

## Метод главных компонент

- один из основных методов уменьшения размерности данных с сохранением максимальной информации

Постановка задачи:

У нас есть  $n$  признаков объектов  $f_j(x)$ ,  $j = \overline{1, n}$   
 объект  $i$ -й выборки будем считать объектом с  $n$  признаковыми описаниями:  $x_i = (f_1(x_i), \dots, f_n(x_i))$ ,  $i = \overline{1, L}$

Рассмотрим матрицу  $F$ :

$$F_{L \times n} = \begin{bmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \vdots & & \vdots \\ f_1(x_L) & \dots & f_n(x_L) \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_L \end{bmatrix}$$

У  $z_i = (g_1(x_i), \dots, g_m(x_i))$  - признаковые описания тех же объектов  $z \in \mathbb{R}^m$ ,  $m < n$

$$G_{L \times m} = \begin{bmatrix} g_1(x_1) & \dots & g_m(x_1) \\ \vdots & & \vdots \\ g_1(x_L) & \dots & g_m(x_L) \end{bmatrix} = \begin{bmatrix} z_1 \\ \vdots \\ z_L \end{bmatrix}$$

У  $U: \hat{x} = zU^T$

$$\Delta^2(G, U) = \sum_{i=1}^L \|\hat{x}_i - x_i\|^2 = \sum_{i=1}^L \|z_i U^T - x_i\|^2 = \|GU^T - F\|^2 \rightarrow \min_{G, U}$$

$$\text{rg } G = \text{rg } U = m \leq \text{rg } F$$

Т Если  $m \leq \text{rg } F$ , то минимум  $\Delta^2(G, U)$  достигается, когда столбцы матрицы  $U$  есть с.л.  $F^T F$ , соответствующие  $m$  максимальным с.л.  
 При этом  $G = FU$ ,  $U$  и  $G$  - ортогональны

Решение:

$$\begin{cases} \frac{\partial \Delta^2}{\partial G} = (GU^T - F)U = 0 \\ \frac{\partial \Delta^2}{\partial U} = G^T(GU^T - F) = 0 \end{cases} \quad \text{т.к. } G \text{ и } U \text{ - независимы} \Rightarrow \begin{cases} G = FU(U^T U)^{-1} \\ U = F^T G(G^T G)^{-1} \end{cases}$$

$R: GU^T = (GR)(R^T U^T)$ ,  $R$  обратимая, это  $U^T U$  и  $G^T G$  - диагональные  
 $G^T G = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$   
 $U^T U = I_m$

$\Rightarrow \begin{cases} G = FU \\ U\Lambda = F^T F U \end{cases} \Rightarrow U\Lambda = F^T F U$ , т.е. столбцы  $U$  - с.л. матрицы  $F^T F$ ,  
 значения  $\lambda_1, \dots, \lambda_m$  - соответствующие им с.л.  
 (при этом соответствующие значения  $\lambda_j$ )

$$\begin{aligned} \Rightarrow \Delta^2(G, U) &= \|F - GU^T\|^2 = \text{tr}(F^T - UG^T)(F - GU^T) = \text{tr} F^T(F - GU^T) = \text{tr} F^T F - \text{tr} F^T G U^T = \\ &= \|F\|^2 - \text{tr} U\Lambda U^T = \|F\|^2 - \text{tr } \Lambda = \sum_{j=1}^n \lambda_j - \sum_{j=1}^m \lambda_j = \sum_{j=m+1}^n \lambda_j \end{aligned}$$

где  $\lambda_1, \dots, \lambda_n$  - все с.л. матрицы  $F^T F$

минимум  $\Delta^2$  достигается, когда  $\lambda_1, \dots, \lambda_m$  - наибольшие из  $n$  с.л.

с.л.  $u_1, \dots, u_m$ , соответствующие этим с.л. - главные компоненты

Геометрический смысл:

Метод аппроксимирует  $n$ -мерное облако точек  $n$ -мерным эллипсоидом, получая которого будем считать матрицей компонент.



## SVD - разложение

- рекурсивный функциональный оператор с целью преобразовать ее к каноническому виду.

I  $\exists A$  - операторная матрица из  $M_{m,n}$

$\exists$  ОНБ  $\{e^k\}_{k=1}^n \subset \mathbb{C}^n, \{q^k\}_{k=1}^m \subset \mathbb{C}^m$  и  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r, 0 \leq r \leq \min\{m, n\}, \sigma_i \geq 0, i = \overline{1, r}$

$$Ae^k = \begin{cases} \sigma_k q^k, & k \leq r \\ 0, & k > r \end{cases} \quad (1)$$

числа  $\sigma_1, \sigma_2, \dots, \sigma_r$  - сингулярные числа матрицы  $A$ .

базисы  $\{e^k\}_{k=1}^n, \{q^k\}_{k=1}^m$  - сингулярные базисы матрицы  $A$

$r$  - ранг матрицы  $\text{Rg}(A)$ , где  $r = \text{rg}(A)$

л-то:

Матрица  $A^*A$  самосопр. и неотрицательна

$$((A^*A)^* = A^*A, (A^*Ax, x) = (Ax, Ax) \geq 0 \quad \forall x \in \mathbb{C}^n)$$

$\Rightarrow \exists$  ОНБ с.л.  $\{e^k\}_{k=1}^n$  матрицы  $A^*A$ . Все ее собственные числа  $\geq 0$

$$A^*Ae^k = \sigma_k^2 e^k, \quad k = \overline{1, n}, \quad \sigma_k^2 \geq 0, \quad k = \overline{1, n}$$

$$\exists \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0, \quad \sigma_{r+1} = \dots = \sigma_n = 0$$

$$\exists z^k = Ae^k, \quad k = \overline{1, r}$$

$$(z^p, z^q) = (Ae^p, Ae^q) = (A^*A e^p, e^q) = \sigma_p^2 (e^p, e^q)$$

$$\Rightarrow (z^p, z^q) = \begin{cases} 0 & p \neq q \\ \sigma_p^2 & p = q \end{cases}$$

$\Rightarrow$  канонич.  $q^k = \sigma_k^{-1} Ae^k, \quad k = \overline{1, r}$  образуют ОНБ системы  $\mathcal{L} \subset \mathbb{C}^m$

(если  $r < m$ , дополним ее функциональными кан. форм.  $q^k, k = \overline{r+1, m}$  до ОНБ  $\mathbb{C}^m$ )

Матричное представление:

представит (1) в виде  $A = V \Sigma W^*$

$V$  - матрица со столбцами  $\{q^k\}_{k=1}^m$

$W$  - - - -  $\{e^k\}_{k=1}^n$

$\Sigma$  - диаг. матрица:

$$\Sigma = \begin{bmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ 0 & & \sigma_r & \\ & & & \ddots & \\ & 0 & & 0 & \dots & 0 \end{bmatrix}$$

Геометрический смысл:

$\exists A$  - линейный оператор. Мы можем разбить его на более простые операторы

сжатия, растяжения. Компоненты сингулярного разложения показывают для преобразования