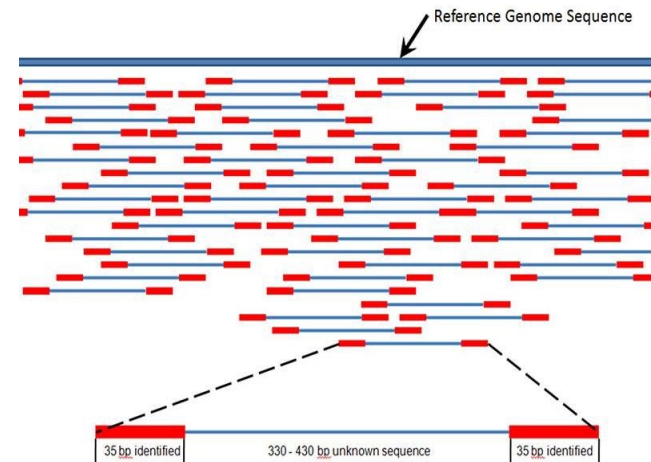


# Понижение размерности



# Биоинформатика

- Задачи анализа генома человека
- Признаки: характеристики генов (более 20.000)
- Маленькие выборки (расшифровка генома — сложная и дорогостоящая процедура)
- Признаков существенно больше, чем объектов!



# Категориальные признаки

- Пример: предсказать, понравится ли пользователю фильм
- Объект: пара «пользователь-фильм»
- Признаки: ID пользователя, ID фильма, ID жанра, ID режиссёра, ID главных актёров, ID композитора, ...
- Как много фильмов снято за всю историю?
- IMDB: >330 тысяч
- После бинарного кодирования получим миллионы признаков

# Анализ текстов

- Пример: предсказание популярности фильма по тексту его сценария
- Признаки: количество вхождений каждого слова из словаря
- Сколько слов в словаре?
- Сотни тысяч признаков
- Если учитывать  $n$ -граммы, то десятки миллионов признаков

# Анализ данных ЭЭГ

- Энцефалограф: 64 датчика, частота сигнала 256 Гц
- Объект: результаты измерений для одного пациента
- За 5 секунд измерений:  $64 * 256 * 5 = 81\,920$  признаков



# UCI Machine Learning Repository



**UCI Machine Learning Repository**  
Center for Machine Learning and Intelligent Systems






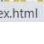
[About](#) [Citation Policy](#) [Donate a Data Set](#) [Contact](#)

☒ Repository ☐ Web

[View ALL Data Sets](#)

Browse Through: 557 Data Sets

Table View List View

Default Task	Name	Data Types	Default Task	Attribute Types	# Instances	# Attributes	Year
<b>Classification (417)</b> Regression (129) Clustering (112) Other (56)	 URL Reputation	Multivariate, Time-Series	Classification	Integer, Real	2396130	3231961	2009
<b>Attribute Type</b> Categorical (38) Numerical (375) Mixed (55)	 KASANDR	Multivariate	Causal-Discovery	Integer	17764280	2158859	2017
<b>Data Type</b> Multivariate (434) Univariate (27) Sequential (55) Time-Series (112) Text (63) Domain-Theory (23) Other (21)	 Gas sensor arrays in open sampling settings	Multivariate, Time-Series	Classification	Real	18000	1950000	2013
<b>Area</b> Life Sciences (131) Physical Sciences (55) CS / Engineering (205) Social Sciences (31) Business (40) Game (10) Other (80)	 YouTube Multiview Video Games Dataset	Multivariate, Text	Classification, Clustering	Integer, Real	120000	1000000	2013
<b># Attributes</b> Less than 10 (141) 10 to 100 (252) Greater than 100 (99)	 Twin gas sensor arrays	Multivariate, Time-Series, Domain-Theory	Classification, Regression	Real	640	480000	2016
	 Deepfakes: Medical Image Tamper Detection	Multivariate	Classification	Real	20000	200000	2020
	 Gas sensor array exposed to turbulent gas mixtures	Multivariate, Time-Series	Classification, Regression	Real	180	150000	2014
	 ElectricityLoadDiagrams20112014	Time-Series	Regression, Clustering	Real	370	140256	2015

<https://archive.ics.uci.edu/ml/index.html>

<https://archive.ics.uci.edu/ml/datasets.php>

# Задача понижения размерности

- Дано: матрица «объекты-признаки»  $X$  размера  $\ell \times D$
- Найти: новую матрицу «объекты-признаки»  $Z$  размера  $\ell \times d$
- $d < D$

# Но зачем?

- Проклятие размерности
- Шумовые признаки
- Переобучение
- Интерпретируемость модели
- Скорость работы модели
- Визуализация данных

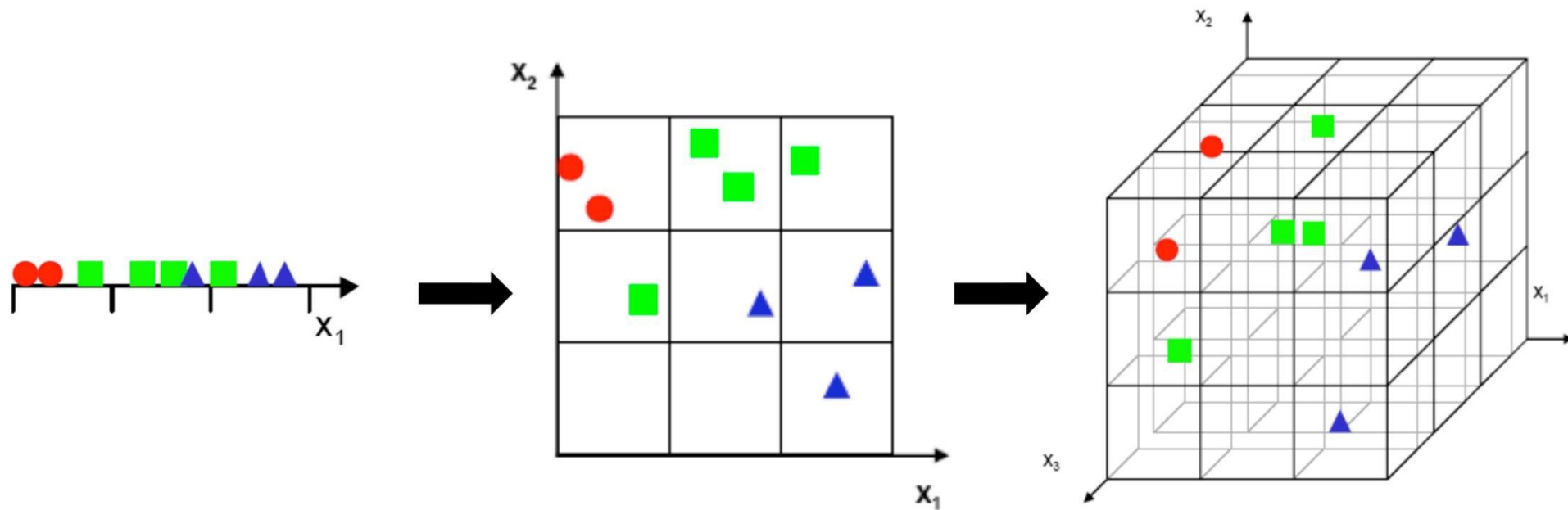


# Проклятие размерности

- Задача: классификация пончиков на вкусные и невкусные
- 100 объектов
- Цвет: 10 вариантов
- Цвет + размер:  $10 * 4 = 40$  вариантов
- Цвет + размер + форма:  $10 * 4 * 4 = 160$  вариантов

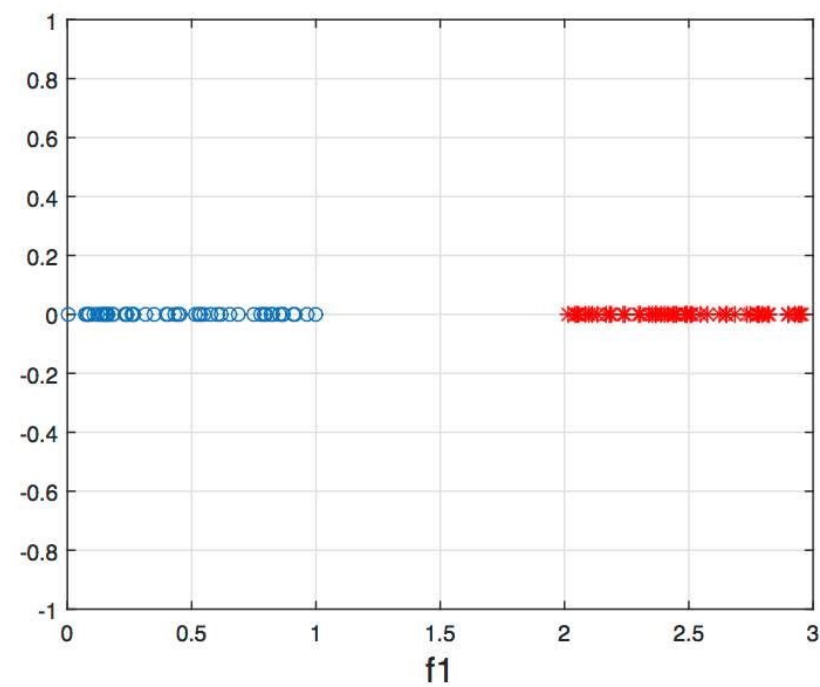


# Проклятие размерности



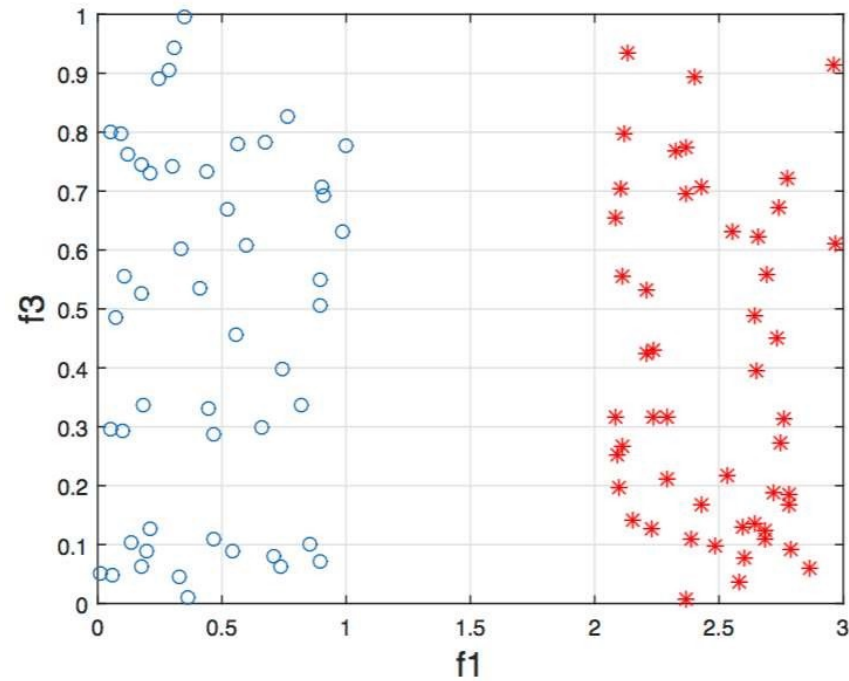
# Плохие признаки

Информативный  
признак



# Плохие признаки

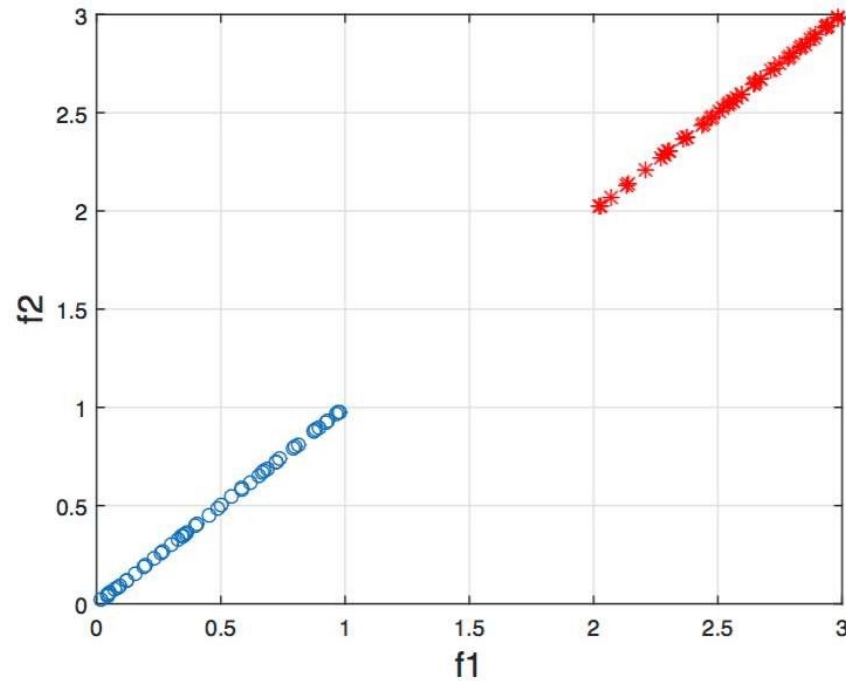
f3 — шумовой  
признак



# Плохие признаки

Коррелирующие  
признаки

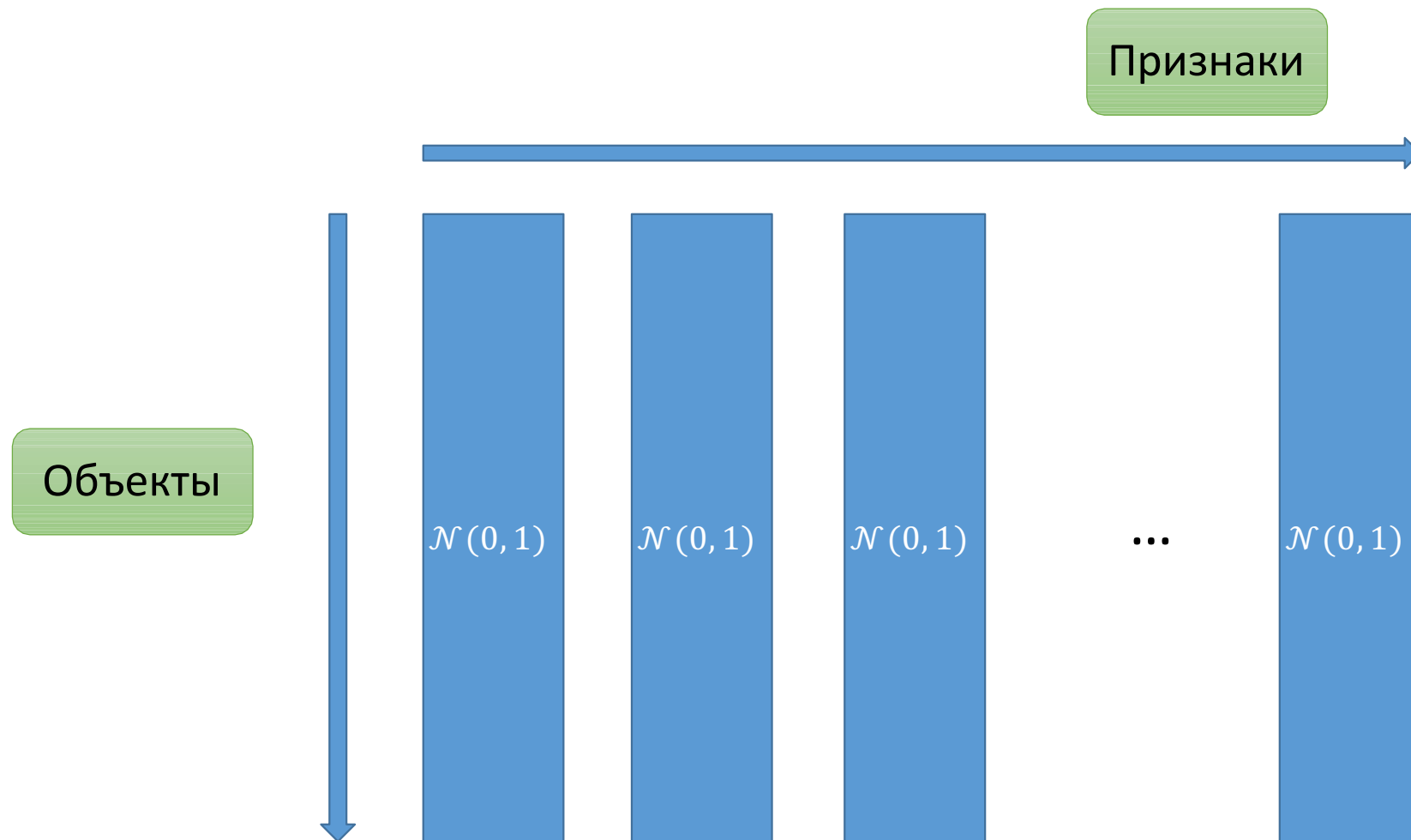
$f_2$  — избыточный  
признак



# Шумовые признаки

- Признаки, которые никак не связаны с целевой переменной
- Но по обучающей выборке это не всегда можно понять

# Шумовые признаки



# Шумовые признаки

- Генерируем случайные признаки
- Если их много, то некоторые будут хорошо коррелировать с ответами

$y$	$x_1$	$x_2$	$x_3$	$x_4$
-1	1.11	<b>-0.5</b>	0.42	0.33
-1	1.22	<b>-0.46</b>	-1.98	-0.55
1	-1.56	<b>0.04</b>	0.39	-1.67
1	-0.48	<b>1.32</b>	0.88	-0.27



# Ускорение моделей

- Чем больше признаков, тем дольше обучаются модели
- Чем дольше обучаются модели, тем меньше экспериментов удаётся провести
- Чем сложнее модели, тем дольше они вычисляют прогнозы
- Могут быть жёсткие ограничения на скорость
- Пример: рекомендательные системы

# Методы понижения размерности

- Отбор признаков (feature selection)
  - Выбрать  $d$  самых важных признаков
- **Извлечение признаков (feature extraction)**
  - Найти  $d$  новых признаков, выражающихся через исходные

# МЕТОД ГЛАВНЫХ КОМПОНЕНТ (PRINCIPAL COMPONENT ANALYSIS, PCA)

# МЕТОД ГЛАВНЫХ КОМПОНЕНТ

Цель: *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
- $z_1, \dots, z_d$  - новые числовые признаки,  $d \leq n$

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

# МЕТОД ГЛАВНЫХ КОМПОНЕНТ

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

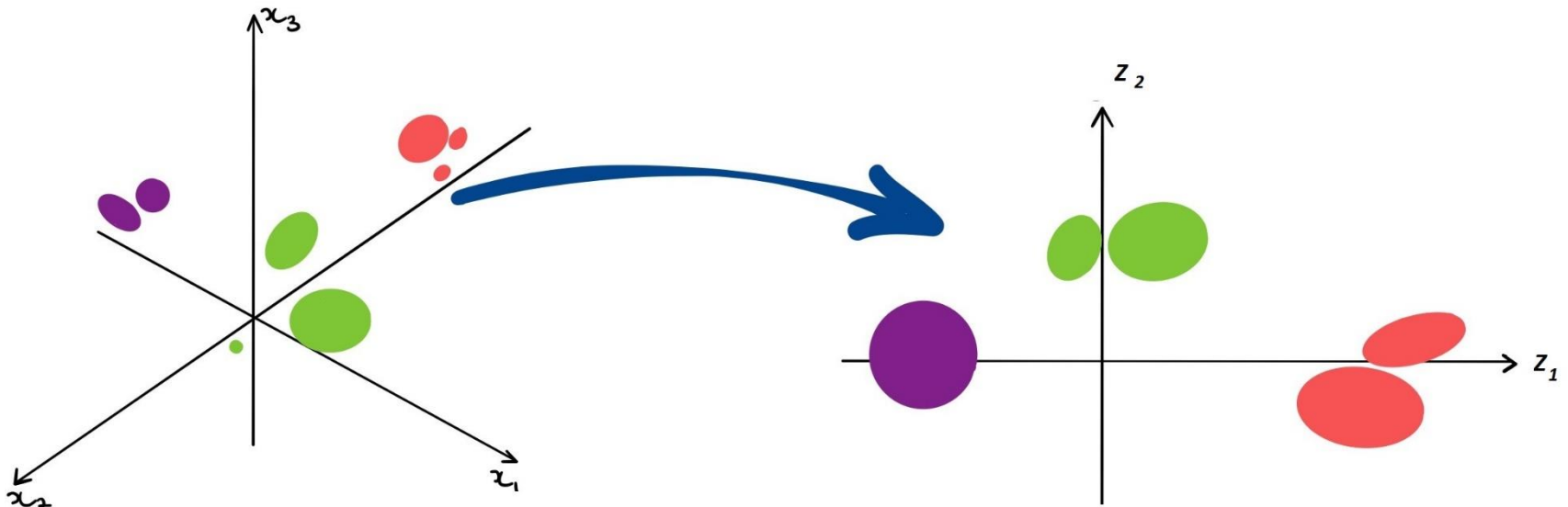
$$\begin{cases} z_1 = u_{11}x_1 + \dots + u_{1n}x_n \\ z_2 = u_{21}x_1 + \dots + u_{2n}x_n \\ \dots \\ z_d = u_{d1}x_1 + \dots + u_{dn}x_n \end{cases}$$

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

# МЕТОД ГЛАВНЫХ КОМПОНЕНТ

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

Геометрически это означает, что мы проецируем пространство признаков размерности  $n$  на некоторое линейное подпространство размерности  $d$ :



# ПОЯСНЕНИЕ: ПРОЕКЦИЯ

- Проекция вектора  $x$  на вектор (компоненту)  $u_i$ :  
$$(x, u_i)$$

- Проекция выборки  $X$  на компоненту  $u_i$ :  
$$Xu_i$$



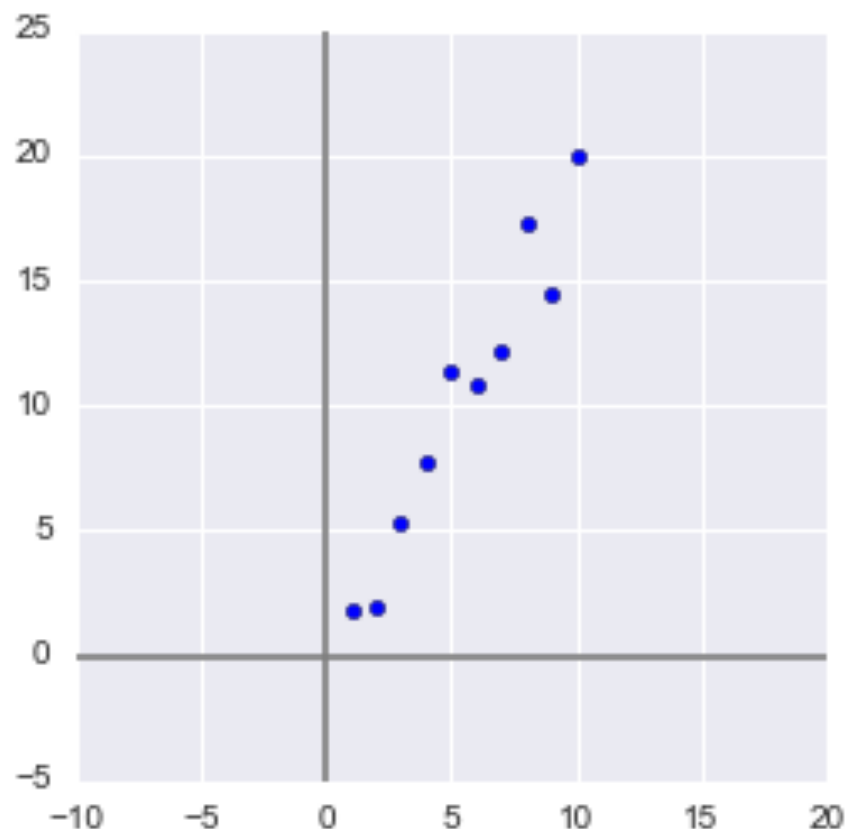
# МЕТОД ГЛАВНЫХ КОМПОНЕНТ

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.**

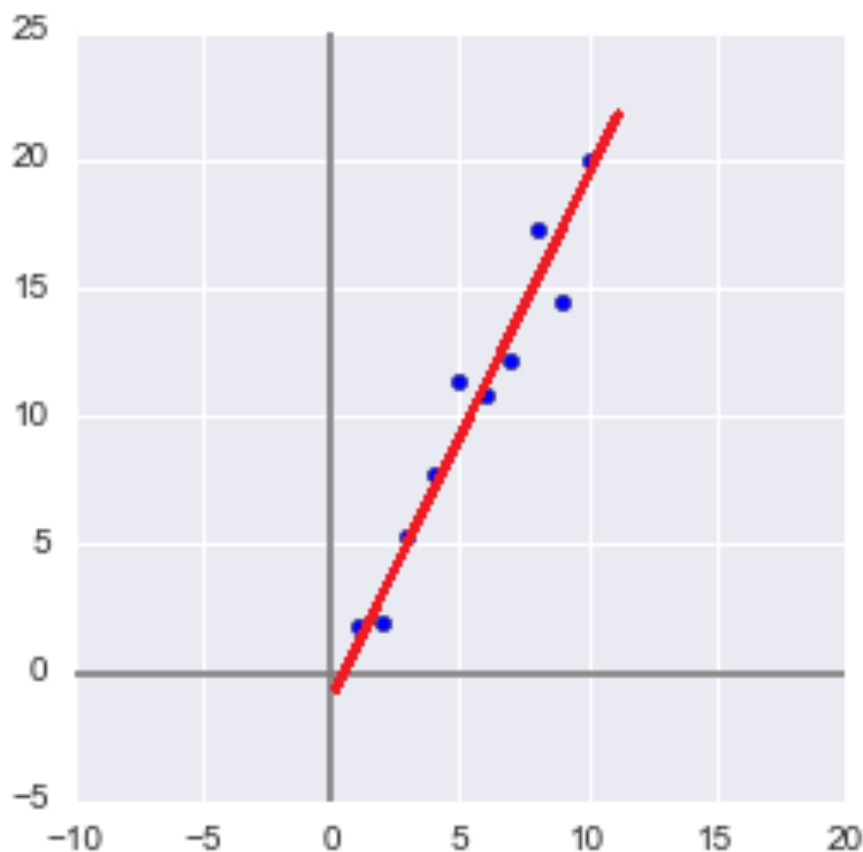
# ПРИМЕР

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:



# ПОСТАНОВКА ЗАДАЧИ

Будем искать такие компоненты  $u_1, u_2, \dots, u_d$ , что:

- 1) Они ортогональны, т.е.  $(u_i, u_j) = 0$
- 2) Они нормированы, т.е.  $\|u_i\| = 1$
- 3) дисперсия проекции выборки на них максимальна:

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

# ВАЖНОЕ ДЕЙСТВИЕ

*Центрируем исходные данные, то есть вычтем из каждого признака его среднее значение.*

# ДИСПЕРСИЯ ПРОЕКЦИИ

- Мы уже выяснили, что проекция выборки  $X$  на компоненту  $u_i$ :

$$Xu_i$$

- Тогда проекция выборки на первые  $d$  компонент, задаваемых столбцами матрицы  $U_d$ :

$$XU_d$$

# ДИСПЕРСИЯ ПРОЕКЦИИ

- Мы уже выяснили, что проекция выборки  $X$  на компоненту  $u_i$ :

$$Xu_i$$

- Тогда проекция выборки на первые  $d$  компонент, задаваемых столбцами матрицы  $U_d$ :

$$XU_d$$

- Тогда дисперсия проекции – это след ковариационной матрицы:

$$\text{tr}((XU_d)^T(XU_d)) = \sum_{i=1}^d ||Xu_i||^2 \rightarrow \max_u$$

# ПЕРВЫЙ ШАГ

- Будем искать первую компоненту,  $u_1$ :

$$\begin{cases} ||Xu_1||^2 \rightarrow \max_{u_1} \\ ||u_1||^2 = 1 \end{cases}$$



# ПЕРВЫЙ ШАГ

- Будем искать первую компоненту,  $u_1$ :

$$\begin{cases} ||Xu_1||^2 \rightarrow \max_{u_1} \\ ||u_1||^2 = 1 \end{cases}$$

Решение:

- Запишем лагранжиан

$$L(u_1, \lambda) = ||Xu_1||^2 + \lambda(||u_1||^2 - 1)$$

# ПЕРВЫЙ ШАГ

- Будем искать первую компоненту,  $u_1$ :

$$\begin{cases} ||Xu_1||^2 \rightarrow \max_{u_1} \\ ||u_1||^2 = 1 \end{cases}$$

Решение:

- Запишем лагранжиан

$$L(u_1, \lambda) = ||Xu_1||^2 + \lambda(||u_1||^2 - 1)$$

- $\frac{\partial L}{\partial u_1} = ?$

# ПЕРВЫЙ ШАГ

- Будем искать первую компоненту,  $u_1$ :

$$\begin{cases} ||Xu_1||^2 \rightarrow \max_{u_1} \\ ||u_1||^2 = 1 \end{cases}$$

Решение:

- Запишем лагранжиан

$$L(u_1, \lambda) = ||Xu_1||^2 + \lambda(||u_1||^2 - 1)$$

- $\frac{\partial L}{\partial u_1} = 2X^T Xu_1 + 2\lambda u_1 = 0 \Rightarrow X^T Xu_1 = -\lambda u_1$  - собств.в-р.

# ПЕРВЫЙ ШАГ

- Будем искать первую компоненту,  $u_1$ :

$$\begin{cases} ||Xu_1||^2 \rightarrow \max_{u_1} \\ ||u_1||^2 = 1 \end{cases}$$

Решение:

- Запишем лагранжиан

$$L(u_1, \lambda) = ||Xu_1||^2 + \lambda(||u_1||^2 - 1)$$

- $\frac{\partial L}{\partial u_1} = 2X^T Xu_1 + 2\lambda u_1 = 0 \Rightarrow X^T Xu_1 = -\lambda u_1$  - собств.в-р.
- $||Xu_1||^2 = u_1^T X^T Xu_1 = \lambda u_1^T u_1 = \lambda \rightarrow \max_{u_1}$  - max  
собств. значение.

# ПЕРВЫЙ ШАГ

- Будем искать первую компоненту,  $u_1$ :

$$\begin{cases} ||Xu_1||^2 \rightarrow \max_{u_1} \\ ||u_1||^2 = 1 \end{cases}$$

Ответ:

$u_1$  - собственный вектор матрицы ковариаций  $X^T X$  с максимальным собственным значением.

# 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.

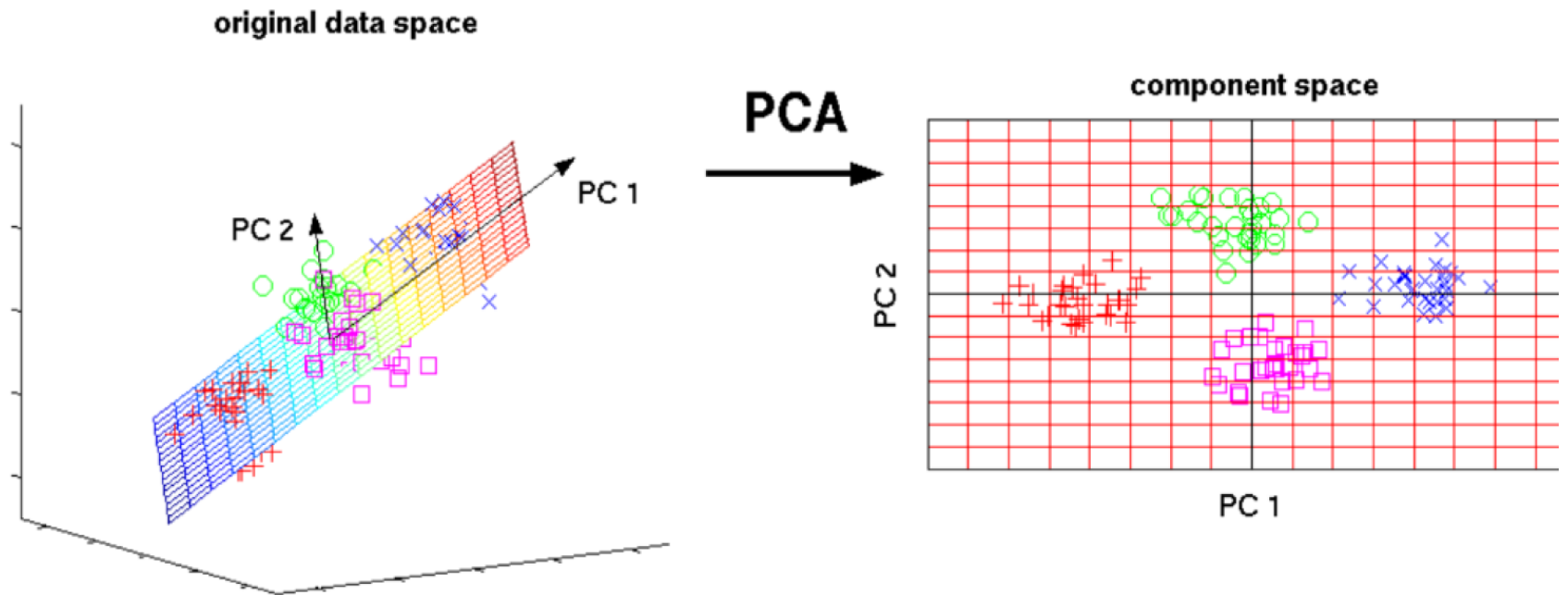
# КОНСТРУКТИВНОЕ ПОСТРОЕНИЕ БАЗИСА В РСА

- 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  $\{u_1, u_2, \dots, u_d\}$ .*

# ПРОЕКЦИЯ НА ГИПЕРПЛОСКОСТЬ





# ПРИМЕНЕНИЕ МЕТОДА

- Когда главные компоненты найдены, можно проецировать на них и новые данные:

$$Z' = X'U_d.$$

# 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:

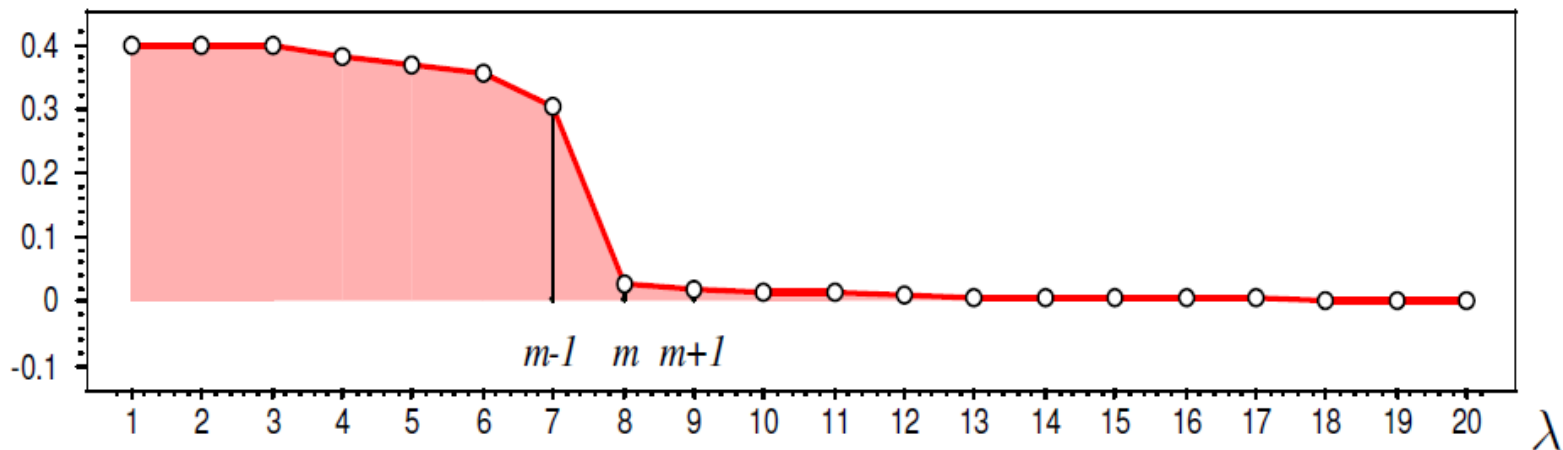
$$\delta = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_n} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n \lambda_i}$$

# 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

#efaces=1, res=57.804

#efaces=2, res=57.611

#efaces=5, res=54.054

#efaces=10, res=52.01

#efaces=20, res=45.897



#efaces=40, res=35.868

#efaces=60, res=29.624

#efaces=80, res=24.103

#efaces=100, res=20.317

#efaces=150, res=16.154



#efaces=200, res=13.257

#efaces=300, res=9.581

#efaces=400, res=6.908

#efaces=1000, res=0.924

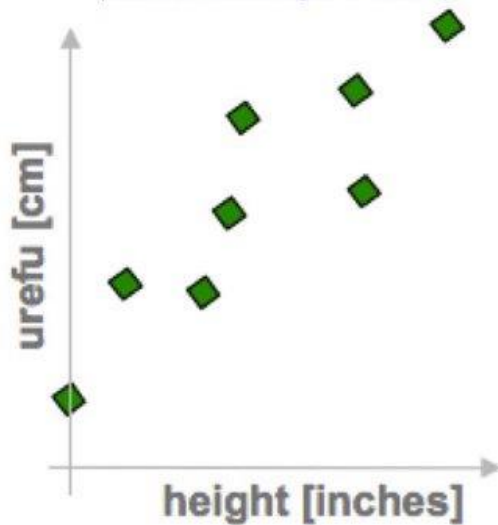
#efaces=1071, res=0.653



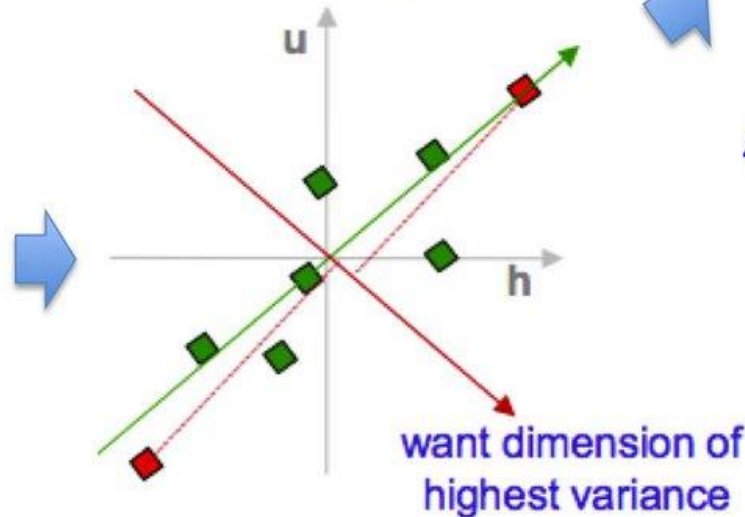
# PCA in a nutshell

## 1. correlated hi-d data

("urefu" means "height" in Swahili)



## 2. center the points



## 3. compute covariance matrix

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

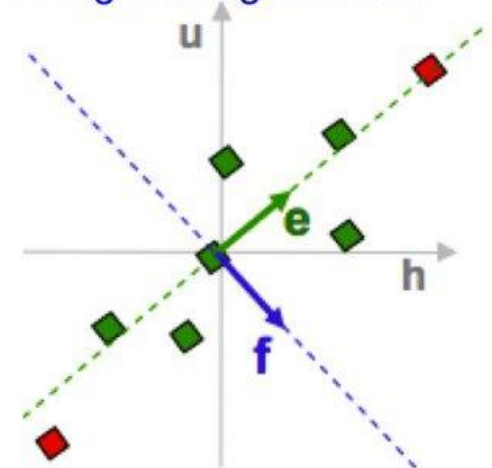
## 4. eigenvectors + eigenvalues

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

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

$\text{eig}(\text{cov}(\text{data}))$

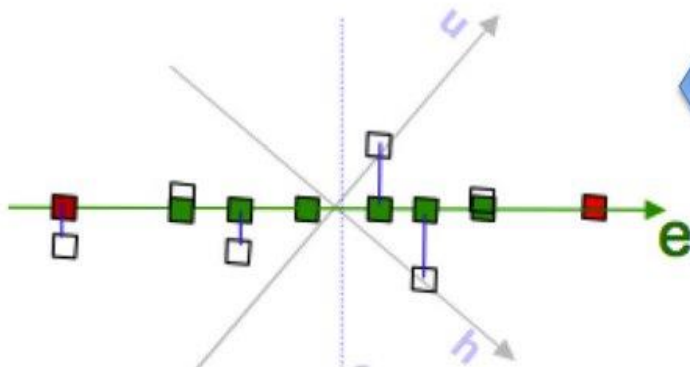
## 5. pick $m < d$ eigenvectors w. highest eigenvalues



## 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  $\lambda_i$  - собственные значения матрицы  $A^T A$ .



# СИНГУЛЯРНОЕ РАЗЛОЖЕНИЕ МАТРИЦЫ (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  $\lambda_i$  - 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  $m \leq n$ :

$$\begin{array}{c} m \times n \\ \boxed{A} \end{array} = \begin{array}{c} m \times m \\ \boxed{U} \end{array} \cdot \begin{array}{c} m \times n \\ \boxed{\begin{array}{c} \sigma_1 \sigma_2 \dots \sigma_m \\ \Sigma \end{array}} \end{array} \cdot \begin{array}{c} n \times n \\ \boxed{V^T} \end{array}$$

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

- При  $m > n$ :

$$\begin{array}{c} m \times n \\ \boxed{A} \end{array} = \begin{array}{c} m \times m \\ \boxed{U} \end{array} \cdot \begin{array}{c} m \times n \\ \boxed{\begin{array}{c} \sigma_1 \sigma_2 \dots \sigma_n \\ \Sigma \end{array}} \end{array} \cdot \begin{array}{c} n \times n \\ \boxed{V^T} \end{array}$$

# 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  $\Sigma$  – 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  $\Sigma$  – 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  $\Sigma$ , тогда матрица  $U_k \Sigma_k$  содержит  $k$  новых признаков, соответствующих первым  $k$  главным компонентам.

# ЧТО ЛУЧШЕ: PCA ИЛИ SVD?

- Существуют вычислительные трудности с нахождением собственных значений, в этом недостаток PCA.
- Существует итерационный алгоритм для нахождения SVD (без нахождения собственных значений)

[http://www.machinelearning.ru/wiki/index.php?title=Простой\\_итерационный\\_алгоритм\\_сингулярного\\_разложения.](http://www.machinelearning.ru/wiki/index.php?title=Простой_итерационный_алгоритм_сингулярного_разложения)

Поэтому вычислительно эффективнее использовать SVD при прочих равных.