

Санкт-Петербургский государственный университет  
Прикладная математика и информатика  
Статистическое моделирование

Зенкова Наталья Валентиновна

## МНОГОМЕРНЫЙ АНАЛИЗ ДАННЫХ

Лекции Голяндиной Нины Эдуардовны

Санкт-Петербург

2018

# Оглавление

<b>Введение</b>	3
1. Генеральный язык	3
2. Выборочный язык	3
3. Примеры данных до стандартизации и после	4
4. Многомерное нормальное распределение	5
5. Переход к новым признакам	8
5.1. Примеры	8
5.2. Две интерпретации новых признаков. Факторные нагрузки	10
5.3. Важность новых признаков	11
6. Пара фактов из линейной алгебры	12
 <b>Глава 1. Сингулярное разложение</b>	
<b>(Singular Value Decomposition)</b>	14
1.1. Как строится сингулярное разложение	14
1.1.1. Матричный вид сингулярного разложения	15
1.2. Единственность SVD	16
 <b>Глава 2. Анализ главных компонент</b>	
<b>(PCA — principal component analysis)</b>	18
2.1. Анализ главных компонент на выборочном языке	18
2.2. Связь между SVD и АГК	18
2.3. Чем отличается SVD от АГК?	19
2.3.1. Пример	19
2.4. Чему АГК соответствует на статистическом языке?	20
2.5. АГК на выборочном языке (продолжение)	20

# Введение

## 1. Генеральный язык

Пусть случайная величина  $\xi = (\xi_1, \dots, \xi_p)^T \in \mathbb{R}^p$ , где  $p$  — количество признаков. Обозначим

- *вектор средних*  $\mu = (E\xi_1, \dots, E\xi_p)^T \in \mathbb{R}^p$ ;
- *ковариационную матрицу*  $\Sigma = \text{Cov}\xi = E(\xi - E\xi)(\xi - E\xi)^T \in \mathbb{R}^{p \times p}$ .

Рассмотрим два случая:

- *Центрированные данные*:  $\xi^{(c)} = \xi - E\xi$ . Тогда ковариационная матрица будет иметь вид:  $\text{Cov}\xi^{(c)} = E\xi^{(c)}(\xi^{(c)})^T$ .
- *Стандартизованные данные*:  $\xi_i^{(s)} = \frac{\xi_i - E\xi_i}{\sqrt{D\xi_i}}$ <sup>1</sup>.

Запишем теперь это в матричном виде: введём  $\Delta$  диагональную матрицу, состоящую из элементов  $\sqrt{D\xi_i}$  для  $i = 1, \dots, p$ . Тогда  $\xi^{(s)} = \Delta^{-1}\xi^{(c)} \Rightarrow \text{Corr}\xi = \text{Cov}\xi^{(s)}$ , то есть  $\text{Cov}\xi^{(s)}$  — это корреляционная матрица до стандартизации.

## 2. Выборочный язык

Перейдём теперь на **выборочный язык**.

Пусть  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  — выборка объёма  $n$ . Можем задать эмпирическое распределение как  $\{\mathbf{x}_i$  с вероятностью  $1/n\}$ . Возникает вопрос: как задать *матрицу данных*? Удобно было бы ее задать как матрицу размера  $p \times n$ :

$$\begin{array}{|c|c|c|} \hline & & \\ \hline \mathbf{x}_1 & \dots & \mathbf{x}_n \\ \hline \end{array}.$$

Но в книгах встречается следующий вариант *матрицы данных*:

$$\mathbb{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix} = [X_1 : \dots : X_p] \in \mathbb{R}^{n \times p}.$$

---

<sup>1</sup>  $\xi_i$  —  $i$ -ая компонента вектора  $\xi$ .

Теперь знаем, как написать все характеристики, упомянутые выше, на выборочном языке.

- $\bar{X}_j = \frac{1}{n} \sum_{i=1}^n (X_j)_i$  — выборочное среднее.

Центрированная матрица данных  $\mathbb{X}^{(c)}$  <sup>2</sup>:

$$\mathbb{X}^{(c)} = [X_1 - \bar{X}_1 : \dots : X_p - \bar{X}_p].$$

- $s_j^2 = \frac{1}{n} \|\mathbb{X}_j^{(c)}\|^2$  — выборочная дисперсия.

Матрица из стандартизованных признаков  $\mathbb{X}^{(s)}$ :

$$\mathbb{X}_j^{(s)} = \frac{X_j - \bar{X}_j}{s_j} = \frac{\mathbb{X}_j^{(c)}}{s_j}, \quad \frac{1}{n} \|\mathbb{X}_j^{(s)}\|^2 = 1 \quad \forall 1, \dots, p.$$

- $\mathbb{S} = \frac{1}{n} (\mathbb{X}^{(c)})^T \mathbb{X}^{(c)}$  — выборочная ковариационная матрица. <sup>3</sup>
- $\widehat{\text{Corr}}\xi = \frac{1}{n} (\mathbb{X}^{(s)})^T \mathbb{X}^{(s)}$  — выборочная корреляционная матрица.

### 3. Примеры данных до стандартизации и после

Зададимся вопросом, как выглядят данные до стандартизации и после неё?

1. Данные с ненулевой корреляцией: <sup>4</sup>

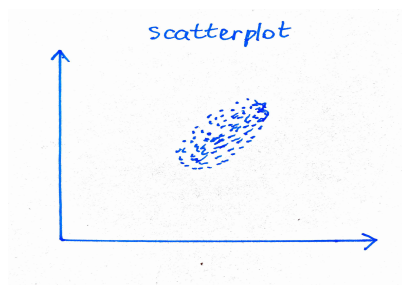


Рис. 1. Данные до стандартизации

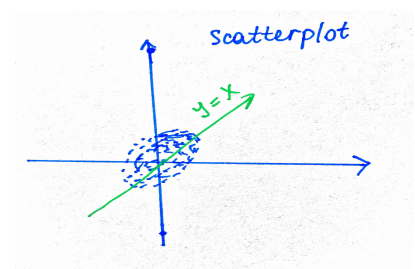


Рис. 2. Данные после стандартизации

<sup>2</sup> Посчитали среднее по каждому признаку и вычли. Таким образом, среднее по каждому признаку нулевое.

<sup>3</sup> Выписать несмещённую оценку ковариационной матрицы.

<sup>4</sup> Давайте подумаем, как будет выглядеть линия регрессии в первом и во втором случаях. Очевидно, что во втором случае она будет проходить через точку (0;0) и совпадать с прямой  $y = x$ .

## 2. Данные с нулевой корреляцией:

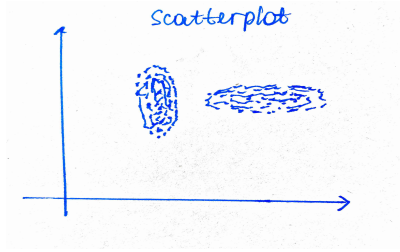


Рис. 3. Данные до стандартизации

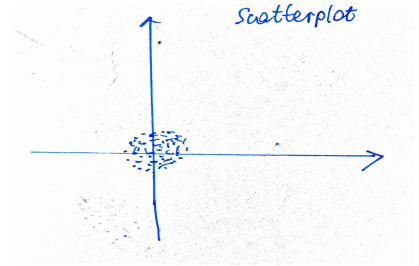


Рис. 4. Данные после стандартизации

## 4. Многомерное нормальное распределение

Рассмотрим  $\xi \sim N(\mu, \Sigma)$ , где  $\mu = E\xi$  — среднее значение  $\xi$ ,  $\Sigma = \text{Cov}\xi$  — ковариационная матрица. Предположим, что  $\Sigma$  — невырожденная матрица.<sup>5</sup> Плотность многомерного нормального распределения задаётся формулой:

$$p(\mathbf{x}, \mu, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right\} \text{ для } \forall \mathbf{x} \in \mathbb{R}^p, \text{ где}$$

$p$  — количество признаков. Очевидно, что в одномерном случае плотность будет иметь следующий вид:

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\}, \text{ где}$$

$\mu, \sigma$  — числа. Введём следующее понятие — *расстояние Махаланобиса*:

$$r_M^2(\mathbf{x}, \mu, \Sigma) = (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu).$$

Отметим, что  $\Sigma$  в данном случае является неотрицательно определённой симметричной матрицей. Но так как ковариационная матрица удовлетворяет данному условию, то получаем, что если расстояния Махаланобиса совпадают, то и плотности многомерного нормального распределения тоже. Посмотрим, что будет являться расстоянием Махаланобиса в одномерном случае ( $p = 1$ ):

$$\sqrt{\frac{(x - \mu)^2}{\sigma^2}} = \frac{|x - \mu|}{\sigma} \text{ — расстояние, измеренное в } \sigma.$$

---

<sup>5</sup>  $|\Sigma| \neq 0$ .

На выборочном языке:

$$\frac{|x_i - \bar{x}|}{s_i}.$$

Таким образом, если у нас есть нормальное распределение, то в одномерном случае мы можем мерить расстояние в сигмах, а в случае многомерного нормального распределения мы смотрим на линии уровня.<sup>6</sup>

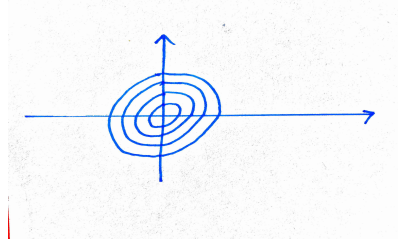


Рис. 5. Зависимые случайные величины

Ответим на следующий вопрос: *какое преобразование необходимо сделать со случайной величиной, имеющей нормальное распределение, чтобы из зависимости получить независимость?*

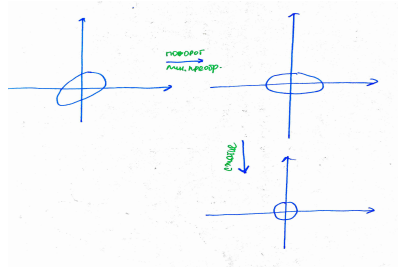


Рис. 6. Отбеливание

Пусть  $\zeta = \Sigma^{-1/2}(\xi - \mu)$  — новая случайная величина, полученная из исходной линейными преобразованиями. Можем считать, что  $E\zeta = \mathbb{O}$ , так как

$$E\zeta = \Sigma^{-1/2}(E\xi - \mu) = \mathbb{O}.$$

---

<sup>6</sup> Находимся на одном расстоянии от центра

Посчитаем ковариацию  $\zeta$ :<sup>7</sup>

$$\mathbb{E}\zeta\zeta^T = \Sigma^{-1/2} \underbrace{\Sigma^1}_{\mathbb{E}\xi\xi^T} \Sigma^{-1/2} = \mathbb{I}.$$

Вспомним несколько свойств, которые характерны для нормально распределённых случайных величин:

1. Если  $\xi \sim N$ , то и всякое линейное преобразование  $\xi$  тоже имеет нормальное распределение.
2. Знаем, что из независимости следует некоррелированность. В случае нормального распределения верно и обратное. В общем случае это неверно. Приведём пример, когда зависимость есть, но корреляция равна 0.

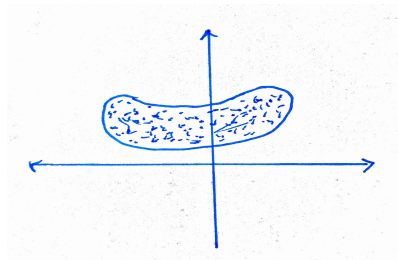


Рис. 7. Пример данных, когда зависимость есть, но корреляция нулевая

3. Условное математическое ожидание является линейной функцией.
4. Ковариационная матрица  $\Sigma$  может быть вырожденной:<sup>8</sup>

---

<sup>7</sup> В данном доказательстве воспользовались симметричностью и неотрицательной определённости ковариационной матрицы  $\Sigma$ .

<sup>8</sup> Наблюдения лежат на одной прямой.

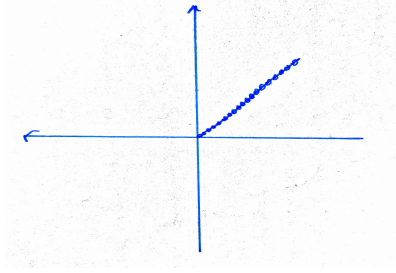


Рис. 8. Ковариационная матрица вырождена

## 5. Переход к новым признакам

Пусть  $\xi$  —  $p$ -мерная случайная величина.<sup>9</sup> Хотим образовать новый признак, который является линейной комбинацией старых признаков.<sup>10</sup>

$$\eta = \mathbf{A}^T \xi = \xi^T \mathbf{A} = a_1 \xi_1 + \dots + a_p \xi_p, \quad \mathbf{A} \in \mathbb{R}^p,$$

$$\text{D}\eta = \mathbf{A}^T \Sigma \mathbf{A}.$$

Теперь пусть  $\eta = (\eta_1, \dots, \eta_d)^T \in \mathbb{R}^d$  —  $d$ -мерная случайная величина,<sup>11</sup> которая получается из исходной  $\eta = \mathbf{A}^T \xi$ , где  $\mathbf{A} \in \mathbb{R}^{p \times d}$ . Соответствующая ковариационная матрица будет иметь вид  $\text{Cov}\eta = \mathbf{A}^T \Sigma \mathbf{A}$ .

**На выборочном языке:**  $Z = \sum_{j=1}^p a_j X_j = \mathbb{X} \mathbf{A} \in \mathbb{R}^n$  или  $\mathbb{Z} = \mathbb{X} \mathbf{A} \in \mathbb{R}^{n \times d}$ .<sup>12</sup>

### 5.1. Примеры

1. Возьмём в качестве индивида школьников, а в качестве признаков — оценки по четырём школьным предметам:

- $X_1$  — оценка по математике;
- $X_2$  — оценка по физике;
- $X_3$  — оценка по русскому;
- $X_4$  — оценка по литературе;

<sup>9</sup>  $p$  — количество признаков.

<sup>10</sup> Новый признак будет лучше, чем старые.

<sup>11</sup>  $d$  новых признаков.

<sup>12</sup>  $\mathbb{X} \in \mathbb{R}^{n \times p}$  — матрица данных.



Сами по себе оценки мало несут информации про общее положение в школе или классе, поэтому давайте создадим новые признаки, которые будут более информативными. Пусть это будут

- $Z_1 = X_1 + X_2 + X_3 + X_4$  — сумма оценок по всем предметам;
- $Z_2 = X_1 + X_2 - (X_3 + X_4)$  — разность между оценками по естественным и гуманитарным наукам.

Выясним, как будет выглядеть матрица  $\mathbb{A}$  в этом случае:

$$\mathbb{A} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & -1 \\ 1 & -1 \end{pmatrix} \in \mathbb{R}^{4 \times 2}.$$

2. Рассмотрим два признака: оценки по математике  $X_1$  и физике  $X_2$  (рисунок 10).

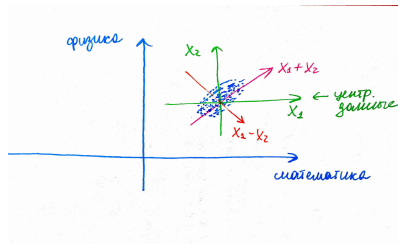


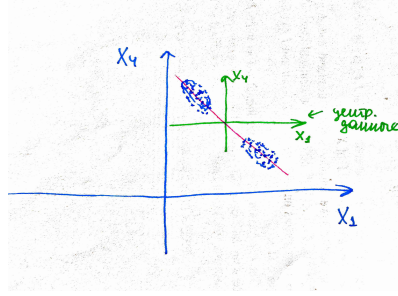
Рис. 9. Scatterplot:  $X_1$  and  $X_2$

В данной ситуации получается, что новый признак  $X_1 + X_2$  лучше всего характеризует данные, то есть это та характеристика, по которой ученики максимально различаются.<sup>1314</sup>

3. Рассмотрим два признака: оценки по математике  $X_1$  и литературе  $X_4$ . Максимальное отличие получаем по разности данных признаков.

<sup>13</sup>  $X_1 - X_2$  — характеристика, по которой ученики минимально отличаются друг от друга.

<sup>14</sup> Необходимо выбрать признак, который бы наилучшим образом характеризовал данные. Такой характеристикой будет тот признак, по чему данные максимально различаются. Таким образом, приходим к анализу главных компонент.

Рис. 10. Scatterplot:  $X_1$  and  $X_4$ 

## 5.2. Две интерпретации новых признаков. Факторные нагрузки

Обсудили, зачем необходимы новые признаки. Также до этого мы требовали, чтобы новые признаки были линейной комбинацией старых. Теперь предположим, что:

- на выборочном языке:  $z_i \perp z_j$ ,<sup>15</sup>
- на генеральном языке:  $\rho(\eta_i, \eta_j) = 0$ .<sup>16</sup>

Зависимость признаков означает, что кусок одного признака входит в другой, то есть получается некоторая избыточность.

Пусть  $\text{rank}(\mathbb{X}) = d$ ,  $(z_1, \dots, z_d)$  — ортогональный базис в  $d$ -мерном подпространстве  $\text{span}(X_1, \dots, X_p) = \text{colspan}(\mathbb{X})$ . Превратим базис в *ортонормированный*:  
 $Q_i = \frac{z_i}{\|z_i\|}$ .<sup>17</sup>

*Новые признаки на самом деле интерпретируются двумя способами:*

---

<sup>15</sup> Признаки ортогональны.

<sup>16</sup> Признаки хотя бы некоррелированы.

<sup>17</sup> Рассмотрим пример: пространство  $\mathbb{R}^2$ , соответствующий базис  $(1, 1)^T/\sqrt{2}$ ,  $(1, -1)^T/\sqrt{2}$ . Это ортонормированный базис по определению.<sup>18</sup> Как вычислить координаты вектора  $(5, 4)^T$  в данном базисе?

Если  $U_1, \dots, U_d$  — ортогональный базис в  $\mathbb{R}^d$ , тогда  $\forall A \in \mathbb{R}^d$  раскладывается по ортонормированному базису:

$$A = \sum_{i=1}^d \langle A, U_i \rangle U_i, \text{ где } \langle A, U_i \rangle - i\text{-ая координата вектора } A \text{ в базисе } \{U_j\}_{j=1}^d.$$

- Новые признаки задаются равенством:<sup>19</sup>

$$\mathbb{X}\mathbb{A} = \mathbb{Z}.$$

- Исходные признаки выражаются через новые:<sup>20</sup>

$$\forall i = 1, \dots, p \quad X_i = \sum_{j=1}^d \langle X_i, Q_j \rangle Q_j. \quad (1)$$

**Определение 1.**  $f_{ij} = \langle X_i, Q_j \rangle$  — факторные леса, факторные нагрузки.

Составим матрицу  $\mathbb{F} = \{f_{ij}\}_{i=1, j=1}^{p, d} = [F_1, \dots : F_d] \in \mathbb{R}^{p \times d}$  — коэффициенты разложения по базису, тогда можем записать разложение 1 в матричном виде:

$$\mathbb{X} = \sum_{j=1}^d Q_j F_j^T = \mathbb{Q}\mathbb{F}^T. \quad (2)$$

$\|Q_j\| = 1$ ,  $\|F_j\| \neq 1$  в разложении 2. Давайте нормируем  $F$ . Пусть  $\sigma_j = \|F_j\|$ ,  $P_j = \frac{F_j}{\|F_j\|}$ , тогда  $\|Q_j\| = \|P_j\| = 1$  и

$$\mathbb{X} = \sum_{j=1}^d \sigma_j Q_j P_j^T = \mathbb{Q}\Sigma\mathbb{P}^T, \text{ где } \Sigma = \text{diag}\{\sigma_1, \dots, \sigma_d\}. \quad (3)$$

Что за матрица  $Q_j P_j^T$ ?<sup>21</sup>  $Q_j P_j^T \in \mathbb{R}^{n \times p}$  и  $\text{rank}(Q_j P_j^T) = 1$ .<sup>22</sup> Таким образом, у нас была исходная матрица ранга  $d$ , а мы превратили её в сумму элементарных матриц ранга 1.

**Предложение 1.**  $\Sigma\mathbb{P}^T = \mathbb{F}^T$ .

*Доказательство.* И это бы надо бы пояснить бы, гыыыы. □

Пусть  $\mathbb{X}^{(j)} = \sigma_j Q_j P_j^T$ , тогда  $\mathbb{X} = \sum_{j=1}^d \mathbb{X}^{(j)}$ .

### 5.3. Важность новых признаков

Ввели новый признак, теперь хотелось бы понять, какую важную часть информации он в себя включает. Введём скалярное произведение двух матриц:<sup>23</sup>

$$(\mathbb{Y}, \mathbb{Z})_F = \sum_{i,j} y_{ij} z_{ij}.$$

<sup>19</sup> Новые признаки есть линейная комбинация старых.

<sup>20</sup> Есть пространство, в нём базис, каждый вектор этого пространства раскладывается по базису.

<sup>21</sup> Вектор умножаем на вектор.

<sup>22</sup> Столбцы пропорциональны.

<sup>23</sup> Чтобы уметь измерять вклад признака, необходимо уметь измерять одним числом, то есть ввести норму матрицы. Скалярное произведение по Фробениусу ( $\|\mathbb{A}\|^2 = \sum_{ij} a_{ij}^2$ ).

Если верно  $\|\mathbb{X}\|^2 = \|\sum_{j=1}^d \mathbb{X}^{(j)}\|_F^2 \stackrel{?}{=} \sum_{j=1}^d \|\mathbb{X}^{(j)}\|_F^2$ , то вклад  $j$ -ого признака —  $\frac{\|\mathbb{X}^{(j)}\|^2}{\|\mathbb{X}\|^2}$ .<sup>24</sup>

**Предложение 2.**  $\|\sum_{j=1}^d \mathbb{X}^{(j)}\|_F^2 = \sum_{j=1}^d \|\mathbb{X}^{(j)}\|_F^2$ , если  $\mathbb{X}^{(j)}$  ортогональны.

*Доказательство.* Доказать как было в тесте. Домножить там на транспонированную матрицу и все дела. Мне лениво сейчас это делать. Скалярное произведение 0. Должны быть ортогональны либо  $Q$ , либо  $P$ . □

Так как  $Q_j$  нормированы (НАДО БЫ И ЭТО ПОЯСНИТЬ, ИБО ЩА Я ВАЩЕ НЕ СООБРАЖАЮ), то можем продолжить равенство:

$$\|\mathbb{X}\|^2 = \left\| \sum_{j=1}^d \mathbb{X}^{(j)} \right\|_F^2 = \sum_{j=1}^d \|\mathbb{X}^{(j)}\|_F^2 = \sum_{j=1}^d \sigma_j^2.$$

**Замечание 1.** Пусть старые признаки  $X_i$  центрированные, тогда новые признаки  $Z_i$  и  $Q_i$  тоже центрированные.

*Доказательство.* АНАААААААЛООООГИИИИИИЧНООО, записать по-человечески. Новые — линейная комбинация старых. Среднее линейной комбинации есть линейная комбинация средних по признакам. Средние по признакам равны 0, получаем 0. □

## 6. Пара фактов из линейной алгебры

1. Унитарная матрица  $\mathbb{U}$  — ортогональная матрица в комплексном случае.

- $\mathbb{U}$  — квадратная матрица:  $\mathbb{U}^T = \mathbb{U}^{-1}$ .
- Столбцы  $\mathbb{U}$  ортонормированы.
- Строки  $\mathbb{U}$  ортонормированы.<sup>25</sup>

Умножение на матрицу  $\mathbb{U}$  означает *поворот* или *отражение*. Пусть есть вектора  $Y, Z$ , после умножения на матрицу  $\mathbb{U}$  получим  $\tilde{Y} = \mathbb{U}Y$  (ВСТАВИТЬ РИСУНОК).<sup>26</sup>

<sup>24</sup> Если отношение равно  $1/2$ , то значит, что  $j$ -ый признак измеряет 50% всей информации.

<sup>25</sup> 2 пункт эквивалентен 3. *Почему?* Если матрица ортогональная, то и транспонированная к ней тоже ортогональная (следует из пункта 1).

<sup>26</sup> Поворачиваем вектора  $Y$  и  $Z$  на какой-то угол, а также при домножении на матрицу  $\mathbb{U}$  не меняются нормы векторов.

**Пример 1.**  $\mathbb{U} = \begin{pmatrix} \cos \phi & \sin \phi \\ -\sin \phi & \cos \phi \end{pmatrix}.$

2. Пусть  $\{P_i\}_{i=1}^r$  — система независимых векторов, рассмотрим линейную оболочку  $\mathcal{L}_r = \text{span}\{P_1, \dots, P_r\}$  в  $\mathbb{R}^L$ ,  $\Pi : \mathbb{R}^L \rightarrow \mathcal{L}_r$  — проектор на  $\mathcal{L}_r$ .<sup>27</sup> Матрица  $\Pi$ :

$$\Pi = \mathbb{P}(\mathbb{P}^T \mathbb{P})^{-1} \mathbb{P}^T.$$

Пусть  $\{P_i\}_{i=1}^r$  — ортонормированный базис  $\mathcal{L}_r$ , тогда

$$\mathbb{P}^T \mathbb{P} = \begin{pmatrix} 1 & \cdots & 0 \\ & \ddots & \\ 0 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{r \times r} \text{ и } \Pi = \mathbb{P} \mathbb{P}^T.$$

---

<sup>27</sup> После проекции размерность вектора не меняется.

## Глава 1

# Сингулярное разложение (Singular Value Decomposition)

## 1.1. Как строится сингулярное разложение

Пусть  $L$  — число признаков,  $K$  — количество индивидов,<sup>1</sup>  $\mathbb{Y} = \mathbb{X}^T \in \mathbb{R}^{L \times K}$  — ненулевая матрица. Обозначим  $\mathbb{S} = \mathbb{Y}\mathbb{Y}^T \in \mathbb{R}^{L \times L}$  — симметричная неотрицательно определённая матрица.

$$\mathbb{S}U_i = \lambda_i U_i, \text{ где}$$

$\{U_i\}_{i=1}^L$  — ортонормированный набор из собственных векторов матрицы  $\mathbb{S}$ ,  
 $\lambda_1 \geq \dots \geq \lambda_L \geq 0$  — собственные числа матрицы  $\mathbb{S}$ .<sup>2</sup>

Пусть  $d = \text{rank} \mathbb{Y} \stackrel{\text{def}}{=} \text{colrank} \mathbb{Y} = \text{rowrank} \mathbb{Y}$ . Знаем, что  $d \leq \min(L, K)$ .

**Предложение 3.** 1.  $d = \text{rank} \mathbb{Y}\mathbb{Y}^T$ .

2.  $\lambda_d > 0, \lambda_i = 0, i > d$ .<sup>3</sup>

3.  $\{U_i\}_{i=1}^d$  образуют ортонормированный базис  $\text{colspan} \mathbb{Y}$ .<sup>4</sup>

Введём вектор<sup>5</sup>

$$V_i \stackrel{\text{def}}{=} \frac{\mathbb{Y}^T U_i}{\sqrt{\lambda_i}} \in \mathbb{R}^k, \quad i = 1, \dots, d.$$

**Предложение 4.** 1.  $\{V_i\}_{i=1}^d$  — ортонормированная система векторов.

2.  $V_i$  — собственные вектора  $\mathbb{Y}\mathbb{Y}^T$ , соответствующие тем же самым собственным числам  $\lambda_i$ . Остальные собственные вектора  $\mathbb{Y}\mathbb{Y}^T$  соответствуют нулевому собственному числу.

---

<sup>1</sup> Для понимания пусть будет так.

<sup>2</sup> Неотрицательная определённость матрицы  $\mathbb{S}$  означает именно это.

<sup>3</sup> Упорядочили собственные числа: первые  $d$  строго положительные, а остальные все нули.

<sup>4</sup> Как понимаю, данное утверждение достаточно элементарное, поэтому все нужно будет как-то на экзамене пояснять, хоть оно и идёт без доказательства. И снова не сейчас мне что-то комментировать подробнее...

<sup>5</sup> «Длинный вектор».  $i = 1, \dots, d$  — на 0 делить не можем же.

3.  $U_i = \frac{\mathbb{Y}V_i}{\sqrt{\lambda_i}}$  — верно и в обратную сторону.

4.  $\mathbb{Y} = \sum_{i=1}^d \sqrt{\lambda_i} U_i V_i^T$  — SVD.<sup>6,7</sup>

Что здесь считать новыми признаками, если  $\mathbb{Y} = \mathbb{X}^T$ ?  $V_i$ ,<sup>8</sup> так как  $U_i \in \mathbb{R}^L$ ,  $V_i \in \mathbb{R}^K$ .

- $U_i$  — ортонормированный базис в пространстве столбцов.
- $V_i$  — ортонормированный базис в пространстве строк.<sup>9</sup>
- $\frac{\lambda_i}{\sum_i \lambda_i}$  — вклад нового признака.

**Определение 2.**  $\sqrt{\lambda_i}$  — сингулярные числа матрицы  $\mathbb{Y}$ ,  $U_i$  — левый сингулярный вектор,  $V_i$  — правый сингулярный вектор.

### 1.1.1. Матричный вид сингулярного разложения

Можно записать двумя способами:

1. Введём  $\mathbb{U}_d = [U_1 : \dots : U_d]$ ,  $\mathbb{V}_d = [V_1 : \dots : V_d]$ ,  $\Lambda_d = \text{diag}(\lambda_1, \dots, \lambda_d)$ . Тогда

$$\mathbb{Y} = \mathbb{U}_d \Lambda_d^{1/2} \mathbb{V}_d^T.$$

2. Возьмём  $\mathbb{U} = [U_1 : \dots : U_d : U_{d+1} : \dots : U_L]$  — ортонормированный базис в  $\mathbb{R}^L$ .<sup>10</sup>

$\mathbb{V} = [U_1 : \dots : U_d : U_{d+1} : \dots : U_K]$  — ортонормированный базис в  $\mathbb{R}^K$ .

$$\Lambda = \begin{pmatrix} \lambda_1 & 0 & \dots & \dots & 0 \\ 0 & \ddots & & & 0 \\ 0 & & \lambda_d & & 0 \\ 0 & & & \ddots & 0 \\ 0 & 0 & \dots & & 0 \end{pmatrix} \in \mathbb{R}^{L \times K}. \text{ Тогда }^{11}$$

$$\mathbb{Y} = \mathbb{U} \Lambda^{1/2} \mathbb{V}^T.$$

<sup>6</sup> Разложение в сумму элементарных матриц.

<sup>7</sup> Самый важный пункт утверждения. Ради него все и вводилось.

<sup>8</sup> Они длинные :)

<sup>9</sup> Можем  $\mathbb{X}$  транспонировать, проделать всё то же самое, а поменяются местами только  $U_i$  и  $V_i$ .

<sup>10</sup>  $U_{d+1}, \dots, U_L$  соответствуют нулевому собственному числу матрицы.

<sup>11</sup>  $\mathbb{U}, \mathbb{V}$  — ортогональные матрицы.

## 1.2. Единственность SVD

Насколько единственно разложение SVD (оно одно существует для матрицы или нет)? Можно подумать, что разложение не единственное, так как

1. Собственные вектора не единственные, то есть если  $U_i$  — собственный вектор, то  $-U_i$  — собственный вектор.<sup>12</sup>

$$\mathbb{Y} = \sum_{i=1}^d \sqrt{\lambda_i} U_i V_i^T = \sum_{i=1}^d \sqrt{\lambda_i} (-U_i) (-V_i)^T.$$

2. Пусть есть два одинаковых собственных числа  $\lambda = \lambda_1 = \lambda_2$ ,  $U_1$  и  $U_2$  — два ортонормированных вектора, соответствующих собственному числу  $\lambda$ . Тогда любая линейная комбинация  $U_1$  и  $U_2$  будет являться также собственным вектором и будет соответствовать тому же собственному числу, то есть  $\forall \alpha, \beta \alpha U_1 + \beta U_2$  — с.в. с с.ч.  $\lambda$ . Таким образом, если у нас есть два одинаковых собственных числа, то они порождают подпространство размерности 2, и любой ортонормированный базис в этом подпространстве подходит нам в качестве собственных векторов.<sup>13</sup>

Получаем, что единственности в буквальном смысле не получается. Поэтому сформулируем необходимо нам утверждение.

**Предложение 5** (Единственность SVD). Пусть  $\mathbb{Y} = \sum_{i=1}^L c_i P_i Q_i^T$  — некоторое разложение в сумму элементарных матрицы (биортогональное разложение), такое что:

1.  $c_1 \geq \dots \geq c_L \geq 0$ ;
2.  $\{P_i\}_{i=1}^L$  — ортонормированные,  $\{Q_i\}_{i=1}^L$  — ортонормированные.

Тогда  $\mathbb{Y} = \sum_{i=1}^L c_i P_i Q_i^T$  — SVD, то есть любое биортогональное разложение является сингулярным.

**Замечание 2.** В частности:

- $c_d > 0$ ,  $c_{d+1} = \dots = c_L = 0$ ,

---

<sup>12</sup> В вещественном случае собственный вектор определяется единственным образом с точностью до константы, модуль которой равен 1 (это всего 1 и -1). Но сингулярное разложение обобщается в комплексном случае, а здесь уже констант, по модулю равных 1, много.

<sup>13</sup> Если у нас есть два одинаковых собственных числа, то мы можем брать любой базис, но сумма двух матриц постоянна, то есть она не меняется от выбора базиса. Можно так-то даже в этом убедиться.



- $c_i^2 = \lambda_i$  — собственные числа  $\mathbb{Y}\mathbb{Y}^T$ ,
- $P_i$  — собственные вектора  $\mathbb{Y}\mathbb{Y}^T$ ,
- $Q_i$  — собственные вектора  $\mathbb{Y}^T\mathbb{Y}$ ,
- $Q_i = \frac{\mathbb{Y}^T P_i}{\sqrt{\lambda_i}}$ ,  $i = 1, \dots, d$  ( $d = \text{rank} \mathbb{Y}$ ).

ДОПИСАТЬ SVD, пропущена часть

## Глава 2

## Анализ главных компонент (PCA — principal component analysis)

### 2.1. Анализ главных компонент на выборочном языке

Вспомним, как выглядит SVD на выборочном языке:

$$D_1 = \{1, \dots, L\}, D_2 = \{1, \dots, K\}, \quad \mu_1, \mu_2 — \text{считающие меры},$$

$$\mathbb{Y} = \sum_{i=1}^d \sqrt{\lambda_i} U_i V_i^T.$$

Теперь перейдём на выборочный язык АГК. Помним, что  $\mathbb{Y} = \mathbb{X}^T$ , где столбцы — индивиды, строки — признаки; индивидов  $K$ , а признаков  $L$ .<sup>1</sup>

$$D_1 = \{1, \dots, L\}, D_2 = \{1, \dots, K\},$$

$$\mu_1 — \text{считающая мера}, \mu_2(\{i\}) = \frac{1}{K} — \text{вероятностная мера, где } i — \text{номер индивида}.$$

Предполагаем, что  $\mathbb{Y}$  — центрированная по строчкам, то есть среднее по признакам равно 0, и тогда получаем:

$$\mathbb{Y} = \sum_{i=1}^d \sqrt{\tilde{\lambda}_i} \tilde{U}_i \tilde{V}_i^T.$$

### 2.2. Связь между SVD и АГК

Необходимо найти связь между  $\tilde{\lambda}_i, \tilde{U}_i, \tilde{V}_i$  и  $\lambda_i, U_i, V_i$  соответственно.<sup>2</sup> Знаем, что

$$U_i — \text{о.н.с с.в. матрицы } \mathbb{Y}\mathbb{Y}^T = \mathbb{X}^T\mathbb{X},$$

$$\tilde{U}_i — \text{о.н.с с.в. матрицы } \frac{1}{K}\mathbb{Y}\mathbb{Y}^T = \frac{1}{K}\mathbb{X}^T\mathbb{X}.$$

То есть получили, что  $U_i$  и  $\tilde{U}_i$  совпадают с точностью до коэффициента.<sup>3</sup>

Таким образом, получаем следующие соотношения:<sup>4</sup>

---

<sup>1</sup> Визуально:  $\mathbb{Y}$  — горизонтальная матрица, а  $\mathbb{X}$  — вертикальная.

<sup>2</sup> В SVD веса 1, а здесь 1 и  $1/K$ .

<sup>3</sup> Если в одном и том же пространстве есть два нормированных вектора: один нормирован с одними весами, другой с другими, то они не должны совпадать. Но здесь у нас понятие нормированности одинаковое за счёт того, что у  $U_i$  вес один и тот же — 1.

<sup>4</sup> **Хм, считаю, что в конспекте была опечатка и так верно.**

- $U_i = \tilde{U}_i$ ,
- $\lambda_i = \frac{\tilde{\lambda}_i}{K}$ ,
- $V_i = \sqrt{K}\tilde{V}_i$ .

Также заметим, что<sup>5</sup>

$$\|\mathbb{Y}\|_{1,2}^2 = \sum_{ij} \frac{1}{K} Y_{ij}^2 = \frac{\|\mathbb{Y}\|_F^2}{K}.$$

## 2.3. Чем отличается SVD от АГК?

1. Столбцы в матрице  $\mathbb{Y}$  не равноправны, то есть SVD полностью симметрично, а АГК нет. Если формально, то получаем, что разные нормы в пространстве признаков.
2. В АГК предполагается, что признаки центрированы, а индивиды нет.
3. В АГК рекомендована нормировка.

*Когда нормируем признаки?* Когда признаки измерены в разных шкалах.<sup>6</sup>

### 2.3.1. Пример

- Нормируем признаки, если есть, например, данные в сантиметрах и метрах.
- Не нормируем признаки, если есть, например, баллы за задачи и хотим, чтобы главная компонента отражала уровень по результатам задач. Пусть есть сложные (от 1 до 10) и простые (от 0 до 5) задачи. Ясно, что получить 2.5 балла за простую задачу и 5 баллов за сложную — это разные вещи, поэтому если мы нормируем данные, то мы сравниваем эти две вещи.

---

<sup>5</sup> **Записать через интегралы.**

<sup>6</sup> Если что-то измерено в шкале от 0 до 1 млн, а что-то от 0 до 100, то результат АГК будет странным :) Всегда первая главная компонента будет там, где миллионы.

## 2.4. Чему АГК соответствует на статистическом языке?

Перейдём к  $\mathbb{X} \in \mathbb{R}^{n \times p}$  ( $L \rightarrow p$  — количество индивидов,  $K \rightarrow n$  — число признаков). Пусть  $\mathbb{X}$  — центрированы. Берём признак, какую норму нужно рассматривать? Вероятностную норму.<sup>7</sup> Что означает характеристика  $\|X_i\|_2^2$  на статистическом языке?

$$\|X_i\|_2^2 = \frac{1}{n} \sum_{j=1}^n ((X_i)_j - \bar{X}_i)^2 = S^2(X_i) — \text{выборочная дисперсия вектора } X_i.^8$$

Что означает матрица данных  $\|X_i\|_{1,2}^2$  на статистическом языке? Используем верхнюю строчку и предполагаем, что  $\mathbb{X}$  центрированы.

$$\|X_i\|_{1,2}^2 = \sum_{i=1}^p \|X_i\|_2^2 = \sum_{i=1}^p S^2(X_i) — \text{total variance}.^9$$

Посчитаем норму вектора главных компонент (считаем, что АГК:  $\mathbb{Y} = \sum_{i=1}^p \sqrt{\lambda_i} U_i V_i^T$ ):

$$Z_i = \mathbb{X} U_i = \sqrt{\lambda_i} V_i, \text{ где } Z_i — \text{проекция на } i\text{-ое направление и } i = 1 \dots, n.$$

Знаем, что  $\|Z_i\|_2^2 = S^2(Z_i)$ . Учитывая, что  $V_i$  нормированы, то

$$\|Z_i\|_2^2 = S^2(Z_i) = \lambda_i.$$

## 2.5. АГК на выборочном языке (продолжение)

---

<sup>7</sup> У длинных векторов вторая норма.