

Partially separable structure (PSS) May 20, 2021

Paul Raynaud

May 20, 2021

Polytechnique Montréal, Grenoble INP

Problem of interest:

$$\min_{x \in \mathbb{R}^n} f(x) \quad f(x) := \sum_{i=1}^N f_i(x)$$

- large problems $n > 10^3$
- $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ does not depend on all of x
- $f_i \in \mathcal{C}^2, i = 1, \dots, N$

Example:

$$\min_{x \in \mathbb{R}^n} f_1(x_1, x_2) + f_n(x_{n-1}, x_n) + \sum_{i=2}^{n-1} f_i(x_{i-1}, x_i, x_{i+1})$$

Definition

The linear operator U_i gives the (linear combination of) variables used by f_i :

$$\begin{aligned}f(x) &= \sum_{i=1}^N \widehat{f}_i(U_i x) \\ \nabla f(x) &= \sum_{i=1}^N U_i^\top \nabla \widehat{f}_i(U_i x) \\ \nabla^2 f(x) &= \sum_{i=1}^N U_i^\top \nabla^2 \widehat{f}_i(U_i x) U_i \\ \nabla^2 f(x) \approx B &= \sum_{i=1}^N U_i^\top \widehat{B}_i U_i\end{aligned}$$

- $\widehat{f}_i: \mathbb{R}^{n_i} \rightarrow \mathbb{R}$ an element function
- $U_i \in \mathbb{R}^{n_i \times n}$ usually a linear operator far more efficient than a matrix
- $\widehat{B}_i \in \mathbb{R}^{n_i \times n_i}, i = 1, \dots, N$
- If $\max_{i=\{1, \dots, N\}} n_i \ll n$, store $\{\widehat{B}_i\}_{i=1}^N$ requires (much) less memory than B

Theorem (Griewank and Toint [1982a])

Every problem having a sparse hessian is partially separable.

Motivation for studying the PSS

The PSS allows partitioned QN updates (PQN) (Griewank and Toint [1982b])

$$B = \sum_{i=1}^N B_i = \sum_{i=1}^N U_i^\top \hat{B}_i U_i$$

- Apply QN update to each \hat{B}_i using $U_i s$ and $\nabla \hat{f}_i(U_i(x+s)) - \nabla \hat{f}_i(U_i x)$
- $\sum_{i=1}^N B_i$ still satisfies secant equation
- Advantages:
 - does not increase memory requirements $\{\hat{B}_i\}_{i=1}^N$ (\neq standard QN)
 - keep the sparsity of B (\neq L-BFGS)
 - fully parallelizable: each \hat{B}_i update is independent: $(U_i s, \hat{y}_i)$ such that $\hat{y}_i = \nabla \hat{f}_i(U_i(x+s)) - \nabla \hat{f}_i(U_i x)$
 - rank update \gg 1 or 2

A trust-region method or a linesearch framework around the PQN update leads us to solve a partitioned linear system at every iteration:

- Conjugate gradient (CG)
 - require matrix-vector products: $Bv = \left(\sum U_i^\top \hat{B}_i U_i\right)v$
 - can compute $\hat{B}_i U_i v$ in parallel and assemble with U_i^\top
- (multi-)frontal factorization (Conn et al. [1994])
 - Cholesky factorization dedicated to partitioned matrix
 - $\{U_i\}_{i=1}^N$ provide the sparsity of B
 - the permutation applied to the matrix is **critical**: front size, filling, parallelizable blocs
- partitioned trust-region method (Conn et al. [1996])

Efficient derivatives computation

Reduce $f(x) = \sum_{i=1}^N \widehat{f}_i(U_i x)$ evaluation required to compute ∇f from $\{\nabla \widehat{f}_i\}_{i=1}^N$ in case every $\widehat{f}_i(x)$ are evaluated at once and by using the structure $\{U_i\}_{i=1}^N$.

$$f(x) = \sum_{i=1}^5 f_i(x) = 1^\top F(x) = \widehat{f}_1(x_1, x_3) + \widehat{f}_2(x_1, x_4) + \widehat{f}_3(x_2, x_3) + \widehat{f}_4(x_2, x_4) + \widehat{f}_5(x_3)$$

$$F(x) = \begin{pmatrix} f_1(x) \\ f_2(x) \\ f_3(x) \\ f_4(x) \\ f_5(x) \end{pmatrix}, \quad \nabla F = \begin{pmatrix} \square & & \Delta & & \\ \square & & & & \diamond \\ & \diamond & \Delta & & \\ & \diamond & & & \diamond \\ & & & \Delta & \end{pmatrix}, \quad S_c = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$

If ∇F is dense the seed S to compute ∇F is I_4 implying 4 f evaluations. The PSS $\{U_i\}_{i=1}^N$ induce graph structure whose a proper coloring define the compressed seed S_c implying 2 f evaluations.

Automatic differentiation (AD)

Compute the derivative of a numerical procedure $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Forward mode
 - more efficient than reverse if $m > n$
 - memoryless method
- Reverse mode
 - more efficient than forward if $m < n$
 - must build a tape of the numerical procedure.

- If every \hat{f}_i is available and evaluate at once a similar procedure to compressed seed may be used (Bischof et al. [1997])
- If each \hat{f}_i is available individually:
 - forward mode is more efficient since $n_i \ll n$
 - the tape of each \hat{f}_i is much smaller than f (smaller expression tree)
 - in practice, $\hat{f}_i = \hat{f}_j$ allowing to reduce the number of tapes needed
- Hessian-vector products $\nabla^2 \hat{f}_i(U_i x) U_i v$ combine both approaches and their properties in PSS. It allows a complete parallel procedure to compute $\nabla^2 f(x) v$ from $\sum_{i=1}^N U_i^\top \nabla^2 \hat{f}_i(U_i x) U_i v$

- Dedicated crossover operator, the key of genetic algorithms (Durand and Alliot [1998])
- Specific to DFO:
 - Interpolations based on the knowledge of $\{\hat{f}_i\}_{i=1}^N$
 - By interpolating each \hat{f}_i , $\approx n_i^2$ points instead of $\approx n^2$
 - Reducing the $\{\hat{f}_i\}_{i=1}^N$ evaluations depending the structure to obtain those n_i^2 points
 - Dedicated efficient procedure to recompute $f, \nabla f$ if $x_{k+1} - x_k$ is sparse, only the $\hat{f}_i, \nabla \hat{f}_i$ impact must be recompute
- Brute Force Optimizer (BFO): (Porcelli and Toint [2021])

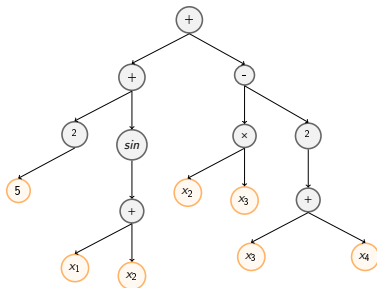
- Problem structure must be explicit by the modeler
- \triangle if $\sum n_i^2 \geq n^2$: not applicable in large scale, require more space and computation than BFGS, ex: $f(x) = \sum_{i=1}^n \hat{f}_i(x_1, x_2, \dots, x_i)$
 - Method to find a new basis to increase the sparsity of the problem
Kim et al. [2009]

- Study of PSS is about 40 years old Griewank and Toint [1982a]
- During the last 40 years, work mainly done by Conn, Gould and Toint
- Resulting LANCELOT a Fortran software using the SIF format
- AMPL (commercial software) also uses the PSS and detects it automatically.

- Provide modern software to detect PSS automatically:
 - Assess convexity of the f_i automatically
 - Construct new optimization methods that exploit PSS
 - 4 julia modules
 - Make it easily usable (\neq LANCELOT)
- Survey on PSS

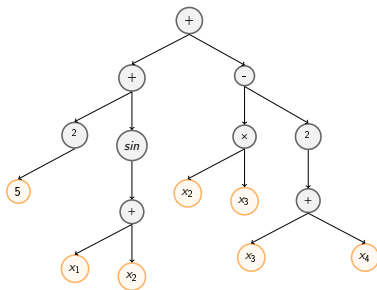
- Detect PSS $(\{\hat{f}_i\}_{i=1}^N, \{U_i\}_{i=1}^N)$ automatically from f
- Automatic strict convexity detection and bounds propagation
- Interfaced to JuMP, NLPModelJuMP, ADNLPModel

Example



$f_1(x) = 5^2$	$[5]$	<i>constant</i>	<i>non strictly convex</i>
$f_2(x) = \sin(x_1 + x_2)$	$[-1, 1]$	<i>nonlinear</i>	<i>non strictly convex</i>
$f_3(x) = x_2 \times x_3$	$[-\infty, \infty]$	<i>quadratic</i>	<i>non strictly convex</i>
$f_4(x) = -(x_3 + x_4)^2$	$[0, \infty]$	<i>quadratic</i>	<i>non strictly convex</i>

Example



$$U_1 = 0 \quad U_2 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \text{ or } U_2 = \begin{pmatrix} 1 & 1 & 0 & 0 \end{pmatrix} \quad U_3 = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

$$U_4 = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \text{ or } U_4 = \begin{pmatrix} 0 & 0 & 1 & 1 \end{pmatrix}$$

- define the algorithm structures around PSS
- Test problem Rosenbrock function ($\mathbb{R}^n \rightarrow \mathbb{R}$)

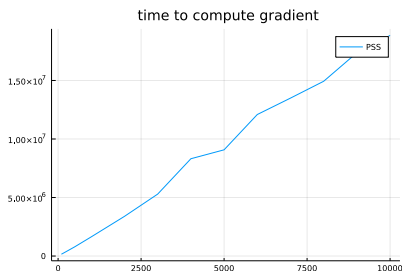


Figure 1: PSF gradient t/n

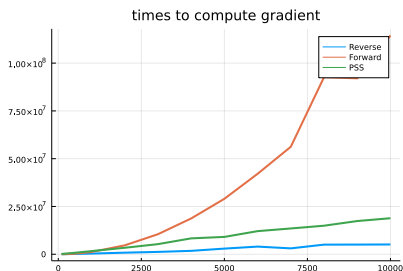


Figure 2: comparison ton AD t/n

- Trust-region methods using partitioned-QN solved by CG
- 40 PS problems of size $n = 1000$

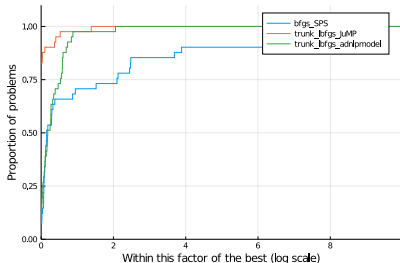


Figure 3: time

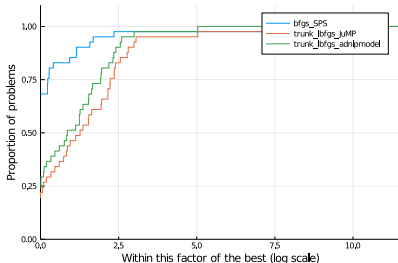


Figure 4: obj+5grad

- Dedicated to partitioned structured: vectors/matrices. Also some specificity about PSS.
- Multi-frontal factorization implementation

Currently a trust region using a P-BFGS update must solve at each iterate the partitioned linear system:

$$\begin{aligned} Ax &= b \\ \sum_i U_i^T \hat{A}_i U_i x &= \sum_i U_i^T \hat{b}_i \\ \hat{A}_i &= \hat{B}_i, \hat{b}_i = -\nabla \hat{f}_i, x = s \end{aligned}$$

The complexity of the whole method:

- PQN: update $\{\hat{B}_i\}_{i=1}^N$ (fully parallelizable, depend of n_i)
- TR management is constant
- Solving the partitioned linear system: CG is state of the art

The following properties must hold:

- Do not form B
- Be parallel
- To use it with a trust region an approximate solution is enough.
- The solution must be a descent direction

The ideal would be an iterative method that iteratively check TR constraint (similar to CG).

In **completely** separable case solving each $\widehat{A}_i \widehat{x}^i = \widehat{b}_i$, $x^i \in \mathbb{R}^{n_i}$ solve $Ax = b$.

$$\begin{pmatrix} \widehat{A}_1 & 0 \\ 0 & \widehat{A}_2 \end{pmatrix} \begin{pmatrix} x \\ \end{pmatrix} = \begin{pmatrix} b^1 \\ b^2 \end{pmatrix}$$

- Plan to form a solution x from $\{\widehat{x}^i\}_{i=1}^2$ such that $\widehat{A}_i \widehat{x}^i = \widehat{b}_i$
- Consequently each $\widehat{A}_i x^i = \widehat{b}_i$ may be solve in parallel.

$$x = \begin{pmatrix} \widehat{x}^1 \\ \widehat{x}^2 \end{pmatrix} \tag{1}$$

The cas of two overlapping blocs

Suppose $\hat{A}_1, \hat{A}_2 \in \mathbb{R}^{n_1 \times n_1}, \mathbb{R}^{n_2 \times n_2}$ such that $A \in \mathbb{R}^{n \times n}$.

$$Ax = \begin{pmatrix} \hat{A}_{1,1} & \hat{A}_{1,2} & 0 \\ \hat{A}_{1,2,1} & \hat{A}_{1,2,2} + \hat{A}_{2,1,1} & \hat{A}_{2,1,2} \\ 0 & \hat{A}_{2,2,1} & \hat{A}_{2,2,2} \end{pmatrix} x = \begin{pmatrix} \hat{b}_{1,1} \\ \hat{b}_{1,2} + \hat{b}_{2,1} \\ \hat{b}_{2,2} \end{pmatrix} = b$$

Suppose \hat{x}^1, \hat{x}^2 such that $\hat{A}_p \hat{x}^p = \hat{b}_p$, $p=1,2$ and an approximation $x^?$ of x^* such that:

$$x^? = \begin{pmatrix} \hat{x}_1^1 \\ \hat{x}^? \\ \hat{x}_2^2 \end{pmatrix} \quad x^1 = \begin{pmatrix} x_1^1 \\ x_2^1 \end{pmatrix} \quad x^2 = \begin{pmatrix} x_1^2 \\ x_2^2 \end{pmatrix}$$

Consequently $Ax^? - b$

$$\begin{aligned}\widehat{A}_{p1,1} \widehat{x}_1^p + \widehat{A}_{p1,2} \widehat{x}_2^p &= \widehat{b}_{p1} \\ \widehat{A}_{p2,1} \widehat{x}_1^p + \widehat{A}_{p2,2} \widehat{x}_2^p &= \widehat{b}_{p2}\end{aligned}$$

Replace \widehat{x}^p by $U_i x^?$:

$$\begin{aligned}\widehat{A}_{11,1} \widehat{x}_1^1 + \widehat{A}_{11,2} \widehat{x}_2^1 + \widehat{A}_{11,2} (\widehat{x}^? - \widehat{x}_2^1) &= \widehat{b}_{11} \\ \underbrace{\widehat{A}_{11,1} \widehat{x}_1^1 + \widehat{A}_{11,2} \widehat{x}_2^1 - \widehat{b}_{11}}_{=0} + \widehat{A}_{11,2} (\widehat{x}^? - \widehat{x}_2^1) &= 0 \\ \widehat{A}_{12,1} \widehat{x}_1^1 + \widehat{A}_{12,2} \widehat{x}_2^1 + \widehat{A}_{12,2} (\widehat{x}^? - \widehat{x}_2^1) &= \widehat{b}_{12} \\ \underbrace{\widehat{A}_{12,1} \widehat{x}_1^1 + \widehat{A}_{12,2} \widehat{x}_2^1 - \widehat{b}_{12}}_{=0} + \widehat{A}_{12,2} (\widehat{x}^? - \widehat{x}_2^1) &= 0 \\ \widehat{A}_{21,1} \widehat{x}_1^2 + \widehat{A}_{21,2} \widehat{x}_2^2 + \widehat{A}_{21,1} (\widehat{x}^? - \widehat{x}_1^2) &= \widehat{b}_{21} \\ \underbrace{\widehat{A}_{21,1} \widehat{x}_1^2 + \widehat{A}_{21,2} \widehat{x}_2^2 - \widehat{b}_{21}}_{=0} + \widehat{A}_{21,1} (\widehat{x}^? - \widehat{x}_1^2) &= 0 \\ \widehat{A}_{22,1} \widehat{x}_1^2 + \widehat{A}_{22,2} \widehat{x}_2^2 + \widehat{A}_{22,1} (\widehat{x}^? - \widehat{x}_1^2) &= \widehat{b}_{22} \\ \underbrace{\widehat{A}_{22,1} \widehat{x}_1^2 + \widehat{A}_{22,2} \widehat{x}_2^2 - \widehat{b}_{22}}_{=0} + \widehat{A}_{22,1} (\widehat{x}^? - \widehat{x}_1^2) &= 0\end{aligned}$$

The residual $Ax^? - b$ is the following:

$$Ax^? - b = - \begin{pmatrix} \hat{A}_{1,2}(\hat{x}^? - \hat{x}_2^1) \\ \hat{A}_{1,2}(\hat{x}^? - \hat{x}_2^1) + \hat{A}_{2,1}(\hat{x}^? - \hat{x}_1^2) \\ \hat{A}_{2,1}(\hat{x}^? - \hat{x}_1^2) \end{pmatrix}$$

This equation link $x^?, \hat{x}^?$ and a approximate solution $Ax = b$ (ie $Bs = -g$). We would like to minimize the residual $Ax^? - b$; depending only of $\hat{x}^?$.

Remark: The optimum of this problem $Ax^? - b$ may be not null since the approximate $x^?$ is arbitrary.

$$\min_{\hat{x}^? \in \mathbb{R}^{n_{inter}}} \left\| \underbrace{\begin{pmatrix} \hat{A}_{11,2}(\hat{x}^? - \hat{x}_2^1) \\ \hat{A}_{12,2}(\hat{x}^? - \hat{x}_2^1) + \hat{A}_{21,1}(\hat{x}^? - \hat{x}_1^2) \\ \hat{A}_{22,1}(\hat{x}^? - \hat{x}_1^2) \end{pmatrix}}_{\in \mathbb{R}^n} \right\|$$

- Problem dimension: $n = n_1 + n_2 - n_{inter}$
- Variable dimension: n_{inter}
- Directional derivative of $\hat{x}^?$ are combination of $\hat{A}_{1,2}$ and $\hat{A}_{2,1}$

An ongoing work:

- Still don't know how to solve this new problem
- May be extends to more than 2 blocs
- Litterature review about bloc matrix resolution (ADMM)

References

- A Griewank and Ph Toint. On the unconstrained optimization of partially separable functions. In M. J. D Powell, editor, *Nonlinear Optimization 1981*, pages 301–312. Academic press, 1982a. Publication editors : M.J.D. Powell.
- A Griewank and Ph Toint. Partitioned variable metric updates for large structured optimization problems. 39:119–137, 1982b. doi: 10.1007/BF01399316.
- A Conn, N Gould, M Lescrenier, and Ph Toint. Performance of a multifrontal scheme for partially separable optimization. In *Advances in Optimization and Numerical Analysis*, pages 79–96. Springer, 1994. doi: 10.1007/978-94-015-8330-5_6.
- A Conn, N Gould, A Sartenaer, and Ph Toint. Convergence properties of minimization algorithms for convex constraints using a structured trust region. 6(4):1059–1086, 1996. doi: <https://doi.org/10.1137/S1052623492236481>.
- C Bischof, A Bouaricha, P Khademi, and J Moré. Computing gradients in large-scale optimization using automatic differentiation. 9:185–194, 1997. doi: <https://doi.org/10.1287/ijoc.9.2.185>.

N Durand and J-M Alliot. Genetic crossover operator for partially separable functions. In *GP 1998, 3rd annual conference on Genetic Programming*, Madison, United States, 1998. URL <https://hal-enac.archives-ouvertes.fr/hal-00937718>.

M. Porcelli and Ph Toint. Exploiting problem structure in derivative free optimization, 2021.

S Kim, M Kojima, and Ph Toint. Recognizing underlying sparsity in optimization. 119(2):273–303, Jul 2009. ISSN 1436-4646. doi: 10.1007/s10107-008-0210-4. URL <https://doi.org/10.1007/s10107-008-0210-4>.