

Cours ANN 203 – Examen écrit, vendredi 6 mai 2022

NOTE: Les matrices et les vecteurs considérés dans tous les exercices sont supposés à valeurs réelles. La notation $\det(X)$ désigne le déterminant de la matrice carrée X .

Exercice E22-1 Calculer la décomposition QR de la matrice

$$A = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

Le réflecteur de Householder pour la première (et unique) colonne de A est

$$H = I - 2 \frac{vv^H}{\|v\|^2} = \frac{1}{5} \begin{bmatrix} -3 & -4 \\ -4 & 3 \end{bmatrix} \quad (v = \begin{Bmatrix} 8 \\ 4 \end{Bmatrix})$$

On trouve alors

$$A = QR, \quad \text{avec } Q = H, R = HA = \begin{bmatrix} -5 \\ 0 \end{bmatrix}.$$

Exercice E22-2 (itérations de Jacobi pour le problème aux valeurs propres symétrique) *Cet exercice utilise la norme de Frobenius $\|A\|_F$ d'une matrice A . On rappelle (pour une matrice réelle A) que*

$$\|A\|_F^2 := \sum_{i,j} a_{ij}^2 = \text{Tr}(AA^T) = \text{Tr}(A^T A).$$

Soit $A \in \mathbb{R}^{n \times n}$ une matrice réelle symétrique. On considère dans cet exercice une approche alternative à celles enseignées dans le cours, appelée itérations de Jacobi, pour la détermination des valeurs propres de A . Cette méthode utilise la décomposition additive $A = A_D + A_N$ de A , où A_D est la partie diagonale de A (et donc A_N est obtenue en remplaçant la diagonale de A par zéro). L'idée générale de l'algorithme est de réduire $\|A_N\|_F$ à chaque itération (dans le but de converger vers une matrice diagonale). Chaque itération est de la forme $A^{(k+1)} = Q^{(k)} A^{(k)} Q^{(k)T}$ où $Q^{(k)}$ est une matrice orthogonale choisie de façon à annuler certains termes de $A^{(k)}$.

- (a) Pour toute matrice carrée A et toute matrice orthogonale $Q \in \mathbb{R}^{n \times n}$, montrer que (i) on a $\|QAQ^T\|_F = \|A\|_F$, et (ii) les matrices A et QAQ^T ont les mêmes valeurs propres.
- (b) Cas $n = 2$: soit $A = \begin{bmatrix} a & b \\ b & d \end{bmatrix}$ une matrice symétrique avec $b \neq 0$ (on a donc $A_D = \begin{bmatrix} a & 0 \\ 0 & d \end{bmatrix}$). Montrer qu'il existe une matrice orthogonale $Q = \begin{bmatrix} c & -s \\ s & c \end{bmatrix}$, définie par deux réels c, s vérifiant $c^2 + s^2 = 1$, telle que $B := QAQ^T$ est diagonale. On déterminera les réels c, s (deux solutions possibles) en fonction de a, b, d , et pourra pour cela introduire le paramètre $\tau = (d - a)/2b$ et l'inconnue auxiliaire $t = s/c$. Ce procédé permet donc (sans surprise) de diagonaliser $A \in \mathbb{R}^{2 \times 2}$ symétrique en une seule itération.
- (c) On extrapole maintenant le procédé précédent à $n \geq 2$ quelconque. On choisit deux lignes k, ℓ de A (pour lesquelles $a_{k\ell} \neq 0$). Définir une matrice orthogonale $Q(k, \ell) \in \mathbb{R}^{n \times n}$ telle que $B := Q(k, \ell)AQ^T(k, \ell)$ laisse les lignes et colonnes de A autres que k, ℓ inchangées et introduise un zéro aux positions $a_{k\ell}$ et $a_{\ell k}$ de A .

Montrer alors, posant $B = B_D + B_N$, que

$$\|B_N\|_F^2 = \|A_N\|_F^2 - 2a_{k\ell}^2,$$

et interpréter ce résultat (utilité de la méthode de calcul de B ci-dessus, meilleur choix de k, ℓ).

Estimer (à l'ordre principal en n) le nombre d'opérations arithmétiques nécessaires à l'évaluation de $Q(k, \ell)AQ^T(k, \ell)$.

- (d) Proposer sur la base des questions précédentes un algorithme itératif susceptible de fournir une approximation des valeurs et vecteurs propres de A , et expliquer comment et à quelle condition ces approximations sont obtenues. Proposer un critère d'arrêt des itérations de Jacobi.

Eléments de solution:

- (a) Pour montrer (i), on écrit (en exploitant l'orthogonalité de Q via $Q^T Q = I$)

$$\begin{aligned} \|QAQ^T\|_F^2 &= \text{Tr}\{QAQ^T(QAQ^T)^T\} = \text{Tr}(QAQ^TQA^TQ^T) = \text{Tr}(QAA^TQ^T) \\ &= \text{Tr}(AA^TQ^TQ) = \text{Tr}(AA^T) = \|A\|_F^2. \end{aligned}$$

La propriété (ii) découle (à nouveau par orthogonalité de Q) de

$$QAQ^T - \lambda I = QAQ^T - \lambda QQ^T = Q(A - \lambda I)Q^T,$$

qui permet de déduire

$$\det(QAQ^T - \lambda I) = \det\{Q(A - \lambda I)Q^T\} = \det\{Q^T Q(A - \lambda I)\} = \det(A - \lambda I).$$

Les matrices $QAQ^T - \lambda I$ et $A - \lambda I$ ont donc le même polynôme caractéristique, et ainsi les mêmes valeurs propres.

- (b) Un calcul élémentaire donne, en utilisant les quantités suggérées

$$(QAQ^T)_{12} = (QAQ^T)_{21} = (a-d)cs + (c^2 - s^2)b \quad , \quad (QAQ^T)_{12} = 0 \Leftrightarrow t^2 + 2\tau t - 1 = 0.$$

La condition d'annulation donne deux valeurs $t = -\tau \pm \sqrt{\tau^2 + 1}$, conduisant à quatre couples (c, s) possibles:

$$c = \pm \frac{1}{1+t^2}, \quad s = ct, \quad t = -\tau \pm \sqrt{\tau^2 + 1}$$

- (c) On définit $Q(k, \ell) \in \mathbb{R}^{n \times n}$ comme la matrice identité en-dehors des lignes et colonnes k, ℓ , dont les intersections forment une matrice 2×2 de rotation:

$$R(k, \ell, p) = \begin{bmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & c & \dots & -s & \dots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \dots & s & \dots & c & \dots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{bmatrix}$$

De plus, les coefficients c, s sont calculés comme en (b) pour $\begin{bmatrix} a & b \\ b & d \end{bmatrix} = \begin{bmatrix} a_{kk} & a_{k\ell} \\ a_{\ell k} & a_{\ell\ell} \end{bmatrix}$. Le calcul de $B := Q(k, \ell)AQ^T(k, \ell)$ ne modifie que les lignes et colonnes k, ℓ de A .

D'après (a), on a $\|B\|_F = \|A\|_F$. D'autre part, pour toute matrice $X \in \mathbb{R}^{n \times n}$, on a $\|X\|^2 = \|X_D\|_F^2 + \|X_N\|_F^2$. Par conséquent:

$$\|B_N\|_F^2 = \|A_N\|_F^2 + \|A_D\|_F^2 - \|B_D\|_F^2 = \|A_N\|_F^2 + \|A_D^{k\ell}\|_F^2 - \|B_D^{k\ell}\|_F^2,$$

où $A^{k\ell}, B^{k\ell} \in \mathbb{R}^{2 \times 2}$ sont les sous-matrices obtenues en prenant les lignes et colonnes k, ℓ de A, B . Les résultats de (a), (b) appliqués à $A^{k\ell}, B^{k\ell}$ montrant que $\|A^{k\ell}\|_F^2 = \|B^{k\ell}\|_F^2$ et $\|B^{k\ell}\|_F^2 = \|B_D^{k\ell}\|_F^2$ (puisque $B_N^{k\ell} = 0$ par construction de $B^{k\ell}$), on déduit

$$\|A_D^{k\ell}\|_F^2 - \|B_D^{k\ell}\|_F^2 = \|A_D^{k\ell}\|_F^2 - \|B^{k\ell}\|_F^2 = \|A_D^{k\ell}\|_F^2 - \|A^{k\ell}\|_F^2 = -\|A_N^{k\ell}\|_F^2 = -2a_{k\ell}^2$$

et on obtient donc $\|B_N\|_F^2 = \|A_N\|_F^2 - 2a_{k\ell}^2$: chaque itération de Jacobi réduit $\|A_N\|_F$.

Exercice E22-3 (choix du paramètre de relaxation dans la méthode itérative SOR) Soit $A \in \mathbb{R}^{n \times n}$ inversible, que l'on décompose additivement comme $A = D - L - U$ (où D , $-L$, $-U$ sont la diagonale et les parties triangulaires inférieure et supérieure strictes de A , respectivement). La méthode itérative SOR (successive over-relaxation) pour le système linéaire $Ax = b$ repose sur les itérations définies par

$$[D - \eta L]x_{k+1} = [\eta U + (1 - \eta)D]x_k + \eta b, \quad k = 0, 1, 2, \dots \quad (x_0 \in \mathbb{R}^n: \text{initialisation arbitraire}),$$

ou $\eta \in \mathbb{R}$ est le paramètre de relaxation de la méthode. On observe que chaque itération demande la résolution d'un système triangulaire. Le choix de η est important, et l'objet de cet exercice est de montrer que certaines valeurs de η ne peuvent pas convenir, quels que soient A, b .

- Expliquer pourquoi les itérations SOR sont susceptibles de fournir la solution de $Ax = b$.
- Mettre l'itération SOR générique sous la forme $x_{k+1} = R(\eta)x_k + f(\eta)$, avec $R \in \mathbb{R}^{n \times n}$ et $f \in \mathbb{R}^n$. Quelle est la condition nécessaire et suffisante sur la matrice d'itération $R(\eta)$ pour que les itérations SOR convergent quelle que soit l'initialisation x_0 ?
- Montrer que $\det(R(\eta)) = (1 - \eta)^n$. En déduire un ensemble de valeurs de η pour lesquelles la satisfaction de la condition nécessaire et suffisante du (b) est impossible.

L'algorithme SOR n'est ainsi susceptible de converger que pour les valeurs de η non exclues par le résultat de (b). On rappelle que $\det(AB) = \det(A)\det(B)$ quelles que soient $A, B \in \mathbb{R}^{n \times n}$.

Eléments de solution:

- Supposons que la suite x_k converge vers une limite x . En passant à la limite dans l'équation définissant les itérations SOR, x vérifie alors

$$[D - \eta L]x = [\eta U + (1 - \eta)D]x + \eta b \implies [\eta D - \eta L - \eta U]x = \eta b,$$

et est donc solution du système linéaire $Ax = b$ à résoudre. La convergence de la suite x_k n'est cependant pas *a priori* garantie.

- L'itération SOR générique s'écrit (multiplication à gauche par $[D - \eta L]^{-1}$)

$$x_{k+1} = R(\eta)x_k + f(\eta), \quad R(\eta) := [D - \eta L]^{-1}[\eta U + (1 - \eta)D], \quad f(\eta) := \eta[D - \eta L]^{-1}b.$$

La condition nécessaire et suffisante de convergence est $\rho(R(\eta)) < 1$, où $\rho(X)$ est le rayon spectral de X , égal au maximum du module de ses valeurs propres (théorème 4.1).

- Par les propriétés classiques du déterminant, on a

$$\det(R(\eta)) = \frac{\det([\eta U + (1 - \eta)D])}{\det([D - \eta L])}$$

De plus, les matrices $[\eta U + (1 - \eta)D]$ et $[D - \eta L]$ étant toutes deux triangulaires, leur déterminant est égal au produit de leurs termes diagonaux, et ainsi

$$\det(R(\eta)) = \frac{(1 - \eta)^n \det(D)}{\det(D)} = (1 - \eta)^n.$$

La condition de convergence $\rho(R(\eta)) < 1$ ne peut pas être vérifiée si $|\det(R(\eta))| \geq 1$, puisque cette dernière inégalité implique qu'une au moins des valeurs propres de $R(\eta)$ est de module égal ou supérieur à 1. La convergence des itérations SOR est donc impossible pour $\eta \notin]0, 2[$. Pour $0 < \eta < 2$, la convergence des itérations SOR est possible mais nullement garantie *a priori*. Cette convergence est établie pour certaines familles de matrices A , par exemple les matrices A SPD.

Exercice E22-4 (modification de rang 1 de systèmes linéaires) Soit $M \in \mathbb{R}^{n \times n}$ une matrice dense réelle inversible, et soit $f \in \mathbb{R}^n$. Le système $Mx = f$ a une solution unique $x \in \mathbb{R}^n$, que l'on suppose déjà calculée (par exemple par une méthode directe adaptée aux caractéristiques de M).

Etant donnés deux vecteurs $u, v \in \mathbb{R}^n$ (considérés comme vecteurs colonne, selon l'usage habituel), on considère la matrice modifiée $M_1 := M + uv^T$ et le système modifié $M_1 x_1 = f$. La modification $M_1 - M = uv^T$ de M est de rang 1 (expliquer pourquoi), et est donc "petite" en ce sens. On cherche à calculer la solution modifiée x_1 économiquement, connaissant M, x et u, v .

- (a) Montrer que $1 + v^T M^{-1} u \neq 0$ est une condition nécessaire et suffisante d'inversibilité de M_1 , et que l'inverse de M_1 est alors donné par la formule de Sherman-Morrison

$$M_1^{-1} = (M + uv^T)^{-1} = M^{-1} - \frac{1}{1 + v^T M^{-1} u} M^{-1} uv^T M^{-1}$$

(on pourra utiliser la propriété suivante: $\det(I + pq^T) = 1 + p^T q$ pour tous vecteurs $p, q \in \mathbb{R}^n$). Que peut-on dire de la modification $M_1^{-1} - M^{-1}$ de l'inverse de M ?

- (b) A l'aide de la formule ci-dessus, définir une méthode aussi économique que possible en opérations arithmétiques pour calculer la solution x_1 de $M_1 x_1 = f$ connaissant M, x et u, v . Estimer (à l'ordre principal en n) le surcoût d'opérations arithmétiques requis pour obtenir x_1 , et comparer au coût de la résolution préalable du système initial $Mx = f$ par une méthode directe.
- (c) On considère maintenant le problème de moindres carrés

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2, \quad A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m \quad (m \geq n),$$

pour lequel on suppose l'unicité de la solution x . Quelle condition doit vérifier A pour que cela soit le cas? Montrer que la solution x du problème ci-dessus vérifie le système linéaire

$$A^T A x = A^T b$$

("équations normales" du problème aux moindres carrés). Quelles sont les principales propriétés de ce système?

- (d) On envisage alors une version modifiée du problème aux moindres carrés ci-dessus, par ajout d'une ligne à A et b :

$$\min_{x_1 \in \mathbb{R}^n} \|A_1 x_1 - b_1\|_2^2, \quad A_1 = \begin{bmatrix} A \\ \alpha^T \end{bmatrix} \in \mathbb{R}^{(m+1) \times n}, \quad b_1 = \begin{Bmatrix} b \\ \beta \end{Bmatrix} \in \mathbb{R}^{m+1}$$

avec $\alpha \in \mathbb{R}^n, \beta \in \mathbb{R} \quad (m \geq n)$,

Former les équations normales de ce problème aux moindres carrés modifié. Montrer que la matrice de ce nouveau système est une perturbation de rang 1 de la matrice initiale $M = A^T A$, et qu'elle est inversible. En déduire une méthode économique de calcul de x_1 par actualisation de x , que l'on détaillera.

Eléments de solution:

- (a) On a $M_1 := M(I + M^{-1} uv^T)$, ce qui entraîne

$$\det(M_1) = \det(M) \det(I + M^{-1} uv^T) = (1 + v^T M^{-1} u) \det(M).$$

La condition nécessaire et suffisante d'inversibilité de M_1 est $\det(M_1) \neq 0$, et donc $1 + v^T M^{-1} u \neq 0$ ($\det(M) \neq 0$ puisque M est par hypothèse inversible). Pour montrer que M_1^{-1} est donné par la formule de Sherman-Morrison, il suffit de vérifier que la multiplication (à gauche ou à droite) par M_1 donne l'identité. Par exemple:

$$\begin{aligned} (M + uv^T) \left(M^{-1} - \frac{1}{1 + v^T M^{-1} u} M^{-1} uv^T M^{-1} \right) \\ = I + \left(1 - \frac{1}{1 + v^T M^{-1} u} - \frac{v^T M^{-1} u}{1 + v^T M^{-1} u} \right) uv^T M^{-1} = I \end{aligned}$$

Dans la formule de Sherman-Morrison, on a $M^{-1}uv^T M^{-1} = (M^{-1}u)(M^{-T}v)^T$, qui est une matrice de rang 1 (produit d'un vecteur colonne par un vecteur ligne, interprétable comme le produit tensoriel de deux vecteurs).

- (b) La solution de $M_1 x_1 = f$ est algébriquement donnée (en utilisant la formule de Sherman-Morrison) par

$$x_1 = M^{-1}f = M^{-1}f - \frac{1}{1 + v^T M^{-1}u} M^{-1}u v^T M^{-1}f$$

Cette expression est généralement inefficace en pratique car elle nécessiterait le calcul (de coût excessif, et parfois risqué sur le plan de la stabilité ou de la précision) de M^{-1} . Cet obstacle est contourné en remarquant que (i) $M^{-1}f = x$ et (ii) $z := M^{-1}u$ s'obtient en résolvant le système $Mz = u$ (une tâche raisonnable en réutilisant la factorisation LU, Cholesky... de M déjà employée pour obtenir x). On obtient ainsi

$$x_1 = x - \frac{1}{1 + v^T z} (v^T x) z$$

Supposons pour fixer les idées M inversible mais non symétrique. Le coût initial principal de la résolution de $Mx = f$ réside dans la factorisation LU de M ($\sim 2n^3/3$ opérations), les deux systèmes triangulaires demandant ensuite $\sim n^2$ opérations chacun. Le surcoût lié à l'évaluation subséquente de x_1 par la formule ci-dessus consiste en (i) $\sim 2n^2$ opérations pour les deux systèmes triangulaires fournissant la solution z de $Mz = u$ (réutilisation de la factorisation LU), (ii) calcul des produits scalaires $v^T x$ et $v^T z$ ($\sim 2n$ opérations chacun), (iii) calcul final de x_1 ($\sim 2n$ opérations). Les coûts de calcul de x puis x_1 sont donc, à l'ordre principal, $\sim 2n^3/3$ et $\sim 2n^2$ respectivement, et le surcoût du calcul de x_1 est donc faible relativement au coût du calcul de x . Un argument similaire s'applique au cas où M est SPD, avec utilisation d'une factorisation initiale de Cholesky.

- (c) La solution du problème de moindres carrés (qui existe toujours) est unique si A est de rang n (rang colonne maximal). Cette solution x doit par exemple vérifier la condition nécessaire d'optimalité $\nabla_x (\|Ax - b\|_2^2) = 0$. On a :

$$\|Ax - b\|_2^2 = x^T (A^T A) x - 2x^T A^T b + b^T b \implies \nabla_x (\|Ax - b\|_2^2) = 2(A^T A)x - 2A^T b.$$

La condition nécessaire d'optimalité est donc le système des équations normales proposé.

La matrice $A^T A$ est clairement carrée, symétrique et positive ($x^T (A^T A)x = \|Ax\|_2^2$ pour tout $x \in \mathbb{R}^n$). La matrice A étant de rang colonne maximal, son noyau est trivial ($\mathcal{N}(A^T A) = \{0\}$), ce qui implique $x^T (A^T A)x > 0$ pour tout $x \neq 0$. La matrice $A^T A$ est donc SPD. La résolution du problème initial de moindres carrés se ramène à celle du système (carré SPD) des équations normales.

- (d) Le problème de moindres carrés modifié conduit aux équations normales $A_1^T A_1 x_1 = A_1^T b_1$, pour lesquelles on trouve

$$A_1^T A_1 = A^T A + \alpha \alpha^T, \quad A_1^T b_1 = A^T b + \beta \alpha$$

La matrice $A_1^T A_1$ est donc une modification de rang 1 de la matrice $A^T A$, à laquelle on peut appliquer la formule de Sherman-Morrison ($M = A^T A$, $u = v = \alpha$), et on obtient

$$(A_1^T A_1)^{-1} = (A^T A)^{-1} - \frac{1}{1 + \alpha^T (A^T A)^{-1} \alpha} (A^T A)^{-1} \alpha \alpha^T (A^T A)^{-1}.$$

Procédant ensuite comme en (b), la solution du problème de moindres carrés modifié est *a priori* donnée par

$$x_1 = (A^T A)^{-1} (A^T b) - \frac{1}{1 + \alpha^T (A^T A)^{-1} \alpha} (A^T A)^{-1} \alpha \alpha^T (A^T A)^{-1} A^T b + \beta \left[(A^T A)^{-1} \alpha - \frac{1}{1 + \alpha^T (A^T A)^{-1} \alpha} (A^T A)^{-1} \alpha \alpha^T (A^T A)^{-1} \alpha \right].$$

En utilisant $(A^T A)^{-1}(A^T b) = x$ et définissant $z := (A^T A)^{-1}\alpha$ par résolution du système $(A^T A)z = \alpha$, on obtient alors

$$x_1 = x - \frac{1}{1 + \alpha^T x}(\alpha^T x)z + \beta \left(z - \frac{\alpha^T z}{1 + \alpha^T z} z \right) = x + \frac{\beta - \alpha^T x}{1 + \alpha^T z} z.$$

L'évaluation de la formule ci-dessus, une fois x connu, demande donc (i) le calcul de z par résolution de $(A^T A)z = \alpha$ (on réutilise pour cela la factorisation LDL^T ou de Cholesky de $A^T A$), (ii) l'évaluation des produits scalaires $\alpha^T x$ et $\alpha^T z$, (iii) l'évaluation finale de x_1 à l'aide de ces éléments. Des remarques similaires à celles du (c) s'appliquent quant à l'évaluation du surcoût (faible en termes relatifs) entraîné par l'évaluation de x_1 .