

§ 1.5. Approximation polynomiale/rationnelle

En vue du § 1.5., avec Ω K -spectrale,
 on a envie d'approcher $f(A)$ par $p(A)$
 avec $\|f(A) - p(A)\| \leq K \cdot \|f - p\|_e$

Q1 comment évaluer $p(A)$?

Q2 comment choisir p ?

La réponse de Q1 dépend de la base utilisée
 par exprimer p . Base canonique

$$p(A) = \sum_{j=0}^m p_j A^j$$

travail principal : Produit de deux matrices $\in \mathbb{C}^{n \times n}$,
 disons, de complexité $M(n)$.

Évaluation successive $A^{j+1} = A^j \cdot A \Rightarrow$ complexité $n \cdot M(n)$
 même chose par Horner

1.5.1. Algo de Paterson - Stockmeyer

On écrit $m = r \cdot s$ ($m \in [r \cdot s, (r+1) \cdot s - 1]$), $p_i = 0 \forall i > m$

et $p(A) = \sum_{j=0}^r \mathcal{B}_j \cdot (A^s)^j$, $\mathcal{B}_j = \sum_{k=0}^{s-1} p_{sj+k} \cdot A^k$
 évaluation Horner $(s-1)M(n)$

total complexité $(r+2) \cdot M(n)$, minimum pas $r = s = \sqrt{n}$

Nécessite stockage de s^2 matrices (au lieu d'une).
reste attractive par PCA. b!

Not. On peut montrer qu'il existe une constante c
 de sorte que, pour la partie S de 1.5.1. en virgule
 flottante: $\|p(A) - S\| \leq \frac{c \varepsilon \cdot n^2}{1 - c \varepsilon \cdot n^2} \cdot \left\| \sum |p_j| \cdot |A|^j \right\|$

le second membre pourrait être bien plus grand que
 $\|p(A)\|$! Cancellation, avec évaluation des séries partielles
 de $\exp(z)$ en $z = -\delta$! $\rightarrow TP$

L'exo suivant nous donne une idée de l'erreur théorique si on utilise des cones partielles d'un développement de Taylor.

Exo)
 1.5.2 Soit f analytique dans $\{ |z| \leq \rho(A) \}$, alors

$$E_n := f(A) - \sum_{j=0}^n \frac{f^{(j)}(0)}{j!} A^j = \int_0^1 f^{(n+1)}(t \cdot A) \cdot \frac{(1-t)^n}{n!} A^{n+1} dt$$

$$\text{En déduire que } \|E_n\| \leq \frac{\|A\|^{n+1}}{(n+1)!} \cdot \max_{t \in [0,1]} \|f^{(n+1)}(tA)\|.$$

Dans l'exo 1.5.2., avec $f(z) = \cos(z)$ et n impair, il reste à estimer $\max_{t \in [0,1]} \|\cos(tA)\| \leq \cosh(\|A\|)$

de même en vue de 1.5.1 : $\left\| \sum_{j=0}^n \frac{f^{(j)}(0)}{j!} A^j \right\| \leq \cosh(\|A\|)$.

Pour obtenir alors une erreur relative petite, on souhaite alors que $\cosh(\|A\|) / \|\cos(A)\|$ ne soit pas trop grand, ce qui est vrai si $\|A\| \leq 1$ mais faux si $\|A\|$ grand. Par conséquent, il vaut mieux appliquer plusieurs fois des formules de réduction d'argument $\cos(A) = 2 \cos^2\left(\frac{A}{2}\right) - 1$ avant d'utiliser Taylor (même chose valable par l'exponentialité).

Au lieu de travailler avec les séries de Taylor, on peut également calculer des interpolants dans une base de Newton (\rightarrow choisir l'ordre de Leja) et l'évaluer avec Horner (ou 1.5.1. si répétition cyclique des points d'interpolation).

Un choix approprié des points d'interpolation dans \mathbb{R} permet de répondre à Q2, voir plus tard.

Pour évaluer $R(A)$ par une fraction rationnelle $R = P/Q$, ^{deg P, Q ≤ m} on dispose de au moins 3 possibilités

(a) évaluation $P(A), Q(A) \rightarrow$ une inversion de complexité $\mathcal{O}(m)$

(b) décomposition en termes simples

$$R(A) = \sum a_j (A - z_j I)^{-1} \text{ par exemple } \Rightarrow m \cdot \mathcal{O}(m)$$

(c) Manipulation d'une FC

attention, les résidus a_j devraient être de taille modérée

1.5.3. Fractions continues (FC):

avec α_n, β_n polynômes

L'écriture $C(z) = \beta_0(z) + \frac{\alpha_1(z)}{\beta_1(z)} + \dots$ formelle désigne une suite de convergents

$$C_m(z) = \beta_0 + \frac{\alpha_1}{\beta_1} + \dots + \frac{\alpha_m}{\beta_m} = \frac{P_m}{Q_m} \text{ calculé par}$$

$$P_{-1} = 1, P_0 = \beta_0, Q_{-1} = 0, Q_0 = 1 \text{ et pour } n \geq 0$$

$$\begin{Bmatrix} P_{m+1} \\ Q_{m+1} \end{Bmatrix} = \beta_{m+1} \begin{Bmatrix} P_m \\ Q_m \end{Bmatrix} + \alpha_{m+1} \begin{Bmatrix} P_{m-1} \\ Q_{m-1} \end{Bmatrix}.$$

Avec les choix

$$(i) \alpha_{m+1} = \frac{C_{m+1} - C_m}{C_m - C_{m-1}}, \beta_{m+1} = 1 - \alpha_{m+1} \quad (C_n = \infty)$$

$$(ii) \alpha_{m+1} = \frac{P_m Q_{m+1} - Q_m P_{m+1}}{P_m Q_{m-1} - Q_m P_{m-1}}, \beta_{m+1} = \frac{P_{m+1} Q_{m-1} - P_{m-1} P_{m+1}}{P_m Q_{m-1} - Q_m P_{m-1}}$$

ou construit une FC avec convergents $C_n = P_n/Q_n$.

Une évaluation numérique stable se fait à l'aide d'une remontée:

$$C_m^{(n)} = \beta_m, \quad k < m \quad C_k^{(n)} = \beta_{k+1} + \frac{\alpha_{k+1}}{C_{k+1}^{(n)}}, \quad C_0^{(n)} = C_m^{(n)}$$

\rightarrow complexité revient à $m \cdot \mathcal{O}(m) + \mathcal{O}(m) \cdot \mathcal{O}(m)$?

1.5.4. Interpolants à poles prescrits

Pour Q polynôme et $z_1, \dots, z_m \in \mathbb{C}, Q(z_j) \neq 0$, on cherche P polynôme de degré $< m$ de sorte que

$R_{m,Q} = P/Q$ interpole f aux points z_1, \dots, z_m .

Formule de Lagrange: si z_1, \dots, z_m distincts $\frac{m}{\prod_{k=1}^m (z-z_k)}$

$$R_{m,Q}(z) = B(z) \cdot \frac{\sum_{j=1}^m \frac{f(z_j)}{(z-z_j) B'(z_j)}}{Q(z)}, \quad B(z) = \prod_{k=1}^m (z-z_k)$$

\Rightarrow existence/unicité.

Lien avec interpolation polynomiale

$P = Q \circ R_{m,Q}$ polynôme d'interp. de f aux points z_1, \dots, z_m

donc $f(z) - R_{m,Q}(z) = B(z) \cdot [z_1, \dots, z_m, z](Q \circ f)$

Formule d'Hermite

si f analytique dans Ω et $z_1, \dots, z_m \in \text{Int}(\Omega)$

$$f(z) - R_{m,Q}(z) = \frac{B(z)}{2\pi i} \int_{\partial\Omega} \frac{f(\zeta)}{B(\zeta)} \cdot \frac{d\zeta}{\zeta-z}$$

Choix des points d'interpolation? Idées: Q simple les s.g. de f .

si on se fixe Q et on souhaite erreur petite sur compact $E \subset \Omega$, il faudrait que

$\|B\|_E = \|\frac{1}{B}\|_{\Omega}$ soit le plus petit possible dit "problème de Zolotarev" (si on cherche z_j et Q).

1.5.5. Interpolants à poles libres/rationnelles
ici de type $[m-1/m]$

si on cherche

Jci pour $z_1, \dots, z_m \in \mathbb{C}$ on cherche à trouver $R_m = P_m/Q_m$ avec $P_m \in \mathbb{P}_{m-1}, Q_m \in \mathbb{P}_m$ de sorte $f Q_m - P_m$ s'annule aux z_1, \dots, z_m comptant multiplicités \Rightarrow si $Q_m(z_j) \neq 0$ (ce qui générique est vrai)

R_m interpole f aux z_1, \dots, z_m .

Idée: poles de R_m dépendent singulièrement de f .

Lemme (i) Un tel interpolant existe.

(ii) La fraction P_m/Q_m est unique

(iii) $\forall Q$ de degré $\leq m$: $R_m = R_{2m,Q_m} \cdot Q$ (avec cas générique) $Q \neq Q$

Preuve ad (i) système homogène à $2m$ équations et $2m+1$ inco.

ad (ii) Soient P_m/Q_m et \tilde{P}_m/\tilde{Q}_m deux solutions $\omega = P_m \tilde{Q}_m - \tilde{P}_m Q_m = (P_m - f Q_m) \tilde{Q}_m - (\tilde{P}_m - f Q_m) Q_m$

Si $z_1 = z_2 = \dots = z_m = \infty$ on parle d'un approximant de Padé en z_0 : $f(z) Q(z) - P(z) = O((z-z_0)^{2m})$ de type $[m-1/m]$

Suivant la preuve du lemme (ii), on arrive à montrer dans le cas "générique" (ou on a égalité ^{part} des degrés) que

$$P_{n+j} Q_n - P_n Q_{n+j+1} = \gamma_{n,j} \cdot (z-z_1) \cdots (z-z_{2n})$$

avec $\gamma_{n,j}$ des polynômes de degré $\leq j$, et donc d'après 1.5.3 (ii) $\alpha_{n+1} = - \underbrace{\gamma_{n,0}}_{\in \mathbb{R}} \cdot (z-z_{2n+1}) \cdots (z-z_{2n})$

et $\beta_{n+1} \in \mathbb{P}_1$ dans la CF ayant les ^{à poles les} convergents P_{2n+1} .

L'idée des interpolants rationnelles est que les poles trouvent eux-même la singularité de la fonction, par exemple $P_n \rightarrow f(z) = \log(z+1)$ pour les approximations de Padé en 0 uniformément dans tout compact de $\mathbb{C} \setminus [-\infty, -1]$, l'intervalle $(-\infty, -1]$ comportant tous les pôles. Ceci est une convergence bien plus intéressante que celle des séries partielles de Taylor.

Notons que l'on connaît explicitement des FC pour les approximations de Padé en \mathbb{Q} pour une grande quantité de fonctions [Baker, Gravesland]

Pour une fonction de Markov on peut être plus précis.

1.5.6. Lemme : ^{pas une fonction rationnelle}

Soient $f(z) = \int_a^b \frac{d\mu(x)}{z-x}$ (et $z_1, \dots, z_{2m} \in \mathbb{C} \setminus [a, b]$ de sorte que $w_m(z) = \prod_{j=1}^{2m} (z-z_j)$ soit réel sur \mathbb{R}). Alors le dénominateur Q_m de l'interpolant rationnel $R_m = P_m / Q_m$ aux points z_1, \dots, z_{2m} vérifie la relation d'orthogonalité

$$(x) \quad \int_a^b \frac{Q_m(x)}{w_m(x)} x^j d\mu(x) = 0, \quad j=0, \dots, m-1$$

en particulier, tous les racines r_1, \dots, r_m de Q_m sont simples et appartiennent à $[a, b]$, et $\dots, a_m > 0$

$$P_m(z) / Q_m(z) = \sum_{j=1}^m \frac{a_j}{z-r_j} \quad \text{avec } a_j \in \mathbb{R} \text{ et } a_j > 0.$$

Exemple : $f(z) = \frac{1}{z} \log(1+z) = \int_{-\infty}^{-1} \frac{(-1/x)}{z-x} dx \quad z_1 = \dots = z_{2m} = 0$

(*) devient $\int_{-\infty}^{-1} Q_m(x) \cdot x^{-2m-1+j} dx = 0 \quad j=0, \dots, m-1$

ou encore $\int_0^1 Q_m(-\frac{1}{y}) \cdot y^m \cdot y^k dy = 0 \quad k=0, 1, \dots, m-1$
 $= \text{const} \cdot L_m(-1+2y)$, L_m polynôme de Legendre
 \perp sur $[-1, 1]$

Preuve : $\forall j=0, \dots, m-1 : [z_1, \dots, z_{2m}]_z z^j [Q_m(z)f(z) - P_m(z)] = 0$
 $= [z_1, \dots, z_{2m}]_z \left(\int \frac{x^j Q_m(x) dx}{z-x} - \int \frac{x^j Q_m(x) - z^j Q_m(z)}{z-x} dx + z^j P_m(z) \right)$
 polynôme en z de degré $\leq m+j \leq 2m-1$
 $= [z_1, \dots, z_{2m}]_z \int \frac{x^j Q_m(x) dx}{z-x} \stackrel{\text{res}}{=} - \int \frac{x^j Q_m(x)}{w_m(x)} dx \Rightarrow (*)$

Pour conclure la deuxième partie observons que w_m est réel et ne change pas de signe sur $[a, b]$ par hypothèse.

(*) donne un système homogène à m inconnues et $m-1$ équations à coefficients réels, avec

$$\det \left(\int \frac{x^{j+k}}{w_m(x)} dx \right)_{j,k=0, \dots, m-1} = \int \frac{dx(x_1)}{w_m(x_1)} \dots \int \frac{dx(x_m)}{w_m(x_m)} \cdot x_1^0 x_2^1 \dots x_m^{m-1} \cdot \det \begin{bmatrix} 1 & \dots & 1 \\ x_1 & & x_m \\ \vdots & & \vdots \\ x_1^{m-1} & & x_m^{m-1} \end{bmatrix}$$

$$= \frac{1}{m!} \det \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_m \\ \vdots & & \vdots \\ x_1^{m-1} & \dots & x_m^{m-1} \end{bmatrix}^2 > 0$$

$\Rightarrow Q_m$ de degré $= m$, et Q_m à coeff. réels. Si la propriété sur les racines serait fautive alors il existerait P de degré $< m$ de sorte que $Q_m \cdot P$ est réel et ne change pas de signe sur $[a, b] \Rightarrow \int \frac{Q_m(x)P(x)}{w_m(x)} dx > 0$ en contradiction avec (*). Pour trouver le signe des a_j dans la décomposition en fractions simples, on remarque que, d'après 1.5.4. et le lemme (iii),

(**) $f(z) - \frac{P_m(z)}{Q_m(z)} = \frac{w_m(z)}{Q_m(z)^2} [z_1, \dots, z_{2m}, z] \cdot (Q_m \cdot (Q_m f - P_m))$
 c'est-à-dire $\frac{w_m(z)}{Q_m(z)^2} \int \frac{Q_m(x)^2}{w_m(x)} \cdot \frac{dx(x)}{z-x}$

En écrivant $\tilde{Q}_m(z) = Q_m(z)/(z-z_j) :$

$$\frac{P_m(z)}{Q_m(z)} \stackrel{(**)}{=} \int \left[1 - \frac{(x - \beta_j)^2}{(z - \beta_j)^2} \right] \frac{W_m(z)}{\tilde{Q}_m(z)^2} \cdot \frac{\tilde{Q}_m(x)^2}{W_m(x)} \frac{d\mu(x)}{z-x}$$

$$+ \int \left[1 - \frac{W_m(z)}{\tilde{Q}_m(z)^2} \cdot \frac{\tilde{Q}_m(x)^2}{W_m(x)} \right] \frac{d\mu(x)}{z-x}$$

fonction rationnelle en z n'admettant pas de pôle en z_j

avec $\frac{1}{z-x} \left[1 - \frac{(x - \beta_j)^2}{(z - \beta_j)^2} \right] = \frac{z+x - 2\beta_j}{(z - \beta_j)^2} = \frac{1}{z - \beta_j} + \frac{x - \beta_j}{(z - \beta_j)^2}$

donc $\frac{P_m(z)}{Q_m(z)} = \frac{W_m(z)}{\tilde{Q}_m(z)^2} \left[\frac{1}{z - \beta_j} \int \frac{\tilde{Q}_m(x)^2}{W_m(x)} d\mu(x) + \frac{1}{(z - \beta_j)^2} \int \frac{\tilde{Q}_m(x) \tilde{Q}_m(x)}{W_m(x)} d\mu(x) \right]$

+ analytique en z_j

= $\frac{W_m(\beta_j)}{\tilde{Q}_m(\beta_j)} \cdot \int \frac{\tilde{Q}_m(x)^2}{W_m(x)} \cdot \frac{1}{z - \beta_j} + \text{analytique en } z_j$

= $\alpha_j \geq 0$.

1.5.7. Théorème de M. Riesz :

Sous les hypothèses 1.5.6., si $\lim_{j \rightarrow \infty} \text{dist}(z_j, [a, b]) > 0$
 alors $R_m \rightarrow f$ uniformement sur tout compact $C \subset]a, b[$.

Preuve : M. R.

(***) $\exists C > 0 : \forall m \geq 1, \forall z \in C \subset]a, b[: |R_m(z)| \leq \frac{C}{\text{dist}(z, [a, b])}$
 (livre de Schiffer)
 $\Rightarrow (R_m)$ est normale sur $C \subset]a, b[$

Preuve de (***) :

$$|f(z_n)| = \left| \frac{P_m(z_n)}{Q_m(z_n)} \right| = \left| \sum_{j=1}^m \frac{a_j}{z_n - \beta_j} \right|$$

$$\geq \left| \sum_{j=1}^m \frac{a_j}{z_n - \beta_j} \right| = \sum_{j=1}^m \frac{|a_j| |z_n - \beta_j|}{|z_n - \beta_j|^2}$$

\Rightarrow si $z_n \notin \mathbb{R} : \sum_{j=1}^m |a_j| \leq \frac{|f(z_n)|}{\text{dist}(z_n, [a, b])} =: C$

sinon, par exemple $z_n > b : |f(z_n)| \geq \sum_{j=1}^m \frac{|a_j|}{|z_n - \beta_j|} \geq \sum_{j=1}^m \frac{|a_j|}{|z_n - b|}$

et donc $\sum |a_j| \leq \frac{|f(z_n)|}{\text{dist}(z_n, [a, b])} =: C$ (car si $z_n \rightarrow a$) (28)

et alors $|R_m(z)| \leq \sum \frac{|a_j|}{|z - r_j|} \leq \frac{\sum_{j=1}^m |a_j|}{\text{dist}(z, [a, b])} \leq \frac{c}{\text{dist}(z, [a, b])}$

si Thm. est faux alors $\exists \Lambda$ infini $\exists \hat{z}_0 \in \mathbb{C} \setminus [a, b]$
avec $R_m(\hat{z}_0) \xrightarrow{\Lambda} c \neq f(\hat{z}_0)$. D'après Moulé, $\exists \hat{\Lambda} \subset \Lambda$ infini $\exists g$ analytique dans $\mathbb{C} \setminus [a, b]$:

$R_m \xrightarrow{\hat{\Lambda}} g$ ^{unif.} sur tout sous-ensemble fermé de $\mathbb{C} \setminus [a, b]$.

Notons que $g(z_j) = f(z_j)$ pour $j=1, 2, 3, \dots$ (comptant multiplicité), et les z_j accumulent sur $\mathbb{C} \setminus [a, b]$ par hypothèse $\Rightarrow g = f \Rightarrow g(\hat{z}_0) = f(\hat{z}_0)$ contradictoire. \square

1.5.8. Corollaire :

Sous les hypothèses du 1.5.6.,
si $z_1 = z_2 = \dots \in \mathbb{R}$ (par exemple Padé)
et $r > 0$ tel q. $z_1 - r > b$ alors

$$\|f - R_m\|_{\{|z - z_1| \leq r\}} = |(f - R_m)(z_1 - r)|$$

en particulier, si $\|A - z_1 I\| < z_1 - b$

$$\|f(A) - R_m(A)\| \leq |(f - R_m)(z_1 - \|A - z_1 I\|)|.$$

Preuve $(f - R_m)(z) = \frac{(z - z_1)^{2m}}{Q_m(z)^2} \int \frac{Q_m(x)}{(x - z_1)^{2m}} \cdot \frac{d\mu(x)}{z - x}$

chaque terme atteint son max dans $\{|z - z_1| \leq r\}$
en $z = z_1 - r$. Pour la deuxième partie,
voir théorie 1.4.3. de Neuman \square