

Secrets of Matrix Factorization: Supplementary Document

Je Hyeong Hong
University of Cambridge
jhh37@cam.ac.uk

Andrew Fitzgibbon
Microsoft, Cambridge, UK
awf@microsoft.com

Abstract

This document consists of details to which the original submission article has made references. These include a list of symbols, some matrix identities and calculus rules, a list of VarPro derivatives, symmetry analysis of the un-regularized VarPro Gauss-Newton matrix and MTSS computation details.

Symbol	Meaning
$\mathbb{R}^{p \times q}$	The space of real matrices of size $p \times q$
\mathbb{S}^p	The space of symmetric real matrices of size $p \times p$
\mathbf{M}	The measurement matrix $\in \mathbb{R}^{m \times n}$ (may contain noise and missing data)
\mathbf{W}	The weight matrix $\in \mathbb{R}^{m \times n}$
\mathbf{U}	The first decomposition matrix $\in \mathbb{R}^{m \times r}$
\mathbf{V}	The second decomposition matrix $\in \mathbb{R}^{n \times r}$
\mathbf{u}	$\text{vec}(\mathbf{U}) \in \mathbb{R}^{mr}$
\mathbf{v}	$\text{vec}(\mathbf{V}^\top) \in \mathbb{R}^{nr}$
m	The number of rows in \mathbf{M}
n	The number of columns in \mathbf{M}
p	The number of visible entries in \mathbf{M}
p_j	The number of visible entries in the j -th column of \mathbf{M}
f	The objective $\in \mathbb{R}$
$\boldsymbol{\varepsilon}$	The vector objective $\in \mathbb{R}^p$
\mathbf{J}	Jacobian
\mathbf{H}	Hessian
\mathbf{V}^*	$\arg \min_{\mathbf{V}} f(\mathbf{U}, \mathbf{V}) \in \mathbb{R}^{n \times r}$
Π	Projection matrix $\in \mathbb{R}^{p \times mn}$ selecting nonzero entries of $\text{vec}(\mathbf{W})$
$\tilde{\mathbf{W}}$	$\Pi \text{diag}(\text{vec}(\mathbf{W})) \in \mathbb{R}^{p \times mn}$
$\tilde{\mathbf{m}}$	$\tilde{\mathbf{W}} \text{vec}(\mathbf{M}) \in \mathbb{R}^p$
$\tilde{\mathbf{U}}$	$\tilde{\mathbf{W}}(\mathbf{I}_n \otimes \mathbf{U}) \in \mathbb{R}^{p \times nr}$
$\tilde{\mathbf{V}}$	$\tilde{\mathbf{W}}(\mathbf{V} \otimes \mathbf{I}_m) \in \mathbb{R}^{p \times mr}$
\mathbf{K}_{mr}	The permutation matrix satisfying $\text{vec}(\mathbf{U}^\top) = \mathbf{K}_{mr} \text{vec}(\mathbf{U}) \in \mathbb{R}^{mr \times mr}$
\mathbf{K}_{nr}	The permutation matrix satisfying $\text{vec}(\mathbf{V}^\top) = \mathbf{K}_{nr} \text{vec}(\mathbf{V}) \in \mathbb{R}^{nr \times nr}$
\mathbf{R}	$\mathbf{W} \odot (\mathbf{U}\mathbf{V}^\top - \mathbf{M}) \in \mathbb{R}^{m \times n}$
\mathbf{Z}	$(\mathbf{W} \odot \mathbf{R}) \otimes \mathbf{I}_r \in \mathbb{R}^{mr \times nr}$

Table 1: A list of symbols used in [1].

1. Symbols

A list of symbols used in [1] can be found in Table 1.

2. Matrix identities and calculus rules

Some of the key matrix identities and calculus rules in [6, 5] are illustrated and extended below for convenience.

2.1. Vec and Kronecker product identities

Given that each small letter in bold is an arbitrary column vector of appropriate size, we have

$$\text{vec}(\mathbf{A} \odot \mathbf{B}) = \text{diag}(\text{vec } \mathbf{A}) \text{vec}(\mathbf{B}), \quad (1)$$

$$\text{vec}(\mathbf{A}\mathbf{X}\mathbf{B}) = (\mathbf{B}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{X}), \quad (2)$$

$$\text{vec}(\mathbf{A}\mathbf{B}) = (\mathbf{B}^\top \otimes \mathbf{I}) \text{vec}(\mathbf{A}) = (\mathbf{I}^\top \otimes \mathbf{A}) \text{vec}(\mathbf{B}), \quad (3)$$

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{A}\mathbf{C} \otimes \mathbf{B}\mathbf{D}, \quad (4)$$

$$(\mathbf{A} \otimes \mathbf{b})\mathbf{C} = \mathbf{A}\mathbf{C} \otimes \mathbf{b}, \quad (5)$$

$$\text{vec}(\mathbf{A}^\top) = \mathbf{K}_\mathbf{A} \text{vec}(\mathbf{A}), \quad (6)$$

$$(\mathbf{A} \otimes \mathbf{B}) = \mathbf{K}_1^\top (\mathbf{B} \otimes \mathbf{A}) \mathbf{K}_2, \text{ and} \quad (7)$$

$$(\mathbf{A} \otimes \mathbf{b}) = \mathbf{K}_1^\top (\mathbf{b} \otimes \mathbf{A}). \quad (8)$$

2.2. Differentiation of μ -pseudo inverse

Let us first define μ -pseudo inverse $\mathbf{X}^{-\mu} := (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \mathbf{X}^\top$. The derivative of μ -pseudo inverse is

$$\partial[\mathbf{X}^{-\mu}] = \partial[(\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1}] \mathbf{X}^\top + (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \partial[\mathbf{X}]^\top \quad (9)$$

$$= -(\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \partial[\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I}] (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \mathbf{X}^\top + (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \partial[\mathbf{X}]^\top \quad (10)$$

$$= -(\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} (\partial[\mathbf{X}^\top] \mathbf{X} + \mathbf{X}^\top \partial[\mathbf{X}]) (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \mathbf{X}^\top + (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \partial[\mathbf{X}]^\top \quad (11)$$

$$= -(\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \partial[\mathbf{X}^\top] \mathbf{X} \mathbf{X}^{-\mu} - (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \mathbf{X}^\top \partial[\mathbf{X}] \mathbf{X}^{-\mu} + (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \partial[\mathbf{X}]^\top \quad (12)$$

$$= -\mathbf{X}^{-\mu} \partial[\mathbf{X}] \mathbf{X}^{-\mu} + (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \partial[\mathbf{X}]^\top (\mathbf{I} - \mathbf{X} \mathbf{X}^{-\mu}). \quad (13)$$

The derivative of pseudo-inverse $\mathbf{X}^\dagger := (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is then simply the above for $\mu = 0$.

3. Derivatives for variable projection (VarPro)

A full list of derivatives for un-regularized VarPro can be found in Table 2. Derivations are illustrated in [2]. The regularized setting is also discussed in [2].

3.1. Differentiating $\mathbf{v}^*(\mathbf{u})$

From [1], we have

$$\mathbf{v}^*(\mathbf{u}) := \arg \min_{\mathbf{v}} f(\mathbf{u}, \mathbf{v}) = \tilde{\mathbf{U}}^{-\mu} \tilde{\mathbf{m}}. \quad (14)$$

Substituting (14) to the original objective yields

$$\varepsilon_1^*(\mathbf{u}, \mathbf{v}^*(\mathbf{u})) = \tilde{\mathbf{U}} \mathbf{v}^* - \tilde{\mathbf{m}} = -(\mathbf{I} - \tilde{\mathbf{U}} \tilde{\mathbf{U}}^{-\mu}) \tilde{\mathbf{m}}. \quad (15)$$

Taking $\partial[\mathbf{v}^*]$ and applying (13) gives us

$$\partial[\mathbf{v}^*] = -\tilde{\mathbf{U}}^{-\mu} \partial[\tilde{\mathbf{U}}] \tilde{\mathbf{U}}^{-\mu} \tilde{\mathbf{m}} + (\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} + \mu \mathbf{I})^{-1} \partial[\tilde{\mathbf{U}}]^\top (\mathbf{I} - \tilde{\mathbf{U}} \tilde{\mathbf{U}}^{-\mu}) \tilde{\mathbf{m}} \quad (16)$$

$$= -\tilde{\mathbf{U}}^{-\mu} \partial[\tilde{\mathbf{U}}] \mathbf{v}^* - (\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} + \mu \mathbf{I})^{-1} \partial[\tilde{\mathbf{U}}]^\top \varepsilon_1^* \quad // \text{ noting } \mathbf{v}^*(\mathbf{u}) = \tilde{\mathbf{U}}^{-\mu} \tilde{\mathbf{m}} \text{ and (15)} \quad (17)$$

$$= -\tilde{\mathbf{U}}^{-\mu} \tilde{\mathbf{V}}^* \partial \mathbf{u} - (\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} + \mu \mathbf{I})^{-1} \mathbf{Z}^{*\top} \mathbf{K}_{mr} \partial \mathbf{u}, \quad // \text{ noting bilinearity and the results in [2]} \quad (18)$$

and hence

$$\frac{d\mathbf{v}^*}{d\mathbf{u}} = -(\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} + \mu \mathbf{I})^{-1} (\tilde{\mathbf{U}}^\top \tilde{\mathbf{V}}^* + \mathbf{Z}^{*\top} \mathbf{K}_{mr}). \quad (19)$$

Quantity	Definition	ALS (RW3) + [Approx. Gauss-Newton (RW2)] _{RW2} + [Full Gauss-Newton (RW1)] _{RW1} + [Full Newton] _{FN} w/o (Damping) w/o (Projection constraint) _P
$\mathbf{v}^*(\mathbf{u}) \in \mathbb{R}^{nr}$	$\mathbf{v}^* := \arg \min_{\mathbf{v}} f(\mathbf{u}, \mathbf{v})$	$\mathbf{v}^* = \tilde{\mathbf{U}}^\dagger \tilde{\mathbf{m}}$ ($\tilde{\mathbf{m}} \in \mathbb{R}^p$ consists of non-zero elements of \mathbf{m} .)
$\frac{d\mathbf{v}^*}{d\mathbf{u}} \in \mathbb{R}^{nr \times mr}$	$\frac{d\mathbf{v}^*}{d\mathbf{u}} := \frac{d \text{vec}(\mathbf{V}^{*\top})}{d \text{vec}(\mathbf{U})}$	$\frac{d\mathbf{v}^*}{d\mathbf{u}} = -[\tilde{\mathbf{U}}^\dagger \tilde{\mathbf{V}}^*]_{RW2} - [(\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}})^{-1} \mathbf{Z}^{*\top} \mathbf{K}_{mr}]_{RW1}$ ($\mathbf{Z}^* := (\mathbf{W} \odot \mathbf{R}^*) \otimes \mathbf{I}_r$)
Cost vector $\boldsymbol{\varepsilon}_1^* \in \mathbb{R}^p$	$\boldsymbol{\varepsilon}_1^*$	$\boldsymbol{\varepsilon}_1^* := \tilde{\mathbf{U}} \mathbf{v}^* - \mathbf{m}$
Jacobian $\mathbf{J}_1^* \in \mathbb{R}^{p \times mr}$	$\mathbf{J}_1^* := \frac{d\boldsymbol{\varepsilon}_1^*}{d\mathbf{u}}$	$\mathbf{J}_1^* = (\mathbf{I}_p - [\tilde{\mathbf{U}} \tilde{\mathbf{U}}^\dagger]_{RW2}) \tilde{\mathbf{V}}^* - [\tilde{\mathbf{U}}^\dagger \mathbf{Z}^{*\top} \mathbf{K}_{mr}]_{RW1}$
Cost function $f_1^* \in \mathbb{R}$	$f_1^* := \ \boldsymbol{\varepsilon}_1^*\ _2^2$	$f_1^* = \ \mathbf{W} \odot (\mathbf{U} \mathbf{V}^{*\top} - \mathbf{M})\ _F$
Ordinary gradient $\mathbf{g}_1^* \in \mathbb{R}^{mr}$	$\mathbf{g}_1^* := \frac{df_1^*}{d\mathbf{u}}$	$\frac{1}{2} \mathbf{g}_1^* = \tilde{\mathbf{V}}^{*\top} \boldsymbol{\varepsilon}_1^* = \text{vec}((\mathbf{W} \odot \mathbf{R}^*) \mathbf{V}^*)$
Projected gradient $\mathbf{g}_p^* \in \mathbb{R}^{mr}$	$\mathbf{g}_p^* := \mathbf{P}_p \mathbf{g}_1^*$ ($\mathbf{P}_p := \mathbf{I}_{mr} - \mathbf{I}_r \otimes \mathbf{U} \mathbf{U}^\top$)	$\frac{1}{2} \mathbf{g}_p^* = \frac{1}{2} \mathbf{g}_1^*$
Reduced gradient $\mathbf{g}_r^* \in \mathbb{R}^{mr-r^2}$	$\mathbf{g}_r^* := \mathbf{P}_r \mathbf{g}_1^*$ ($\mathbf{P}_r := \mathbf{I}_r \otimes \mathbf{U}_\perp^\top$)	$\frac{1}{2} \mathbf{g}_r^* = \frac{1}{2} \mathbf{P}_r \mathbf{g}_1^*$
Ordinary Hessian $\mathbf{H}_{1GN}^*, \mathbf{H}_1^* \in \mathbb{S}^{mr}$	$\frac{1}{2} \mathbf{H}_{1GN}^* := \mathbf{J}_1^{*\top} \mathbf{J}_1^* + \langle \lambda \mathbf{I} \rangle$ $\frac{1}{2} \mathbf{H}_1^* := \frac{1}{2} \frac{d\mathbf{g}_1^*}{d\mathbf{u}} + \langle \lambda \mathbf{I} \rangle$	$\frac{1}{2} \mathbf{H}_{1GN}^* = \tilde{\mathbf{V}}^{*\top} (\mathbf{I}_p - [\tilde{\mathbf{U}} \tilde{\mathbf{U}}^\dagger]_{RW2}) \tilde{\mathbf{V}}^*$ $+ [\mathbf{K}_{mr}^\top \mathbf{Z}^* (\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}})^{-1} \mathbf{Z}^{*\top} \mathbf{K}_{mr}]_{RW1} + \langle \lambda \mathbf{I}_{mr} \rangle$ $\frac{1}{2} \mathbf{H}_1^* = \frac{1}{2} \mathbf{H}_{1GN}^* - [2\mathbf{K}_{mr}^\top \mathbf{Z}^* (\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}})^{-1} \mathbf{Z}^{*\top} \mathbf{K}_{mr}]_{FN}$ $- [\mathbf{K}_{mr}^\top \mathbf{Z}^* \tilde{\mathbf{U}}^\dagger \tilde{\mathbf{V}}^* + \tilde{\mathbf{V}}^{*\top} \tilde{\mathbf{U}}^\dagger \mathbf{Z}^{*\top} \mathbf{K}_{mr}]_{FN} + \langle \lambda \mathbf{I}_{mr} \rangle$
Projected Hessian $\mathbf{H}_{pGN}^*, \mathbf{H}_p^* \in \mathbb{S}^{mr}$	$\frac{1}{2} \mathbf{H}_{pGN}^* := \mathbf{P}_p \mathbf{J}_1^{*\top} \mathbf{J}_1^* \mathbf{P}_p + \langle \lambda \mathbf{I} \rangle$ $\frac{1}{2} \mathbf{H}_p^* := \frac{1}{2} \mathbf{P}_p \frac{d\mathbf{g}_p^*}{d\mathbf{u}} \mathbf{P}_p + \langle \lambda \mathbf{I} \rangle$ ($\mathbf{P}_p := \mathbf{I}_r \otimes (\mathbf{I}_m - \mathbf{U} \mathbf{U}^\top)$)	$\frac{1}{2} \mathbf{H}_{pGN}^* = \frac{1}{2} \mathbf{H}_{1GN}^*$ $\frac{1}{2} \mathbf{H}_p^* = \frac{1}{2} \mathbf{H}_{pGN}^* - [2\mathbf{K}_{mr}^\top \mathbf{Z}^* (\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}})^{-1} \mathbf{Z}^{*\top} \mathbf{K}_{mr}]_{FN}$ $- [\mathbf{K}_{mr}^\top \mathbf{Z}^* \tilde{\mathbf{U}}^\dagger \tilde{\mathbf{V}}^* \mathbf{P}_p + \mathbf{P}_p \tilde{\mathbf{V}}^{*\top} \tilde{\mathbf{U}}^\dagger \mathbf{Z}^{*\top} \mathbf{K}_{mr}]_{FN}$ $+ \langle \alpha \mathbf{I}_r \otimes \mathbf{U} \mathbf{U}^\top \rangle_P + \langle \lambda \mathbf{I}_{mr} \rangle$
Reduced Hessian $\mathbf{H}_r^* \in \mathbb{S}^{(mr-r^2)}$	$\frac{1}{2} \mathbf{H}_r^* := \frac{1}{2} \mathbf{P}_r \mathbf{H}_p^* \mathbf{P}_r^\top$ ($\mathbf{P}_r := \mathbf{I}_r \otimes \mathbf{U}_\perp^\top$)	$\frac{1}{2} \mathbf{H}_r^* = \frac{1}{2} \mathbf{P}_r \mathbf{H}_p^* \mathbf{P}_r^\top = \frac{1}{2} \mathbf{P}_r \mathbf{H}_1^* \mathbf{P}_r^\top = \frac{1}{2} \mathbf{P}_r \frac{d\mathbf{g}_1^*}{d\mathbf{u}} \mathbf{P}_r^\top + \langle \lambda \mathbf{I}_{mr} \rangle$ $\begin{cases} \mathbf{K}_{mr} \text{vec}(\mathbf{U}) = \text{vec}(\mathbf{U}^\top) \\ \mathbf{R}^* := \mathbf{W} \odot (\mathbf{U} \mathbf{V}^{*\top} - \mathbf{M}) \\ \mathbf{Z}^* := (\mathbf{W} \odot \mathbf{R}^*) \otimes \mathbf{I}_r \end{cases}$
Ordinary update	Ordinary update for \mathbf{U}	$\mathbf{U} \leftarrow \mathbf{U} - \text{unvec}(\mathbf{H}_1^{*-1} \mathbf{g}_1^*)$
Projected update	Grassmann manifold retraction	$\mathbf{U} \leftarrow qf(\mathbf{U} - \text{unvec}(\mathbf{H}_p^{*-1} \mathbf{g}_p^*))$
Reduced update	Grassmann manifold retraction	$\mathbf{U} \leftarrow qf(\mathbf{U} - \mathbf{P}_r^\top \text{unvec}(\mathbf{H}_r^{*-1} \mathbf{g}_r^*))$

Table 2: List of derivatives and extensions for un-regularized variable projection on low-rank matrix factorization with missing data. This includes notions to Ruhe and Wedin's algorithms.

4. Additional symmetry of the Gauss-Newton matrix for un-regularized VarPro

To observe the additional symmetry, we first need to know how the un-regularized VarPro problem can be reformulated in terms of column-wise derivatives. Below is a summary of these derivatives found in [2].

4.1. Summary of column-wise derivatives for VarPro

The unregularized objective $f_1^*(\mathbf{u}, \mathbf{v}^*(\mathbf{u}))$ can be defined as the sum of squared norm of each column in $\mathbf{R}^*(\mathbf{u}, \mathbf{v}^*(\mathbf{u}))$. i.e.

$$f_1^*(\mathbf{u}, \mathbf{v}^*(\mathbf{u})) := \sum_{j=1}^n \|\tilde{\mathbf{w}}_j(\mathbf{U}\mathbf{v}_j^* - \mathbf{m}_j)\|_2^2, \quad (20)$$

where \mathbf{v}_j^* is the j -th row of $\mathbf{V}^*(\mathbf{U})$, \mathbf{m}_j is the j -th column of \mathbf{M} and $\tilde{\mathbf{w}}_j$ is the non-zero rows of $\text{diag vec}(\mathbf{w}_j)$, where $\mathbf{w}_j \in \mathbb{R}^m$ is the j -th column of \mathbf{W} . Hence, $\tilde{\mathbf{w}}_j$ is a $p_j \times m$ weight-projection matrix, where p_j is the number of non-missing entries in column j .

Now define the followings:

$$\tilde{\mathbf{U}}_j := \tilde{\mathbf{w}}_j \mathbf{U} \quad (21)$$

$$\tilde{\mathbf{m}}_j := \tilde{\mathbf{w}}_j \mathbf{m}_j, \text{ and} \quad (22)$$

$$\boldsymbol{\epsilon}_{1j}^* = \tilde{\mathbf{U}}_j \mathbf{v}_j^* - \tilde{\mathbf{m}}_j. \quad (23)$$

Hence,

$$\mathbf{v}_j^* := \arg \min_{\mathbf{v}_j} \sum_{j=1}^n \|\tilde{\mathbf{U}}_j \mathbf{v}_j^* - \tilde{\mathbf{m}}_j\|_2^2 \quad (24)$$

$$= \arg \min_{\mathbf{v}_j} \|\tilde{\mathbf{U}}_j \mathbf{v}_j^* - \tilde{\mathbf{m}}_j\|_2^2 = \tilde{\mathbf{U}}_j^\dagger \tilde{\mathbf{m}}_j, \quad (25)$$

which shows that each row of \mathbf{V}^* is independent of other rows.

4.1.1 Some definitions

We define

$$\tilde{\mathbf{V}}_j^* := \tilde{\mathbf{w}}_j (\mathbf{v}_j^{*\top} \otimes \mathbf{I}_m) \in \mathbb{R}^{p_j \times mr}, \text{ and} \quad (26)$$

$$\mathbf{Z}_j^* := (\mathbf{w}_j \odot \mathbf{r}_j^*) \otimes \mathbf{I}_r \in \mathbb{R}^{mr \times r}. \quad (27)$$

4.1.2 Column-wise Jacobian \mathbf{J}_{1j}^*

The Jacobian matrix is also independent for each column. From [2], the Jacobian matrix for column j , \mathbf{J}_{1j}^* , is

$$\mathbf{J}_{1j}^* = (\mathbf{I} - \tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^\dagger) \tilde{\mathbf{V}}_j^* - \tilde{\mathbf{U}}_j^{\dagger\top} \mathbf{Z}_j^{*\top} \mathbf{K}_{mr} \quad (28)$$

$$= (\mathbf{I} - \tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^\dagger) \tilde{\mathbf{V}}_j^* - \tilde{\mathbf{U}}_j^\top ((\tilde{\mathbf{U}}_j^\top \tilde{\mathbf{U}}_j)^{-1} \otimes (\mathbf{w}_j \odot \mathbf{r}_j^*)^\top). \quad (29)$$

4.1.3 The column-wise form of Hessian \mathbf{H}_1^*

From [2], the column-wise form of \mathbf{H}_1^* is

$$\begin{aligned} \frac{1}{2} \mathbf{H}_1^* = & \sum_{j=1}^n \left[\tilde{\mathbf{V}}_j^{*\top} (\mathbf{I} - [\tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^\dagger]_{RW2}) \tilde{\mathbf{V}}_j^* + [-1]_{FN} \times [\mathbf{K}_{mr}^\top \mathbf{Z}_j^* (\tilde{\mathbf{U}}_j^\top \tilde{\mathbf{U}}_j)^{-1} \mathbf{Z}_j^{*\top} \mathbf{K}_{mr}]_{RW1} \right. \\ & \left. - [\tilde{\mathbf{V}}_j^{*\top} \tilde{\mathbf{U}}_j^{\dagger\top} \mathbf{Z}_j^{*\top} \mathbf{K}_{mr} + \mathbf{K}_{mr}^\top \mathbf{Z}_j^* \tilde{\mathbf{U}}_j^\dagger \tilde{\mathbf{V}}_j^*]_{FN} \right] + \langle \lambda \mathbf{I} \rangle. \end{aligned} \quad (30)$$

All $\mathbf{Z}_j^{*\top} \mathbf{K}_{mr}$ can be replaced with $\mathbf{I} \otimes (\mathbf{w}_j \odot \mathbf{r}_j^*)^\top$ using (8).

4.2. Symmetry analysis

As derived in [2], the Gauss-Newton matrix for un-regularized VarPro is identical to its projection to the tangent space of \mathbf{U} . The column-wise representation of this is

$$\frac{1}{2}\mathbf{H}_{GN1}^* = \sum_{j=1}^n \left[\tilde{\mathbf{V}}_j^{*\top} (\mathbf{I} - [\tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^\dagger]_{RW2}) \tilde{\mathbf{V}}_j^* + [\mathbf{K}_{mr}^\top \mathbf{Z}_j^* (\tilde{\mathbf{U}}_j^\top \tilde{\mathbf{U}}_j)^{-1} \mathbf{Z}_j^{*\top} \mathbf{K}_{mr}]_{RW1} \right] + \langle \lambda \mathbf{I} \rangle. \quad (31)$$

Simplifying the first term yields

$$\tilde{\mathbf{V}}_j^{*\top} \tilde{\mathbf{V}}_j^* = (\mathbf{v}_j^* \otimes \mathbf{I}) \tilde{\mathbf{W}}_j^\top \tilde{\mathbf{W}}_j (\mathbf{v}_j^{*\top} \otimes \mathbf{I}) \quad (32)$$

$$= (\mathbf{v}_j^* \otimes \tilde{\mathbf{W}}_j^\top) (\mathbf{v}_j^{*\top} \otimes \tilde{\mathbf{W}}_j) = \mathbf{v}_j^* \mathbf{v}_j^{*\top} \otimes \tilde{\mathbf{W}}_j^\top \tilde{\mathbf{W}}_j, \quad // \text{ noting (5)} \quad (33)$$

which is the Kronecker product of two symmetric matrices.

Now defining $\tilde{\mathbf{U}}_{jQ} := qf(\tilde{\mathbf{U}}_j)$ gives $\tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^\dagger = \tilde{\mathbf{U}}_{jQ} \tilde{\mathbf{U}}_{jQ}^\top$. Hence, simplifying the second term yields

$$-\tilde{\mathbf{V}}_j^{*\top} (\tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^\dagger) \tilde{\mathbf{V}}_j^* = -(\mathbf{v}_j^* \otimes \tilde{\mathbf{W}}_j^\top) (\tilde{\mathbf{U}}_{jQ} \tilde{\mathbf{U}}_{jQ}^\top) (\mathbf{v}_j^{*\top} \otimes \tilde{\mathbf{W}}_j) \quad (34)$$

$$= -\mathbf{v}_j^* \mathbf{v}_j^{*\top} \otimes \tilde{\mathbf{W}}_j^\top \tilde{\mathbf{U}}_{jQ} \tilde{\mathbf{U}}_{jQ}^\top \tilde{\mathbf{W}}_j. \quad // \text{ noting (5)} \quad (35)$$

Given that $\tilde{\mathbf{U}}_j = \tilde{\mathbf{U}}_{jQ} \tilde{\mathbf{U}}_{jR}$, we can transform $(\tilde{\mathbf{U}}_j^\top \tilde{\mathbf{U}}_j)^{-1}$ to $\tilde{\mathbf{U}}_{jR}^{-1} \tilde{\mathbf{U}}_{jR}^{-\top}$. The last term then becomes

$$\mathbf{K}_{mr}^\top \mathbf{Z}_j^* (\tilde{\mathbf{U}}_j^\top \tilde{\mathbf{U}}_j)^{-1} \mathbf{Z}_j^{*\top} \mathbf{K}_{mr} = (\mathbf{I} \otimes (\mathbf{w}_j \odot \mathbf{r}_j^*)) \tilde{\mathbf{U}}_{jR}^{-1} \tilde{\mathbf{U}}_{jR}^{-\top} (\mathbf{I} \otimes (\mathbf{w}_j \odot \mathbf{r}_j^*))^\top \quad // \text{ noting (8)} \quad (36)$$

$$= \tilde{\mathbf{U}}_{jR}^{-1} \tilde{\mathbf{U}}_{jR}^{-\top} \otimes (\mathbf{w}_j \odot \mathbf{r}_j^*) (\mathbf{w}_j \odot \mathbf{r}_j^*)^\top. \quad // \text{ noting (5)} \quad (37)$$

Hence, all the terms inside the square brackets in (31) can be represented as the Kronecker product of two symmetric matrices, which has an additional internal symmetry that can be exploited. This means that computing the Gauss-Newton matrix for un-regularized VarPro can be sped up by factor of 4 at most.

5. Computing Mean time to second success (MTSS)

As mentioned in [1], computing MTSS is based on the assumption that there exists a known best optimum on a dataset to which thousands or millions of past runs have converged.

Given that we have such dataset, we first draw a set of random starting points, usually between 20 to 100, from an isotropic Normal distribution (i.e. $\mathbf{U} = \text{randn}(m, r)$ with some predefined seed). For each sample of starting points, we run each algorithm and observe two things,

1. whether the algorithm has successfully converged to the best-known optimum, and
2. the elapsed time.

Since the samples of starting points have been drawn in a pseudo-random manner with a monotonically-increasing seed number, we can now estimate how long it would have taken each algorithm to observe two successes starting from each sample number. The average of these yields MTSS.

An example is included in the next section to demonstrate this procedures more clearly.

5.1. An example illustration

Table 3 shows some experimental data produced by an algorithm on a dataset of which the best-known optimum is 1.225. For each sample of starting points, the final cost and the elapsed time have been recorded.

We first determine which samples have yielded *success*, which is defined as convergence to the best-known optimum (1.225). Then, for each sample number, we can calculate up to which sample number we should have run to obtain two successes. e.g. when starting from sample number 5, we would have needed to run the algorithm up to sample number 7. This allows us to estimate the times required for two successes starting from each sample number (see Table 3). MTSS is simply the average of these times, $432.4/7 = 61.8$.

Sample no.	1	2	3	4	5	6	7	8	9	10
Final cost	1.523	1.225	1.647	1.225	1.52	1.225	1.225	1.647	1.225	1.774
Elapsed time (s)	23.2	15.1	24.7	19.5	25.4	16.3	15.5	21.2	17.8	21.0
Success (S)		S		S		S	S		S	
Sample numbers of next 2 successes	2 & 4	2 & 4	4 & 6	4 & 6	6 & 7	6 & 7	7 & 9	N/A	N/A	N/A
Time required for 2 successes (s)	82.5	59.3	85.9	61.2	57.2	31.8	54.5	N/A	N/A	N/A

Table 3: Above example data is used to compute MTSS of an algorithm. We assume that each sample starts from randomly-drawn initial points and that the best-known optimum for this dataset is 1.225 after millions of trials. In this case, MTSS is the average of the bottom row values which is 61.8.

6. Results

The supplementary package includes a set of csv files that were used to generate the results figures in [1]. They can be found in the `<main>/Results` directory.

6.1. Limitations

Some results are incomplete for the following reasons:

1. RTRMC [3] fails to run on FAC and Scu datasets. The algorithm has also failed on several occasions when running on other datasets. The suspected cause of such failure is the algorithm’s use of Cholesky decomposition, which is unable to decompose non-positive definite matrices.
2. The original Damped Wiberg code (TO_DW) [7] fails to run on FAC and Scu datasets. The way in which the code implements QR-decomposition requires datasets to be *trimmed* in advance so that each column in the measurement matrix has at least the number of visible entries equal to the estimated rank of the factorization model.
3. CSF [4] runs extremely slowly on UGb dataset.
4. Our Ceres-implemented algorithms are not yet able to take in sparse measurement matrices, which is required to handle large datasets such as NET.

References

- [1] Authors. Secrets of matrix factorization: Approximations, numerics and manifold optimization. 2015. 1, 2, 5, 6
- [2] Authors. Secrets of matrix factorization: Further derivations and comparisons. Technical report, further.pdf, 2015. 2, 4, 5
- [3] N. Boumal and P.-A. Absil. RTRMC: A Riemannian trust-region method for low-rank matrix completion. In *Advances in Neural Information Processing Systems 24 (NIPS 2011)*, pages 406–414. 2011. 6
- [4] P. F. Gotardo and A. M. Martinez. Computing smooth time trajectories for camera and deformable shape in structure from motion with occlusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 33(10):2051–2065, Oct 2011. 6
- [5] J. R. Magnus and H. Neudecker. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. 3rd edition, 2007. 2
- [6] T. P. Minka. Old and new matrix algebra useful for statistics. Technical report, Microsoft Research, 2000. 2
- [7] T. Okatani, T. Yoshida, and K. Deguchi. Efficient algorithm for low-rank matrix factorization with missing components and performance comparison of latest algorithms. In *Proceedings of the 2011 IEEE International Conference on Computer Vision (ICCV)*, pages 842–849, 2011. 6