

## 19 Réduction dimensionnelle : cas de l'Analyse en Composantes Principales, ACP (Principal Components Analysis, PCA)

L'analyse en composantes principales (en anglais, *Principal Components Analysis*, PCA) peut-être vue comme un des nombreux pré-traitements disponibles dans la panoplie du Machine-Learner. Mais c'est un outil en fait beaucoup plus général et puissant, qui est plus ancien que l'avènement des méthodes de ML modernes.

En dépit de son âge, cette méthode n'a pas tellement vieilli: de par son élégance, ou son côté "bien fondée formellement", elle reste d'actualité.

La PCA est une méthode d'**apprentissage non supervisé**. Elle fait partie des méthodes dites de **réduction dimensionnelle** (*dimensional reduction*).

La PCA peut aussi être vue comme une méthode de visualisation de données en grande dimension. Cependant, ce serait ignorer sa puissance de la limiter à ce rôle. D'autres méthodes, moins contrôlées mais assez puissantes, comme *t-SNE*, sont bien adaptées pour la visualisation de données en grande  $D$ .

La PCA peut aussi être vue comme une méthode de compression des données. La reconstruction (décompression) est très facile.

### 19.1 Intuition sur : réduction dimensionnelle

De façon générale (et non seulement pour la PCA), les méthodes de réduction dimensionnelle consistent à réduire la dimension de l'espace d'entrée (l'espace des features, de dimension  $D$ ), vers une autre dimension, plus petite,  $D'$  (ou  $p$ ). Concrètement, cela consiste à chercher une application non injective  $\Phi$ :

$$\Phi : \mathbb{R}^D \longrightarrow \mathbb{R}^{D'}, \text{ avec } D' < D \quad (38)$$

Pourquoi est-ce intéressant? Il y a différentes raisons, qui se recoupent entre elles. Réduire la dimension des données:

- Permet de supprimer les redondances dans l'information d'entrée. Par exemple, lorsque 2 attributs sont identiques en valeur, ou de valeur proportionnelle, il n'y a pas d'intérêt à les garder en doublon (c'est même nuisible, a priori).
- Permet de réduire le "bruit" (c'est le nom qu'on donne aux informations "inutiles"... c'est une vaste question de savoir ce que ça veut dire, précisément, et il est très difficile de distinguer le bruit du signal).
- On évite ainsi la malédiction des grandes dimensions. En effet, en général (mais pas nécessairement), plus les dimensions sont grandes, plus on a de paramètres à ajuster et on se retrouve facilement dans des cas de sur-apprentissage.
- On évite donc le sur-apprentissage: en effet, qui dit moins de dimensions en entrée dit, en général, moins de paramètres dans le modèle.
- Cela aide donc l'apprentissage: en termes de vitesse (a priori, c'est le cas, on peut probablement imaginer des cas pathologiques où ça ralentit l'apprentissage).
- Cela aide donc l'apprentissage: en termes de performance: ça dépend. On peut limiter la performance en supprimant les données, mais on peut aussi l'améliorer par la diminution du sur-apprentissage.

### 19.2 PCA: idée

Ici on ne verra qu'une réduction dimensionnelle: la PCA.

Dans ce cas particulier, l'application  $\Phi$  est une projection (application linéaire) de l'espace  $\mathbb{R}^D$  dans le sous-espace  $\mathbb{R}^{D'}$ . On projette les données sur un hyper-plan de dimension  $D'$ , qui a une certaine orientation dans l'espace parent  $\mathbb{R}^D$ . Dans la figure 8, on voit ce que cela peut signifier, dans un cas à petite dimension, ou on peut encore dessiner. On voit aussi, ou bien on rappelle que, une projection est aussi le choix d'un nombre  $D'$  de combinaisons linéaires des  $D$  attributs d'entrée. Chaque dimension/nouvel attribut, après PCA, sera une combinaison linéaire des attributs d'entrée.

La position relative de l'hyper-plan de projection (espace à  $D'$  dimensions) par rapport à l'origine du repère (dans  $\mathbb{R}^D$ ) est encodée dans le vecteur représentant la moyenne des données,  $\langle \vec{x} \rangle$ . On a donc  $\vec{x}' = (\vec{x} - \langle \vec{x} \rangle)P$  les données après projection ( $P_{(D,D')}$  étant la matrice de projection, ou la matrice des coefficients des combinaisons linéaires, si vous voulez). Dans le cas particulier où  $D' = 1$ ,  $P'$  est un vecteur à  $D$  entrée, et  $x'$  est juste un nombre, qui se calcule ainsi:  $x' = \vec{P} \cdot (\vec{x} - \langle \vec{x} \rangle) = (\vec{x} - \langle \vec{x} \rangle) \cdot \vec{P}$ , où  $\cdot$  représente le produit scalaire.

La décompression s'obtient en appliquant la transposée de la matrice de projection, puis en restaurant la moyenne des données:  $\vec{x}_{decompressed} = (\vec{x}' \cdot P^T) + \langle \vec{x} \rangle$

Pourquoi soustrait-on à chaque attribut sa moyenne sur les données ? La réponse se trouve dans la compression vers  $D' = 0$  dimension: si on souhaite garder 0 attribut, et simplement représenter toutes les données par le même vecteur d'attribut, qui serait comme un descripteur de l'ensemble de jeu de données, alors on prend la moyenne (ça peut se justifier plus mathématiquement). Le fait d'annoter chaque point de donnée avec quelques attributs (il y en aura  $D'$  après compression) permet de tenir compte du fait que les différents points de données ont divers écart à la moyenne.. mais la moyenne est le point de départ, le centre de ces fluctuations. Si les données sont standardisées, leur moyenne est nulle et donc cette remarque perd de son intérêt.

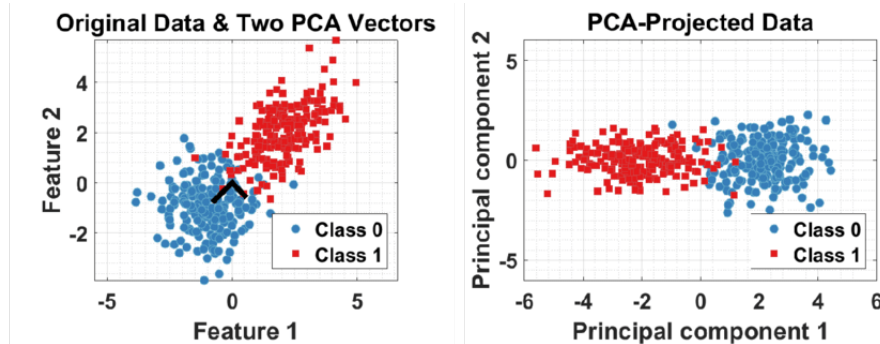


Figure 8: PCA de données à  $D \geq 2$  dimensions (on n'en voit que 2, les 2 premières, mais on peut imaginer qu'il y en a d'autres, orthogonales aux premières) vers  $D' = 2$  dimensions. Si on imagine le cas  $D = 2$ , on fait alors juste une rotation. Si on avait choisi  $D' = 1$ , il ne resterait pour représenter les données que l'axe des abscisses, celui de la première composante: dans l'image de droite, il faudrait imaginer que l'axe des ordonnées a été écrasé, qu'il n'en reste que la valeur moyenne (ou 0).

La question qui se pose maintenant est: comment choisir un bon hyper-plan sur lequel on voudrait projeter ? Le choix fait, dans le cadre de la PCA, est de prendre l'hyperplan **tel que la variance des données, après projection, est maximale**. On favorise donc les dimensions de grande variance (et leurs combinaisons linéaires). On verra que **ce choix est équivalent à un autre choix**, qui serait celui de vouloir **minimiser l'erreur de reconstruction** (en norme L2, c'est-à-dire de minimiser l'erreur quadratique moyenne de reconstruction).

L'idée (naïve) ici est donc de **garder les attributs qui varient le plus** (car se sont ceux qui portent le plus d'information). Par ailleurs, si les données sont très hétérogènes, on peut les standardiser un minimum avant d'effectuer la PCA.

## 19.3 PCA: démonstration du bien-fondé de la méthode

### 19.3.1 Variance des données après projection

Calculons donc la variance des données, après projection. Pour commencer, on ne s'intéresse qu'à une direction de projection. Soit  $\vec{u}_i$  le vecteur qui caractérise l'axe de la projection. Les données projetées deviennent donc de simples nombres, qu'on peut noter  $x'_n = \vec{u}_i^T \vec{x}_n = \sum_d u_{id} x_{nd}$ .

On va prouver que la variance des données le long de cet axe,  $\sigma_i'^2$ , est donnée par:  $\sigma_i'^2 = \vec{u}_i^T C \vec{u}_i$ , où  $C$  est la matrice de covariance des données. (L'idée de la démonstration est de se concentrer sur une direction, ici  $\vec{u}_i$ ). On a déjà parlé de la covariance, mais on est sympa, on rappelle la formule:

$$C = \frac{1}{N} \sum_n (\vec{x}_n - \langle \vec{x} \rangle)(\vec{x}_n - \langle \vec{x} \rangle)^T \quad (39)$$

Pour rappel,  $C$  est clairement une matrice de taille  $(D, D)$  (somme de  $N$  produits matriciels de matrices  $(D, 1)$  par matrices  $(1, D)$ ) et  $\langle \vec{x} \rangle$  représente le vecteur moyen de  $\vec{x}$  (la moyenne empirique). On calcule la variance

des données après projection sur  $\vec{u}_i$ , notée  $\sigma_i'^2$ :

$$\sigma_i'^2 = \frac{1}{N} \sum_n (x'_n - \langle x' \rangle)^2 \quad = \text{Définition de la variance} \quad (40)$$

$$= \frac{1}{N} \sum_n (\vec{u}_i^T \vec{x}_n - \langle \vec{u}_i^T \vec{x} \rangle)^2 \quad (41)$$

$$= \frac{1}{N} \sum_n (\vec{u}_i^T \vec{x}_n - \langle \vec{u}_i^T \vec{x} \rangle)(\vec{u}_i^T \vec{x}_n - \langle \vec{u}_i^T \vec{x} \rangle)^T \quad (42)$$

Où on a utilisé que  $a_i^2 = a_i a_i^T$ , même si  $a_i$  est de taille (1,1).

$$= \frac{1}{N} \sum_n (\vec{u}_i^T \vec{x}_n - \langle \vec{u}_i^T \vec{x} \rangle)(\vec{x}_n - \langle \vec{x} \rangle)^T (\vec{u}_i^T)^T \quad (43)$$

On a utilisé que  $(AB)^T = B^T A^T$ .

$$= \frac{1}{N} \sum_n \vec{u}_i^T (\vec{x}_n - \langle \vec{x} \rangle)(\vec{x}_n - \langle \vec{x} \rangle)^T \vec{u}_i \quad (44)$$

$$= \vec{u}_i^T \frac{1}{N} \sum_n (\vec{x}_n - \langle \vec{x} \rangle)(\vec{x}_n - \langle \vec{x} \rangle)^T \vec{u}_i \quad (45)$$

$$\sigma_i'^2 = \vec{u}_i^T \quad C \quad \vec{u}_i \quad (46)$$

CQFD.

### 19.3.2 Calcul de la meilleure direction

Comme on l'a dit, **le choix fait dans la méthode de la PCA est celui de maximiser la variance des données projetées, donc trouver le  $\vec{u}_i$  qui maximise  $\sigma_i'^2$** . A priori, on pourrait écrire que le meilleur vecteur de projection est donc :

$$(?) \quad \vec{u}_i^* = \operatorname{argmax}_{\vec{u}_i \in \mathbb{R}^D} (\vec{u}_i^T C \vec{u}_i) \quad (?) \quad (47)$$

Cependant, si on résout ce problème d'optimisation, on obtient la solution triviale et non intéressante, qu'il faut prendre un vecteur  $\vec{u}_i$  de la norme la plus grande possible (une matrice de Covariance est par définition semi-définie positive, c.a.d. que ses valeurs propres sont positives ou nulles). Ce n'est pas intéressant, car on fait une projection, ce qui nous intéresse c'est de trouver le bon angle de projection, pas de multiplier nos données par  $\infty$  pour augmenter leur variance.

On décide donc de se restreindre aux vecteurs de norme 1, c'est-à-dire de ne regarder que l'orientation de  $\vec{u}_i$ : on impose  $\vec{u}_i^2 = 1$  On traite donc du problème d'optimisation sous contrainte suivant:

$$\vec{u}_i^* = \operatorname{argmax}_{\vec{u}_i \in \mathbb{R}^D, \vec{u}_i^2=1} (\vec{u}_i^T C \vec{u}_i) \quad (48)$$

Qui peut se réécrire, en attribuant à la contrainte  $\vec{u}_i^2 = 1 \Leftrightarrow (1 - \vec{u}_i^2) = 0$  le multiplicateur<sup>9</sup>  $\lambda$ :

$$\vec{u}_i^* = \operatorname{argmax}_{\vec{u}_i \in \mathbb{R}^D} (\vec{u}_i^T C \vec{u}_i + \lambda(1 - \vec{u}_i^2)) \quad (49)$$

On note:

$$F(\vec{u}_i) = \vec{u}_i^T C \vec{u}_i + \lambda(1 - \vec{u}_i^2) \quad (50)$$

$$\vec{u}_i^* = \operatorname{argmax}_{\vec{u}_i \in \mathbb{R}^D} (F(\vec{u}_i)) \quad (51)$$

Ce problème de maximisation est assez simple, et il se trouve qu'on peut le résoudre exactement. La méthode est classique: on cherche le zéro de la fonction que est dans le argmax, c'est-à-dire on cherche à résoudre  $\vec{\nabla}_{\vec{u}_i} F(\vec{u}_i) = \vec{0}$ . Ce n'est pas évident, mais le gradient de  $\vec{u}^T C \vec{u}$  se calcule, pour une matrice  $C$  symétrique (et donc telle que  $C^T = C$ ):  $\vec{\nabla}_{\vec{u}} (\vec{u}^T C \vec{u}) = 2C\vec{u}$ . En fait, pour faire ce calcul, on repart de la définition initiale,

<sup>9</sup>Ici, le signe à mettre devant la contrainte s'obtient en réfléchissant sur la direction qu'on veut contenir: en fait, ce qu'on veut, c'est que  $(1 - \vec{u}_i^2) \geq 0$ , mais cette inégalité va être saturée car les plus grand  $\vec{u}_i^2$  sont favorisés. D'où l'ajout de  $\lambda(1 - \vec{u}_i^2)$  à la fonction, et non pas son opposé,  $\lambda(\vec{u}_i^2 - 1)$ .

c'est-à-dire de  $\vec{u}^T C \vec{u} = \frac{1}{N} \sum_n (\vec{u}^T \vec{x} - \langle \vec{u}^T \vec{x} \rangle)^2$ . Pour alléger la notation, on va supposer, temporairement, que les données sont centrées, c'est-à-dire que  $\langle \vec{x} \rangle = \vec{0}$ . Ça ne change pas la valeur de la matrice de covariance, et ça allège l'écriture.

$$\vec{\nabla}_{\vec{u}}(\vec{u}^T C \vec{u}) = \vec{\nabla}_{\vec{u}} \frac{1}{N} \sum_n (\vec{u}^T \vec{x}_n - \langle \vec{u}^T \vec{x}_n \rangle)^2 \quad (52)$$

$$= \frac{1}{N} \sum_n \vec{\nabla}_{\vec{u}} (\vec{u}^T \vec{x}_n)^2 \quad (53)$$

$$= \frac{1}{N} \sum_n 2\vec{x}(\vec{u}^T \vec{x}_n)^1 \quad \text{car } \frac{\partial}{\partial s} f^2(s) = 2 \frac{\partial f}{\partial s} f(s) \quad (54)$$

$$\text{Puisque } \vec{u}^T \vec{x} \text{ est un nombre, } \vec{u}^T \vec{x}_n = (\vec{u}^T \vec{x})^T = \vec{x}_n^T \vec{u} \quad (55)$$

$$= \frac{1}{N} \sum_n 2\vec{x}_n(\vec{x}_n^T \vec{u}) \quad (56)$$

$$= 2 \frac{1}{N} \sum_n (\vec{x}_n \vec{x}_n^T) \vec{u} \quad (57)$$

$$= 2C\vec{u} \quad (58)$$

Par ailleurs, il est facile de vérifier que  $\vec{\nabla}_{\vec{u}} \lambda \vec{u}^2 = 2\lambda \vec{u}$ . On peut donc terminer la recherche du  $\text{argmax}(F(\vec{u}_i))$ :

$$\vec{\nabla}_{\vec{u}_i} F(\vec{u}_i) = \vec{0} \Leftrightarrow \quad (59)$$

$$\vec{0} = 2C\vec{u} - 2\lambda \vec{u} \quad (60)$$

$$C\vec{u} = \lambda \vec{u} \quad (61)$$

**C'est l'équation aux valeurs propres !** Comme c'est joli ! C'est un résultat assez profond. Il y a donc **autant de maxima locaux de la fonction  $F(\vec{u}_i)$  qu'il y a de couples  $(\lambda_i, \vec{u}_i)$**  (couples de valeur propre-vecteur propre). La norme des vecteurs propre est imposée par notre contrainte  $\vec{u}^2 = 1$ . La direction (le signe devant  $\vec{u}$ ) est arbitraire, ce qui est normal pour une projection, car une projection dépend de la droite sur laquelle on projette, mais pas de son orientation. Par contre, il y a plein de valeurs propres... laquelle faut-il prendre ? Hé bien, il faut prendre celle qui réalise le maximum global de  $F(\vec{u}_i)$ .

Maintenant qu'on sait qu'on s'intéresse aux vecteurs  $\vec{u}$  de norme 1 et tels que  $C\vec{u} = \lambda \vec{u}$  (c'est-à-dire les vecteurs propres), on peut réécrire  $F$ , en prenant un couple  $(\lambda_i, \vec{u}_i)$  arbitraire:

$$F(\vec{u}_i)_{|\vec{\nabla}_{\vec{u}_i} F(\vec{u}_i) = \vec{0}} = \vec{u}_i^T C \vec{u}_i + \lambda(1 - \vec{u}_i^2) \quad (62)$$

$$= \vec{u}_i^T \lambda_i \vec{u}_i + 0 \quad (63)$$

$$= \lambda_i \vec{u}_i^T \vec{u}_i \quad (64)$$

$$= \lambda_i \quad (65)$$

Plutôt propre, non ? Donc en fait, **la variance des données projetées est exactement égale à la valeur propre correspondant à la direction propre qu'on choisit pour projeter**. Une autre façon de le voir est de juste calculer la variance projetée, sans passer par  $F$ :  $\sigma_i^2 = \vec{u}_i^T (C \vec{u}_i) = \vec{u}_i^T \lambda_i \vec{u}_i = \lambda_i \vec{u}_i^T \vec{u}_i = \lambda_i$ .

**Conclusion: la meilleure direction, parmi toutes les directions, est celle donnée par le vecteur propre qui est associé à la plus grande valeur propre.**

### 19.3.3 Directions suivantes

En pratique, on ne souhaite pas conserver que 1 composante, mais plutôt projeter depuis un espace de dimension  $D$  vers un espace de dimension  $D'$ ,  $1 < D' < D$ . On projette donc non pas sur une droite (un seul vecteur) mais plutôt sur un hyper-plan, dont l'orientation est définie par les vecteurs qui sont dans ce plan. Pensez par exemple à une paire de vecteurs (non colinéaires) dans un espace à  $D = 3$ : ces deux vecteurs définissent naturellement un plan (à  $D' = 2$  dimensions).

Les données **après projection sur un hyper-plan** de dimension  $D'$ , par exemple, s'écrivent de la façon suivante. Appelons  $\vec{u}_1 \in \mathbb{R}^D$ ,  $\vec{u}_2 \in \mathbb{R}^D$  les deux vecteurs qui définissent ce plan. Appelons  $\vec{e}_1 = (1, 0)$ ,  $\vec{e}_2 = (0, 1)$

les deux vecteurs de la base interne au plan (ces vecteurs définissent un repère interne au plan, qui permet de se promener dedans). Alors la projection s'écrit :

$$\vec{x}' = (\vec{x}^T \vec{u}_1) \vec{e}_1 + (\vec{x}^T \vec{u}_2) \vec{e}_2 \quad (66)$$

L'écriture générale est  $\vec{x}' = \sum_{d'=1}^{D'} (\vec{x} \cdot \vec{u}_{d'}) \vec{e}_{d'}$ .

**Quels sont les directions qui définissent cet hyper-plan ?** On ne refait pas tout le raisonnement, mais au vu de la partie précédente, **la 2ème meilleure direction est celle donnée par le vecteur propre qui est associé à la plus grande valeur propre**. Et ainsi de suite pour la 3ème, la 4ème, etc.

On aboutit alors à la remarque suivante. Si on diagonalise la matrice de covariance,  $C$ , on a :  $C = U \Lambda U^{-1}$ , avec  $\Lambda$  (se lit Lambda majuscule, comme  $\lambda$ ) la matrice diagonale des valeurs propres, et  $U$  la matrice de passage vers l'espace où  $C$  est diagonalisée.

$$\Lambda = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_D \end{pmatrix} \quad U = \left( \begin{pmatrix} \cdot \\ \vec{u}_1 \\ \cdot \end{pmatrix} \quad \begin{pmatrix} \cdot \\ \vec{u}_2 \\ \cdot \end{pmatrix} \quad \dots \quad \begin{pmatrix} \cdot \\ \vec{u}_D \\ \cdot \end{pmatrix} \right) \quad (67)$$

On ne rentre pas dans les détails, mais il est légitime de supposer qu'une matrice de covariance de données "normales" sera diagonalisable<sup>10</sup>, avec tous les vecteurs propres distincts deux à deux. Par définition de  $C$ , elle est semi-définie positive, c'est-à-dire que toutes ses valeurs propres sont positives ou nulles. Comme toujours lors d'une diagonalisation,  **$U$  est constituée des vecteurs propres mis côte à côte**, dans le même ordre que les  $\lambda_i$  dans  $\Lambda$ . Ici, on choisit d'**ordonner les  $\lambda_i$  par valeur décroissante** (elles sont toutes positives et réelles).

La matrice de projection sur l'hyper-plan de dimension  $D'$  s'obtient alors (souvenez vous qu'on a trié les  $\lambda_i$  par ordre décroissant) en gardant les  $D'$  premières directions propres :

$$P = \left( \begin{pmatrix} \cdot \\ \vec{u}_1 \\ \cdot \end{pmatrix} \quad \begin{pmatrix} \cdot \\ \vec{u}_2 \\ \cdot \end{pmatrix} \quad \dots \quad \begin{pmatrix} \cdot \\ \vec{u}_{D'} \\ \cdot \end{pmatrix} \right) \quad (68)$$

La matrice  $P$  n'est pas carrée, elle est de taille  $(D, D')$  ( $D$  lignes,  $D'$  colonnes). Comme on l'a dit plus haut, la donnée  $\vec{x}_n$  projetée s'obtient alors ainsi :  $\vec{x}'_n = \sum_{d'=1}^{D'} (\vec{u}_{d'} \cdot \vec{x}_n) \vec{e}_{d'}$ . Ceci est en fait l'écriture d'une multiplication matricielle :

$$\vec{x}_{n,transformed} = (\vec{x}_n)' = (\vec{x}_n - \langle \vec{x} \rangle) \cdot P \quad (69)$$

Ces vecteurs sont de dimension  $D'$ .

### 19.3.4 Décompression

La transformée inverse s'obtient ainsi :

$$\vec{x}_{n,decompressed} = \vec{x}_{n,transformed} \cdot P^T + \langle \vec{x} \rangle \quad (70)$$

On obtient donc de nouveau des vecteurs de dimension  $D$ .

On voit dans cette expression en quoi c'est une expression avec perte. En effet, on a  $\vec{x}_{n,decompressed} = ((\vec{x}_n - \langle \vec{x} \rangle) \cdot P) \cdot P^T + \langle \vec{x} \rangle$ . L'opérateur  $PP^T$  est une matrice de dimension  $(D, D)$ , mais qui intuitivement, "passe par une dimension intermédiaire de dimension  $D'$ ". Mathématiquement, on pourrait constater que  $PP^T$  n'a que  $D'$  directions indépendantes (vecteurs propres linéairement indépendants).

## 19.4 PCA: résumé de l'algorithme

L'algorithme de la PCA est donc le suivant :

1. Calculer la matrice de covariance  $C$  (cela se fait en un temps linéaire en la taille des données) :  $C = \frac{1}{N} \sum_{n=1}^N (\vec{x}_n - \langle \vec{x} \rangle)(\vec{x}_n - \langle \vec{x} \rangle)^T$

<sup>10</sup>Le sous-espace des matrices diagonalisables dans  $\mathbb{C}$  est dense dans l'espace des matrices (à valeurs dans  $\mathbb{C}$ ). Ceci n'est pas vrai dans  $\mathbb{R}$ , mais à part pour de données pathologiques, une matrice de covariance sera diagonalisable.

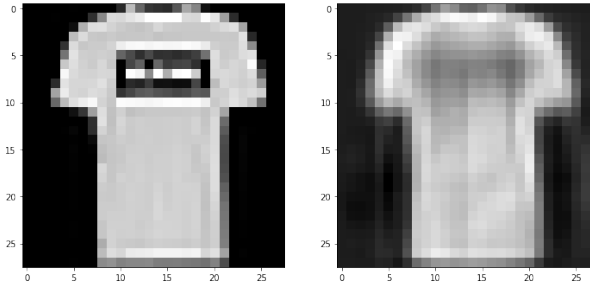


Figure 9: On passe d'une image de dimension  $D = 784$  à  $D' = 30$ , puis on revient à l'espace d'origine pour visualisation. On a  $D' = 30 \sim \sqrt{D}$  soit presque 30 fois moins, et pourtant on peut toujours reconnaître l'image.

2. Diagonaliser  $C$ . Cela se fait en un temps  $O(D^3)$ , a priori (voir remarque plus bas). Cela revient à chercher une matrice de passage  $U$  telle que

$$C = U\Lambda U^{-1}$$

où  $\Lambda$  est une matrice diagonale dont les coefficients diagonaux sont les valeurs propres de  $C$  rangées dans l'ordre décroissant.

3. Garder seulement les  $D'$  premières valeurs propres (on s'intéresse uniquement aux plus grandes valeurs propres qui permettent d'obtenir plus grande variance, que l'on cherche justement à maximiser) :  $P = (\vec{u}_1, \dots, \vec{u}_{D'})_{D, D'}$
4. Projeter (transformer) :  $\vec{x}_{n, transformed} = (\vec{x}_n - \langle \vec{x} \rangle) \cdot P$

Par ailleurs si on ne cherche que les  $D'$  plus grandes valeurs propres, on peut le faire en un temps  $O(D'D^2)$  au lieu de  $O(D^3)$ , ce qui est très appréciable (on part du principe que  $D' \leq D$ ). Il existe aussi des solutions encore approximatives plus rapides avec résultats légèrement aléatoires. D'autant que dans le cas de la projection ce n'est pas grave d'avoir quelques erreurs.

Une fois cet algorithme appliqué, on peut effectuer la transformation inverse pour revenir à l'espace de départ (c'est une compression avec perte).

Attention : Quand on effectue une PCA sur une image de dimension  $D$ , on obtient  $D'$  nombres, mais on n'a plus du tout à faire à une image, la donnée n'est plus visualisable – on peut cependant la décompresser ensuite pour revenir à une image.

## 19.5 Ressources

Concernant ce Chapitre, vous avez:

- (ressource bilingue)  
<https://gitlab.inria.fr/flandes/ias/-/raw/master/CM-slides/exemples-PCA-d%C3%A9mo-minimale-Engl.ipynb?inline=false>
- (ressource bilingue)  
[https://gitlab.inria.fr/flandes/ias/-/raw/master/CM-slides/2022-52-exemples-PCA-d%C3%A9tails+visualisation\\_modes-FR+EN.ipynb?inline=false](https://gitlab.inria.fr/flandes/ias/-/raw/master/CM-slides/2022-52-exemples-PCA-d%C3%A9tails+visualisation_modes-FR+EN.ipynb?inline=false)
- (en anglais)  
<https://plot.ly/ipython-notebooks/principal-component-analysis/>
- Bishop [Bis06] chapter 12, page 561-570 (ou plus si vous êtes intéressé.e.s).

## 20 Modèles Bayésiens

Les techniques décrites ici peuvent être vues comme des *estimations de densité*.

Lorsqu'on souhaite "fitter" les paramètres d'une densité de probabilité sur un ensemble de données, la méthode la plus naturelle consiste à choisir les paramètres tels que les données soient "les plus vraisemblables" ou autrement dit, qu'elles soient "réalistes", c.a.d. qu'il soit crédible que ces données aient été générées par cette distribution, avec ces paramètres.

Cette approche s'appelle l'estimation par le Maximum de Vraisemblance (des données). On cherche ainsi à trouver les paramètres optimaux  $\theta^*$  qui maximisent:

$$\theta^* = \operatorname{argmax}_{\theta} (\mathbb{P}(X|\theta)) \quad (71)$$

Notation: l'exposant  $.$  est souvent utilisé pour désigner la solution d'un problème d'optimisation.

Notation: le chapeau  $\hat{u}$  est souvent utilisé pour désigner l'*estimateur* d'une quantité a priori inconnu. On lira  $\hat{u}$  "u chapeau" ou bien "estimateur de u".

Notation: on dit que la variable aléatoire (v.a.)  $X$  est paramétrée par les paramètres  $\theta$ .

### 20.1 MLE: 1 variable à 1 dimension

Prenons un exemple concret, par exemple la taille en cm d'un certain nombre  $N$  d'étudiants. On dispose de  $N$  points de données, à 1 dimension, c.a.d que ce sont des nombres réels. On suppose que les tirages sont indépendants. Si  $X_n$  est la v.a. associée au  $n$ -ième tirage (indiquée ci dessous), la v.a. des  $N$  tirages est alors la v.a. produit, notée  $X$ :

$$\forall n \in [1, \dots, N], X_n \sim X_1 \quad (72)$$

$$X \sim (X_1, X_2, \dots, X_N) \quad (73)$$

Notations:  $X \sim Y$  signifie que la v.a.  $X$  suit la même loi que la v.a.  $Y$ . L'indépendance entre deux v.a.  $X_1$  et  $X_2$  est notée  $X_1 \perp\!\!\!\perp X_2$ . Ici on a  $X_i \perp\!\!\!\perp X_j, \forall i \neq j$ . Ici, les virgules signalent un "ET" logique.

On dispose de données, qu'on note plutôt  $X_{(N,1)}$  que  $X$ , afin de les distinguer de la variable aléatoire  $X$ . Les données sont la *réalisation* d'une certaine loi (ou bien quelque chose qu'on arrivera jamais à décrire mathématiquement, selon les problèmes). Les données  $X_{(N,1)}$  sont donc une liste:  $X_{(N,1)} = (x_1, x_2, \dots, x_N)$ , avec  $x_n \in \mathbb{R}$ .

On définit alors la probabilité d'avoir obtenu le tirage  $X_{(N,1)}$ , depuis la v.a.  $X$ :

$$\mathbb{P}(X = X_{(N,1)}) = \mathbb{P}(X_1 = x_1, X_2 = x_2, \dots, X_N = x_N) \quad (\text{Définition, toujours vraie}) \quad (74)$$

$$= \prod_n^N \mathbb{P}(X_n = x_n) \quad (\text{Car les } X_n \text{ sont indépendants}) \quad (75)$$

Il faut maintenant rendre les choses plus concrètes en choisissant un *modèle* mathématique, c'est-à-dire une expression **explicite** pour la forme de  $\mathbb{P}(X_n = x_n)$ . Ici, pour aujourd'hui, on suppose que tous les points sont distribués selon la (même) loi **Gaussienne**:

$$\mathbb{P}(X_n \in [x, x + dx]) = \rho(x)dx, \quad (76)$$

$$\rho(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (77)$$

Où  $\theta = (\mu, \sigma)$  sont les 2 paramètres réels à *estimer* (à "fitter", en vocabulaire de Machine Learning). Attention, de façon générale, on aurait pu faire un autre choix de représentation, à la place de la Gaussienne: on aurait pu prendre une loi de Poisson, de Bernoulli, etc. C'est un *choix* de *modélisation*.

Notation:  $\rho(x)$  se lit rho de  $x$ , c'est la *densité* associée à la v.a.  $X$ , qui est une v.a. continue.

L'idée de maximiser la vraisemblance des données revient alors, concrètement, à **maximiser la vraisemblance des données sachant les paramètres**, notée  $\mathcal{L}$  (comme Likelihood, Vraisemblance en anglais):

$$\mathcal{L}(\theta) = \mathbb{P}(X = X_{(N,1)}|\theta) \quad (78)$$

$$\theta^* = \operatorname{argmax}_{\theta} (\mathcal{L}(\theta)) \quad (79)$$

Quel que soit le choix de modélisation, il y aura toujours au moins un paramètre à *estimer* d'après les données: c'est pour cela qu'on parle d'*estimateur par le maximum de vraisemblance* (*Maximum Likelihood Estimate*, en anglais, ou EMV en Français).