

1. **Analiza glavnih komponent (PCA).** Podatke (vrstične vektorje) $\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_n^\top$ predstavimo kot vrstice matrice

$$X = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \in \mathbb{R}^{n \times d}.$$

Komponente vsakega vektorja \mathbf{x}_i^\top si lahko predstavljamo kot (nabor) značilnosti opazovanih objektov. Stolpcem \mathbf{c}_j matrice $X = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_d]$ pogosto pravimo *značilni vektorji* oz. *vektorji značilk*.

Cilj naloge je poiskati t.i. *glavne komponente* $\mathbf{y}_1, \dots, \mathbf{y}_d \in \mathbb{R}^n$, ki so nekorelirane projekcije podatkov \mathbf{x}_i^\top na enotske vektorje $\mathbf{v}_1^\top, \dots, \mathbf{v}_d^\top$, ki maksimizirajo variance $\text{var}(\mathbf{y}_i)$. Oporne točke so:

- *Centralizacija podatkov:* Od vsakega stolpca X odštejemo njegovo srednjo vrednost

$$\bar{X} := X - [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_d],$$

kjer je $\boldsymbol{\mu}_j = \mu_j [1, \dots, 1]^\top$, μ_j pa je srednja vrednost komponent značilnega vektorja \mathbf{c}_j .

- *Izračun singularnega razcepa \bar{X} :* $\bar{X} = USV^\top$, kjer je $U = [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{R}^{n \times n}$, $V = [\mathbf{v}_1, \dots, \mathbf{v}_d] \in \mathbb{R}^{d \times d}$ ter $S \in \mathbb{R}^{n \times d}$ diagonalna matrika s singularnimi vrednostmi $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d$ na diagonali.
- *Glavne komponente X so $\mathbf{y}_1, \dots, \mathbf{y}_d \in \mathbb{R}^n$, ki jih dobimo kot*

$$\mathbf{y}_j = \bar{X} \mathbf{v}_j = \sigma_j \mathbf{u}_j.$$

Odgovori na spodnja vprašanja.

- Naj bo $\Sigma = \frac{1}{n-1} \bar{X}^\top \bar{X}$. Utemelji, da za $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$ velja $\text{cov}(X\mathbf{v}, X\mathbf{w}) = \mathbf{v}^\top \Sigma \mathbf{w}$.
- Kako izrazimo $\text{var}(\mathbf{y}_j) := \text{cov}(\mathbf{y}_j, \mathbf{y}_j)$ s singularnimi vrednostmi matrice \bar{X} ?
- Izračunaj $\text{cov}(\mathbf{y}_j, \mathbf{y}_k)$ za $j \neq k$.

Napiši tri Octave funkcije:

- $[\boldsymbol{\mu}, \mathbf{V}_k, \mathbf{U}_k, \mathbf{D}_k] = \text{pca}(X, k)$, ki za dano matriko s podatki X in za celo število k , $0 \leq k \leq \min(n, d)$, vrne srednje vrednosti $\boldsymbol{\mu}$, matriki \mathbf{V}_k ter \mathbf{U}_k s prvimi k levimi oz. desnimi glavnimi smermi in vektor \mathbf{D}_k prvih k varianc $\text{var}(\mathbf{y}_j)$,
- $Z = \text{proj}(X)$, ki za matriko s podatki X vrne projekcijo $\mathbf{x}_i^\top - [\mu_{i1}, \dots, \mu_{id}]$ na največji dve desni glavni smeri in izriše obe glavni komponenti ter projekciji podatkov na isto sliko,
- $r = \text{threshold}(X, p)$, ki za matriko X ter za število $p \in [0, 1]$ vrne najmanjše celo število r , za katerega velja

$$\frac{\text{var}(\mathbf{y}_1) + \dots + \text{var}(\mathbf{y}_r)}{\text{var}(\mathbf{y}_1) + \dots + \text{var}(\mathbf{y}_d)} \geq p.$$