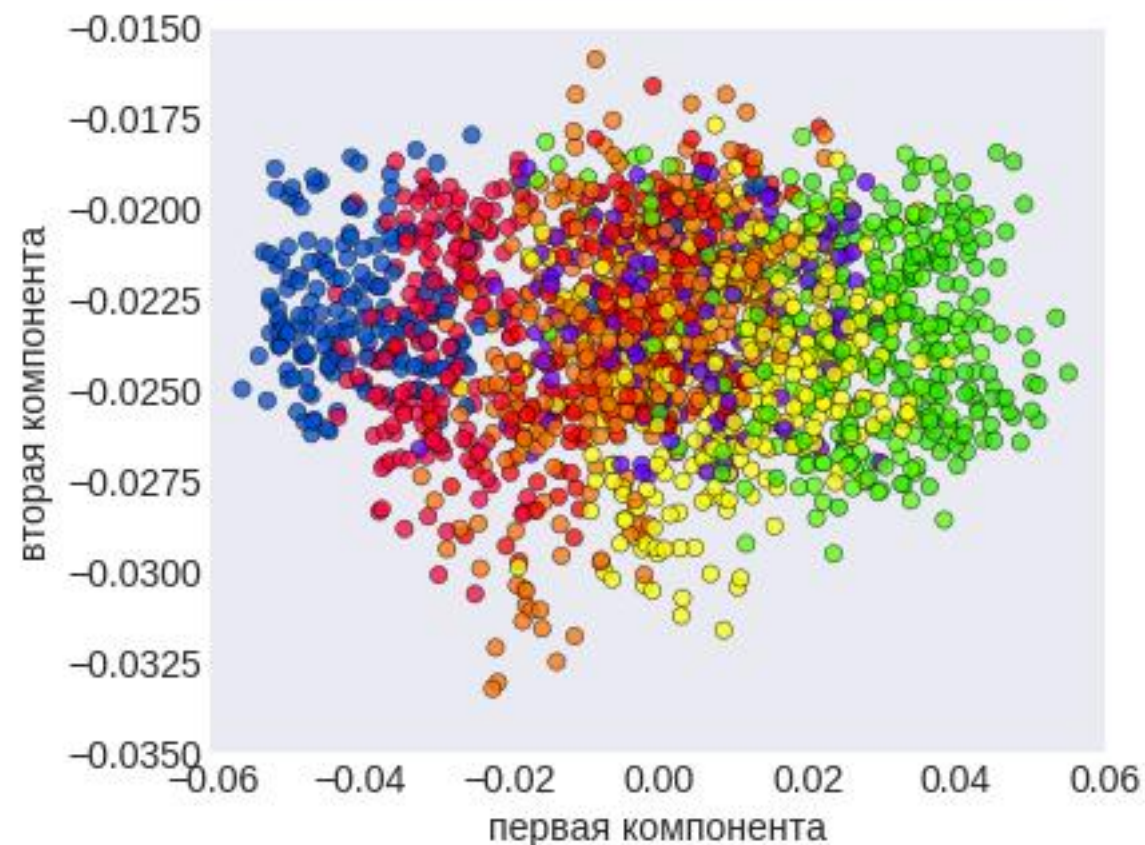
восстановленные по 2 компонентам



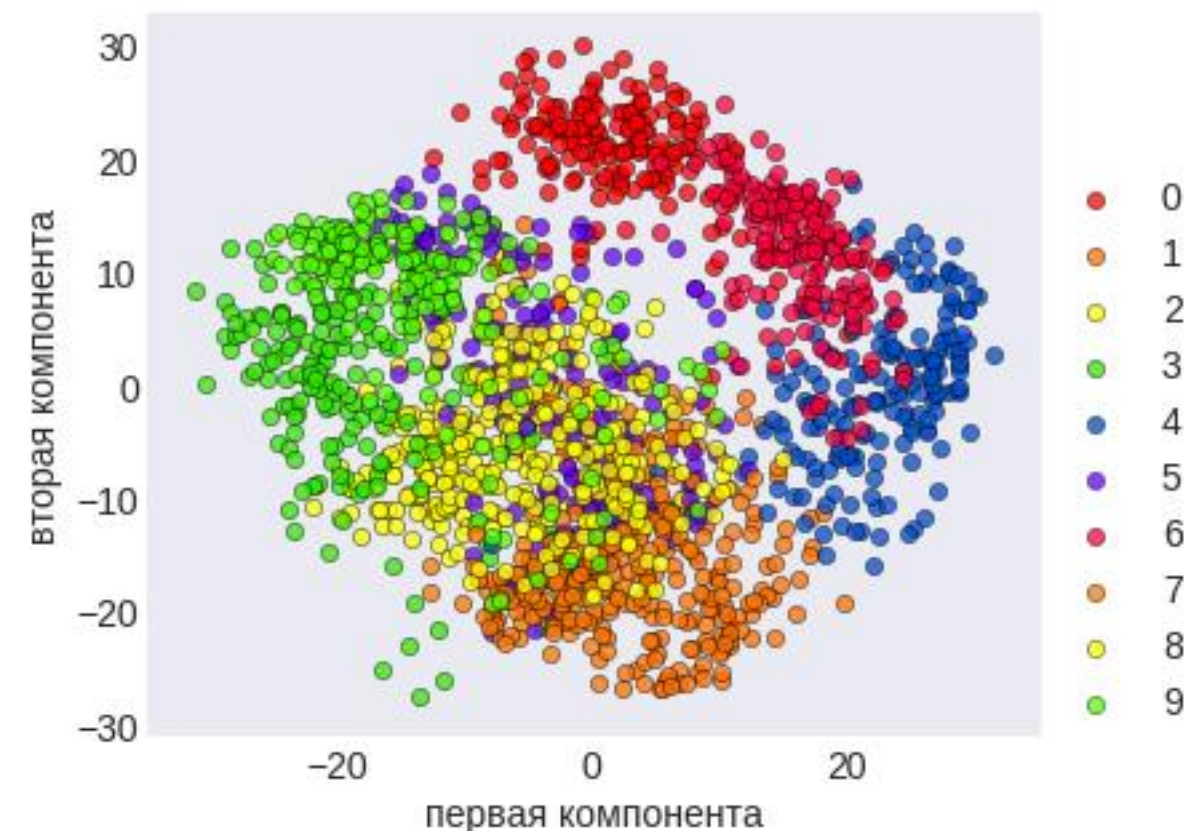


## Эксперименты с лицами «digits»

SVD



PCA



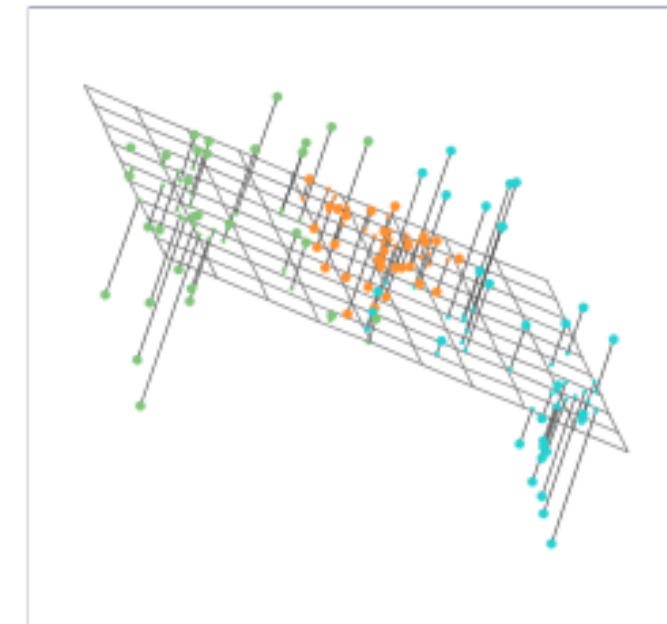
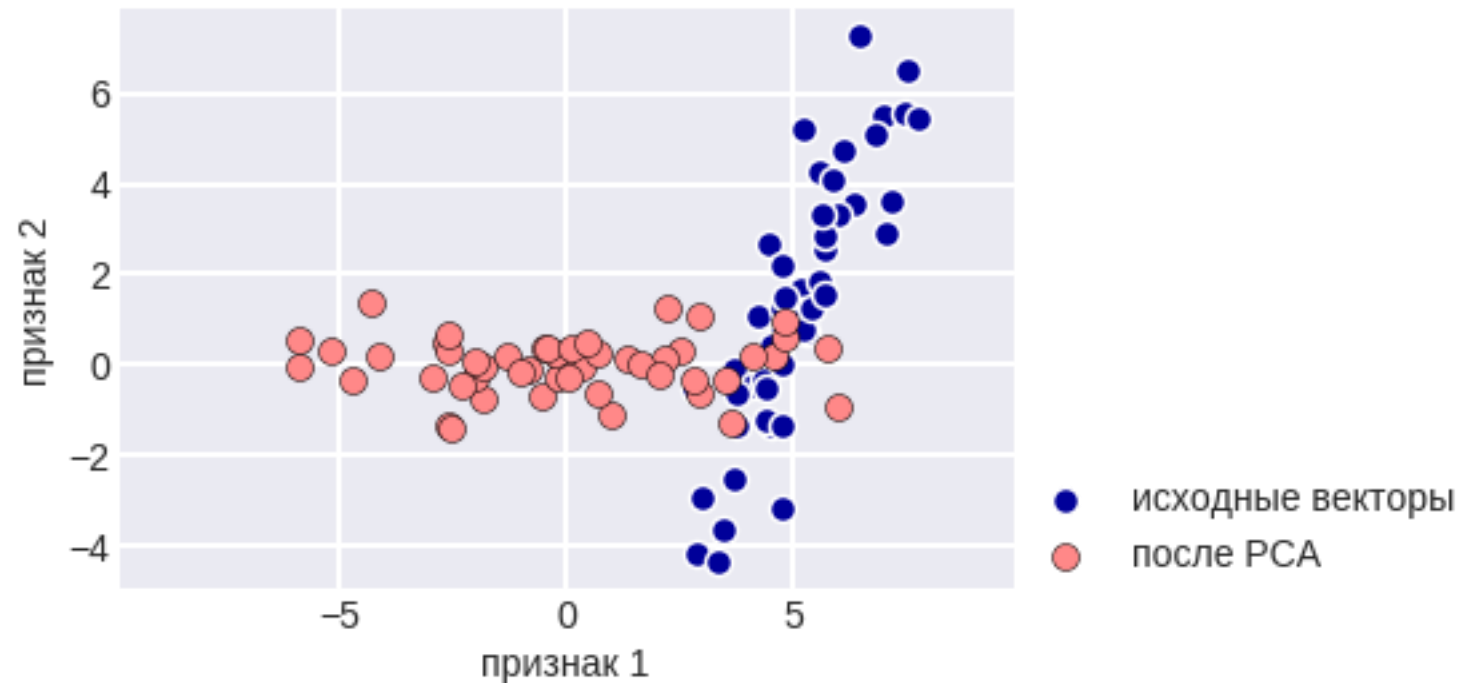
сейчас разберёмся в чём разница

```
from sklearn.decomposition import PCA
estimator = PCA(n_components=10)
X_pca = estimator.fit_transform(X_digits)
```

## Понижение размерности: PCA

### Анализ главных компонент = Principal Component Analysis (PCA)

Представление данных в линейно преобразованном пространстве, если надо – меньшей размерности



Каждый признак нового пространства ищется в виде **линейной комбинации** исходных признаков так, чтобы максимизировать разброс при условии ортогональности (независимости) с уже найденными новыми признаками.

## Понижение размерности: PCA – две интерпретации

ортогональная проекция данных в низкоразмерное пространство, которое

### 1) Maximum Variance Subspace – максимизирует разброс

Находим направление, проекции на которое имеют максимальный разброс

$$z_1^T z_1 = w_1^T X^T X w_1 \rightarrow \max$$

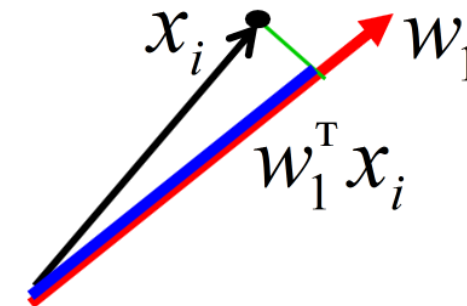
**формулу сейчас поясним**

### 2) Minimum Reconstruction Error – минимизирует MSE (между точками и их проекциями)

Находим направление с минимальной ошибкой восстановления

$$\sum_{i=1}^m \|x_i - (w_1^T x_i) w_1\|^2 \rightarrow \min$$

~ прямая, до которой минимальна сумма  
квадратов расстояний

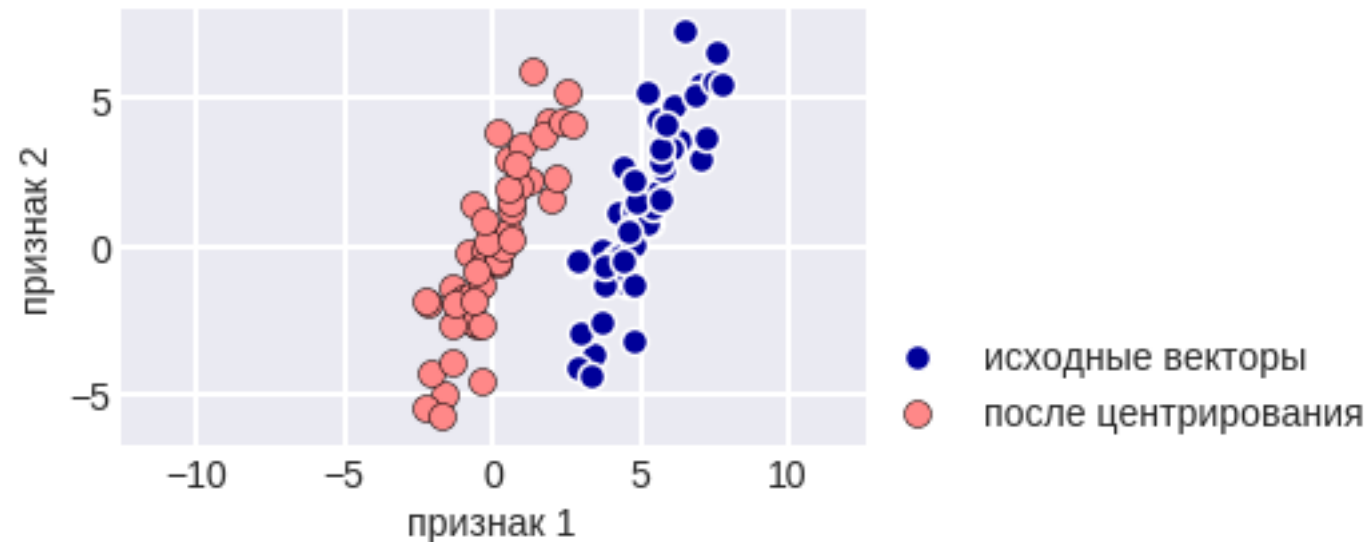


картинка при условии  $\|w_1\| = 1$

## Понижение размерности: PCA

Предполагаем, что все признаки центрированы (**главное отличие от SVD**):

$$\text{mean}(X_i) = 0$$



ищем первый признак в виде (он тоже будет центрированным)

$$Z_1 = w_1 X_1 + \dots + w_n X_n$$

решая задачу (это разброс с учётом центрированности)

$$\frac{1}{m} \sum_{i=1}^m (w_1 x_{i1} + \dots + w_n x_{in})^2 \rightarrow \max, w_1^2 + \dots + w_n^2 = 1$$

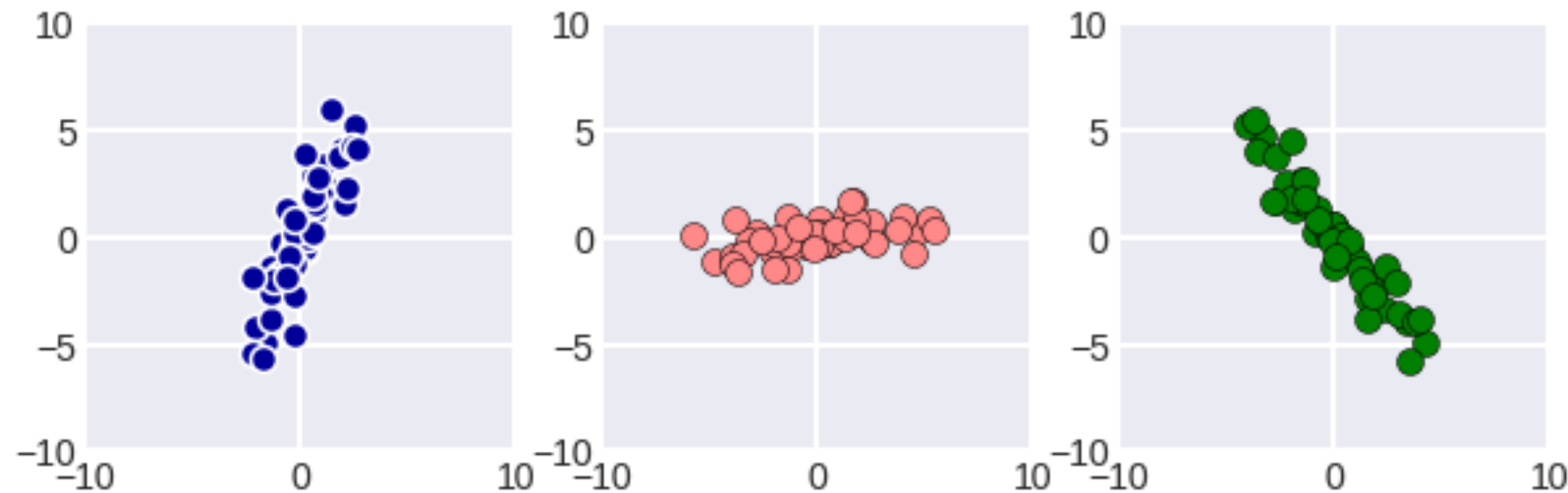
## Понижение размерности: PCA

$$x \rightarrow x^T w_1 = z_1$$

**Матрично... хотим**

$$\|z_1\|^2 \rightarrow \max, \quad z_1 = \frac{1}{\sqrt{m}} X w_1, \quad \|w_1\| = 1$$

$$\begin{cases} z_1^T z_1 = \frac{1}{m} w_1^T X^T X w_1 \rightarrow \max \\ w_1^T w_1 = 1 \end{cases}$$



**ищем удачный поворот... точнее проекцию на ось с max разбросом**  
**«теряется мало информации при переходе к проекции»**

## Понижение размерности: PCA

тут временно избавились от нормирующего множителя

$$J(w) = w^T X^T X w - \lambda(w^T w - 1) \rightarrow \max$$

$$\frac{\partial J}{\partial w} = 2X^T X w - 2\lambda w = 0$$

$$X^T X w = \lambda w$$

**если подставить...**

$$J(w) = \lambda w^T w - \lambda(w^T w - 1) = \lambda$$

**решение – с.в. ~ максимальное с.з. (= разброс в оптимальном решении)**

понятна связь с SVD (**см. дальше**)



## Понижение размерности: PCA

Если искать не один признак, а гиперплоскость, на которую проецируем  
**подробно не доказываем**

$$x^T \rightarrow z^T = x^T V = \begin{pmatrix} v_1^T x \\ \vdots \\ v_k^T x \end{pmatrix}$$

**векторы  $v_1, \dots, v_k$  – с.в., соотв. наибольшим с.з. матрицы ковариаций  $S = \frac{1}{m} X^T X$ :**

$$\lambda_1 \geq \dots \geq \lambda_k \geq \dots$$

**переменные в новом пространстве (координаты вектора  $z$  )  
называются главными компонентами (principal components)**

**а «проекторы»  $v_1, \dots, v_k$  называются главными направлениями / осями  
(principal directions/axes)**

## РСА

- **Вычислить средние**

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

- **Вычислить матрицу ковариаций**

$$S = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^T$$

- **Вычислить  $k$  с.в. соответствующих максимальным  $k$  с.з. матрицы  $S$ :**

$$v_1, \dots, v_k : \lambda_1 \geq \dots \geq \lambda_k$$

- **матрица проекций:  $V = [v_1, \dots, v_k]$**

$$z^T = x^T V = \begin{pmatrix} v_1^T x \\ \vdots \\ v_k^T x \end{pmatrix}$$

**Чаще если**

$$U, L, V = \text{svd}([x_i - \bar{x}]_{i=1}^m)$$

$$X \rightarrow [\lambda_1 u_1, \dots, \lambda_k u_k]$$

эти лямбда корни тех  $\leftarrow$   
(с точностью до константы  $1/m$ )

**Почему:**

Если  $X = ULV^T$  для центрированных данных, то

$$S = X^T X = VL^T U^T ULV^T = VL^2 V^T$$

видим задачу на с.в.:

$$SV = VL^2$$

столбцы  $V$  – главные направления  
столбцы  $UL$  – главные компоненты

**PCA / SVD – другой взгляд: факторизация**

пусть мы не сокращаем размерность,  
а просто проводим линейное преобразование  $\mathbb{R}^n \rightarrow \mathbb{R}^n$ , тогда

$$Z_{m \times n} = X_{m \times n} V_{n \times n}$$

если  $V^T V = V V^T = I$ , то

$$X = Z V^T$$

т.е. это факторизация матрицы  $X$

**Кстати, что означает  $V V^T = I$**

**для вектора  $z^T = x^T V$**

$$\| z \|_2^2 = \| V^T x \|_2^2 = x^T V V^T x = \| x \|_2^2$$

$z = Z[i, :]$                        $x = X[i, :]$

**(тут вектор-строки)**

**т.е. не меняются расстояния  $\Rightarrow$  поворот**

## РСА – другой взгляд: декоррелированность

Рассмотрим центрированные данные, тогда

$$S = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^T = X^T X$$

$$X = \underset{Z}{U} L V^T$$

без сокращения размерности:

$$Z = X V$$

тогда матрица разброса после преобразования (данные тоже будут центрированные)

$$Z^T Z = V^T \textcolor{blue}{X}^T \textcolor{red}{X} V = V^T \textcolor{blue}{V} L^T \textcolor{blue}{U}^T \textcolor{red}{U} L \textcolor{red}{V}^T V = \textcolor{blue}{L}^T L = L^2$$

диагональная матрица, т.е. компоненты полученного вектора не коррелированы  
а разброс ~ диагональные элементы (больше всего у первой координаты и т.д.)

## РСА – другой взгляд: «Minimum Reconstruction Error»

в наших введённых обозначениях, когда мы сначала повернули пространство (с помощью  $V$ ), а потом оставляем  $k$  компонент, ошибка реконструкции

$$\varepsilon = \left\| \underset{Z}{XV - XV|_{\text{zero}}} \right\|_F^2 = \left\| \underset{Z}{XV \mathbf{V}^T - XV|_{\text{zero}} \mathbf{V}^T} \right\|_F^2$$

zero – зануление всех компонент (координат, тут столбцов), начиная с  $k+1$   
 $\mathbf{V}^T$  – от домножения на ортогональную матрицу норма Фробениуса не зависит

тогда

$$\varepsilon = \left\| ULV^T - UL|_{\text{zero}} V^T \right\|_F^2 = \left\| ULV^T - U \text{diag}(\lambda_1, \dots, \lambda_k, 0, \dots, 0) V^T \right\|_F^2$$

$$\varepsilon = \left\| \underbrace{U \text{diag}(0, \dots, 0, \lambda_{k+1}, \dots, \lambda_n)}_{\substack{D \\ H}} V^T \right\|_F^2 = \text{trace}(HH^T) = \text{trace}(UD^2U^T) = \text{trace}(D^2)$$

$$\varepsilon = \lambda_{k+1}^2 + \dots + \lambda_n^2$$

## Понижение размерности: PCA

**Proportion Variance Explained (PVE) / Explained Variance Ratio (EVR)  $k$ -й компоненты –  $\lambda_k^2$**

**Часто смотрят на**

$$\frac{\lambda_1^2 + \dots + \lambda_k^2}{\lambda_1^2 + \dots + \lambda_n^2}$$

**PCA: сколько компонент использовать?**

- можно определить скользящим контролем, если PCA используется для обучения с учителем
  - по графику кумулятивного PVE

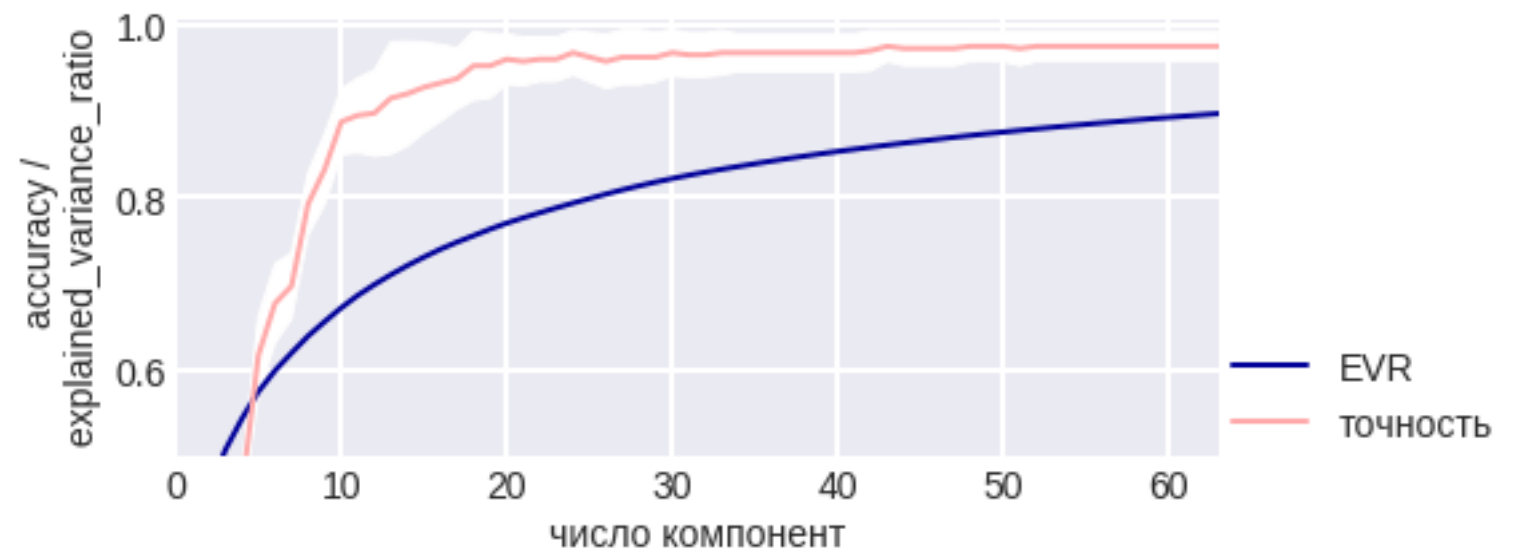


## РСА: сколько компонент использовать?

«digits»



«faces»



показана точность на первых  $k$  компонентах методом логистической регрессии

## РСА: поиск с.в.

### 1. «По определению»

### 2. Итерационный метод для нахождения с.в.

$$w^{(t+1)} = X^T X w^{(t)}$$
$$w^{(t+1)} = w^{(t+1)} / \| w^{(t+1)} \|$$

### 3. EM-алгоритм

см. в [Бишопе] вероятностную трактовку РСА

**РСА: поиск с.в.**

$$S = \frac{1}{m} X^T X$$

**собственные векторы:**  $S = \frac{1}{m} X^T X v_i = \lambda_i v_i$

умножим слева на  $X$

$$\frac{1}{m} \textcolor{red}{X} X^T X v_i = \lambda_i \textcolor{red}{X} v_i$$

$u_i \qquad u_i$

$$\frac{1}{m} X X^T u_i = \lambda_i u_i$$

**Если  $n \gg m$  можно вычислить с.в. и матрицы  $XX^T$**

$$\frac{1}{m} X X^T \sim (\lambda_i, u_i) \quad \Leftrightarrow \quad \frac{1}{m} X^T X \sim (\lambda_i, v_i)$$

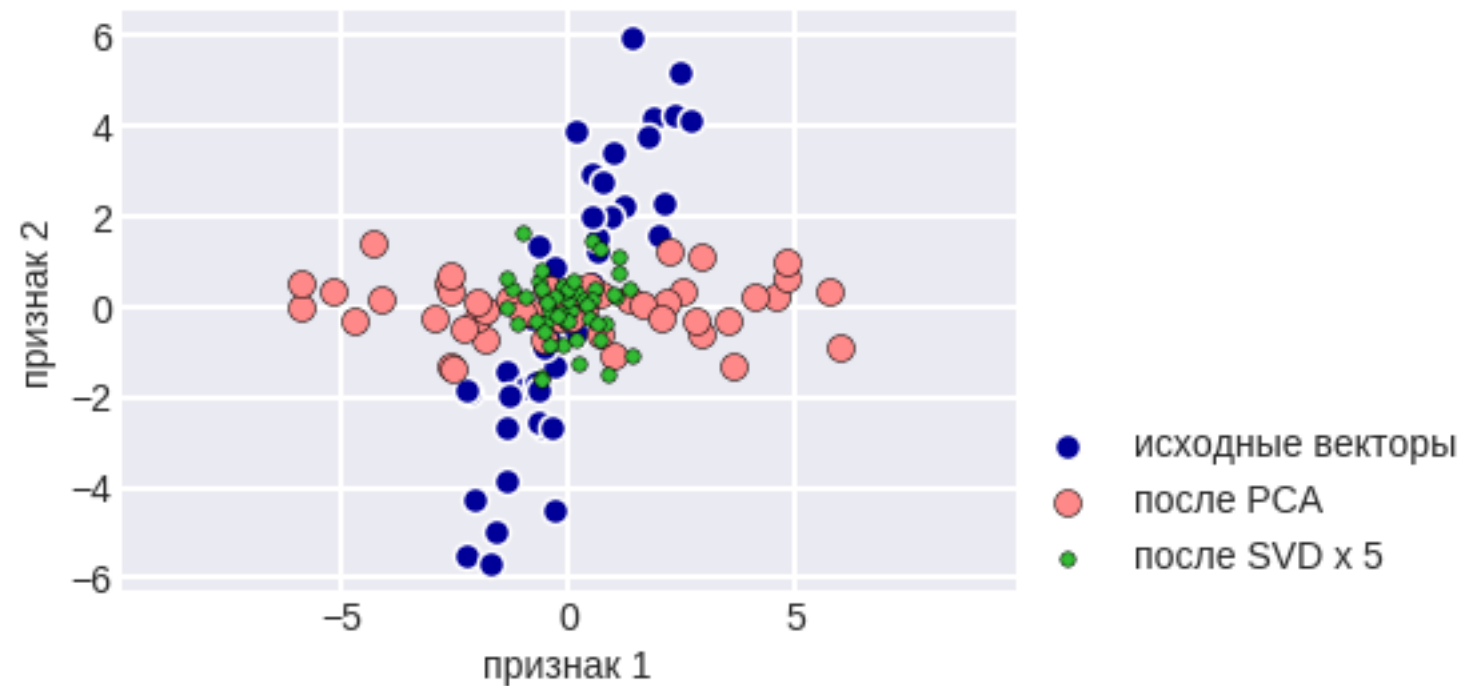
$$v_i = \frac{1}{(m\lambda_i)^{1/2}} X^T u_i$$

**константа из соображения нормировки**  
(подробно не рассказываем)

## Особенности PCA

- PCA зависит от масштаба (стандартизация)
  - PCA чувствителен к выбросам
  - + PCA можно кернализовать...
  - + PCA эквивалентен SVD после централизации данных и  $\lambda$ -масштабирования
    - ⇒ оптимальное линейное преобразование
      - но только линейное
  - + сокращение размерности, можно для больших размерностей
    - 2D может не годиться для интерпретации
  - + новые признаки генерируются с помощью обучения без учителя
    - не подглядываем в целевые значения
  - нет гарантии, что в получаемых признаковых пространствах задача хорошо решается
- ! столбцы в матрицы  $U$  могут быть неоднозначны (с точностью до знака)!**
- могут ли быть другие причины неоднозначности?**

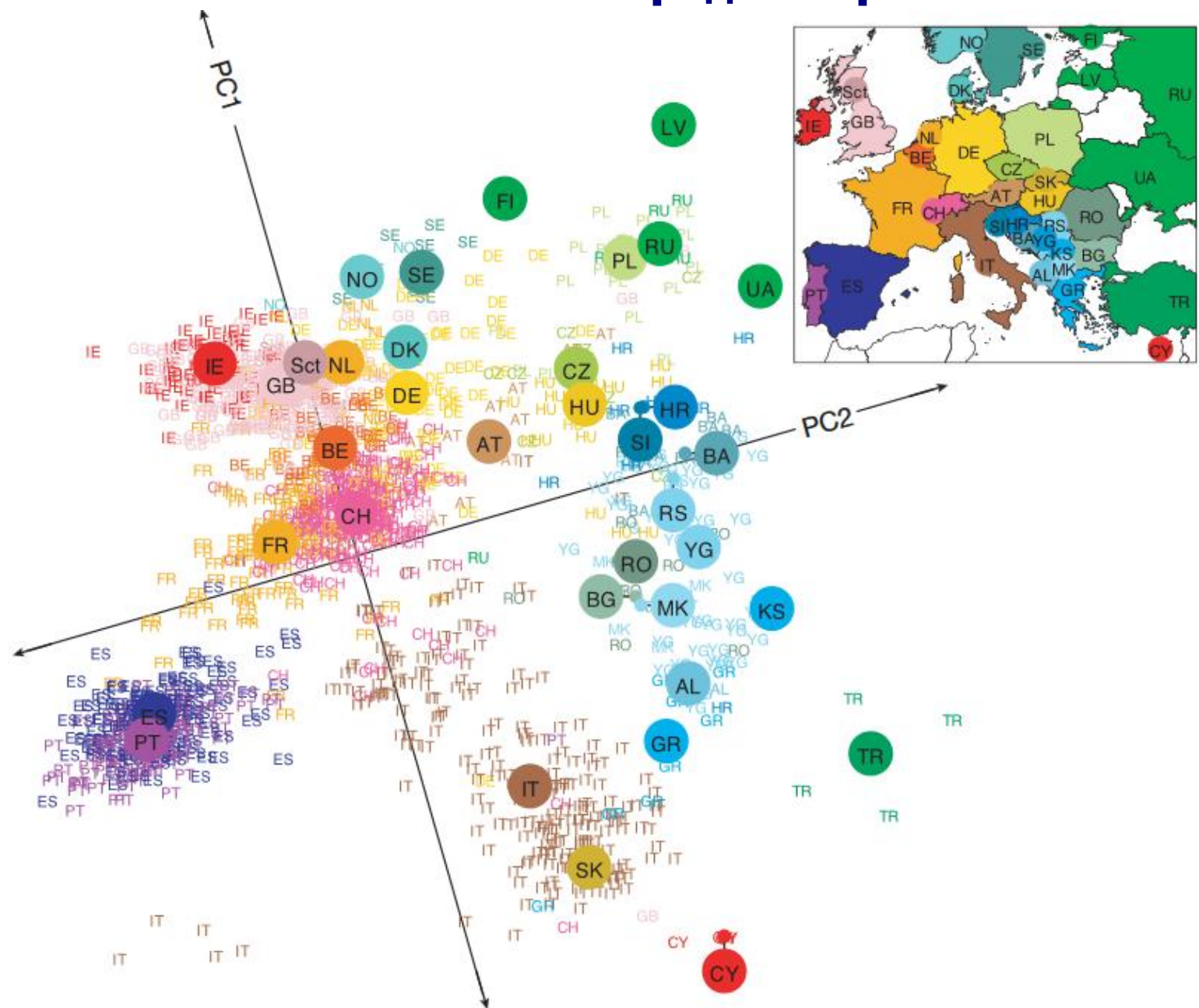
## Минутка кода



```
X = X - X.mean(axis=0)
U, L, V = svd(X)
from sklearn.decomposition import PCA
pca_transformer = PCA()
X2 = pca_transformer.fit_transform(X)
```

```
plt.scatter(X[:, 0], X[:, 1])
plt.scatter(X2[:, 0], X2[:, 1])
# ~ L[0]*U[:, 0], L[1]*U[:, 1]
plt.scatter(5*U[:, 0], 5*U[:, 1])
```

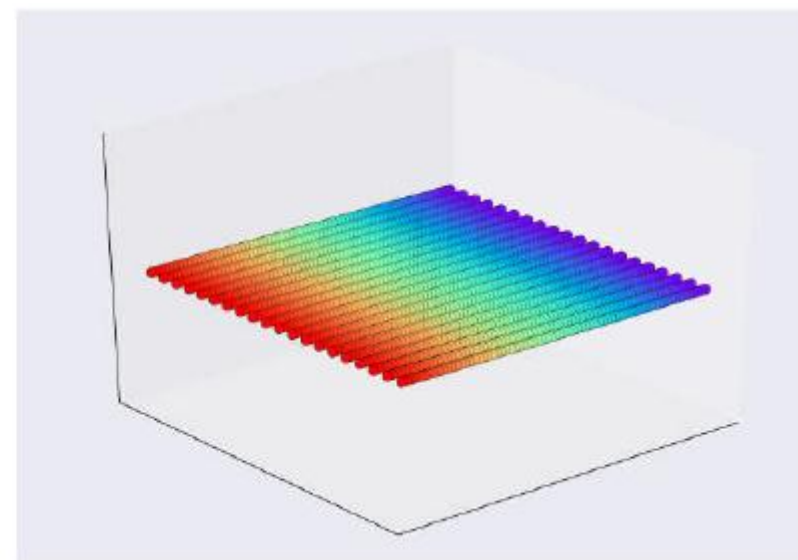
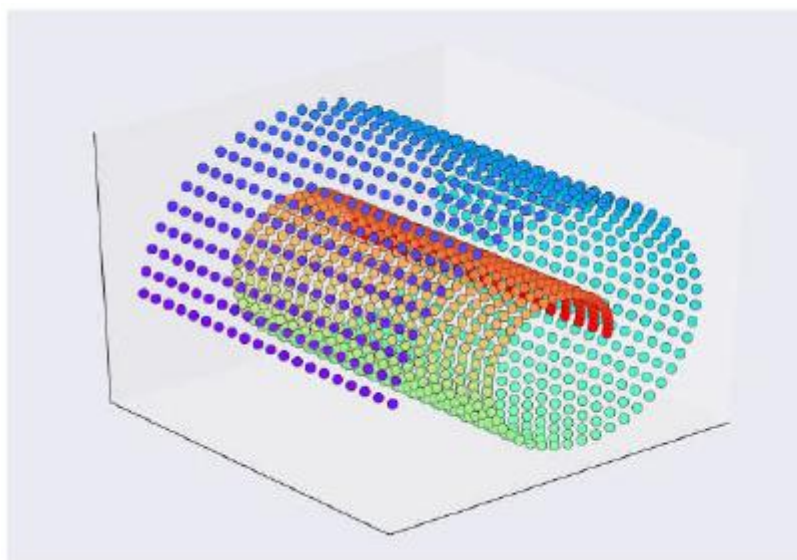
РСА: генотипы народов Европы



<https://www.nature.com/articles/nature07331>



## Нелинейное сокращение размерности



**тема связанная с вложением в поверхности (Manifold Learning)**

**ясно, что метрики здесь не описывают адекватно близость  
мешаются «средние расстояния»**

## Нелинейное сокращение размерности: способы

- kernel PCA

### Manifold based methods

- LLE (Locally Linear Embedding)
- SNE (Stochastic Neighbor Embedding)
- t-SNE (t-distributed Stochastic Neighbor Embedding)
  - IsoMap (Isometric Mapping)
  - MDS (MultiDimensional Scaling)
  - Maximum Variance Unfolding
- Spectral Embedding / Laplacian Eigenmap

### нейросетевые

- autoencoder

## Kernel PCA

**Обычный PCA: после центрирования векторов ищем с.в. матрицы ковариаций (covariance matrix):**

$$S = \frac{1}{m} \sum_{i=1}^m x_i x_i^T$$

**Теперь по аналогии ищем с.в. матрицы:**

$$S = \frac{1}{m} \sum_{i=1}^m \varphi(x_i) \varphi(x_i)^T$$

под суммой стоит не  $K(x_i, x_i)$ , т.к. там не скалярное произведение, а внешнее

$$Su_t = \lambda u_t \qquad \frac{1}{m} \sum_{i=1}^m \varphi(x_i) \varphi(x_i)^T u_t = \lambda u_t$$

**тогда с.в. представимы в виде**

$$u_t = \sum_{j=1}^m a_{jt} \varphi(x_j)$$

## Kernel PCA

**Подставим с.в. в выражение его определения:**

$$\frac{1}{m} \sum_{i=1}^m \varphi(x_i) \varphi(x_i)^T \sum_{j=1}^m a_{jt} \varphi(x_j) = \lambda \sum_{j=1}^m a_{jt} \varphi(x_j)$$

**умножим на  $\varphi(x_s)^T$ , учтём обозначение  $K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$**

$$\frac{1}{m} \sum_{i=1}^m \varphi(x_s)^T \varphi(x_i) \varphi(x_i)^T \sum_{j=1}^m a_{jt} \varphi(x_j) = \lambda \sum_{j=1}^m a_{jt} \varphi(x_s)^T \varphi(x_j)$$

$$\frac{1}{m} \sum_{i=1}^m \varphi(x_s)^T \varphi(x_i) \sum_{j=1}^m a_{jt} \varphi(x_i)^T \varphi(x_j) = \lambda \sum_{j=1}^m a_{jt} \varphi(x_s)^T \varphi(x_j)$$

$$\frac{1}{m} \sum_{i=1}^m K(x_s, x_i) \sum_{j=1}^m a_{jt} K(x_i, x_j) = \lambda \sum_{j=1}^m a_{jt} K(x_s, x_j)$$

## Kernel PCA

Пусть есть матрица  $K = \| K(x_i, x_j) \|_{m \times m}$ , тогда полученное уравнение эквивалентно:

$$K^2 a_t = \lambda m K a_t$$

$$K a_t = \lambda m a_t$$

т.е. получили аналогичную задачу на с.в. ~ kernel PCA

тут не надо вычислять  $\varphi(x_i)$

## Kernel PCA

**Как и в PCA эту матрицу надо нормировать...  
можно показать, что это делается так:**

$$\begin{aligned}\tilde{K} &= (I - E / m) K (I - E / m) = \\ &= K - \left\| \frac{1}{m} \right\| K - K \left\| \frac{1}{m} \right\| + \left\| \frac{1}{m} \right\| K \left\| \frac{1}{m} \right\|\end{aligned}$$

**– центрированная матрица ядра (centered kernel matrix)**

**Пусть  $a_1, \dots, a_k$  – с.в., соответствующие максимальным с.з. матрицы  $\tilde{K}$  тогда**

$$x \rightarrow (z_1, \dots, z_k)$$

$$z_t = \varphi(x)^T u_t$$

$$z_t = \varphi(x)^T u_t = \sum_{j=1}^m a_{tj} \varphi(x)^T \varphi(x_j) = \sum_{j=1}^m a_{tj} K(x, x_j)$$

## kernel PCA

1. Построить  $m \times m$ -матрицу  $K$ .
2. Центрировать, получить матрицу  $\tilde{K}$
3. Найти наибольших с.з. и соответствующие с.в.

$$a_1, \dots, a_k$$

4. Перенормировать их:

$$a'_s = \frac{a_s}{\sqrt{\lambda_s m}}$$

5. Выполнить вложение:

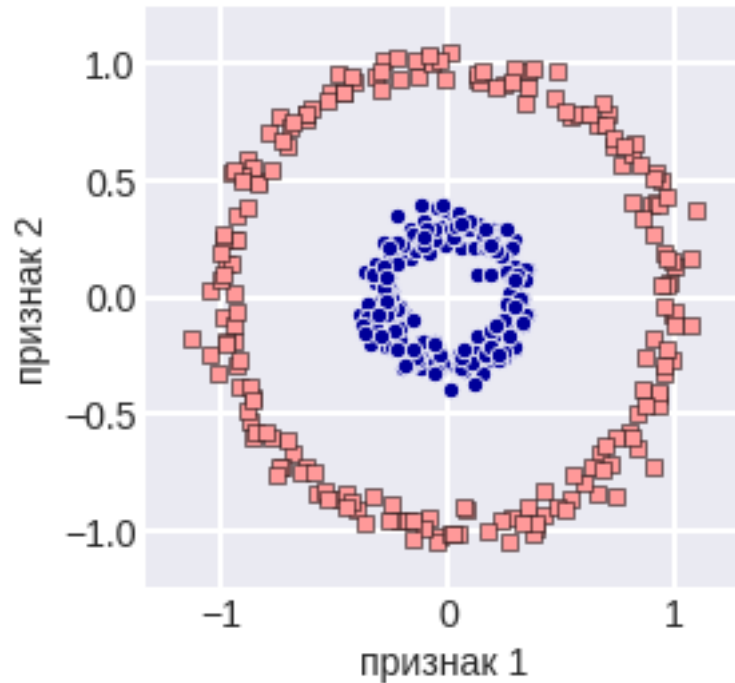
$$x \rightarrow (z_1, \dots, z_k)$$

$$z_t = \sum_{j=1}^m a_{tj} K(x, x_j)$$

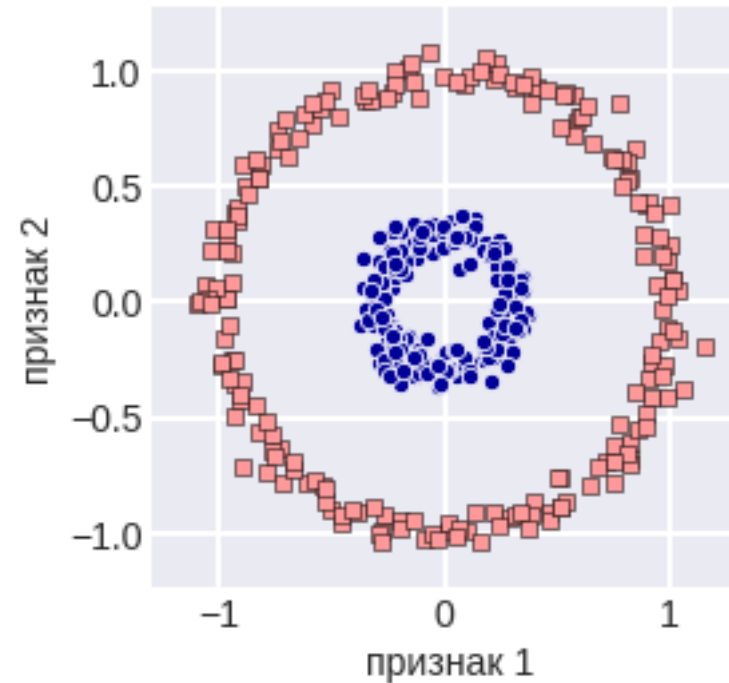


## kernel PCA: примеры

данные



PCA



kernel PCA

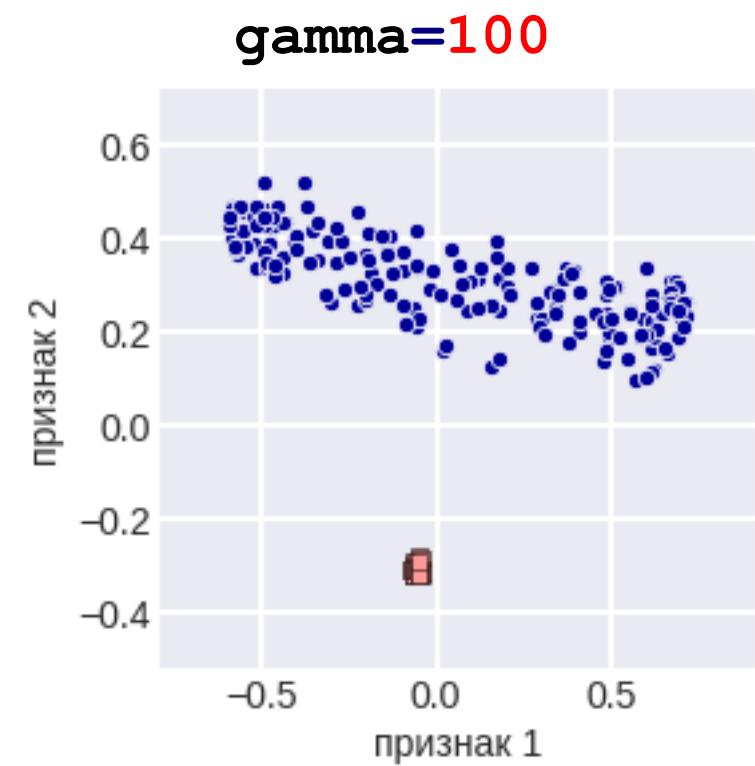
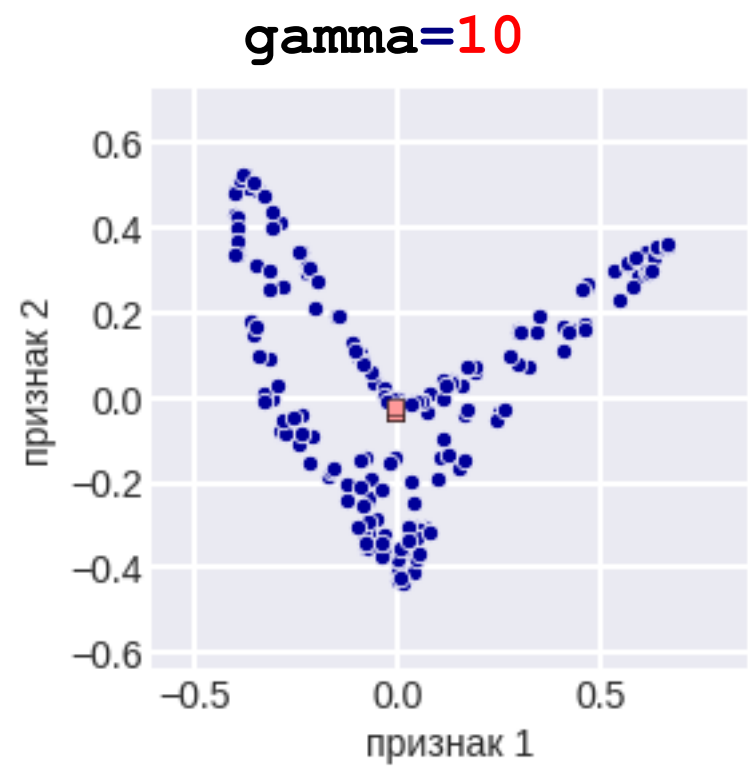
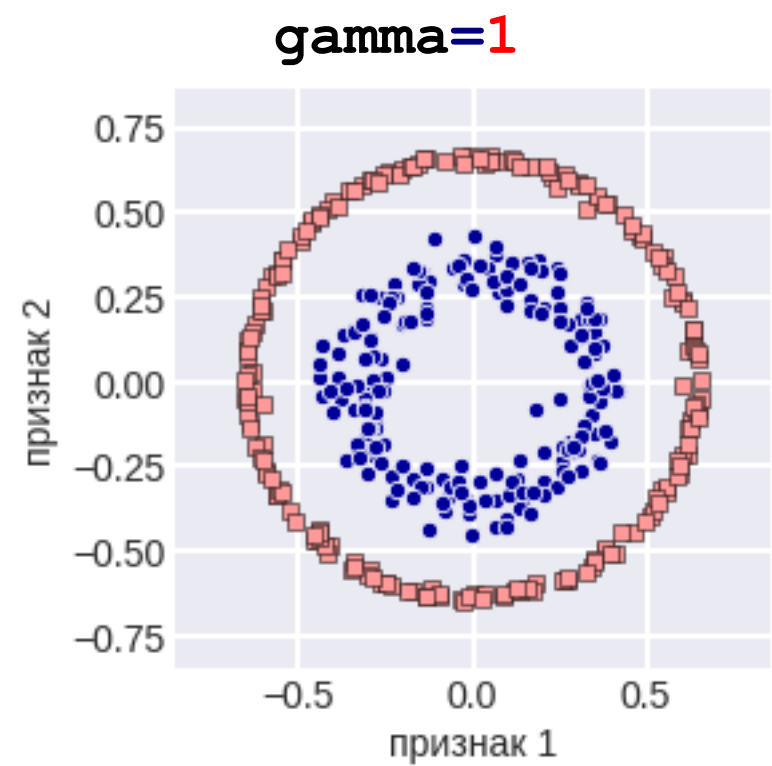


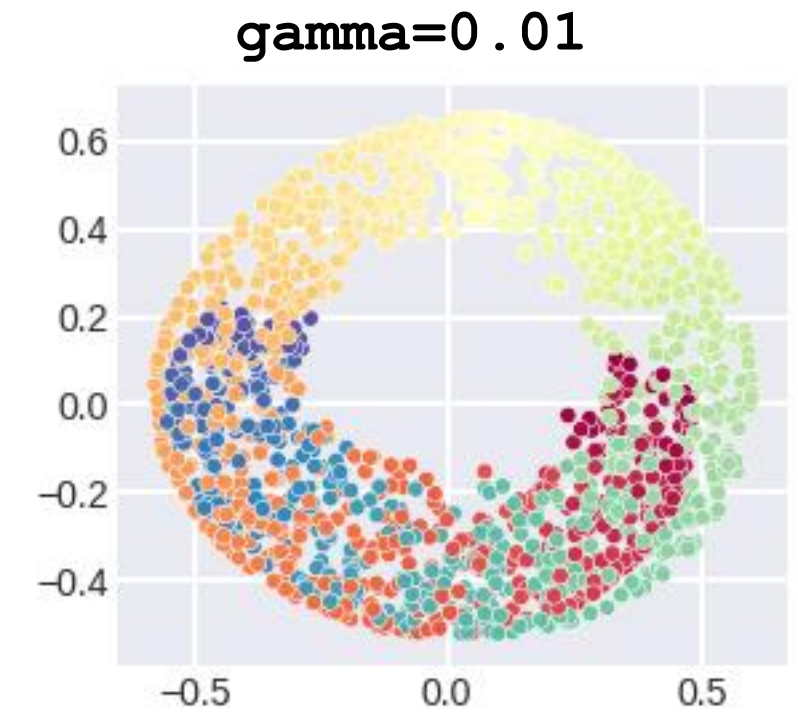
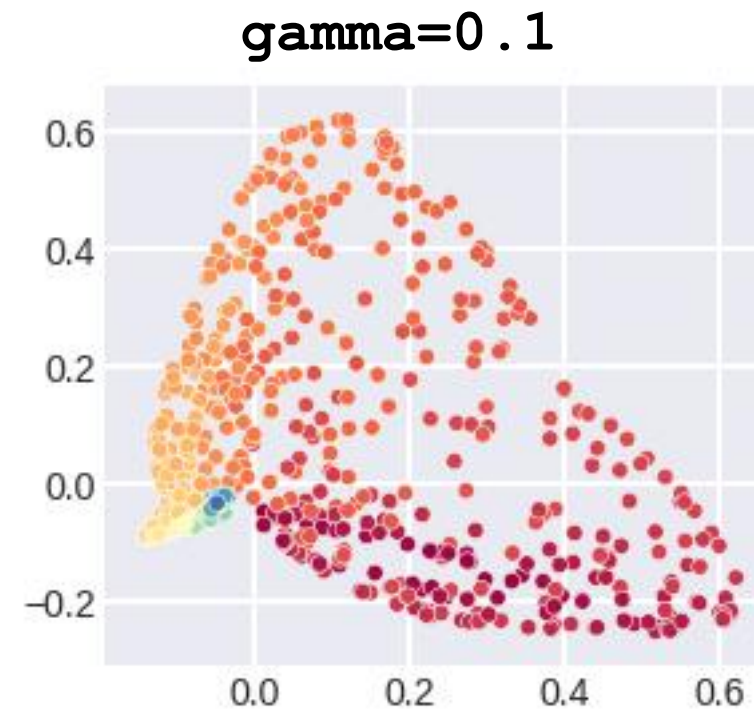
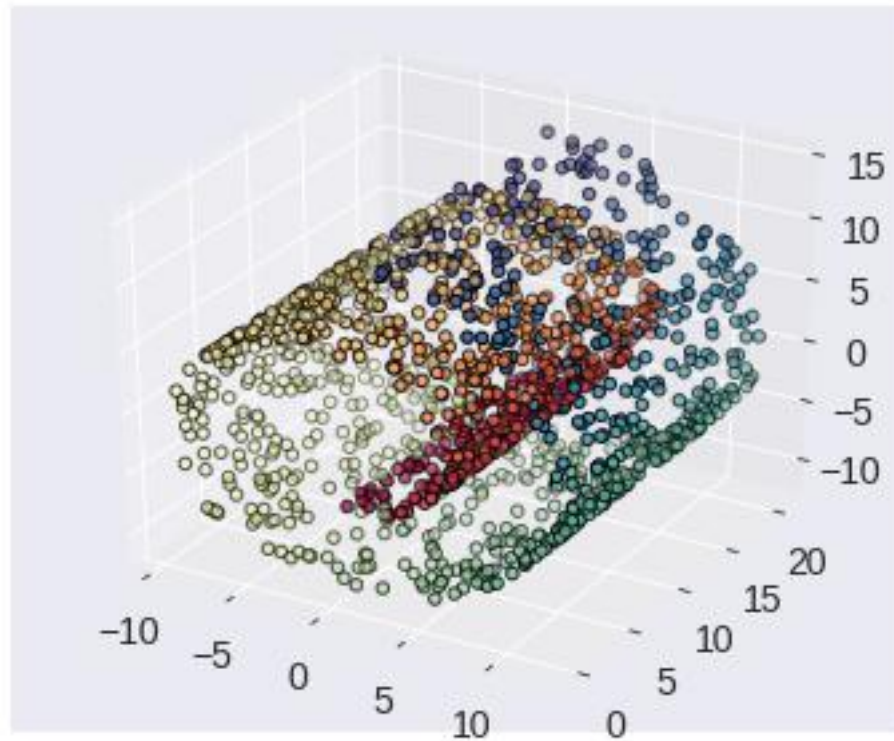
```
from sklearn.decomposition import KernelPCA
kpca = KernelPCA(kernel="rbf", gamma=1, random_state=1)
X_kpca = kpca.fit_transform(X)
```

есть возможность учить и обратное преобразование `fit_inverse_transform=True`

**пример кода:** [https://scikit-learn.org/stable/auto\\_examples/decomposition/plot\\_kernel\\_pca.html](https://scikit-learn.org/stable/auto_examples/decomposition/plot_kernel_pca.html)

kernel PCA: примеры



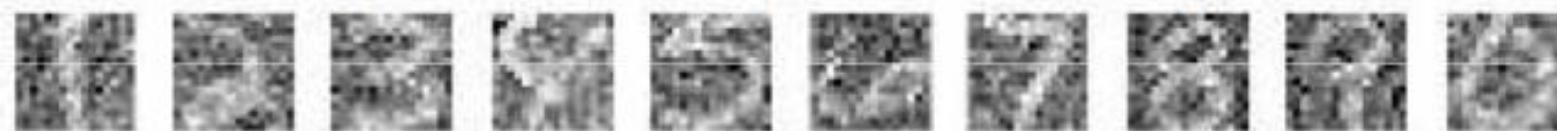
**kernel PCA: примеры**

## kernel PCA: устранение шума

данные:



данные + гауссовский шум



линейный PCA



RBF kernel PCA



[http://www.cs.haifa.ac.il/~rita/uml\\_course/lectures/KPCA.pdf](http://www.cs.haifa.ac.il/~rita/uml_course/lectures/KPCA.pdf)

## **kernel PCA: свойства**

- не всегда удобен для визуализации**
- не всегда получается желаемое...**
- + нелинейное сокращение размерности (и преобразование пространства)**
- + может быть полезен для генерации признаков**

## Итоги

**USL – определение «природы» (структуры) неразмеченных данных**

**Часто удаётся найти «хорошее» маломерное пространство**

**Также нет идеальных методов**  
не забывать про нормировку признаков  
однородность пространства

**интерактивная демка**

<http://colah.github.io/posts/2014-10-Visualizing-MNIST/>

**хорошая презентация по теме**

<https://sites.uclouvain.be/inma/reddot/slides/lee09.pdf>