

— Sparse Rational Approximation —
*Pole filtering by regularization & feature selection
for complex matrix rational approximation*

Pablo Ducru^a

^a*Nuclear Science & Engineering department, Massachusetts Institute of Technology,
77 Massachusetts Ave., Cambridge, MA 02139, USA.*

Abstract

MISSION:

- Meromorphic operators are the fundamental representation of resonance phenomena in physics, and Runge's theorem establishes rational fractions are universal approximators of meromorphic functions on any given region of the complex plane.
- Thus, myriad operators in physics or electrical engineering are best characterized by a complex rational fraction expansion.
- To this day, most methods of rational fraction approximation require to specify the order of the approximant: i.e. the number of roots and poles.
- When methods do seek to systematically find the number of poles, they presently take one of the following approaches:
 - solving with various number of poles and select best performance,
 - *ad hoc* problem specific poles filtering, such as low-band or high-frequency filters.
 - pole clustering and replacement with the number of clusters,
 - pole trimming, often with a hard-threshold manual way.
- In this research project, we propose and investigate a group LASSO regularization method for residues sparsification as a means to eliminate spurious poles.

WHAT HAS BEEN ACHIEVED:

- Kernel exploration: we searched for a RKHS structure to embed our problem in. This quest was inconclusive. The method is thus parametric.

- The major result of this research project is theorem 2, which generalizes the proximal forward backward splitting PFBS algorithm to the case of complex matrices. We demonstrate that in a given gradient sense (c.f. “cartesian differentiation” introduced in definition 1), the real-case PFBS algorithm can be formally extended to encompass complex matrices.
- The derived algorithm has implemented in the Julia Language and in a toy problem it was possible to eliminate a pole when LS was unable to.

SETBACK:

- Cauchy matrices condition numbers render any practical application out-of-reach.

WAY FORWARD:

- Write-up the Kernel exploration inconclusive results, to explain why this problem has no apparent RKHS structure (the “pole-residue” atoms are difficult to provide with a Hilbert structure).
- Implementation and performance of the complex Frobenius PFBS algorithm 1, to show the convergence, and verify theorem 2, should now take care of the ill-conditioned problem.
- Elastic-Net regularization seems the way to go, though QR factorization may also be of great use.

Keywords: Rational approximation, feature selection, Group LASSO, complex variable optimization.

1. Residues sparsification problem

Provided observed complex data points $(z_k)_{k \in \llbracket 1, N_k \rrbracket} \in \mathbb{C}$ with weights $(\rho_k)_{k \in \llbracket 1, N_k \rrbracket} \in \mathbb{R}_+$ (for homoscedastic case $\rho_k = \frac{1}{N_k}$) and their corresponding matrices $\mathbf{Y}_k \in \mathbb{C}^{q \times s}$, and given N_p fixed poles $(p_n)_{n \in \llbracket 1, N_p \rrbracket} \in \mathbb{C}$, we search for the residues and affine coefficients matrices $(\mathbf{R}_n)_{n \in \llbracket 1, N_p \rrbracket}, \mathbf{C}, \mathbf{A} \in \mathbb{C}^{q \times s}$ that minimize the rational approximation regularized weighted least square (LS) problem:

$$\min_{\mathbf{R}_n, \mathbf{C}, \mathbf{A}} \left\{ \sum_{k=1}^{N_k} \rho_k \left\| \sum_{n=1}^{N_p} \frac{\mathbf{R}_n}{z_k - p_n} + \mathbf{C} + z_k \mathbf{A} - \mathbf{Y}_k \right\|_F^2 + \lambda \|\mathbf{R}\|_{L_1 - F} \right\} \quad (1)$$

where $\mathbf{R} := [\mathbf{R}_1, \dots, \mathbf{R}_n, \dots, \mathbf{R}_{N_p}]$ and $\|\mathbf{R}\|_{L_1 - F} := \sum_{n=1}^{N_p} \|\mathbf{R}_n\|_F$ is a sparsity-based group-LASSO penalization, with $\|\cdot\|_F$ being the Frobenius norm associated to the Hilbert-Schmidt hermitian product. This $L_1 - F$ regularization will drive low-incidence residue matrices to zero when solving problem (1). This

Email address: p_ducru@mit.edu (Pablo Ducru)

residues sparsification can then be used to discard superfluous poles by identifying the poles associated to zero-matrix residues.

The sparsity-inducing regularization in problem (1) is in essence a feature extraction by ‘kitchen-sink approach’: we prescribe a number N_p of poles deemed to be higher than the optimal number $N_p^{\text{opt}} \ll N_p$, and filter-out the poles associated to zero residues $\mathbf{R}_n = \mathbf{0}$, as prescribed in the final algorithm 2.

Defining the coefficients matrix \mathbf{K} we seek to find

$$\mathbf{K} := \begin{bmatrix} \mathbf{R} \\ \mathbf{C} \\ \mathbf{A} \end{bmatrix} \quad (2)$$

and the weights matrices \mathbf{W}_k composed of the Cauchy matrices at the poles (p_n) and data points (z_k):

$$\mathbf{W}_k := \left[\frac{\mathbb{I}}{z_k - p_1}, \dots, \frac{\mathbb{I}}{z_k - p_{N_p}}, \mathbb{I}, z_k \mathbb{I} \right] \quad (3)$$

where \mathbb{I} designates the $\mathbb{C}^{q \times q}$ identity matrix, the residues sparsification problem (1) appears as a special case of the more general problem:

$$\min_{\mathbf{K}} \left\{ \sum_{k=1}^{N_k} \rho_k \|\mathbf{W}_k \mathbf{K} - \mathbf{Y}_k\|_F^2 + \lambda \|\mathbf{R}\|_{L_1-F} \right\} \quad (4)$$

where the specific \mathbf{W}_k weights matrices are tailored to the simple poles and residues feature space.

Defining the empirical error

$$\widehat{\mathcal{E}}(\mathbf{R}, \mathbf{C}, \mathbf{A}) := \sum_{k=1}^{N_k} \rho_k \left\| \sum_{n=1}^{N_p} \frac{\mathbf{R}_n}{z_k - p_n} + \mathbf{C} + z_k \mathbf{A} - \mathbf{Y}_k \right\|_F^2 \quad (5)$$

or equivalently

$$\widehat{\mathcal{E}}(\mathbf{K}) := \sum_{k=1}^{N_k} \rho_k \|\mathbf{W}_k \mathbf{K} - \mathbf{Y}_k\|_F^2 \quad (6)$$

we are thus solving the regularization problem:

$$\min_{\mathbf{R}, \mathbf{C}, \mathbf{A}} \left\{ \widehat{\mathcal{E}}(\mathbf{R}, \mathbf{C}, \mathbf{A}) + \lambda \|\mathbf{R}\|_{L_1-F} \right\} \quad (7)$$

2. In quest of a reproducing kernel Hilbert space

We could however seek to perform this search of optimal complexity in an unsupervised way, by invoking a representer theorem and calling the reproducing kernel on the given training points $(z_k)_{k \in \llbracket 1, N_k \rrbracket}$.

In this section, we present some research we undertook to exhibit a reproducing kernel Hilbert space structure to this problem. In effect, we failed to unveil a convincing RKHS structure.

The essential reason for this is that the “pole-residue” atoms are difficult to provide with a Hilbert space nature. The use of contour integrals and the residues theorem has enabled us to exhibit some sesquilinear hermitian operator with reproducing property, but the space is not stable through the kernel. We have not yet transcribed here these hand-written results.

3. Proximal Forward Backward Splitting Iteration for Complex Matrices

In this section we establish that the traditional Proximal Forward-Backward Splitting (PFBS) iteration can be generalized to deal with complex matrices in the case of problem (1). Our main result is theorem 2, where we entirely characterize a PFBS iteration for complex matrices, with a specific definition of the gradient — the cartesian differential — introduced in definition 1.

This is not an obvious result at all, because holomorphic differentiation with respect to the complex variable changes the nature and behavior of complex gradients. We however notice and demonstrate that the quadratic norm structure of problem (1) lends itself well to the use of the usually ill-suited definition of complex differential — here named Cartesian differential and specified in definition 1 — which then enables us to seamlessly extend the expression of a PFBS iteration to complex matrices.

Definition 1. CARTESIAN COMPLEX DIFFERENTIAL.

Let $z \in \mathbb{C}$, $z = x + iy$, $x, y \in \mathbb{R}$. The holomorphic differential is defined as

$$\partial_z := \frac{1}{2} (\partial_x - i\partial_y) \tag{8}$$

Note that $\partial_z z = 1$, and $\partial_z \bar{z} = 0$, where $\bar{z} := x - iy$ designates the complex conjugate.

We here define the Cartesian differential as:

$$\bar{\partial}_z := \partial_x + i\partial_y \tag{9}$$

Note that $\bar{\partial}_z z = 0$, and $\bar{\partial}_z \bar{z} = 2$.

The Cartesian differential is thus *a priori* a poor choice for complex differentials. We however here demonstrate that it is well suited to our specific quadratic problem (1).

Lemma 1. CARTESIAN DIFFERENTIAL OF FROBENIUS NORM.

Let $\mathbf{A} \in \mathbb{C}^{n \times m}$ be a complex matrix, and $\|\cdot\|_F^2$ be the Frobenius norm associated with the Hilbert-Schmidt hermitian product $\langle \mathbf{A}, \mathbf{B} \rangle := \text{Tr}(\mathbf{A}^* \mathbf{B})$, with $\mathbf{A}^* := \overline{\mathbf{A}}^\top$, then the following equalities hold:

$$\partial_{\mathbf{A}} \|\mathbf{A}\|_F^2 = \overline{\mathbf{A}} \quad \text{and} \quad \partial_{\mathbf{A}} \|\mathbf{A}\|_F = \frac{1}{2} \frac{\overline{\mathbf{A}}}{\|\mathbf{A}\|_F} \quad (10)$$

while Cartesian differentiation yields:

$$\bar{\partial}_{\mathbf{A}} \|\mathbf{A}\|_F^2 = 2\mathbf{A} \quad \text{and} \quad \bar{\partial}_{\mathbf{A}} \|\mathbf{A}\|_F = \frac{\mathbf{A}}{\|\mathbf{A}\|_F} \quad (11)$$

Proof. Applying the Cartesian differential definition 1 to each element $a_{ij} \in \mathbb{C}$ of \mathbf{A} in the expression of the norm $\|\mathbf{A}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^m |a_{ij}|^2 = \sum_{i=1}^n \sum_{j=1}^m \Re[a_{ij}]^2 + \Im[a_{ij}]^2$ readily yields the results. \square

Lemma 1 shows that the Cartesian differential symbolically behaves on the norm of complex matrices like a real differential would behave on the norm of real matrices, all the while it does not behave like the real differential on the vector themselves (recall $\bar{\partial}_{\mathbf{A}} \mathbf{A} = \mathbf{0}$). We can now use Lemma 1 to demonstrate the first main result: theorem 1.

Theorem 1. PROXIMAL OPERATOR OF FROBENIUS PENALIZATION FOR COMPLEX MATRICES.

Let $\mathbf{A} \in \mathbb{C}^{n \times m}$ be a complex matrix and $\lambda \in \mathbb{R}_+$, the proximal operator of a Frobenius norm penalization is defined as:

$$\text{prox}_{\lambda \|\cdot\|_F}(\mathbf{A}) := \arg \min_{\mathbf{M} \in \mathbb{C}^{n \times m}} \left\{ \frac{1}{2} \|\mathbf{M} - \mathbf{A}\|_F^2 + \lambda \|\mathbf{M}\|_F \right\} \quad (12)$$

In the Cartesian differential sense, the proximal operator has the closed explicit form:

$$\text{prox}_{\lambda \|\cdot\|_F}(\mathbf{A}) = \left(\mathbf{A} - \lambda \frac{\mathbf{A}}{\|\mathbf{A}\|_F} \right) \mathbf{1}_{\left\{ \|\mathbf{A}\|_F \geq \lambda \right\}} = \mathbf{A} \cdot \max \left\{ 0, 1 - \lambda \frac{1}{\|\mathbf{A}\|_F} \right\} \quad (13)$$

Which is the formally the same expression as for real matrices.

Proof. The optimality condition on the subgradient, valid for both real and imaginary partial differentials, can thus also be expressed with the Cartesian differential from definition 1, i.e. let \mathbf{M}_{opt} be the minimizer, then $\mathbf{0} \in \bar{\partial}_{\mathbf{M}} \left(\frac{1}{2} \|\mathbf{M}_{\text{opt}} - \mathbf{A}\|_F^2 + \lambda \|\mathbf{M}_{\text{opt}}\|_F \right)$. Applying lemma 1 to this Cartesian differential yields:

$\mathbf{A} \in \left(\mathbb{I} + \lambda \frac{\mathbb{I}}{\|\cdot\|_F}\right) \mathbf{M}_{\text{opt}}$, where $\mathbb{I} \in \mathbb{C}^{n \times m}$ indicates the identity matrix. Formally, we are thus seeking to invert: $\mathbf{M}_{\text{opt}} = \left(\mathbb{I} + \lambda \frac{\mathbb{I}}{\|\cdot\|_F}\right)^{-1} \mathbf{A}$. Let us now demonstrate that this inverse is given by: $\mathbf{M}_{\text{opt}} = \left(\mathbf{A} - \lambda \frac{\mathbf{A}}{\|\mathbf{A}\|_F}\right) \mathbb{1}_{\{\|\mathbf{A}\|_F \geq \lambda\}}$. First, let us notice that since $\lambda \geq 0$, the real and imaginary parts of each element of \mathbf{A} and of \mathbf{M}_{opt} must be of same sign, and taking the norm yields $\|\mathbf{A}\|_F \geq \lambda$, unless $\mathbf{M}_{\text{opt}} = \mathbf{0}$ in which case $\mathbf{A} = \mathbf{0}$. In the case $\|\mathbf{A}\|_F \geq \lambda$ we can show (13) holds by noticing that $\mathbf{A} = \left(\mathbb{I} + \lambda \frac{\mathbb{I}}{\|\cdot\|_F}\right) \mathbf{M}_{\text{opt}}$ is equivalent to $\mathbf{A} = \mathbf{A} - \lambda \frac{\mathbf{A}}{\|\mathbf{A}\|_F} + \lambda \frac{\mathbf{A} - \lambda \frac{\mathbf{A}}{\|\mathbf{A}\|_F}}{\left\|\mathbf{A} - \lambda \frac{\mathbf{A}}{\|\mathbf{A}\|_F}\right\|_F}$, i.e. $\frac{\mathbf{A}}{\|\mathbf{A}\|_F} = \frac{\mathbf{A} - \lambda \frac{\mathbf{A}}{\|\mathbf{A}\|_F}}{\left\|\mathbf{A} - \lambda \frac{\mathbf{A}}{\|\mathbf{A}\|_F}\right\|_F}$, where this colinearity result is true iff $\|\mathbf{A}\|_F \geq \lambda$. This terminates the proof of (13). \square

Having established the closed-form of the proximal operator in the Cartesian differential sense in theorem 1, we can now define a Cartesian differential proximal forward backward splitting iteration scheme:

Definition 2. CARTESIAN DIFFERENTIAL PROXIMAL FORWARD BACKWARD SPLITTING ITERATION.

For complex matrices, we define the Cartesian differential proximal forward backward splitting iteration for problem (7) with learning rate $\gamma \in \mathbb{R}_+$ as:

$$\begin{bmatrix} \mathbf{R}_n^{(t+1)} \\ \mathbf{C}^{(t+1)} \\ \mathbf{A}^{(t+1)} \end{bmatrix} := \begin{bmatrix} \text{prox}_{\lambda \|\cdot\|_F} \left(\mathbf{R}_n^{(t)} - \gamma \partial_{\mathbf{R}_n} \widehat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) \right) \\ \mathbf{C}^{(t)} - \gamma \partial_{\mathbf{C}} \widehat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) \\ \mathbf{A}^{(t)} - \gamma \partial_{\mathbf{A}} \widehat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) \end{bmatrix} \quad (14)$$

Where the block expression stems from the block-property of the proximal operator with respect to the group penalization: $\|\mathbf{R}\|_{L_1-F} := \sum_{n=1}^{N_p} \|\mathbf{R}_n\|_F$.

This definition is a direct translation to complex matrices of the real PFBS algorithm, where the real gradient $\partial(\cdot)$ has been replaced with the Cartesian derivative $\partial(\cdot)$.

We are now equipped to establish our main result, theorem 2:

Theorem 2. PROXIMAL FORWARD BACKWARD SPLITTING ALGORITHM FOR FROBENIUS PENALIZATION OF COMPLEX MATRICES.

The Cartesian differential proximal forward backward splitting iteration for problem (1) takes the form:

$$\begin{bmatrix} \mathbf{R}_n^{(t+1)} \\ \mathbf{C}^{(t+1)} \\ \mathbf{A}^{(t+1)} \end{bmatrix} = \begin{bmatrix} \left(\mathbf{R}_n^{(t)} - \gamma \partial_{\mathbf{R}_n} \widehat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) \right) \cdot \max \left\{ 0, 1 - \frac{\lambda}{\left\| \mathbf{R}_n^{(t)} - \gamma \partial_{\mathbf{R}_n} \widehat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) \right\|_F} \right\} \\ \mathbf{C}^{(t)} - \gamma \partial_{\mathbf{C}} \widehat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) \\ \mathbf{A}^{(t)} - \gamma \partial_{\mathbf{A}} \widehat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) \end{bmatrix} \quad (15)$$

where:

$$\bar{\partial}_{\mathbf{R}_n} \hat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) = 2 \sum_{k=1}^{N_k} \rho_k \frac{1}{(z_k - p_n)^*} \left[\sum_{n=1}^{N_p} \frac{\mathbf{R}_n^{(t)}}{z_k - p_n} + \mathbf{C}^{(t)} + z_k \mathbf{A}^{(t)} - \mathbf{Y}_k \right] \quad (16)$$

$$\bar{\partial}_{\mathbf{C}} \hat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) = 2 \sum_{k=1}^{N_k} \rho_k \left[\sum_{n=1}^{N_p} \frac{\mathbf{R}_n^{(t)}}{z_k - p_n} + \mathbf{C}^{(t)} + z_k \mathbf{A}^{(t)} - \mathbf{Y}_k^{(t)} \right] \quad (17)$$

$$\bar{\partial}_{\mathbf{A}} \hat{\mathcal{E}}(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)}) = 2 \sum_{k=1}^{N_k} \rho_k z_k^* \left[\sum_{n=1}^{N_p} \frac{\mathbf{R}_n^{(t)}}{z_k - p_n} + \mathbf{C}^{(t)} + z_k \mathbf{A}^{(t)} - \mathbf{Y}_k^{(t)} \right] \quad (18)$$

Moreover, the Cartesian differential PFBS algorithm, which repeats iteration (15), is guaranteed to converge if the learning rate $\gamma \in \mathbb{R}_+$ respects:

$$\gamma \leq \frac{1}{2 \max \left\{ \mathcal{S}_p \left(\sum_{k=1}^{N_k} \rho_k \mathbf{W}_k^* \mathbf{W}_k \right) \right\}} \quad (19)$$

Proof. Expression (15) is a direct application of theorem 1 to express the proximal operator for complex matrices in the Cartesian differential sense.

To establish expressions (16), (17), and (18), we are going to show that the Cartesian differential satisfies (20), which is equivalent to expressions (16), (17), and (18) combined.

$$\bar{\partial}_{\mathbf{K}} \hat{\mathcal{E}}(\mathbf{K}^{(t)}) = 2 \sum_{k=1}^{N_k} \rho_k \mathbf{W}_k^* [\mathbf{W}_k \mathbf{K}^{(t)} - \mathbf{Y}_k] \quad (20)$$

We drop the (t) iteration superscript for clarity, and decompose into the real and imaginary components: $\hat{\mathcal{E}}(\mathbf{K}) = \sum_{k=1}^{N_k} \rho_k (\Re[\mathbf{W}_k \mathbf{K} - \mathbf{Y}_k]^2 + \Im[\mathbf{W}_k \mathbf{K} - \mathbf{Y}_k]^2)$. Expressing the real and imaginary part $\mathbf{W}_k = \Re[\mathbf{W}_k] + i\Im[\mathbf{W}_k]$ and $\mathbf{K} = \Re[\mathbf{K}] + i\Im[\mathbf{K}]$ yields $\hat{\mathcal{E}}(\mathbf{K}) = \sum_{k=1}^{N_k} \rho_k \left\{ \left(\Re[\mathbf{W}_k] \Re[\mathbf{K}] - \Im[\mathbf{W}_k] \Im[\mathbf{K}] - \Re[\mathbf{Y}_k] \right)^2 + \left(\Im[\mathbf{W}_k] \Re[\mathbf{K}] + \Re[\mathbf{W}_k] \Im[\mathbf{K}] - \Im[\mathbf{Y}_k] \right)^2 \right\}$. The Cartesian differential $\bar{\partial}_{\mathbf{K}} := \partial_{\Re[\mathbf{K}]} + i\partial_{\Im[\mathbf{K}]}$ then yields $\bar{\partial}_{\mathbf{K}} \hat{\mathcal{E}}(\mathbf{K}) = \sum_{k=1}^{N_k} \rho_k \left\{ 2 \left(\Re[\mathbf{W}_k] \Re[\mathbf{K}] - \Im[\mathbf{W}_k] \Im[\mathbf{K}] - \Re[\mathbf{Y}_k] \right) \cdot \left(\Re[\mathbf{W}_k] - i\Im[\mathbf{W}_k] \right) + 2 \left(\Im[\mathbf{W}_k] \Re[\mathbf{K}] + \Re[\mathbf{W}_k] \Im[\mathbf{K}] - \Im[\mathbf{Y}_k] \right) \cdot \left(\Im[\mathbf{W}_k] + i\Re[\mathbf{W}_k] \right) \right\}$ i.e. $\bar{\partial}_{\mathbf{K}} \hat{\mathcal{E}}(\mathbf{K}) = \sum_{k=1}^{N_k} \rho_k \left\{ 2 \Re[\mathbf{W}_k \mathbf{K} - \mathbf{Y}_k] \cdot \left(\Re[\mathbf{W}_k] - i\Im[\mathbf{W}_k] \right) + 2i \Im[\mathbf{W}_k \mathbf{K} - \mathbf{Y}_k] \cdot \left(\Re[\mathbf{W}_k] - i\Im[\mathbf{W}_k] \right) \right\}$ that is, noticing $\mathbf{W}_k^* = \left(\Re[\mathbf{W}_k] - i\Im[\mathbf{W}_k] \right)$, the Cartesian differential takes the form of equation (20), which establishes expressions (16), (17), and (18).

We now seek to establish the learning rate inequality (19), according to the well known property that maximum decrease is achieved with the learning rate $\gamma = \frac{1}{L}$. To do this, one must express the Hessian of

the loss function, and find its maximum eigenvalue, so as to obtain the corresponding Lipschitz constant L , according to $L = \max \left\{ \mathcal{S}p \left(\text{Hess} \left(\widehat{\mathcal{E}} \right) \right) \right\}$. We remind that, since $\partial_z z = 0$, the Cartesian differential of (20) would be zero: $\partial_{\mathbf{K}}^2 \widehat{\mathcal{E}} \left(\mathbf{K}^{(t)} \right) = 2 \sum_{k=1}^{N_k} \rho_k \mathbf{W}_k^* \mathbf{W}_k \left(\partial_{\mathbf{K}} \mathbf{K}^{(t)} \right) = 0$. However, each derivative along the real or imaginary component yields: $\partial_{\Re[\mathbf{K}]} \partial_{\mathbf{K}} \widehat{\mathcal{E}} \left(\mathbf{K}^{(t)} \right) = 2 \sum_{k=1}^{N_k} \rho_k \mathbf{W}_k^* \mathbf{W}_k = -i \partial_{\Im[\mathbf{K}]} \partial_{\mathbf{K}} \widehat{\mathcal{E}} \left(\mathbf{K}^{(t)} \right)$. This means that the Hessian can be interpreted as $2 \sum_{k=1}^{N_k} \rho_k \mathbf{W}_k^* \mathbf{W}_k$, which will yield real positive eigenvalues according to the spectral theorem. maximal eigenvalue of the Hessian yields the Lipschitz constant and thus learning rate γ according to (19), Q.E.D. \square

4. Computational PFBS algorithm

Having established an analytical expression of the PFBS iteration in theorem 2, we here transcribe it into the computational Frobenius penalization complex matrix PFBS algorithm 1, to solve problem (1), for a given λ , provided a tolerance ϵ on the empirical error defined on the training set as in (5) or (6):

Algorithm 1 Complex Frobenius PBFS algorithm

Input: $\mathbf{R}^{(0)}$, $\mathbf{C}^{(0)}$, and $\mathbf{A}^{(0)}$ initial guesses ; λ meta-parameter ; ϵ accuracy ; γ learning rate from (19)
while $\widehat{\mathcal{E}} \left(\mathbf{R}^{(t)}, \mathbf{C}^{(t)}, \mathbf{A}^{(t)} \right) \leq \epsilon$ **do**
 Update $\mathbf{R}_n^{(t+1)}, \mathbf{C}^{(t+1)}, \mathbf{A}^{(t+1)}$ using equation (15)
 $t \rightarrow t + 1$
end while
Output: $\mathbf{R}^{(\lambda)} = \left[\mathbf{R}_n^{(t_{\text{end}})} \right]_{n \in \llbracket 1, N_p \rrbracket}$, $\mathbf{C}^{(\lambda)} = \mathbf{C}^{(t_{\text{end}})}$, and $\mathbf{A}^{(\lambda)} = \mathbf{A}^{(t_{\text{end}})}$.

Defining the cross validation error as the empirical L_2 loss on a cross-validation grid $(z_\ell)_{\ell \in \llbracket 1, N_{\text{CV}} \rrbracket} \in \mathbb{C}$:

$$\widehat{\mathcal{E}}_{\text{CV}} \left(\mathbf{R}^{(\lambda)}, \mathbf{C}^{(\lambda)}, \mathbf{A}^{(\lambda)} \right) := \sum_{\ell=1}^{N_{\text{CV}}} \rho_\ell \left\| \sum_{n=1}^{N_p} \frac{\mathbf{R}_n^{(\lambda)}}{z_\ell - p_n} + \mathbf{C}^{(\lambda)} + z_\ell \mathbf{A}^{(\lambda)} - \mathbf{Y}_\ell \right\|_F^2 \quad (21)$$

where $\mathbf{R}^{(\lambda)}$, $\mathbf{C}^{(\lambda)}$, and $\mathbf{A}^{(\lambda)}$ are the solutions of problem (1) for a given λ , the meta-optimization process of λ selection can then be performed by selecting the λ as that minimizes the empirical L_2 norm on a cross-validation grid $(z_\ell)_{\ell \in \llbracket 1, N_{\text{CV}} \rrbracket} \in \mathbb{C}$: i.e.

$$\lambda_{\text{CV}} := \arg \min_{\lambda \in \mathbb{R}_+} \left\{ \widehat{\mathcal{E}}_{\text{CV}} \left(\mathbf{R}^{(\lambda)}, \mathbf{C}^{(\lambda)}, \mathbf{A}^{(\lambda)} \right) \right\} \quad (22)$$

In practice, we can solve the meta-optimization problem (22) on a logarithmic grid of $\lambda_m \in [\lambda_{\min}, \lambda_{\max}]$

values comprised in the interval

$$0 \leq \lambda_m \leq \frac{1}{2 \max \left\{ \mathcal{S}_p \left(\sum_{k=1}^{N_k} \rho_k \mathbf{W}_k^* \mathbf{W}_k \right) \right\}} \quad (23)$$

The order in which we span this grid is important: we start with λ_{\max} and then solve problem (1) by running algorithm 1, yielding $(\mathbf{R}^{(\lambda_{\max})}, \mathbf{C}^{(\lambda_{\max})}, \mathbf{A}^{(\lambda_{\max})})$, on which we can calculate $\hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_{\max})}, \mathbf{C}^{(\lambda_{\max})}, \mathbf{A}^{(\lambda_{\max})})$. We then solve problem (1) with the next lower $\lambda_{\max-1}$ value, by running algorithm 1 with initial guesses $(\mathbf{R}^{(\lambda_{\max})}, \mathbf{C}^{(\lambda_{\max})}, \mathbf{A}^{(\lambda_{\max})})$. This will yield $(\mathbf{R}_n^{(\lambda_{\max-1})}, \mathbf{C}^{(\lambda_{\max-1})}, \mathbf{A}^{(\lambda_{\max-1})})$, from which we can compute $\hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_{\max-1})}, \mathbf{C}^{(\lambda_{\max-1})}, \mathbf{A}^{(\lambda_{\max-1})})$. We then repeat the process going down the λ_m grid as long as $\hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_{m-1})}, \mathbf{C}^{(\lambda_{m-1})}, \mathbf{A}^{(\lambda_{m-1})}) \leq \hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_m)}, \mathbf{C}^{(\lambda_m)}, \mathbf{A}^{(\lambda_m)})$, until we find a minimum when $\hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_{m-1})}, \mathbf{C}^{(\lambda_{m-1})}, \mathbf{A}^{(\lambda_{m-1})}) > \hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_m)}, \mathbf{C}^{(\lambda_m)}, \mathbf{A}^{(\lambda_m)})$.

At each new λ_m iteration, some \mathbf{R}_n residues will have been drawn to zero, and we can then discard the corresponding poles p_n . We can then re-calculate the learning rate

$$\gamma_m = \frac{1}{2 \max \left\{ \mathcal{S}_p \left(\sum_{k=1}^{N_k} \rho_k \mathbf{W}_k^{(m)*} \mathbf{W}_k^{(m)} \right) \right\}} \quad (24)$$

where $\mathbf{W}_k^{(m)} := \left[\frac{\mathbb{I}}{z_k - p_1}, \dots, \frac{\mathbb{I}}{z_k - p_{N_p^{(m)}}}, \mathbb{I}, z_k \mathbb{I} \right]$ with $N_p^{(m)} \leq N_p^{(m-1)}$, where $(p_j)_{j \in \llbracket 1, N_p^{(m-1)} \rrbracket}$ is the set of remaining poles, from which the spurious poles have been discarded.

With this process, the meta-optimization is made faster, by pre-converging the initial guesses and sparsity exhibited by sub-optimal λ_m values. This yields the complex Frobenius PFBS meta-algorithm 2

Algorithm 2 Complex Frobenius PBFS meta-algorithm

Input: $\mathbf{R}^{(0)}$, $\mathbf{C}^{(0)}$, and $\mathbf{A}^{(0)}$ initial guesses ; $[\lambda_{\max}, \lambda_{\max}]$ meta-parameter grid ; ϵ accuracy
while $\hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_{m-1})}, \mathbf{C}^{(\lambda_{m-1})}, \mathbf{A}^{(\lambda_{m-1})}) \leq \hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_m)}, \mathbf{C}^{(\lambda_m)}, \mathbf{A}^{(\lambda_m)})$ **do**
 run algorithm 1 with initial guesses $(\mathbf{R}^{(\lambda_{m-1})}, \mathbf{C}^{(\lambda_{m-1})}, \mathbf{A}^{(\lambda_{m-1})})$, learning rate γ_m from (24)
 eliminate spurious poles (p_n) corresponding to zero residues $\mathbf{R}_n = \mathbf{0}$
 compute cross validation error $\hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_m)}, \mathbf{C}^{(\lambda_m)}, \mathbf{A}^{(\lambda_m)})$
 $m \rightarrow m - 1$
end while
Output: $(\mathbf{R}^{(\lambda_{\text{opt}})}, \mathbf{C}^{(\lambda_{\text{opt}})}, \mathbf{A}^{(\lambda_{\text{opt}})})$, optimal meta-parameter λ_{opt} , lowest cross validation error $\hat{\mathcal{E}}_{\text{CV}}(\mathbf{R}^{(\lambda_{\text{opt}})}, \mathbf{C}^{(\lambda_{\text{opt}})}, \mathbf{A}^{(\lambda_{\text{opt}})})$.

5. PFBS algorithm : implementation & computational performance

To test the validity of our theorem (2), and the subsequent computational complex Frobenius PFBS algorithm 1 and meta-algorithm 2, we devise a test set of points $(z_k)_{k \in \llbracket 1, N_k \rrbracket} \in \mathbb{C}$ and matrices $\mathbf{Y}_k \in \mathbb{C}_{q \times s}$ with weights $(\rho_k)_{k \in \llbracket 1, N_k \rrbracket} \in \mathbb{R}_+$, which we will have generated from a known reference complex matrix operator $\mathbf{Y}(z) = \frac{\mathbf{R}_n^{\text{ref}}}{z_k - p_n} + \mathbf{C}^{\text{ref}} + z_k \mathbf{A}^{\text{ref}}$, so that $\mathbf{Y}_k := \mathbf{Y}(z_k)$.

We then searched for poles with and without noise: i.e. $p_n := p_n + \epsilon$. When no noise was introduced, the LS method is clearly superior, and manages to find all the residues (including the zero ones). However, when the poles are perturbed by less than 10%, the group-LASSO was able to eliminate the corresponding residue. This shows that our method is theoretically valid, albeit more work has to be done to palliate the issue of extreme condition numbers of the Cauchy matrices. Elastic net or QR seems the way to go.

Note that we could not have introduced the noise in the residues, i.e. $\mathbf{Y}_k := \mathbf{Y}(z_k) + \epsilon_{\mathcal{N}}$, with Gaussian noise $\epsilon_{\mathcal{N}} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$.

6. Conclusions

- We theoretically generalized the PFBS algorithm with L_1 -Frobenius group LASSO penalization to the case of complex matrices.
- We implemented it in Julia and proved it write on the polynomial and residues case.
- We stumbled upon very significant ill-conditioning problems that made it very hard to continue
- Elastic Net regularization to address this issue is our next step.
- De-bug for all cases.
- Make sure the regularization does trim the right poles !
- study the relation between number of points and conditioning.
- Code Elastic net for conditioning