

Fondements du Machine Learning

Florentin Goyens

Notes de cours - L3 IM2D - 2022/2023

- Ces notes sont basées sur une version précédente du cours rédigée par Clément Royer.
- Pour toute remarque, envoyer un mail à florentin.goyens@dauphine.psl.eu.
Merci aux étudiant(e)s ayant signalé des erreurs ou des coquilles.
- **Objectifs d'apprentissage**
À l'issue de ce cours, l'étudiant(e) sera capable de
 - Donner la formule de la décomposition en valeurs singulières, et appliquer cette décomposition à des problèmes matriciels;
 - Reconnaître et formuler des problèmes aux moindres carrés linéaires;
 - Donner des solutions de ces problèmes et expliciter leur lien avec la décomposition en valeurs singulières;
 - Définir des estimateurs statistiques et établir les propriétés de ces estimateurs;
 - Comprendre le compromis biais-variance pour les modèles linéaires;
 - Donner la définition d'une composante principale, et appliquer l'analyse en composantes principales à des données matricielles.

Sommaire

1	Introduction	3
1.1	Derrière le Machine Learning...	3
1.2	Contexte et objectifs du cours	4
1.3	Notations	4
1.3.1	Convention de notation	4
1.3.2	Notations vectorielles	4
1.3.3	Notations matricielles	5
2	Algèbre linéaire et décomposition en valeurs singulières	6
2.1	Rappels d'algèbre linéaire	6
2.2	Valeurs propres et décomposition spectrale	9
2.3	Décomposition en valeurs singulières	9
2.3.1	Principe de la décomposition	10
2.3.2	Preuve constructive de la décomposition	11
2.3.3	Décomposition tronquée et approximation	13
2.4	Exercices	15
3	Premiers pas avec le modèle linéaire	17
3.1	Introduction	17
3.2	Résolution de systèmes non linéaires	18
3.2.1	Cas d'un système carré	18
3.2.2	Cas d'un système rectangulaire	19
3.2.3	Pseudo-inverse et SVD	20
3.3	Moindres carrés linéaires	21
3.3.1	Solution au sens des moindres carrés	21
3.3.2	Résolution du problème aux moindres carrés	23
3.4	Conclusion	23
3.5	Exercices	24

Chapitre 1

Introduction

1.1 Derrière le Machine Learning...

Le terme *machine learning*, dont les traductions varient entre apprentissage machine, apprentissage automatique et apprentissage artificiel, fait partie d'un ensemble de mots-clés qui ont récemment gagné en popularité. Parmi ceux-ci, on trouve également l'analyse de données (*data analysis*), la fouille de données (*data mining*), l'intelligence artificielle (*artificial intelligence*, ou simplement *AI*), les masses de données (*Big Data*), etc. L'utilisation de cette terminologie est parfois hasardeuse : on leur préférera donc la notion de **sciences des données**, ou *data science*.

La notion de donnée est en effet au coeur des différents concepts sus-mentionnés, et représente un enjeu majeur dans de nombreux secteurs d'activités. Pour les entreprises de service telles que les GAFAM ¹, il s'agit de fournir une valeur ajoutée dans leur service autrement gratuit via la façon dont les données des utilisateurs sont exploitées. En recherche et développement, la quantité massive de données générées dans certains domaines (biologie, médecine) pose d'importants défis mathématiques et informatiques. Plus globalement, les approches guidées par les données (*data-driven*) deviennent de plus en plus populaires, car elles permettent de pallier le manque de modèles formels ou implémentables. C'est le cas par exemple pour la modélisation météorologique à grande échelle : nos capacités de calcul ne nous permettent pas de faire évoluer un modèle de prédiction à l'échelle du globe, mais il est possible de collecter un grand volume de données et d'en extraire les tendances majeures.

Exemple 1.1 (Systèmes de recommandation) *Les plate-formes commerciales telles que Netflix ou Youtube suggèrent du contenu pertinent à leurs utilisateurs en fonction de leurs préférences. Pour ce faire, elles disposent d'une matrice d'avis : il s'agit d'un tableau en deux dimensions, l'une représentant les clients et l'autre les produits. Les grandes questions qui se posent sont donc :*

- 1) *Quels sont les éléments principaux de nos préférences ?*
- 2) *Comment gérer un grand nombre d'avis ?*
- 3) *Les avis reflètent-ils vraiment la réalité ?*

Dans ce cours, on considèrera deux approches d'analyse de données. La première, dite fonctionnelle, supposera que les données d'apprentissage suivent une distribution de probabilité connue : il

¹Google, Apple, Facebook, Amazon et Microsoft.

est ainsi possible de construire un modèle de ces données adapté à la distribution sous-jacente. Les approches de **régression linéaire** seront utilisés dans ce contexte d'apprentissage (dit supervisé). La seconde approche, dite prédictive, ne pré-supposera pas de distribution sur les données, et visera à extraire de l'information des données même (on parle ainsi d'apprentissage non supervisé). L'**analyse en composantes principales** (ou ACP, voir Chapitre ??) sera ici l'outil-clé pour obtenir cette information.

1.2 Contexte et objectifs du cours

Dans ce cours, on s'intéressera à développer des techniques visant à extraire de l'information d'un jeu de données. On considèrera que l'on dispose d'une quantité importante de données, non seulement pour qu'il soit intéressant d'en extraire de l'information, mais aussi pour que ces données puissent représenter des tendances. Les techniques que nous emploierons reposent sur des algorithmes, c'est-à-dire des traitements systématiques à appliquer aux données. Comme on le verra, le développement d'un algorithme efficace repose à la fois sur des arguments mathématiques et sur une implémentation bien pensée.

Ce cours se concentre plus spécifiquement sur l'obtention de modèles **linéaires** des relations entre les données; ces données seront de plus vues comme des réalisations de variables aléatoires (généralement gaussiennes). Ce choix se justifie par la pertinence et l'efficacité de ces modèles simples dans la pratique. Il permet également d'utiliser des résultats et algorithmes issus de l'algèbre linéaire, de l'optimisation et des statistiques.

1.3 Notations

1.3.1 Convention de notation

- Les scalaires seront représentés par des lettres minuscules : $a, b, c, \alpha, \beta, \gamma$.
- Les vecteurs seront représentés par des lettres minuscules **en gras** : $\mathbf{a}, \mathbf{b}, \mathbf{c}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}$.
- Les lettres majuscules en gras seront utilisées pour les matrices : $\mathbf{A}, \mathbf{B}, \mathbf{C}$.
- En l'absence d'ambiguïté, on pourra omettre les indices de début et de fin dans une somme finie afin d'alléger les notations. On pourra de même utiliser un seul symbole de notation pour plusieurs indices et ainsi écrire de manière équivalente $\sum_{i=1}^m \sum_{j=1}^n$, $\sum_i \sum_j$ ou $\sum_{i,j}$ si le contexte le permet.

1.3.2 Notations vectorielles

- On notera \mathbb{R}^n l'ensemble des vecteurs à n composantes réelles, et on considèrera toujours que n est un entier supérieur ou égal à 1.
- Un vecteur $\mathbf{x} \in \mathbb{R}^n$ sera pensé (par convention) comme un vecteur colonne. On notera $x_i \in \mathbb{R}$ sa i -ème coordonnée dans la base canonique de \mathbb{R}^n . On aura ainsi $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$, que l'on notera plus succinctement $\mathbf{x} = [x_i]_{1 \leq i \leq n}$.

- Étant donné un vecteur (colonne) $\mathbf{x} \in \mathbb{R}^n$, le vecteur ligne correspondant sera noté \mathbf{x}^T . On aura donc $\mathbf{x}^T = [x_1 \ \cdots \ x_n]$ et $[\mathbf{x}^T]^T = \mathbf{x}$.
- Pour tout $n \geq 1$, les vecteurs $\mathbf{0}_n$ et $\mathbf{1}_n$ représentent les vecteurs colonnes de \mathbb{R}^n dont tous les éléments sont égaux à 0 ou 1, respectivement.

1.3.3 Notations matricielles

- On notera $\mathbb{R}^{m \times n}$ l'ensemble des matrices à m lignes et n colonnes à coefficients réels, où m et n seront des entiers supérieurs ou égaux à 1. Les espaces $\mathbb{R}^{m \times 1}$ et \mathbb{R}^m étant isomorphes (ce que l'on note $\mathbb{R}^{m \times 1} \simeq \mathbb{R}^m$), on pourra considérer un vecteur de \mathbb{R}^m comme une matrice de $\mathbb{R}^{m \times 1}$, et vice versa. Une matrice $\mathbf{A} \in \mathbb{R}^{n \times n}$ est dite carrée (dans le cas général, on parlera de matrice rectangulaire).
- Étant donnée une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$, on notera A_{ij} le coefficient en ligne i et colonne j de la matrice. La notation $[A_{ij}]_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}}$ sera donc équivalente à \mathbf{A} . Sans ambiguïté sur la taille de la matrice, on notera simplement $[A_{ij}]$.
- Selon les besoins, on utilisera \mathbf{a}_i^T pour la i -ème ligne de \mathbf{A} ou \mathbf{a}_j pour la j -ème colonne de \mathbf{A} .
Selon le cas, on aura donc $\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix}$ ou $\mathbf{A} = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_n]$.
- Pour une matrice $\mathbf{A} = [A_{ij}] \in \mathbb{R}^{m \times n}$, la matrice transposée de \mathbf{A} , notée \mathbf{A}^T , est la matrice de $\mathbb{R}^{n \times m}$ telle que

$$\forall i = 1 \dots m, \forall j = 1 \dots n, \quad \mathbf{A}_{ji}^T = A_{ij}.$$

Nota Bene : Cette notation généralise donc la correspondance entre vecteurs lignes et vecteurs colonnes.

- Pour tout $n \geq 1$, la matrice \mathbf{I}_n représentera la matrice identité de $\mathbb{R}^{n \times n}$ (avec des 1 sur la diagonale et des 0 partout ailleurs), tandis que les matrices $\mathbf{0}_n$ et $\mathbf{1}_n$ représenteront les matrices dont tous les éléments sont égaux à 0 ou 1, respectivement. De manière plus générale, les notations $\mathbf{0}_{m,n}$ et $\mathbf{1}_{m,n}$ seront utilisées pour les matrices de $\mathbb{R}^{m \times n}$ ne contenant respectivement que des 0 et des 1.

Chapitre 2

Algèbre linéaire et décomposition en valeurs singulières

2.1 Rappels d'algèbre linéaire

On considère l'espace des vecteurs \mathbb{R}^n muni de sa structure d'espace vectoriel normé de dimension n :

- Pour tous $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, la somme des vecteurs \mathbf{x} et \mathbf{y} est notée $\mathbf{x} + \mathbf{y} = [x_i + y_i]_{1 \leq i \leq n}$;
- Pour tout $\lambda \in \mathbb{R}$, on définit $\lambda \mathbf{x} := \lambda \cdot \mathbf{x} = [\lambda x_i]_{1 \leq i \leq n}$.
- Pour tous vecteurs $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, le produit scalaire entre \mathbf{x} et \mathbf{y} est noté $\mathbf{x}^T \mathbf{y}$, et défini par

$$\mathbf{x}^T \mathbf{y} := \sum_{i=1}^n x_i y_i.$$

Il s'agit d'une forme bilinéaire symétrique définie positive. En particulier, on a $\mathbf{y}^T \mathbf{x} = \mathbf{x}^T \mathbf{y}$.

- La norme Euclidienne $\|\cdot\|$ sur \mathbb{R}^n dérivée du produit scalaire est définie pour tout vecteur $\mathbf{x} \in \mathbb{R}^n$ par

$$\|\mathbf{x}\| := \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}.$$

On dira que $\mathbf{x} \in \mathbb{R}^n$ est unitaire si $\|\mathbf{x}\| = 1$.

- Il existe une famille libre et génératrice de \mathbb{R}^n de taille n . Par exemple, tout vecteur \mathbf{x} de \mathbb{R}^n s'écrit $\mathbf{x} = \sum_{i=1}^n x_i \mathbf{e}_i$, où $\mathbf{e}_i = [0 \cdots 0 \ 1 \ 0 \cdots 0]^T$ est le i -ème vecteur de la base canonique (le coefficient 1 se trouvant en i -ème position).

Définition 2.1 (Sous-espace engendré) Soient $\mathbf{x}_1, \dots, \mathbf{x}_p$ p vecteurs de \mathbb{R}^n . Le *sous-espace engendré* par les vecteurs $\mathbf{x}_1, \dots, \mathbf{x}_p$ est le sous-espace vectoriel

$$\text{vect}(\mathbf{x}_1, \dots, \mathbf{x}_p) := \left\{ \mathbf{x} = \sum_{i=1}^p \alpha_i \mathbf{x}_i \mid \alpha_i \in \mathbb{R} \ \forall i \right\}.$$

Ce sous-espace est de dimension au plus $\min(n, p)$.

Lorsque l'on travaille avec des matrices, on s'intéresse généralement aux sous-espaces définis ci-dessous.

Définition 2.2 (Sous-espaces matriciels) Soit une matrice $A \in \mathbb{R}^{m \times n}$, on définit les deux sous-espaces suivants :

- Le **noyau** (kernel/null space en anglais) de A est le sous-espace vectoriel

$$\ker(A) := \{x \in \mathbb{R}^n \mid Ax = 0_m\}$$

- L'**image** (range space) de A est le sous-espace vectoriel

$$\text{Im}(A) := \{y \in \mathbb{R}^m \mid \exists x \in \mathbb{R}^n, y = Ax\}$$

La dimension de ce sous-espace vectoriel s'appelle le **rang** de A . On la note $\text{rang}(A)$ et on a $\text{rang}(A) \leq \min\{m, n\}$.

Théorème 2.1 (Théorème du rang) Pour toute matrice $A \in \mathbb{R}^{m \times n}$, on a

$$\dim(\ker(A)) + \text{rang}(A) = n.$$

Définition 2.3 (Normes matricielles) On définit sur $\mathbb{R}^{m \times n}$ la norme d'opérateur $\|\cdot\|$ et la norme de Frobenius $\|\cdot\|_F$ par

$$\forall A \in \mathbb{R}^{m \times n}, \begin{cases} \|A\| &:= \max_{\substack{x \in \mathbb{R}^n \\ x \neq 0_n}} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \|Ax\| \\ \|A\|_F &:= \sqrt{\sum_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} A_{ij}^2}. \end{cases}$$

Nous terminons cette section par quelques définitions de sous-ensembles de matrices carrées qui nous seront utiles dans le cours.

Définition 2.4 (Matrice symétrique) Une matrice carrée $A \in \mathbb{R}^{n \times n}$ est dite **symétrique** si elle vérifie $A^T = A$.

Définition 2.5 (Matrice inversible) Une matrice carrée $A \in \mathbb{R}^{n \times n}$ est dite **inversible** s'il existe $B \in \mathbb{R}^{n \times n}$ telle que $BA = AB = I_n$ (où l'on rappelle que I_n désigne la matrice identité de $\mathbb{R}^{n \times n}$).

Si elle existe, une telle matrice B est unique : elle est appelée **l'inverse de A** et on la note A^{-1} .

Définition 2.6 (Matrice (semi-)définie positive) Une matrice symétrique $A \in \mathbb{R}^{n \times n}$ est dite **semi-définie positive** si

$$x^T Ax \geq 0, \quad \text{pour tout } x \in \mathbb{R}^n.$$

Elle est dite **définie positive** lorsque $x^T Ax > 0$ pour tout vecteur x non nul.

Notons qu'il existe des matrices non-symétriques qui satisfont $x^T Ax \geq 0$ pour tout $x \in \mathbb{R}^n$. Mais ces matrices ne sont pas d'un grand intérêt en pratique, et nous restreignons donc la définition de matrice semi-définie positive aux matrices symétriques.

Définition 2.7 (Matrice orthogonale) Une matrice carrée $P \in \mathbb{R}^{n \times n}$ est dite *orthogonale* si $P^T = P^{-1}$. C'est à dire, $PP^T = P^TP = I_n$. Par extension, on dira que $Q \in \mathbb{R}^{m \times n}$ avec $m \geq n$ est orthogonale si $Q^TQ = I_n$ (les colonnes de Q sont donc orthonormées dans \mathbb{R}^m).

Si $Q \in \mathbb{R}^{n \times n}$ est une matrice orthogonale, alors Q^T est également orthogonale.

On utilisera fréquemment la propriété des matrices orthogonales énoncée ci-dessous.

Lemme 2.1 Soit une matrice $A \in \mathbb{R}^{m \times n}$ et soit deux matrices orthogonales $U \in \mathbb{R}^{m \times m}$ et $V \in \mathbb{R}^{n \times n}$. On a

$$\|A\| = \|UA\| = \|AV\| \text{ et } \|A\|_F = \|UA\|_F = \|AV\|_F,$$

c'est-à-dire que la multiplication par une matrice orthogonale ne modifie pas la norme d'opérateur.

Démonstration. On montre tout d'abord que pour tout vecteur $x \in \mathbb{R}^m$, on a $\|Ux\| = \|x\|$. En utilisant la définition de la norme et celle d'une matrice orthogonale, on a :

$$\|Ux\|^2 = x^T U^T U x = x^T x = \|x\|^2,$$

ce qui établit le résultat. Par conséquent, on a également

$$\|UAx\| = \|Ax\|$$

pour tout vecteur x . En revenant à la définition de la norme d'opérateur, on obtient ainsi

$$\|UA\| = \max_{\|x\|=1} \|UAx\| = \max_{\|x\|=1} \|Ax\| = \|A\|,$$

ce qui est bien le résultat recherché. Pour le résultat sur $\|AV\|$, on note que V^T est une matrice orthogonale inversible avec $VV^T = I_n$, donc que pour tout x , il existe z tel que $x = V^T z$ et $\|x\| = \|z\|$ d'après ce qui précède. On a ainsi :

$$\begin{aligned} \|AV\| &= \max_{\|x\|=1} \|AVx\| = \max_{\|V^T z\|=1} \|AVV^T z\| \\ &= \max_{\|z\|=1} \|Az\| = \|A\|, \end{aligned}$$

ce qu'il fallait démontrer.

On démontre maintenant le résultat pour la norme de Frobenius, dont on rappelle qu'elle est définie par $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2}$. On a ainsi

$$\|A\|_F = \sqrt{\|a_1^\top\|^2 + \dots + \|a_m^\top\|^2} = \sqrt{\|a_1\|^2 + \dots + \|a_n\|^2}$$

où $a_1^\top, \dots, a_m^\top$ et a_1, \dots, a_n représentent les lignes et les colonnes de A , respectivement. Ainsi, la norme de Frobenius d'une matrice au carré est égale à la somme des carrés des normes de ses colonnes ou de ses lignes). Comme on a montré que la norme d'un vecteur ne change pas par transformation orthogonale, pour toute matrice $U \in \mathbb{R}^{m \times m}$ orthogonale, on a

$$\sqrt{\|Ua_1\|^2 + \dots + \|Ua_n\|^2} = \sqrt{\|a_1\|^2 + \dots + \|a_n\|^2},$$

d'où $\|UA\|_F = \|A\|_F$. En considérant les lignes de A , on montre de même que $\|A\|_F = \|AV\|_F$.

□

Par corollaire immédiat du lemme précédent, on note qu'une matrice $Q \in \mathbb{R}^{m \times n}$ orthogonale avec $m \geq n$ vérifie nécessairement $\|Q\| = \|I_n\| = 1$ et $\|Q\|_F = \|I_n\|_F = \sqrt{n}$.

2.2 Valeurs propres et décomposition spectrale

Définition 2.8 (Valeur propre) Soit une matrice $A \in \mathbb{R}^{n \times n}$. On dit que $\lambda \in \mathbb{R}$ est une **valeur propre de A** si

$$\exists v \in \mathbb{R}^n, v \neq 0_n, \quad Av = \lambda v.$$

Le vecteur v est appelé un **vecteur propre associé à la valeur propre λ** . L'ensemble des valeurs propres de A s'appelle le **spectre de A** .

Le sous-espace engendré par les vecteurs propres associés à la même valeur propre d'une matrice s'appelle un sous-espace propre. Sa dimension correspond à l'ordre de multiplicité de la valeur propre relativement à la matrice.

Proposition 2.1 Pour toute matrice $A \in \mathbb{R}^{n \times n}$, on a les propriétés suivantes :

- La matrice A possède n valeurs propres complexes mais pas nécessairement réelles.
- Si la matrice A est semi-définie positive (respectivement définie positive), alors ses valeurs propres sont réelles positives (respectivement strictement positives).
- Le noyau de A est engendré par les vecteurs propres associés à la valeur propre 0.

Théorème 2.2 (Théorème spectral) Toute matrice carrée $A \in \mathbb{R}^{n \times n}$ symétrique admet une décomposition dite **spectrale** de la forme :

$$A = P \Lambda P^{-1},$$

où $P \in \mathbb{R}^{n \times n}$ est une matrice orthogonale, dont les colonnes p_1, \dots, p_n forment une base orthonormée de vecteurs propres, et $\Lambda \in \mathbb{R}^{n \times n}$ est une matrice diagonale qui contient les n valeurs propres de A $\lambda_1, \dots, \lambda_n$ sur la diagonale.

Il est à noter que la décomposition spectrale n'est pas unique. En revanche, l'ensemble des valeurs propres est unique, que l'on prenne en compte les ordres de multiplicité ou non.

La décomposition spectrale définie dans le théorème 2.2 est particulièrement importante car elle permet de synthétiser l'information de A par son effet sur les vecteurs p_i . Ainsi, lorsque $|\lambda_i| \gg 1$, on aura $\|Ap_i\| \gg \|p_i\|$, et la matrice aura donc un effet expansif dans la direction de p_i (ou sa direction opposée lorsque $\lambda_i < 0$). De même, si $|\lambda_i| \ll 1$, la matrice aura un effet contractant dans la direction de p_i : le cas extrême est $\lambda_i = 0$, c'est-à-dire que $p_i \in \ker(A)$ et la matrice ne conserve donc pas d'information relative à p_i .

Géométriquement parlant, on voit ainsi que, pour tout vecteur $x \in \mathbb{R}^n$ décomposé dans la base des p_i que l'on multiplie par A , les composantes de ce vecteur associées aux plus grandes valeurs propres¹ seront augmentées, tandis que celles associées aux valeurs propres de petite magnitude seront réduites (voire annihilées dans le cas d'une valeur propre nulle).

2.3 Décomposition en valeurs singulières

La décomposition en valeurs singulières (ou SVD, pour *Singular Value Decomposition*) est une technique fondamentale en analyse et compression de données, particulièrement utile pour compresser des signaux audio, des images, etc.

¹On parle ici de plus grandes valeurs propres en valeur absolue, ou magnitude.

2.3.1 Principe de la décomposition

Soit une matrice rectangulaire $A \in \mathbb{R}^{m \times n}$: dans le cas général, les dimensions de la matrice diffèrent, et on ne peut donc pas parler de valeurs propres de la matrice A . On peut en revanche considérer les deux matrices

$$A^T A \in \mathbb{R}^{n \times n} \quad \text{et} \quad A A^T \in \mathbb{R}^{m \times m}.$$

Ces matrices sont symétriques réelles, et par conséquent diagonalisables. Par ailleurs, elles sont fortement liées à la matrice A . Le lemme ci-dessous illustre quelques-unes des propriétés de $A^T A$; des résultats similaires peuvent être démontrés pour $A A^T$.

Lemme 2.2 *Pour toute matrice $A \in \mathbb{R}^{m \times n}$, les propriétés suivantes sont vérifiées :*

- i) $A^T A$ est symétrique;
- ii) $A^T A$ est semi-définie positive;
- iii) $\ker(A^T A) = \ker(A)$;
- iv) $\text{rang}(A^T A) = \text{rang}(A)$;
- v) $\text{Im}(A^T A) = \text{Im}(A^T)$.

Ces résultats sont à la base de la construction de la décomposition en valeurs singulières, dont on donne l'énoncé ci-dessous.

Théorème 2.3 (Décomposition en valeurs singulières) *Toute matrice $A \in \mathbb{R}^{m \times n}$ admet une décomposition en valeurs singulières (SVD²) de la forme*

$$A = U \Sigma V^T,$$

où $U \in \mathbb{R}^{m \times m}$ est orthogonale ($U^T U = I_m$), $V \in \mathbb{R}^{n \times n}$ est orthogonale ($V^T V = I_n$) et $\Sigma \in \mathbb{R}^{m \times n}$ est telle que $\Sigma_{ij} = 0$ si $i \neq j$ et $\Sigma_{ii} \geq 0$.

L'ensemble des valeurs $\{\Sigma_{ii}\}$ pour $1 \leq i \leq \min\{m, n\}$, noté $\{\sigma_1, \dots, \sigma_{\min\{m, n\}}\}$ est appelé l'ensemble des valeurs singulières de la matrice A . Les colonnes de V (resp. de U) sont appelées les vecteurs singuliers à droite (resp. à gauche) de A .

Remarque 2.1 *Comme dans le cas de la décomposition en valeurs propres, il n'y a pas unicité de la décomposition en valeurs singulières, mais il y a unicité de l'ensemble des valeurs singulières.*

Exemple 2.1 *La décomposition en valeurs singulières d'une matrice de $\mathbb{R}^{3 \times 2}$ est de la forme*

$$A = \underbrace{[u_1 \ u_2 \ u_3]}_U \underbrace{\begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \\ 0 & 0 \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} v_1^T \\ v_2^T \end{bmatrix}}_{V^T}$$

où $\sigma_1 \geq 0, \sigma_2 \geq 0$, les u_i forment une base orthonormée de \mathbb{R}^3 et les v_i forment une base orthonormée de \mathbb{R}^2 .

²Dans la suite, on utilisera fréquemment l'algorithme anglo-saxon SVD pour faire référence à la décomposition en valeurs singulières.

Proposition 2.2 Soit $A \in \mathbb{R}^{m \times n}$ et $U\Sigma V^T$ une décomposition en valeurs singulières de A . Alors :

- i) Les carrés des valeurs singulières sont les valeurs propres de $A^T A$.
- ii) Si $\text{rang}(A) = r$, alors il y a exactement r valeurs singulières non nulles.

2.3.2 Preuve constructive de la décomposition

Il existe plusieurs manières d'établir la décomposition en valeurs singulières. On propose ci-dessous une preuve constructive de la SVD, qui permet d'obtenir une décomposition explicite de la matrice.

Soit une matrice $A \in \mathbb{R}^{m \times n}$ que l'on supposera non nulle (s'il s'agit de la matrice nulle, la décomposition $U = I_m$, $\Sigma = A$ et $V = I_n$ convient). Par définition de $\|A\| = \max_{\|z\|=1} \|Az\|$, il existe un vecteur $v_1 \in \mathbb{R}^n$ de norme 1 tel que $\|Av_1\| = \|A\|$. En posant $u_1 = \frac{Av_1}{\sigma_1}$ avec $\sigma_1 = \|A\|$, on obtient ainsi $u_1 \in \mathbb{R}^m$ et $v_1 \in \mathbb{R}^n$ tels que $\|u_1\| = 1$, $\|v_1\| = 1$ et $Av_1 = \sigma_1 u_1$. Ces vecteurs constitueront nos premiers vecteurs singuliers.

Comme les vecteurs u_1 et v_1 sont unitaires, on peut définir deux matrices orthogonales U_1 et V_1 , respectivement dans $\mathbb{R}^{m \times m}$ et $\mathbb{R}^{n \times n}$, telles que u_1 et v_1 soient les premières colonnes des matrices U_1 et V_1 . Soient $u_1^{(1)}, \dots, u_{m-1}^{(1)}$ et $v_1^{(1)}, \dots, v_{n-1}^{(1)}$ les colonnes restantes de U_1 et V_1 . On peut alors considérer la matrice $U_1^T A V_1$, dont les coefficients sont donnés par :

$$\left[\begin{array}{c|c} u_1^T A v_1 & [u_1^T A v_j^{(1)}]_{j=1, \dots, n-1} \\ \hline [(u_i^{(1)})^T A v_1]_{i=1, \dots, m-1} & [(u_i^{(1)})^T A v_j^{(1)}]_{\substack{i=1, \dots, m-1 \\ j=1, \dots, n-1}} \end{array} \right].$$

Notre but est de montrer que

$$U_1^T A V_1 = \left[\begin{array}{c|c} \sigma_1 & \mathbf{0}_{n-1}^T \\ \hline \mathbf{0}_{m-1} & A_1 \end{array} \right],$$

où $A_1 \in \mathbb{R}^{(m-1) \times (n-1)}$.

On a $u_1^T A v_1 = \sigma_1 v_1^T v_1 = \sigma_1 \|v_1\|^2 = \sigma_1$ d'après la première question.

Par ailleurs, pour tout $i = 1, \dots, m-1$, on a :

$$u_i^{(1)T} A v_1 = \sigma_1 u_i^{(1)T} u_1 = 0$$

car U_1 est une matrice orthogonale.

En posant $w = [u_1^T A v_j^{(1)}]_{j=1, \dots, n-1}$ et $A_1 = [(u_i^{(1)})^T A v_j^{(1)}]_{\substack{i=1, \dots, m-1 \\ j=1, \dots, n-1}}$, on obtient alors

$$A = U_1 B V_1^T, \quad B = \left[\begin{array}{c|c} \sigma_1 & w^T \\ \hline 0 & A_1 \end{array} \right].$$

Pour que cette décomposition corresponde à une SVD, il faut maintenant s'assurer que $w = \mathbf{0}_{n-1}$. Pour cela, on observe que :

$$B \begin{bmatrix} \sigma_1 \\ w \end{bmatrix} = \left[\begin{array}{c|c} \sigma_1 & w^T \\ \hline 0 & A_1 \end{array} \right] \begin{bmatrix} \sigma_1 \\ w \end{bmatrix} = \begin{bmatrix} \sigma_1^2 + w^T w \\ A_1 w \end{bmatrix}.$$

Par conséquent,

$$\left\| B \begin{bmatrix} \sigma_1 \\ \mathbf{w} \end{bmatrix} \right\| = \left\| \begin{bmatrix} \sigma_1^2 + \mathbf{w}^T \mathbf{w} \\ \mathbf{A}_1 \mathbf{w} \end{bmatrix} \right\| \geq \left\| \begin{bmatrix} \sigma_1^2 + \mathbf{w}^T \mathbf{w} \\ 0 \end{bmatrix} \right\| \geq \sigma_1^2 + \mathbf{w}^T \mathbf{w}$$

où l'on a utilisé la définition de la norme d'un vecteur pour obtenir l'avant-dernière égalité. On remarque que $\sigma_1^2 + \mathbf{w}^T \mathbf{w} = \left\| \begin{bmatrix} \sigma_1 \\ \mathbf{w} \end{bmatrix} \right\|^2$.

Il s'agit maintenant d'établir que $\mathbf{w} = 0$. Pour cela, on utilise le fait que $\|B\| = \|A\| = \sigma_1$, car la norme ne change pas par multiplication avec une matrice orthogonale (cf Lemme 2.1). Cela signifie que

$$\begin{aligned} \sigma_1 &= \max_{\substack{\mathbf{x} \in \mathbb{R}^n \\ \mathbf{x} \neq 0}} \frac{\|B\mathbf{x}\|}{\|\mathbf{x}\|} \\ &\geq \frac{\left\| B \begin{bmatrix} \sigma_1 \\ \mathbf{w} \end{bmatrix} \right\|}{\left\| \begin{bmatrix} \sigma_1 \\ \mathbf{w} \end{bmatrix} \right\|} \\ &\geq \frac{\sigma_1^2 + \mathbf{w}^T \mathbf{w}}{\sqrt{\sigma_1^2 + \mathbf{w}^T \mathbf{w}}} \\ &= \sqrt{\sigma_1^2 + \mathbf{w}^T \mathbf{w}}. \end{aligned}$$

Si $\mathbf{w}^T \mathbf{w} > 0$, alors on obtient $\sigma_1 > \sqrt{\sigma_1^2} = \sigma_1$, ce qui est absurde. On en déduit donc que $\mathbf{w}^T \mathbf{w} = \|\mathbf{w}\|^2 = 0$, donc que \mathbf{w} est un vecteur nul.

On a ainsi prouvé que l'on pouvait définir des matrices U_1 et V_1 telles que

$$A = U_1^T \left[\begin{array}{c|c} \sigma_1 & 0 \\ \hline 0 & A_1 \end{array} \right] V_1^T.$$

avec $\sigma_1 = \|A\|$. Si A_1 est la matrice nulle, alors nous avons bien obtenu une décomposition en valeurs singulières de A . Sinon, on applique le même raisonnement que précédemment à la matrice A_1 , de sorte à obtenir des matrices orthogonales $U_2 \in \mathbb{R}^{(m-1) \times (m-1)}$ et $V_2 \in \mathbb{R}^{(n-1) \times (n-1)}$ telles que

$$A_1 = U_2 \left[\begin{array}{c|c} \sigma_2 & 0 \\ \hline 0 & A_2 \end{array} \right] V_2.$$

où $\sigma_2 = \|A_1\|$ et $A_2 \in \mathbb{R}^{(m-2) \times (n-2)}$. Soient

$$U_{12} = U_1 \begin{bmatrix} 1 & \mathbf{0}_{m-1}^T \\ \mathbf{0}_{m-1} & U_2 \end{bmatrix}, \quad V_{12} = V_1 \begin{bmatrix} 1 & \mathbf{0}_{n-1}^T \\ \mathbf{0}_{n-1} & V_2 \end{bmatrix}.$$

En utilisant les propriétés de U_1 , U_2 , V_1 et V_2 , on peut vérifier que ces matrices sont bien orthogonales comme produits de matrices elles-mêmes orthogonales, et que l'on a

$$A = U_{12} \left[\begin{array}{c|c|c} \sigma_1 & 0 & 0 \\ \hline 0 & \sigma_2 & 0 \\ \hline 0 & 0 & A_2 \end{array} \right] V_{12}^T.$$

En répétant cette procédure jusqu'à obtenir un bloc A_i vide, on obtiendra une décomposition qui s'identifie à la décomposition en valeurs singulières. \square

Remarque 2.2 Si la preuve ci-dessus permet de construire une décomposition en valeurs singulières, elle n'est pas forcément aisée à utiliser en pratique. Le calcul des vecteurs \mathbf{u}_i et \mathbf{v}_i (et donc de $\|\mathbf{A}_i\|$) est en effet une étape coûteuse en grande dimension.

Les implémentations modernes de la décomposition en valeurs singulières reposent sur des techniques d'algèbre linéaire (QR avec pivotage, factorisation symétrique). La version la plus utilisée, basée sur l'algorithme de Golub et Kahan [3], possède de très bonnes garanties de stabilité numérique, ce qui est fondamental pour une implémentation efficace. Cela explique en partie le succès de la décomposition en valeurs singulières et son utilisation très répandue dans de nombreuses applications.

2.3.3 Décomposition tronquée et approximation

Le principal intérêt de la décomposition en valeurs singulières est de permettre de compresser la représentation de données matricielles. Dans de nombreuses applications, il est fréquent que les matrices de données présentent peu de valeurs singulières élevées, et beaucoup de petites valeurs singulières. On peut alors se demander quelle est la perte d'information que l'on réalise en omettant ces valeurs singulières.

La première réduction d'information que l'on peut opérer consiste à éliminer les valeurs singulières nulles dans la représentation de la matrice. C'est le sens du résultat ci-dessous.

Théorème 2.4 (SVD réduite) Toute matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$ de rang r admet une SVD réduite de la forme

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T, \quad (2.3.1)$$

où $\mathbf{U} \in \mathbb{R}^{m \times r}$ avec $\mathbf{U}^T \mathbf{U} = \mathbf{I}_r$, $\mathbf{V} \in \mathbb{R}^{n \times r}$ avec $\mathbf{V}^T \mathbf{V} = \mathbf{I}_r$ et $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$ est diagonale à coefficients diagonaux strictement positifs.

La décomposition (2.3.1) est plus compacte que la décomposition originelle. En particulier, il suffit de stocker $(n+m)r$ réels³ pour pouvoir reconstruire la matrice, ce qui peut être plus avantageux que de stocker les mn coefficients de la matrice \mathbf{A} .

Définition 2.9 (SVD tronquée) Soit $\mathbf{A} \in \mathbb{R}^{m \times n}$ une matrice de rang r et $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ sa SVD réduite, avec

$$\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_r], \quad \mathbf{V} = [\mathbf{v}_1 \cdots \mathbf{v}_r], \quad \mathbf{\Sigma} = \begin{bmatrix} \sigma_1 & 0 \cdots & 0 \\ 0 & \ddots & 0 \\ 0 & \cdots & 0 & \sigma_r \end{bmatrix}.$$

On suppose que $\sigma_1 \geq \cdots \geq \sigma_r$. Alors, pour tout $k \leq r$, la décomposition $\mathbf{U}_k \mathbf{\Sigma}_{k,k} \mathbf{V}_k^T$, où

$$\mathbf{U}_k = [\mathbf{u}_1 \cdots \mathbf{u}_k], \quad \mathbf{V}_k = [\mathbf{v}_1 \cdots \mathbf{v}_k], \quad \mathbf{\Sigma}_{k,k} = \begin{bmatrix} \sigma_1 & 0 \cdots & 0 \\ 0 & \ddots & 0 \\ 0 & \cdots & 0 & \sigma_k \end{bmatrix}$$

s'appelle la *décomposition en valeurs singulières tronquée à k valeurs*, ou *k-SVD*.

³Le coût est de $(n+m+1) \cdot r$ si on stocke les valeurs singulières séparément, mais on peut également les incorporer dans \mathbf{U} ou \mathbf{V} , auquel cas le coût de stockage sera de $(n+m) \cdot r$.

On remarque que cette factorisation est encore moins coûteuse en terme de coefficients que la décomposition en valeurs singulière réduite. Contrairement à celle-ci, la k -SVD supprime de l'information issue de la matrice \mathbf{A} dès lors que $k < r$: tout l'enjeu du processus de troncature consiste à préserver la majeure partie de l'information contenue dans la matrice, c'est-à-dire les valeurs singulières les plus importantes.

Application en compression d'images On considère une image 200x320 pixels stockée sous la forme d'une matrice $\mathbf{A} \in \mathbb{R}^{m \times n}$. Si on calcule la décomposition en valeurs singulières de \mathbf{A} , on peut voir que le rang de la matrice est $m = 200$.

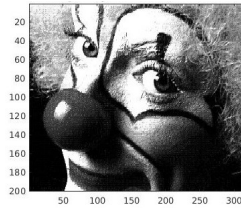


Figure 2.1: Image de clown 200x320 pixels; la matrice correspondante est de rang 200.

On peut cependant se demander si toutes les valeurs singulières sont nécessaires pour encoder l'image. On considère donc plusieurs troncatures de la *SVD*, correspondant à différentes valeurs de k inférieures au rang véritable. Les résultats de la figure 2.2 montrent qu'il n'est pas nécessaire de considérer la décomposition complète pour obtenir une image nette (voire indiscernable de l'image d'origine à l'oeil nu).

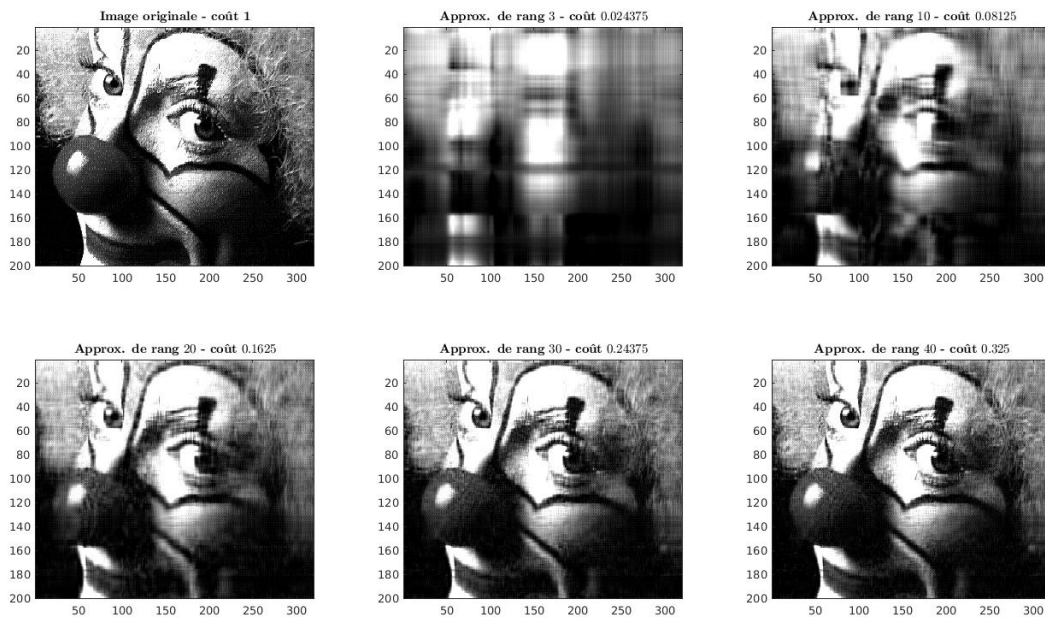


Figure 2.2: Image de clown 200x320 et SVD tronquées avec $k \in \{3, 10, 20, 30, 40\}$; le rang de chaque matrice est indiqué, ainsi que le ratio $\frac{(m+n)*k}{mn}$.

2.4 Exercices

Exercice 0 : Matrices rectangulaires

Soit $A \in \mathbb{R}^{m \times n}$. Démontrer les propriétés suivantes sans utiliser la décomposition en valeurs singulières :

- $\ker(A^T A) = \ker(A)$;
- $\text{rang}(A^T) = \text{rang}(A)$;
- $\text{rang}(A^T A) = \text{rang}(A)$, et en déduire que $\text{Im}(A^T A) = \text{Im}(A^T)$.

Question subsidiaire : Comment utiliser la décomposition en valeurs singulières pour établir ces propriétés ?

Exercice 1 : Valeurs propres et valeurs singulières

Soit $U \Sigma V^T$ une SVD de $A \in \mathbb{R}^{m \times n}$ avec $U = [u_1 \dots u_m]$, $V = [v_1 \dots v_n]$, et $\sigma_i = [\Sigma]_{ii} \forall i \leq \min\{m, n\}$.

- Montrer que

$$\forall i, 1 \leq i \leq \min\{m, n\}, \quad A v_i = \sigma_i u_i \text{ et } A^T u_i = \sigma_i v_i.$$

- En déduire un lien entre les valeurs propres de $A^T A$ et les valeurs singulières de A .
- Si $\text{rang}(A) = r$, que peut-on dire des valeurs singulières ?

Exercice 2 : Normes et valeurs singulières

- a) En utilisant une décomposition en valeurs singulières de \mathbf{A} , montrer que $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^{\min\{n,m\}} \sigma_i^2}$ et $\|\mathbf{A}\| = \sigma_1$, où $\sigma_1 \geq \dots \geq \sigma_{\min\{n,m\}}$ sont les valeurs singulières de \mathbf{A} .
- b) Si \mathbf{A} est carrée et diagonalisable, que vaut $\|\mathbf{A}\|$?

Exercice 3 : SVD tronquée

Soit $\mathbf{X} \in \mathbb{R}^{m \times n}$ de rang r , $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ une SVD réduite de \mathbf{X} avec $\sigma_1 > \dots > \sigma_r > 0$.

- a) On considère la SVD tronquée $\mathbf{X}_k = \mathbf{U}_k \mathbf{\Sigma}_{k,k} \mathbf{V}_k^T$. Evaluer $\|\mathbf{X} - \mathbf{X}_k\|_F^2$.
- b) Application : Supposons que $\sigma_i^2 = \frac{1}{2^i} \forall i = 1, \dots, m$. Pour tout $\epsilon > 0$, déterminer le plus petit k tel que $\|\mathbf{X} - \mathbf{X}_k\|_F^2 \leq \epsilon^2$.

Exercice 4 : Carrés de matrices (CC 2019/2020)

Soit $\mathbf{A} \in \mathbb{R}^{n \times n}$ symétrique et définie positive. Soit $\mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$ une décomposition en valeurs propres de \mathbf{A} , avec \mathbf{P} orthogonale et $\mathbf{\Lambda}$ diagonale à coefficients diagonaux ordonnés par ordre décroissant $\Lambda_{11} \geq \Lambda_{22} \geq \dots \geq \Lambda_{nn}$.

- a) Donner une décomposition en valeurs singulières de \mathbf{A} basée sur la décomposition en valeurs propres ci-dessus.
- b) Soit une matrice symétrique $\mathbf{B} \in \mathbb{R}^{n \times n}$ telle que $\mathbf{B}^2 = \mathbf{A}$ et $\mathbf{B} \succ 0$. Donner une décomposition en valeurs propres et une décomposition en valeurs singulières de \mathbf{B} en fonction de celles de \mathbf{A} .
- c) Que vaut $\mathbf{B}^T \mathbf{B}$? Retrouver alors le lien entre les valeurs propres de $\mathbf{B}^T \mathbf{B}$ et les valeurs singulières de \mathbf{B} .

Chapitre 3

Premiers pas avec le modèle linéaire

Dans le chapitre précédent, nous avons introduit la décomposition en valeurs singulières, une technique visant à extraire de l'information d'un jeu de données : ce paradigme est celui de l'apprentissage *non supervisé*. Ce chapitre aborde un autre paradigme, celui de l'apprentissage *supervisé*, dans lequel il s'agira de déterminer un modèle (une fonction linéaire pour les besoins de ce cours) décrivant une relation entre différents éléments d'un jeu de données, pour potentiellement prédire (on parle également d'inférer) le comportement de données futures.

3.1 Introduction

On considère un jeu de données ayant m éléments ou individus, et on associe à chaque individu n caractéristiques¹ sous la forme d'un vecteur de \mathbb{R}^n . Soient $\mathbf{x}_1, \dots, \mathbf{x}_m$ ces vecteurs : on les regroupe alors sous la forme d'une matrice de données

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix} \in \mathbb{R}^{m \times n}. \quad (3.1.1)$$

Exemple 3.1 • Chaque ligne de \mathbf{X} représente un individu, et les n composantes de \mathbf{x}_i sont des données médicales (âge, poids, taux de cholestérol, etc).

- Chaque ligne de \mathbf{X} est une "vectorisation" d'une image 2D, et les valeurs de \mathbf{x}_i sont celles des pixels, en niveau de gris. Ainsi, une image de taille 480*640 serait transformée (en mettant les lignes bout à bout, par exemple) en un vecteur de taille $n = 480 * 640 = 307200$.

Sans autre information que la matrice elle-même, on peut appliquer des techniques d'algèbre linéaire pour en extraire de l'information : c'est ce que réalisait la SVD dans le chapitre précédent. Dans un contexte d'apprentissage supervisé, on associe chaque vecteur de caractéristiques \mathbf{x}_i à un **label** $y_i \in \mathbb{R}$, qui peut représenter une classe à laquelle l'individu appartient (malade/non malade, image de chien ou de chat, etc). Ces labels sont concaténés pour former un vecteur de labels $\mathbf{y} \in \mathbb{R}^m$.

Par conséquent, notre but n'est plus seulement d'analyser l'information de la matrice \mathbf{X} , mais bien de trouver une relation entre les caractéristiques \mathbf{X} et les labels \mathbf{y} . Pour ce cours, on postulera que

¹Ou *features* en anglais.

cette relation est linéaire : on va donc chercher une fonction $h : \mathbb{R}^n \rightarrow \mathbb{R}$ de la forme $h(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}$. On souhaite que h permette d'obtenir les y_i à partir des x_i , c'est-à-dire que l'on voudrait avoir

$$h(\mathbf{x}_i) = \mathbf{x}_i^T \boldsymbol{\beta} = y_i \quad \forall i = 1, \dots, m,$$

que l'on peut ré-écrire sous forme matricielle comme

$$\mathbf{X} \boldsymbol{\beta} = \mathbf{y}.$$

On se trouve donc en présence d'un système linéaire que l'on va vouloir résoudre. Rien ne garantit a priori que ce système possède une solution, ou que cette solution (si elle existe) est unique. Une étude plus approfondie des systèmes d'équations linéaires semble donc nécessaire.

3.2 Résolution de systèmes non linéaires

Définition 3.1 Un système linéaire de m équations à n inconnues β_1, \dots, β_n est donné par

$$\begin{array}{cccccc} x_{11}x_1 & + & x_{12}\beta_2 & + & \dots & + & x_{1n}\beta_n & = & y_1 \\ x_{21}x_1 & + & x_{22}\beta_2 & + & \dots & + & x_{2n}\beta_n & = & y_2 \\ \vdots & & \vdots & & \vdots & & \vdots & & \vdots \\ x_{m1}x_1 & + & x_{m2}\beta_2 & + & \dots & + & x_{mn}\beta_n & = & y_m \end{array}$$

ou, sous forme compacte,

$$\mathbf{X} \boldsymbol{\beta} = \mathbf{y},$$

avec $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\boldsymbol{\beta} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$.

Dans la suite, nous allons déterminer les conditions d'existence de solutions à ce système linéaire.

3.2.1 Cas d'un système carré

On s'intéresse tout d'abord au cas où $n = m$: on a donc un système dit "carré" de n équations à n inconnues de la forme

$$\mathbf{X} \boldsymbol{\beta} = \mathbf{y}, \tag{3.2.1}$$

où $\mathbf{X} \in \mathbb{R}^{n \times n}$, $\mathbf{y} \in \mathbb{R}^n$, et $\boldsymbol{\beta} \in \mathbb{R}^n$ représente les paramètres inconnus de notre modèle.

L'existence et l'unicité de solutions au système (3.2.1) dépendent des propriétés de la matrice \mathbf{X} et du vecteur \mathbf{y} , comme le montrent les exemples ci-dessous.

Exemple 3.2 a) Le système

$$\begin{cases} \beta_1 + \beta_2 & = & 0, \\ 3\beta_1 + 2\beta_2 & = & 1. \end{cases}$$

possède une unique solution $\beta_1 = 1$, $\beta_2 = -1$.

b) Le système

$$\begin{cases} \beta_1 + \beta_2 & = & 0, \\ 2\beta_1 + 2\beta_2 & = & -1. \end{cases}$$

ne possède pas de solution.

c) Le système

$$\begin{cases} \beta_1 + 2\beta_2 = 2, \\ 2\beta_1 + 4\beta_2 = 4. \end{cases}$$

possède une infinité de solutions.

On peut donc se trouver dans trois cas différents. Les deux derniers vont demander une analyse plus détaillée. En revanche, le premier cas correspond à une matrice \mathbf{X} inversible : dans cette situation, il existe une caractérisation de la solution du système.

Théorème 3.1 (Résolution d'un système carré inversible) Soient $\mathbf{X} \in \mathbb{R}^{n \times n}$ une matrice inversible et $\mathbf{y} \in \mathbb{R}^n$. Le système carré inversible $\mathbf{X}\boldsymbol{\beta} = \mathbf{y}$ possède une unique solution $\boldsymbol{\beta}^*$ donnée par

$$\boldsymbol{\beta}^* = \mathbf{X}^{-1}\mathbf{y}.$$

Lorsque la matrice \mathbf{X} n'est pas inversible en revanche, c'est-à-dire que $\text{rang}(\mathbf{X}) < n$, il existera une infinité de solutions si $\mathbf{y} \in \text{Im}(\mathbf{X})$, et aucune si $\mathbf{y} \notin \text{Im}(\mathbf{X})$.

3.2.2 Cas d'un système rectangulaire

On considère maintenant le cas d'un système linéaire rectangulaire non carré, soit

$$\mathbf{X}\boldsymbol{\beta} = \mathbf{y}, \quad \mathbf{X} \in \mathbb{R}^{m \times n}, \quad m \neq n. \quad (3.2.2)$$

Comme le montrent les exemples ci-dessous, on retrouve les mêmes cas que pour un système carré.

Exemple 3.3 a) Le système

$$\begin{aligned} \beta_1 + \beta_2 &= 0, \\ 3\beta_1 + 2\beta_2 &= 1, \\ 6\beta_1 + 5\beta_2 &= 1. \end{aligned}$$

possède une unique solution ($\beta_1 = 1, \beta_2 = -1$).

b) Le système

$$\beta_1 + 2\beta_2 = 2.$$

possède une infinité de solutions.

c) Le système

$$\begin{aligned} \beta_1 &= 2, \\ \beta_2 &= 3, \\ \beta_1 + \beta_2 &= 0. \end{aligned}$$

ne possède pas de solution.

Il est donc nécessaire d'analyser le système afin de déterminer s'il possède ou non des solutions. On peut être plus spécifique pour certains types de systèmes dont la matrice \mathbf{X} est de rang plein.

Définition 3.2 Une matrice $\mathbf{X} \in \mathbb{R}^{m \times n}$ est de **rang plein** si $\text{rang}(\mathbf{X}) = \min\{m, n\}$.

Toute matrice carrée inversible est de rang plein, mais cette notion est plus générique. On a ainsi les cas particuliers suivants.

Théorème 3.2 Soient $\mathbf{X} \in \mathbb{R}^{m \times n}$ de rang plein et $\mathbf{y} \in \mathbb{R}^m$ avec $m \neq n$. On a :

- a) Si $\text{rang}(\mathbf{X}) = n$, alors $\mathbf{X}^T \mathbf{X}$ est inversible et si $\mathbf{y} \in \text{Im}(\mathbf{X})$, le système (3.2.2) possède une unique solution donnée par $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$;
- b) Si $\text{rang}(\mathbf{X}) = m$, alors $\mathbf{X} \mathbf{X}^T$ est inversible et le vecteur $\mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$ est solution de (3.2.2)

Il est donc toujours possible de déterminer une solution au problème (3.2.2) lorsque \mathbf{X} est de rang plein et que $\mathbf{y} \in \text{Im}(\mathbf{X})$. Lorsque \mathbf{X} n'est pas de rang plein, on retombe en revanche dans les mêmes difficultés que pour un système carré.

Comme on le verra dans la section suivante, les différentes expressions établies dans les théorèmes (3.1) et (3.2) correspondent en fait à une même formule.

3.2.3 Pseudo-inverse et SVD

Le concept d'inverse d'une matrice carrée peut être généralisé au cas d'une matrice rectangulaire : c'est le principe de la pseudo-inverse² que l'on décrit ci-dessous.

Définition 3.3 (Pseudo-inverse d'une matrice) Soit une matrice $\mathbf{X} \in \mathbb{R}^{m \times n}$. Il existe une unique matrice $\mathbf{M} \in \mathbb{R}^{n \times m}$ vérifiant les équations de Penrose :

$$\begin{aligned} \mathbf{X} \mathbf{M} \mathbf{X} &= \mathbf{X} \\ \mathbf{M} \mathbf{X} \mathbf{M} &= \mathbf{M} \\ (\mathbf{X} \mathbf{M})^T &= \mathbf{X} \mathbf{M} \\ (\mathbf{M} \mathbf{X})^T &= \mathbf{M} \mathbf{X} \end{aligned}$$

Cette matrice s'appelle la **pseudo-inverse** de \mathbf{X} et on la note \mathbf{X}^\dagger .

Sous sa forme générale, le calcul de la pseudo-inverse ne semble pas évident. Il existe heureusement une formule explicite de la pseudo-inverse basée sur la décomposition en valeurs singulières.

Proposition 3.1 Soit une matrice diagonale par blocs $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ de la forme

$$\left[\begin{array}{ccc|c} \sigma_1 & 0 \cdots & 0 & 0 \\ 0 & \ddots & 0 & \vdots \\ 0 & \cdots 0 & \sigma_r & 0 \\ \hline 0 & \cdots & \cdots & 0 \end{array} \right]$$

avec $\sigma_1 \geq \cdots \geq \sigma_r > 0$. La pseudo-inverse de la matrice $\mathbf{\Sigma}$ est la matrice $\mathbf{\Sigma}^\dagger \in \mathbb{R}^{n \times m}$ définie par

$$\mathbf{\Sigma}^\dagger = \left[\begin{array}{ccc|c} \frac{1}{\sigma_1} & 0 \cdots & 0 & 0 \\ 0 & \ddots & 0 & \vdots \\ 0 & \cdots 0 & \frac{1}{\sigma_r} & 0 \\ \hline 0 & \cdots & \cdots & 0 \end{array} \right].$$

²On parle aussi d'inverse généralisée ou d'inverse de Moore-Penrose.

Le résultat de la proposition 3.1 illustre bien le concept de pseudo-inverse : on a ainsi “inversé” uniquement le bloc diagonal contenant des valeurs non nulles, le reste étant simplement transposé pour inverser les espaces de départ et d’arrivée.

On a alors le résultat générique suivant.

Théorème 3.3 (Formule de pseudo-inverse) *Soit une matrice $X \in \mathbb{R}^{m \times n}$ et $U\Sigma V^T$ une décomposition en valeurs singulières de cette matrice. Alors, la pseudo-inverse de X est donnée par*

$$X^\dagger = V\Sigma^\dagger U^T. \quad (3.2.3)$$

On peut vérifier que la formule (3.2.3) satisfait bien aux équations de Penrose. Ces dernières suggèrent qu’il est possible d’utiliser la pseudo-inverse d’une manière similaire à l’inverse d’une matrice carrée inversible : le lien entre inverse et pseudo-inverse est encore plus ténu, comme le montre le corollaire suivant.

Corollaire 3.1 *Soit $X \in \mathbb{R}^{m \times n}$ une matrice de rang plein. Alors,*

- i) *Si $\text{rang}(X) = m$, alors $X^\dagger = X^T(XX^T)^{-1}$;*
- ii) *Si $\text{rang}(X) = n$, alors $X^\dagger = (X^T X)^{-1}X^T$;*
- iii) *Si $\text{rang}(X) = n$ et que $n = m$, alors X est carrée inversible, et les deux formules ci-dessus correspondent à $X^\dagger = X^{-1}$.*

La pseudo-inverse est ainsi apparue dans les sections 3.2.1 et 3.2.2 lorsque l’on supposait que la matrice X était de rang plein. En ce sens, il semble que la technique de pseudo-inverse soit adaptée aux problèmes bien posés. L’approche par moindres carrés, développée dans la section suivante, va permettre de formaliser cette propriété, et de mettre en lumière le rôle plus large joué par la pseudo-inverse.

3.3 Moindres carrés linéaires

Comme expliqué dans la section précédente, la notion de solution d’un système linéaire perd de son sens lorsque le système ne possède pas de solution, ou une infinité. Pour cette raison, on définit un autre concept de solution, dite au sens des moindres carrés.

3.3.1 Solution au sens des moindres carrés

Un problème aux moindres carrés linéaires est un problème d’optimisation, et plus précisément de minimisation : étant donnés $X \in \mathbb{R}^{m \times n}$ et $y \in \mathbb{R}^m$, on ne cherche plus à résoudre $X\beta = y$ de manière exacte, mais plutôt à minimiser l’écart entre les vecteurs $X\beta$ et y . Du point de vue mathématique, on évaluera cet écart via la norme euclidienne, et on cherchera donc à résoudre le problème (dit aux moindres carrés linéaires) suivant :

$$\min_{\beta \in \mathbb{R}^n} \frac{1}{2} \|X\beta - y\|^2. \quad (3.3.1)$$

Il s'agit d'un problème de minimisation de la fonction $\beta \mapsto \frac{1}{2}\|\mathbf{X}\beta - \mathbf{y}\|^2$ selon β , sous sa forme standard en optimisation.³ Notons que si l'on décompose la norme, on obtient :

$$\|\mathbf{X}\beta - \mathbf{y}\|^2 = \sum_{i=1}^m (\mathbf{x}_i^T \beta - y_i)^2.$$

Cet objectif que l'on cherche à réduire représente donc bien l'attachement aux données, c'est-à-dire la correspondance entre notre modèle linéaire et les labels de chaque individu.

Définition 3.4 (Solution au sens des moindres carrés) Soient $\mathbf{X} \in \mathbb{R}^{m \times n}$ et $\mathbf{y} \in \mathbb{R}^m$. L'ensemble des vecteurs $\beta \in \mathbb{R}^n$ tels que la valeur de $\frac{1}{2}\|\mathbf{X}\beta - \mathbf{y}\|^2$ soit minimale, que l'on note :

$$\arg \min_{\beta \in \mathbb{R}^n} \frac{1}{2}\|\mathbf{X}\beta - \mathbf{y}\|^2, \quad (3.3.2)$$

s'appelle l'ensemble des **solutions au sens des moindres carrés**.

La notion de solution au sens des moindres carrés est à distinguer de celle d'une solution du système linéaire, pour laquelle nous introduisons la terminologie ci-dessous.

Définition 3.5 (Solution au sens classique) Soient $\mathbf{X} \in \mathbb{R}^{m \times n}$ et $\mathbf{y} \in \mathbb{R}^m$. L'ensemble des vecteurs $\beta \in \mathbb{R}^n$ tels que $\mathbf{X}\beta = \mathbf{y}$ s'appelle l'ensemble des **solutions au sens classique** du système linéaire.

Cet ensemble peut être vide, et il est nécessairement inclus dans l'ensemble des solutions au sens des moindres carrés.

Le concept de solution au sens des moindres carrés nous permet d'introduire différents concepts, liés notamment à la pseudo-inverse.

Théorème 3.4 Pour tous $\mathbf{X} \in \mathbb{R}^{m \times n}$ et $\mathbf{y} \in \mathbb{R}^m$, les propriétés suivantes sont vérifiées :

1. l'ensemble défini par (3.3.2) est toujours non vide;
2. le vecteur $\mathbf{X}^\dagger \mathbf{y}$ est une solution au sens des moindres carrés;
3. parmi toutes les solutions au sens des moindres carrés, le vecteur $\mathbf{X}^\dagger \mathbf{y}$ est la solution de norme minimale :

$$\text{pour tout } \mathbf{v} \in \arg \min_{\beta \in \mathbb{R}^n} \frac{1}{2}\|\mathbf{X}\beta - \mathbf{y}\|^2, \text{ on a } \|\mathbf{v}\| \geq \|\mathbf{X}^\dagger \mathbf{y}\|.$$

Ce théorème ne sera pas démontré dans ce cours (voir à cet égard le cours de *Méthodes numériques : Optimisation* au semestre 2, ou le cours *Mathématiques pour les sciences des données*). En revanche, nous exploiterons fortement ce résultat pour calculer une solution au sens des moindres carrés.

³Pour des raisons de normalisation, on introduit notamment un facteur 1/2.

3.3.2 Résolution du problème aux moindres carrés

Le théorème ci-dessous récapitule l'ensemble des cas à considérer dans le calcul d'un modèle linéaire.

Théorème 3.5 Soient $\mathbf{X} \in \mathbb{R}^{m \times n}$, $\mathbf{y} \in \mathbb{R}^m$ et $\beta^* = \mathbf{X}^\dagger \mathbf{y}$.

1) Si $m = n$, on distingue trois cas :

- a) Si $\text{rang}(\mathbf{X}) = m = n$, on a $\beta^* = \mathbf{X}^{-1} \mathbf{y}$, et il s'agit de l'unique solution au sens classique et au sens des moindres carrés;
- b) Si $\text{rang}(\mathbf{X}) < m$ et $\mathbf{y} \in \text{Im}(\mathbf{X})$, alors β^* est une solution au sens classique et de norme minimale au sens des moindres carrés. Les problèmes (3.2.1) et (3.3.1) admettent chacun une infinité de solutions.
- c) Si $\text{rang}(\mathbf{X}) < m$ et $\mathbf{y} \notin \text{Im}(\mathbf{X})$, alors il n'existe pas de solution au sens classique; en revanche, β^* est la solution de norme minimale au sens des moindres carrés.

2) Si $m < n$ (système sous-déterminé), on distingue trois cas :

- b) Si $\text{rang}(\mathbf{X}) = m$, alors $\beta^* = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$ et il s'agit à la fois d'une solution classique et de la solution de norme minimale au sens des moindres carrés. Les problèmes (3.2.2) et (3.3.1) admettent chacun une infinité de solutions.
- b) Si $\text{rang}(\mathbf{X}) < m$ et $\mathbf{y} \in \text{Im}(\mathbf{X})$, alors β^* est une solution au sens classique et de norme minimale au sens des moindres carrés. Les problèmes (3.2.1) et (3.3.1) admettent chacun une infinité de solutions.
- b) Si $\text{rang}(\mathbf{X}) < m$ et $\mathbf{y} \notin \text{Im}(\mathbf{X})$, alors il n'existe pas de solution au sens classique; en revanche, β^* est la solution de norme minimale au sens des moindres carrés.

3) Si $m > n$ (système sur-déterminé), on distingue trois cas :

- a) Si $\text{rang}(\mathbf{X}) = n$, alors $\beta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ et il s'agit de l'unique solution au sens des moindres carrés. C'est une solution au sens classique lorsque $\mathbf{y} \in \text{Im}(\mathbf{X})$.
- b) Si $\text{rang}(\mathbf{X}) < n$ et $\mathbf{y} \in \text{Im}(\mathbf{X})$, alors β^* est une solution au sens classique et de norme minimale au sens des moindres carrés. Les problèmes (3.2.1) et (3.3.1) admettent chacun une infinité de solutions.
- c) Si $\text{rang}(\mathbf{X}) < n$ et $\mathbf{y} \notin \text{Im}(\mathbf{X})$, alors il n'existe pas de solution au sens classique; en revanche, β^* est la solution de norme minimale au sens des moindres carrés.

3.4 Conclusion

Dans le contexte de l'apprentissage supervisé, on peut être amené à vouloir expliquer nos données par un modèle linéaire. Ce choix de modélisation fait naturellement apparaître des systèmes linéaires, dont l'étude repose sur des propriétés issues de l'algèbre linéaire. Il apparaît alors que le problème peut être bien ou mal posé, selon que le système possède une, des ou même aucune solution(s).

On a ainsi introduit le concept de *problème aux moindres carrés associé à un système linéaire*, qui a permis de formuler la procédure d'apprentissage du modèle linéaire en prenant en compte les cas où le modèle ne peut pas expliquer les données de manière unique. Grâce à la pseudo-inverse, nous avons pu caractériser une solution du problème qui fournit la meilleure erreur d'approximation au sens des moindres carrés, tout en permettant d'avoir un modèle plus simple au sens de la norme.

3.5 Exercices

Exercice 5 : Modélisation par moindres carrés

On dispose de m notes y_1, \dots, y_m données par un même voyageur durant m séjours d'une nuit dans un hôtel. On cherche à déterminer une note globale du voyageur sur cet hôtel qui soit la plus cohérente possible avec l'ensemble des m notes.

- Modéliser le problème sous la forme d'un système linéaire et d'un problème aux moindres carrés.
- Pour les deux formulations obtenues, décrire si le problème possède ou non une solution, et donner l'ensemble des solutions s'il existe.
- Si les séjours n'étaient pas identiques, comment pourrait-on tenir compte de cela dans la modélisation ? Quel serait alors l'impact sur la résolution des problèmes associés ?

Exercice 6 : Projections et moindres carrés

Soient $\mathbf{y} \in \mathbb{R}^m$ et \mathcal{S} un sous-espace vectoriel de dimension $n < m$ de \mathbb{R}^m . La *projection orthogonale de \mathbf{y} sur \mathcal{S}* est définie comme la solution du problèmes aux moindres carrés avec contraintes :

$$\min_{\gamma \in \mathcal{S}} \frac{1}{2} \|\gamma - \mathbf{y}\|^2.$$

- En utilisant une base orthonormée de \mathcal{S} , former un problème sans contraintes équivalent au problème de départ;
- Quels sont les deux cas à considérer pour étudier l'ensemble des solutions ? Décrire et interpréter l'ensemble des solutions dans chaque cas.

Exercice 7 : Matrices de rang 1

Soit une matrice $\mathbf{X} \in \mathbb{R}^{m \times 1}$ de rang 1.

- Donner une décomposition en valeurs singulières de \mathbf{X} .
- Soit $\mathbf{y} \in \mathbb{R}^m$. Résoudre le problème aux moindres carrés

$$\min_{\beta \in \mathbb{R}} \frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|^2.$$

- On suppose que \mathbf{X} est une matrice orthogonale. D'après l'exercice 6, que représente alors $\mathbf{X}\beta^*$, où β^* est la solution du problème aux moindres carrés ?

Bibliographie

- [1] S. Boyd and L. Vandenberghe. *Introduction to Applied Linear Algebra - Vectors, Matrices and Least Squares*. Cambridge University Press, Cambridge, United Kingdom, 2018.
- [2] S. L. Brunton and J. N. Kutz. *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press, Cambridge, United Kingdom, 2019.
- [3] G. H. Golub and C. F. Van Loan. *Matrix computations*. The Johns Hopkins University Press, Baltimore, fourth edition, 2013.