

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 Σ_X et Σ_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 Σ_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 Σ_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 ?

Solutions

Solution de l'exercice 2.1: Double ACP

a) Question de cours. Pour le iii), il s'agit des axes principaux.

b) On a

$$\bar{\mathbf{v}} = \frac{1}{m_1 + m_2} \left(\sum_{i=1}^{m_1} \mathbf{x}_i + \sum_{i=1}^{m_2} \mathbf{w}_i \right) = \frac{m_1}{m_1 + m_2} \bar{\mathbf{x}} + \frac{m_2}{m_1 + m_2} \bar{\mathbf{w}}.$$

c)

$$\Sigma_V = \frac{1}{m_1 + m_2 - 1} \left(\sum_{i=1}^{m_1} (\mathbf{x}_i - \bar{\mathbf{v}})(\mathbf{x}_i - \bar{\mathbf{v}})^T + \sum_{i=1}^{m_2} (\mathbf{w}_i - \bar{\mathbf{v}})(\mathbf{w}_i - \bar{\mathbf{v}})^T \right)$$

d) Pour le cas i), on obtient que $\bar{\mathbf{v}} = \mathbf{0}$, car $m_1 = m_2$ et $\mathbf{x}_i = -\mathbf{w}_i$ pour tout i . Le jeu de données est centré en l'origine, ce qui est logique puisqu'il ne consiste que de paires de vecteurs opposés. Pour ii), on obtient que $\Sigma_V = \frac{m-1}{2m-1} \Sigma_X$, ce qui veut notamment dire que les variances seront plus faibles. Là encore, c'est logique parce qu'on a dupliqué le jeu de données en combinant \mathbf{X} et \mathbf{W} .

Solution de l'exercice 2.2: SVD et ACP (Examen 2020/2021)

a) $\mathbf{U} \in \mathbb{R}^{m \times m}$ orthogonale, $\Sigma \in \mathbb{R}^{m \times n}$ est nulle sauf potentiellement sur la diagonale, où les coefficients sont $\sigma_1 \geq \dots \geq \sigma_{\min(m,n)}$.

b) $\bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$.

c) En centrant les données, on définit un nouveau jeu de données $\mathbf{X}^c = \mathbf{X} - \mathbf{1}_m \bar{\mathbf{x}}^T$.

d) On suppose que $\bar{\mathbf{x}} = \mathbf{0}_{\mathbb{R}^n}$.

i) Le jeu de données est centré en $\mathbf{0}_{\mathbb{R}^n}$.

ii) Par définition, la matrice de covariance empirique est égale à $\frac{1}{m-1} \mathbf{X}^c (\mathbf{X}^c)^T$, où \mathbf{X}^c est le jeu de données centrées. Comme on a ici $\mathbf{X}^c = \mathbf{X}$, on en conclut que $\frac{1}{m-1} \mathbf{X}^c (\mathbf{X}^c)^T = \frac{1}{m-1} \mathbf{X} \mathbf{X}^T$.

iii) Immédiat, la seule différence avec le cas général étant que les deux dimensions sont égales.

iv) Σ_X est symétrique semi-définie positive, donc il existe une base orthonormée $\mathbf{P} \in \mathbb{R}^{n \times n}$ telle que $\Sigma_X = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^T$, où $\mathbf{\Lambda}$ est une matrice diagonale à coefficients positifs ou nuls (ordonnés par ordre décroissant sans perte de généralité). Cette décomposition définit aussi une SVD de Σ_X .

v) Les coefficients diagonaux de Σ_X représentent les variances empiriques des attributs.

Solution de l'exercice 2.3: ACP et données augmentées (Examen 2019/2020)

- a) Cours
- b) Cours
- c) Au lieu de considérer $\Sigma = \frac{1}{m-1} \mathbf{X}^c (\mathbf{X}^c)^T$, on peut considérer la matrice $\frac{1}{m-1} (\mathbf{X}^c)^T \mathbf{X}^c$, qui est une matrice de taille $m \times m$, donc avec beaucoup moins de coefficients. Pour autant, les deux matrices ont les mêmes valeurs propres.
- d) On trouve que rajouter la moyenne réduit les valeurs propres de la matrice de covariance, et biaise les données vers $\bar{\mathbf{x}}$. L'individu moyen ne change pas quand on ajoute l'individu moyen. En effet, l'individu moyen associé à $\bar{\mathbf{X}}$ est

$$\frac{1}{m+1} \left(\sum_{i=1}^m \mathbf{x}_i + \bar{\mathbf{x}} \right) = \frac{m}{m+1} \bar{\mathbf{x}} + \frac{1}{m+1} \bar{\mathbf{x}} = \bar{\mathbf{x}}.$$

Par ailleurs, la matrice de covariance associée à $\bar{\mathbf{X}}$ est

$$\begin{aligned} \Sigma_{\bar{\mathbf{X}}} &= \frac{1}{(m+1)-1} \left(\sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T + (\bar{\mathbf{x}} - \bar{\mathbf{x}})(\bar{\mathbf{x}} - \bar{\mathbf{x}})^T \right) \\ &= \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \\ &= \frac{m-1}{m} \Sigma_{\mathbf{X}}. \end{aligned}$$

Si l'on considère les valeurs propres de $\Sigma_{\bar{\mathbf{X}}}$, on voit que celles-ci sont plus petites (et ≥ 0) que celles de $\Sigma_{\mathbf{X}}$, et donc l'inertie du nuage de points est réduite.