

PCA

On se donne une matrice $\mathbf{X} = (x_{i,j})$ avec $1 \leq i \leq n$ le nombre d'observations et $1 \leq j \leq p$ le nombre de variables.

Matrice de pondération des observations : D_w une matrice diagonale telle que $\text{tr}(D_w) = 1$. Généralement, chaque observation possède le même poids, on a donc $D_w = \frac{1}{n} I_n$.

Il est nécessaire de centrer \mathbf{X} : $\mathbf{Y} = (y_{i,j}) = (x_{i,j} - \bar{x}_j)$ avec $\bar{x}_j = \sum_{i=1}^n w_i x_{i,j}$.

Matrice des covariances empiriques : $\mathbf{V} = (s_{i,j})$ avec $s_{i,j} = \sum_{k=1}^n w_k (x_{k,i} - \bar{x}_i)(x_{k,j} - \bar{x}_j) = \sum_{k=1}^n w_k y_{k,i} y_{k,j}$.

Matrice des données réduites et standardisées : $\mathbf{Z} = (z_{i,j})$ avec $(z_{i,j}) = \frac{y_{i,j}}{s_j}$.

Matrice des coefficients de corrélation : $\mathbf{R} = (r_{i,j})$ avec $(r_{i,j}) = \frac{s_{i,j}}{s_i s_j}$.

Remarque : \mathbf{V} et \mathbf{R} sont des matrices carrées de dimension p symétriques et semi-définies positives. Il y a deux types d'APC. L'APC canonique qui utilise \mathbf{Y} et qui effectue son analyse sur \mathbf{V} , elle est utilisée généralement quand les variables sont dans la même unité. L'APC standardisé qui utilise \mathbf{Z} et qui effectue son analyse sur $\mathbf{V} = \mathbf{R}$.

Produit scalaire sur \mathbb{R}^p : on se donne une matrice \mathbf{M} symétrique et définie positive. Alors $\langle u, v \rangle = u^T M v$.

Pour un nuage d'observations (\mathbb{R}^p), on prend $\mathbf{M} = I_p$. Pour un nuage de variables (\mathbb{R}^n), on prends $\mathbf{M} = D_w$. C'est ce dernier cas qui nous interesse. On a alors : $\langle y_j, y_k \rangle = s_{j,k}$, $\|y_j\|^2 = s_j^2$, $d^2(y_j, y_k) = \text{var}(y_j - y_k)$ et $\cos(y_j, y_k) = r_{j,k}$.

Inertie totale du nuage d'observation : $I(\mathbf{X}) = \sum_{i=1}^n w_i d^2(x_i, g) = \sum_{i=1}^n w_i \|x_i - g\|^2 = \sum_{i=1}^n w_i \|y_i\|^2 = I(\mathbf{Y})$.

Remarques : l'inertie correspond à la variance au cas multivarié, $I(\mathbf{Z}) \neq I(\mathbf{Y})$ et $\mathbf{I} = \text{tr}(\mathbf{V})$.

Inertie expliquée par un sous espace vectoriel F (de dimension $q < p$) de \mathbb{R}^p : $I_F(\mathbf{Y}) = \sum_{i=1}^n w_i \|P(y_i)\|^2$ avec P la projection orthogonale sur le sous espace.

Cas particulier en dimension 1 : si $F = \text{vect}(u)$ avec $\|u\| = 1$, on a $I_F(\mathbf{Y}) = u^T V u$.

On construit par récurrence E_q notre sous espace vectoriel : $E_q = \Delta u_1 \oplus \Delta u_2 \oplus \dots \oplus \Delta u_q$, avec u_i le vecteur normalisé associé au i -ème plus grand vecteur propre de \mathbf{V} λ_i de la manière suivante : $I_{\Delta u_i}(\mathbf{Y}) = \lambda_i$.

Nombre d'axes à retenir :

Critère d'inertie : $\frac{I_{E_q}}{I} = \frac{\sum_{i=1}^q \lambda_i}{\sum_{i=1}^p \lambda_i} \geq \alpha$ avec en général $\alpha = 0.8$.

Règle de Kaiser : $\lambda_q \geq \frac{I}{p}$ et $\lambda_{q+1} < \frac{I}{p}$.

Recherche d'un coude en faisant un graphique des valeurs propres dans l'ordre décroissant.

On note $c_j = Y u_j$ la projection de \mathbf{Y} sur u_j . On nomme ce vecteur une composante principale. Sa moyenne est nulle et sa variance est λ_j .

Clustering

On a la même matrice \mathbf{X} et on essaie de partitionner les observations $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n\}$.

Distance de Minkosky : distance en utilisant la norme L_q . Distance Euclidienne quand $q = 2$ et distance de Manhattan quand $q = 1$. Distance de Chebyshev en utilisant la norme infinie. En utilisant la norme définie par le produit scalaire, on a la distance de Mahalanobis en prenant $M = V^{-1}$.

Pour les distances discrètes, la distance de Hamming est le nombre de différence, la distance de Levenshtein est le nombre d'ajout, de suppression et de remplacement. Enfin, en notant $a_{i,j}$ le nombre d'éléments en communs entre x_i et x_j , $b_{i,j}$ ceux qui apparaissent dans x_i mais pas dans x_j , ... On a le coefficient de simple matching $\frac{a_{i,j} + d_{i,j}}{p}$ et l'indice de Jaccard

$$\frac{a_{i,j}}{a_{i,j} + b_{i,j} + c_{i,j}}.$$

Méthodes combinatoires

On cherche une partition de \mathcal{X} en K clusters en minimisant une fonction de coût définie à partir d'un critère de similarité/dissimilarité. On notera g_k le centroïde du cluster C_k .

On a : $I = \sum_{i=1}^n \|x_i - g\|^2 = \sum_{k=1}^K \sum_{x_i \in C_k} d^2(x_i, g_k) + \sum_{k=1}^K \#C_k d^2(g_k, g) = I_w + I_B$ avec I_w l'inertie intra-classes qu'on cherche à minimiser et I_B l'inertie inter-classes qu'on cherche à maximiser.

Algorithme des K moyennes : convergence garantie mais seulement vers un min local, faible complexité, facile à interpréter, sensible à l'initialisation, adapté aux clusters de formes convexes et de tailles équilibrées. Par contre, il est sensible aux valeurs aberrantes, ne trouve que des clusters convexes, et on doit fixer K en avance.

Clustering hiérarchique agglomératif (AHC) : bottom up (ce qui est utilisé en général) ou top down. Il faut définir la distance inter classe pour l'utiliser. **Simple linkage** : distance entre les deux points les plus proches. Tendence à produire des clusters larges, sensible au bruit et aux valeurs aberrantes. **Complete linkage** : distance entre les deux points les plus éloignés. Tendence à produire des clusters compacts, sensible au bruit et aux valeurs aberrantes. **Average linkage** : moyenne des distances entre les points de C_a et de C_b . Tendence à construire des clusters plus homogènes, moins sensible aux valeurs aberrantes. **Centroïde linkage** : $d(C_a, C_b) = d(g_a, g_b)$. Moins sensible aux valeurs aberrantes.

Le nombre K de clusters n'est pas requis à l'avance, il suffit de faire un dendrogramme. L'algorithme est totalement déterministe. La complexité est élevée.

Pour trouver K , on peut tracer I_w en fonction de K et trouver le coude. On peut aussi chercher le score silhouette le plus élevé : $s(x_i) = \frac{b(x_i) - a(x_i)}{\max(a(x_i), b(x_i))}$ avec $a(x_i)$ la distance moyenne du point à son groupe et $b(x_i)$ la distance moyenne du point au groupe voisin.

Méthode basée sur la densité On veut trouver des clusters de formes complexes et qu'ils ne soient pas influencés par les valeurs aberrantes. On pose ϵ et n_{min} . On cherche les points centraux : ceux dont le nombre de points dans leur ϵ -voisinage est $\geq n_{min}$. tous les points dans ces ϵ -voisinages appartiennent aux clusters définis par les points centraux. On relie les clusters dont l'intersection n'est pas nulle. En général, on prend $n_{min} = 2p$. Pour trouver ϵ , on regarde pour tous les points la distance avec leur n_{min} i-ème voisin, on trace dans l'ordre décroissant et on prend le coude.

Classification

Classification discriminante : on cherche une ou plusieurs frontières, $P(y|x)$.

Classification générative : on cherche les clusters, $P(x|y)$.

Techniques d'encodage : one hot encoding, label encoding et ordinal encoding (on préserve un ordre).

Méthode des ensembles

Bagging (Bootstrap Aggregating) : on entraîne k modèles sur une partition du jeu d'entraînement de cardinal k . Ensuite, on utilise ces modèles pour déterminer la classe finale d'une observation (on utilise souvent le vote majoritaire).

Boosting : à chaque itération, on entraîne un modèle en augmentant le poids des erreurs du modèle précédent.

Stacking : on utilise un méta-modèle qui prend en entrée la sortie de plusieurs modèles simples et qui détermine la classe finale.

Interprétabilité

SHAP (SHapley Additive exPlanation) : attribue une valeur d'importance à chaque caractéristique d'entrée, indiquant son impact sur la prédiction finale.

LIME (Local Interpretable Model-agnostic Explanation) : approxime localement le comportement d'un modèle complexe par un modèle interpretable.

Optimisation des hyper-paramètres

Grid search : brut force

Random search : on regarde aléatoirement les combinaisons d'hyper-paramètres et on choisit la meilleure.

Bayesian optimisation : construit un modèle probabiliste pour trouver la meilleure combinaison d'hyper-paramètres.

Modèles non linéaires Arbre de décision : on choisit la variable v qu'on veut utiliser en label. Pour chaque variable v_i , on la décompose en ses sous ensembles se_j . On isole ces sous ensembles et on calcule l'indice de Gini pour les sous ensembles associés de v : $1 - \sum p_k^2$. On pondère ces indices par la répartition des sous ensembles se_j . On choisit le v_i avec le score le plus faible.

t-SNE : modèle de réduction de dimension utilisé pour la visualisation.

Modèle linéaire

Régression logistique : donne la probabilité d'appartenir à une certaine classe. La fonction d'activation est la fonction sigmoïde : $\sigma(z) = \frac{1}{1+e^{-z}}$ avec $z = \beta_0 + \sum \beta_i x_i$. On doit maximiser la fonction de vraisemblance :

$$L(\beta) = \sum ((y_i \log(\sigma(z_i)) + (1 - y_i) \log(1 - \sigma(z_i))).$$

Perceptron : $w_i \leftarrow w_i + \alpha(y - \hat{y})x_i$.

Large Margin Classifier : on cherche à maximiser la marge entre les classes. Equation de l'hyperplan séparateur : $w^T x + b = 0$.

On minimise $\|w\|$ sous contraintes $y_i(w^T x_i + b) \geq 1$ pour tous $i \in \llbracket 1 ; n \rrbracket$.

SVM : Large Margin Classifier + kernel trick pour créer des frontières non linéaires. Possibilité de classification multiclass.

Régression

Modèle linéaire simple

$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$, avec $\epsilon_i \sim N(0, \sigma^2)$. On veut minimiser $L(\beta_0, \beta_1) = \sum \epsilon_i^2 = \sum (y_i - \beta_0 - \beta_1 x_i)^2$. Solution : $\hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$ et $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$.

Modèle linéaire multiple

$\mathbf{Y} = \mathbf{X}\beta + \epsilon$. On veut minimiser $L(\beta) = \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|^2 = \frac{1}{2} \sum_i (y_i - \beta_0 - \sum_j \beta_j x_{i,j})^2$. Si $\mathbf{X}^T \mathbf{X}$ est inversible, alors la solution est $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{X}^\dagger \mathbf{Y}$ où \mathbf{X}^\dagger est la pseudo-inverse de Moore-Penrose de \mathbf{X} . On a alors $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} = \mathbf{X}\mathbf{X}^\dagger \mathbf{Y} = \mathbf{P}\mathbf{Y}$ avec \mathbf{P} le projecteur sur l'espace généré par les variables de \mathbf{X} .

Modèle non linéaire

$\mathbf{Y} = \Phi\beta + \epsilon$ avec Φ la matrice où les observations \mathbf{x}_i sont remplacées par $\phi(\mathbf{x}_i)$.

Approximation polynomiale : $\phi(\mathbf{x}) = [1, \mathbf{x}, \mathbf{x}^2, \dots, \mathbf{x}^{d^*}]$ avec d^* le degré du polynôme voulu.

Approximation gaussienne : $\phi_j(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mu_j\|^2}{2s^2}\right)$, où μ_j sont les centres et s est l'écart-type.

Régression régularisée

L'objectif est de lisser la régression pour éviter l'overfitting. La fonction objectif est transformée de la manière suivante :

$$\mathcal{L}_q(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \alpha \|\beta\|_q^q \text{ où } \alpha \text{ est le paramètre de régularisation et } \|\beta\|_q = \left(\sum_{i=1}^d |\beta_i|^q\right)^{1/q}.$$

Lasso regression : on prend $q = 1$. Ce lissage change les poids des coefficients de manière à ce que certains deviennent négligeables. Il permet donc de savoir quelles variables ne sont pas très utiles, et ainsi pouvoir réduire la dimension pour optimiser la complexité.

Ridge regression : on prend $q = 2$. Ce lissage ne va pas rendre en général des variables négligeables devant d'autres. Cependant, le résultat sera un peu plus proche du résultat sans lissage, ce qui permet de mieux tendre vers la solution.