

TD 02 : Analyse en composantes principales

Fondements du Machine Learning, L3 IM2D

Octobre 2024



Exercice 2.1: Double ACP (Examen 2020/2021)

Dans cet exercice, on considère deux tableaux de données $\mathbf{X} \in \mathbb{R}^{m_1 \times n}$ et $\mathbf{W} \in \mathbb{R}^{m_2 \times n}$ représentant des exemples, ou individus, dans un espace à n dimensions. On note ainsi

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_{m_1}^T \end{bmatrix} \quad \text{et} \quad \mathbf{W} = \begin{bmatrix} \mathbf{w}_1^T \\ \vdots \\ \mathbf{w}_{m_2}^T \end{bmatrix}.$$

- a) On considère d'abord les matrices de données individuellement.
- Donner la formule de définition des individus moyens des tableaux de données \mathbf{X} et \mathbf{W} , que l'on notera $\bar{\mathbf{x}}$ et $\bar{\mathbf{w}}$.
 - Donner les matrices de covariance empirique respectives des tableaux \mathbf{X} et \mathbf{W} , notées $\Sigma_{\mathbf{X}}$ et $\Sigma_{\mathbf{W}}$.
 - Que représentent les vecteurs propres d'une matrice de covariance empirique dans le contexte de l'analyse en composantes principales ?
- b) On considère maintenant le tableau de données $\mathbf{V} = \begin{bmatrix} \mathbf{X} \\ \mathbf{W} \end{bmatrix} \in \mathbb{R}^{(m_1+m_2) \times n}$.
- Exprimer l'individu moyen de ce tableau, noté $\bar{\mathbf{v}}$, en fonction des formules de $\bar{\mathbf{x}}$ et $\bar{\mathbf{w}}$ établies en question a-i).
 - Donner une expression de la matrice de covariance $\Sigma_{\mathbf{V}}$ de la matrice \mathbf{V} .
- c) On suppose enfin qu'il existe $\alpha \in \mathbb{R}$ tel que $\mathbf{X} = \alpha \mathbf{W}$. Proposer une interprétation des valeurs suivantes :
- Valeur de $\bar{\mathbf{v}}$ lorsque $\alpha = -1$;
 - Valeur de $\Sigma_{\mathbf{V}}$ lorsque $\alpha = 1$.

Exercice 2.2: SVD et ACP (Examen 2020/2021)

Soit une matrice $\mathbf{X} \in \mathbb{R}^{m \times n}$ de rang $n < m$. On supposera que \mathbf{X} est la concaténation de m vecteurs de \mathbb{R}^n notés $\mathbf{x}_1, \dots, \mathbf{x}_m$, de sorte que $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix}$.

- On considère une décomposition en valeurs singulières réduite $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, où les valeurs singulières sont ordonnées par ordre décroissant. Donner alors les tailles des matrices \mathbf{U} , $\mathbf{\Sigma}$, \mathbf{V} ainsi que leurs caractéristiques.
- Rappeler la définition de l'individu moyen $\bar{\mathbf{x}}$ correspondant au nuage de points $\{\mathbf{x}_i\}_{i=1}^m$.
- On suppose d'abord que $\bar{\mathbf{x}} \neq \mathbf{0}_{\mathbb{R}^n}$. Expliquer comment modifier \mathbf{X} de sorte à obtenir un nouveau nuage de points (noté \mathbf{X}^c) d'individu moyen $\mathbf{0}_{\mathbb{R}^n}$.
- On suppose maintenant que $\bar{\mathbf{x}} = \mathbf{0}_{\mathbb{R}^n}$.
 - Quelle est l'implication de ce résultat pour le nuage de points ?
 - Donner la définition de la matrice de covariance empirique associée à \mathbf{X} , notée $\mathbf{\Sigma}_X$, et justifier que $\mathbf{\Sigma}_X = \frac{1}{m-1}\mathbf{X}^T\mathbf{X}$.
 - En utilisant la SVD de \mathbf{X} de la question a), écrire une SVD de $\mathbf{\Sigma}_X$.
 - Justifier que toute décomposition en valeurs propres de $\mathbf{\Sigma}_X$ en est également une SVD.
 - Que représentent les coefficients diagonaux de $\mathbf{\Sigma}_X$?

Exercice 2.3: ACP et données augmentées (Examen 2019/2020)

Soit $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix} \in \mathbb{R}^{m \times n}$ représentant n caractéristiques de m individus $\mathbf{x}_1, \dots, \mathbf{x}_m$.

- Rappeler la formule de calcul de l'individu moyen, noté $\bar{\mathbf{x}}$, ainsi que de la matrice de covariance empirique, notée $\mathbf{\Sigma}$, pour le tableau de données \mathbf{X} .
- Rappeler le lien entre la première composante principale de \mathbf{X} et la matrice $\mathbf{\Sigma}$. Comment cette composante peut-elle être utilisée pour approcher le tableau de données \mathbf{X} ?
- Dans le cas où $1 \leq m \ll n$, la matrice $\mathbf{\Sigma}$ peut être de taille trop grande pour pouvoir être stockée et manipulée numériquement. Quelle technique vue en TP peut-on alors utiliser pour obtenir (entre autres) la première composante principale ?
- On considère maintenant le tableau de données augmenté avec l'individu moyen, c'est-à-dire la matrice $\bar{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \bar{\mathbf{x}}^T \end{bmatrix} \in \mathbb{R}^{(m+1) \times n}$. Calculer l'individu moyen ainsi que la matrice de covariance de ce nouveau tableau en fonction de $\bar{\mathbf{x}}$ et $\mathbf{\Sigma}$. Comment interprétez-vous ces résultats ?