

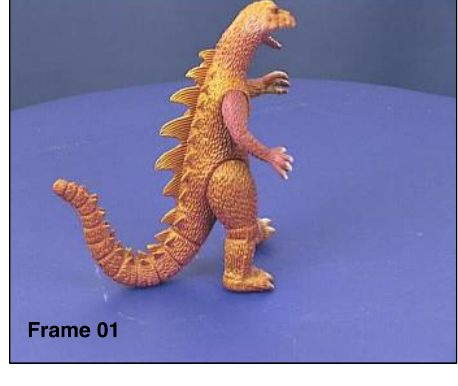
1. Matrix factorization with missing data

Given a noisy measurement matrix $M \in \mathbb{R}^{m \times n}$, where some of the values are unobserved or weighted by another matrix $W \in \mathbb{R}^{m \times n}$, we are interested in finding a set of low-rank matrices $U \in \mathbb{R}^{m \times r}$ and $V \in \mathbb{R}^{n \times r}$ that minimizes the error function

$$f(U, V) = \|\epsilon(U, V)\|_2^2 = \|W \odot (UV^T - M)\|_F^2 + \mu(\|U\|_F^2 + \|V\|_F^2) \quad (1.1)$$

where the operator \odot is Hadamard or element-wise product, and the norms are Frobenius.

- This type of low-rank models appears in rigid and non-rigid structure-from-motion (SfM), photometric stereo and recommender systems.
- For computer vision problems, the regularization parameter μ is often set to 0.



(a) Rigid SfM [3]



(b) Non-rigid SfM [3]



(c) Photometric Stereo [4]

2. Known strategies in the literature

A. Block coordinate-descent (ALS)

- Fix U and optimize over V , and vice versa.
- Closed form updates for both U and V .

inputs: M, W, r, U, μ
 $\tilde{W} \leftarrow \text{diag}(\text{vec}(W))$ with zero-rows removed.
 $\tilde{m} \leftarrow \tilde{W} \text{vec}(M)$
 $V \leftarrow \text{unvec}(\tilde{U}^{-\mu} \tilde{m})$
repeat
 $U \leftarrow \text{unvec}(\tilde{V}^{-\mu} \tilde{m})$
 $V \leftarrow \text{unvec}(\tilde{U}^{-\mu} \tilde{m})$
until convergence
outputs: U, V

- $X^{-\mu} := (X^T X + \mu I)^{-1} X^T$
(When $\mu = 0$, $X^{-\mu} = X^\dagger$.)
- $\tilde{U} := \tilde{W}(I \otimes U)$
- $\tilde{V} := \tilde{W}(V \otimes I)$
- $V^*(U) = \text{unvec}(\tilde{U}^{-\mu} \tilde{m})$

B. Joint optimization [3]

- Optimize over U and V simultaneously using a Newton-like solver.
- Bilinear structure \rightarrow sparse Hessian (or JTT) matrix.

inputs: $M, W, r, U, V, \mu, \lambda_0$
 $\lambda \leftarrow \lambda_0$
 $\tilde{W} \leftarrow \text{diag}(\text{vec}(W))$ with zero-rows removed.
 $\tilde{m} \leftarrow \tilde{W} \text{vec}(M)$
repeat
 $\mathbf{z} \leftarrow [\text{vec}(U) \ ; \ \text{vec}(V^T)]$
 Compute the gradient $\mathbf{g} := df(U, V)/d\mathbf{z}$.
 Compute the Hessian $H := d^2 f(U, V)/d\mathbf{z}^2$ or its approximation.
repeat
 $\Delta \mathbf{z} \leftarrow (H + \lambda I)^{-1} \mathbf{g}$
 Retrieve ΔU and ΔV from $\Delta \mathbf{z}$.
 $\lambda \leftarrow \lambda * 10$
until $f(U + \Delta U, V + \Delta V) < f(U, V)$
 $U \leftarrow U + \Delta U$
 $V \leftarrow V + \Delta V$
 $\lambda \leftarrow \lambda/100$
until convergence
outputs: U, V

C. Variable projection (VarPro) [6]

- Replace V by its closed form minimizer $V^*(U) := \text{argmin}_V f(U, V) = \text{unvec}(\tilde{U}^{-\mu} \tilde{m})$.
- Optimize the reduced problem over U using a Newton-like solver.
- Nonlinear structure \rightarrow denser Hessian (or JTT) matrix.

inputs: $M, W, r, U, \mu, \lambda_0$
 $\lambda \leftarrow \lambda_0$
 $\tilde{W} \leftarrow \text{diag}(\text{vec}(W))$ with zero-rows removed.
 $\tilde{m} \leftarrow \tilde{W} \text{vec}(M)$
repeat
 Compute the gradient $\mathbf{g} := df(U, V^*(U))/d\text{vec}(U)$.
 Compute the Hessian $H := d^2 f(U, V^*(U))/d\text{vec}(U)^2$ or its approximation.
repeat
 $\Delta U \leftarrow \text{unvec}((H + \lambda I)^{-1} \mathbf{g})$
 $\lambda \leftarrow \lambda * 10$
until $f(U + \Delta U, V^*(U + \Delta U)) < f(U, V^*(U))$
 $U \leftarrow U + \Delta U$
 $V \leftarrow V^*(U)$
 $\lambda \leftarrow \lambda/100$
until convergence
outputs: U, V

- $X^{-\mu} := (X^T X + \mu I)^{-1} X^T$
(When $\mu = 0$, $X^{-\mu} = X^\dagger$.)
- $\tilde{U} := \tilde{W}(I \otimes U)$
- $\tilde{V} := \tilde{W}(V \otimes I)$
- $V^*(U) = \text{unvec}(\tilde{U}^{-\mu} \tilde{m})$

3. Random restarts

- Any algorithm started from a random starting point (i.e. random U and/or V) is improved by starting from multiple points.
- To find a strong optimum, just run from many random starting points. Report optimum that is seen twice and is lower than all the others. Hence,

A new meta algorithm, RUSSO-X

RUSSO- X is literally “run algorithm- X until seen second optimum”, where this optimum is the best one observed so far. It is a simple meta-algorithm which wraps any existing low-rank matrix factorization algorithm and improves its success rate.

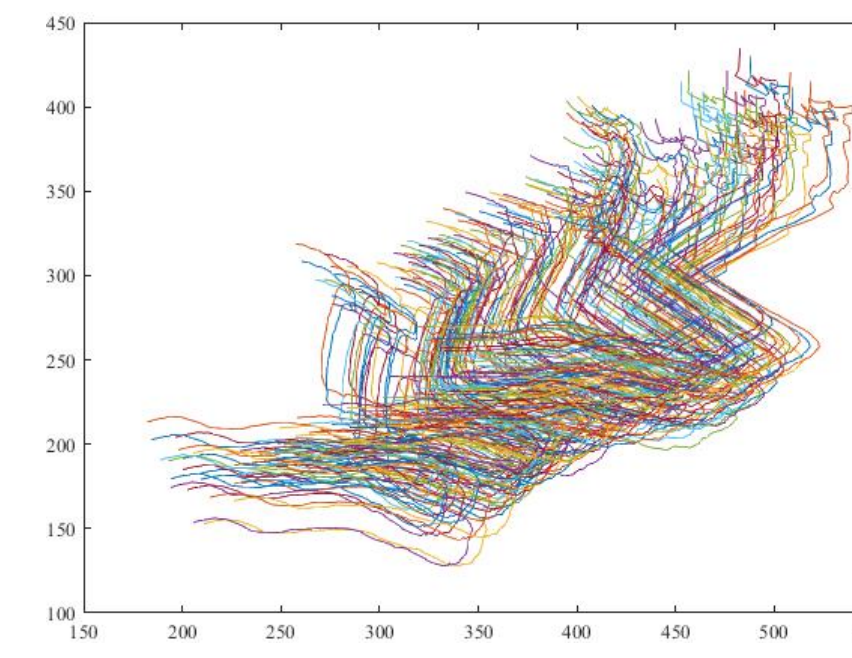
inputs: M, W, r, X, N
 $k \leftarrow 0$
 $f_{\text{best}} \leftarrow 0$
repeat
 $U \leftarrow \text{randn}(m, r)$ // m is the row size of M .
 $f_{\text{current}} \leftarrow \min_{U, V} f(U, V)$ using algorithm- X
if $|f_{\text{current}} - f_{\text{best}}| < 10^{-6}$ **do**
break // output the RUSSO- X min. value.
else if $f_{\text{current}} < f_{\text{best}}$ **do**
 $f_{\text{best}} \leftarrow f_{\text{current}}$ // update the current candidate for RUSSO- X min.
end
 $k \leftarrow k + 1$
until $k = N$
output: f_{best}

- An example of 10 runs from random starting points illustrates how RUSSO- X works in practice. (NB: RUSSO- X still returns the best optimum observed so far even when it is not observed twice.)

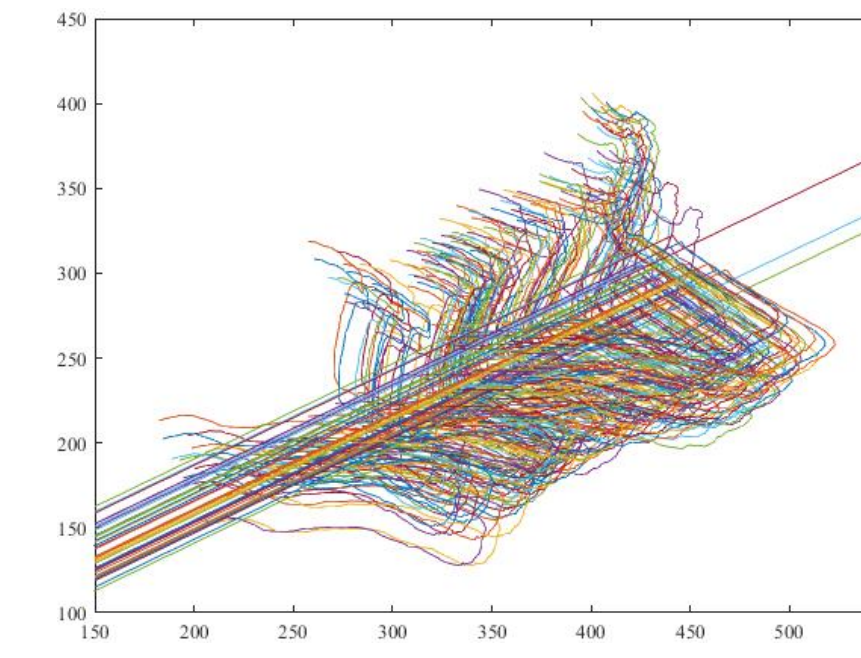
Run sequence	1	2	3	4	5	6	7	8	9	10	Success rate
Minimum	1.42	1.58	1.42	1.14	1.31	1.02	2.04	1.02	1.02	1.28	3 / 10
RUSSO- X min.	1.42	1.02	1.02	1.02	1.02	1.02	1.02	1.02	1.02	1.28	8 / 10

A side note: Why do we care about the best optimum?

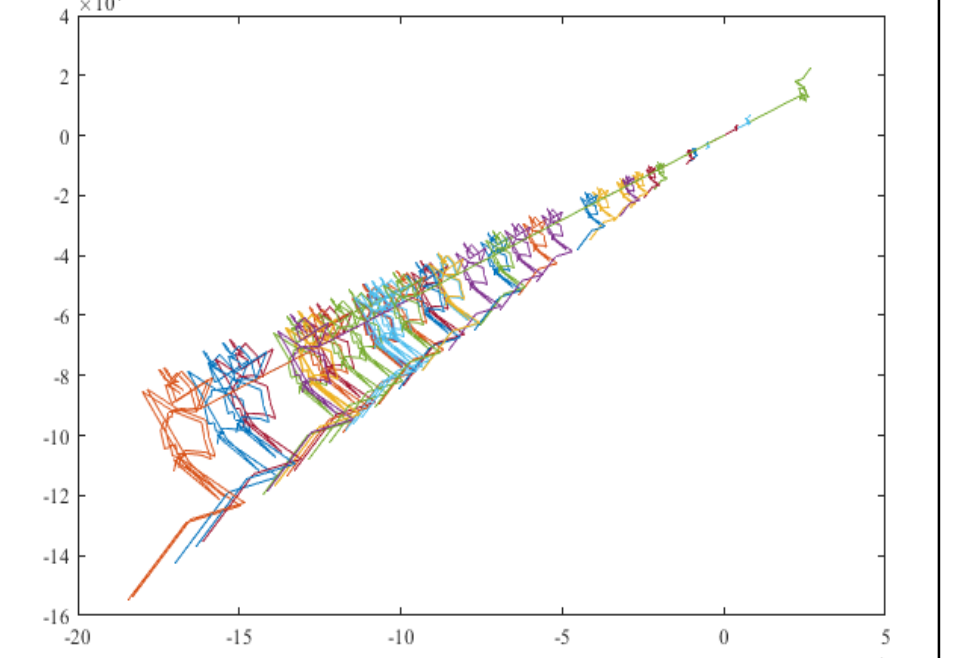
The reconstruction of point trajectories (GIR in §8) shown below illustrates that a solution with function value just .06% above the optimum can have significantly worse extrapolation properties.



(a) Best known minimum (0.3228)



(b) Second best solution (0.3230)



(c) Second best, zoomed to image

4. Approximations

- All top matrix factorization algorithms in computer vision [5, 7, 9] benefit from the use of the **unregularized** VarPro, which is illustrated in §2.C.
- These just differ in how much approximation each makes in computing the Hessian. We use the RWY convention where Y is the algorithm number mentioned in [10].

$$\frac{1}{2} H = \underbrace{\tilde{V}^*{}^T \tilde{V}^*}_{0.5 H_{RW3}} - \underbrace{\tilde{V}^*{}^T \tilde{U} \tilde{U}^\dagger \tilde{V}^*}_{G_{RW2}} - \underbrace{K_{mr}^T Z^* (\tilde{U}^\dagger \tilde{U})^{-1} Z^{*T} K_{mr}}_{G_{RW1}} + \underbrace{K_{mr}^T Z^* \tilde{U}^\dagger \tilde{V}^*}_{G_{FN}} + \underbrace{\tilde{V}^*{}^T \tilde{U}^\dagger Z^{*T} K_{mr}}_{G_{FN}^T}$$

$$\frac{1}{2} H_{RW2} = \frac{1}{2} H_{RW3} - G_{RW2}$$

$$\frac{1}{2} H_{RW1} = \frac{1}{2} H_{RW2} + G_{RW1}$$

$$\frac{1}{2} H = \frac{1}{2} H_{RW2} - G_{RW1} + G_{FN} + G_{FN}^T$$

- $K_{mr} \text{vec}(U) = \text{vec}(U^T)$
- $\tilde{U} := \tilde{W}(I \otimes U)$
- $V^*(U) = \text{unvec}(\tilde{U}^{-\mu} \tilde{m})$
- $\tilde{V}^*(U) := \tilde{W}(V^*(U) \otimes I)$
- $Z^*(U) := (W \odot W \odot (M - UV^T(U))) \otimes I$

- The most influential term is $-\tilde{V}^* \tilde{U} \tilde{U}^\dagger \tilde{V}^*$, which brings accuracy from $O(0)$ to $O(100)\%$.

5. Manifold optimization

- As noted in [3], the unregularized matrix factorization problem have an innate gauge freedom. $(UV^T = (UA^{-1})(VA^T)^T$ for any invertible matrix A .)
- For unregularized VarPro, this translates to U residing in the Grassmann manifold. i.e. U is a solution \rightarrow entire $\text{col}(U)$ is a solution space.
- To incorporate the manifold structure, we need to include the following steps in addition to the pseudocode illustrated in §2.C. (Background theory in [1].)

At each iteration,

- $\mathbf{g} \leftarrow P^T \mathbf{g}$ // project the gradient and the Hessian to the tangent space of U .
- $H \leftarrow PHP^T$ // project the Hessian to the tangent space of U .
- $U \leftarrow \text{unvec}(qf(U + \Delta U))$ // retract updated U back to the Grassmann manifold.

NB: For this problem, P can be $I \otimes U_\perp^T \in \mathbb{R}^{(m-r)r \times mr}$ or $I \otimes (I - UU^T) \in \mathbb{R}^{mr \times mr}$.

Acknowledgement

The work was supported by Microsoft and Toshiba Research Europe. We thank Roberto Cipolla, Christopher Zach, Bamdev Mishra and the anonymous reviewers.

The conference travel was supported by the British Machine Vision Association (BMVA), the Cambridge University Engineering Department (Rex Moir Fund), Christ's College (University of Cambridge) and the Cambridge Philosophical Society.

Project repository

<https://github.com/jhh37/matrix-factorization/>

Approximations, Numerics, Manifold Optimization and Random Restarts

6. Unified analysis of algorithms

A. Alternation and VarPro algorithms for the unregularized problem ($\mu \leftarrow 0$)

	ALS	VH PF [8]	TW WB [11]	TO DW [9]	DRW1	DRW1P	DRW2	DRW2P	PG CSF [7]	CH LM S GN [5]	CH LM S RW2	CH LM M [5]	CH LM M GN [5]	CH LM M RW2	NB RTRMC [2]	
1	1	1	2	2	1	1	1	1	1	1	1	1	1	1	1	1:
																2:
																3:
																4:
																5:
																6:
																7:
																8:
																9:
																10:
																11:
																12:
																13:
																14:
																15:
																16:
																17:
																18:

inputs: M, W, r, U, λ_0
 $\lambda \leftarrow \lambda_0$
 $\tilde{W} \leftarrow \text{diag}(\text{vec}(W))$ with zero-rows removed.
 $\tilde{\mathbf{m}} \leftarrow \tilde{W} \text{vec}(M)$
repeat
 $\mathbf{g} \leftarrow V^*(U)^T \text{vec}(UV^*(U)^T - M)$ // for all algorithms listed here.
 $H \leftarrow V^{*T} V^*$
 $H \leftarrow H - \tilde{V}^{*T} \tilde{U} \tilde{U}^\dagger \tilde{V}^*$
 $H \leftarrow H + K_{mr}^T Z^* (\tilde{U}^T \tilde{U})^{-1} Z^{*T} K_{mr}$
 $H \leftarrow H - K_{mr}^T Z^* (\tilde{U}^T \tilde{U})^{-1} Z^{*T} K_{mr} + K_{mr}^T Z^* \tilde{U}^\dagger \tilde{V}^* + \tilde{V}^{*T} \tilde{U}^\dagger Z^{*T} K_{mr}$
 $P \leftarrow I \otimes (I - UU^T) \in \mathbb{R}^{mr \times mr}$ // (5) is also a no-operation if 7 and 8 are.
 $P \leftarrow I \otimes U_\perp^T \in \mathbb{R}^{(m-r)r \times mr}$
 $\mathbf{g} \leftarrow P\mathbf{g}$ // (7) is a no-operation since $\mathbf{g} = P\mathbf{g}$.
 $H \leftarrow PHP^T$ // (8) is a no-operation since $H = PH$.
 $H \leftarrow H + I \otimes UU^T$ // relaxed constraint to promote $U^T \Delta U = 0$.
repeat
 $\Delta U \leftarrow \text{unvec}(H^{-1} \mathbf{g})$
 $\Delta U \leftarrow \text{unvec}((H + \lambda I)^{-1} \mathbf{g})$
 $\lambda \leftarrow \lambda * 10$
until $f(U + \Delta U, V^*(U + \Delta U)) < f(U, V^*(U))$
 $U \leftarrow U + \Delta U$
 $U \leftarrow qf(U)$ // [U,~]=qr(U,0) in MATLAB.
 $V \leftarrow \text{unvec}(\tilde{U}^\dagger \tilde{\mathbf{m}})$
 $\lambda \leftarrow \lambda / 100$
until convergence
outputs: U, V

B. Joint optimization algorithms

	CE LM	AB DN [3]	
1	1	1	inputs: M, W, r, U, V, μ , λ_0 $\lambda \leftarrow \lambda_0$ $\tilde{W} \leftarrow \text{diag}(\text{vec}(W))$ with zero-rows removed. $\tilde{\mathbf{m}} \leftarrow \tilde{W} \text{vec}(M)$ repeat $\mathbf{g} \leftarrow \begin{bmatrix} \tilde{V}^T \text{vec}(UV^T - M) + \mu \text{vec}(U) \\ \tilde{U}^T \text{vec}(UV^T - M) + \mu \text{vec}(V^T) \end{bmatrix}$ $H \leftarrow \begin{bmatrix} \tilde{V}^T \tilde{V} + \mu I & \tilde{V}^T \tilde{U} \\ \tilde{U}^T \tilde{V} & \tilde{U}^T \tilde{U} + \mu I \end{bmatrix}$ $H \leftarrow H + \begin{bmatrix} 0 & K_{mr}^T Z \\ Z^T K_{mr} & 0 \end{bmatrix}$ repeat $\Delta \mathbf{z} \leftarrow (H + \lambda I)^{-1} \mathbf{g}$ Retrieve ΔU and ΔV from $\Delta \mathbf{z}$. $\lambda \leftarrow \lambda * 10$ until $f(U + \Delta U, V + \Delta V) < f(U, V)$ $U \leftarrow U + \Delta U$ $V \leftarrow V + \Delta V$ $\lambda \leftarrow \lambda / 100$ until convergence outputs: U, V
2	2	2	

- $K_{mr} \text{vec}(U) = \text{vec}(U^T)$
- $\tilde{U} := \tilde{W}(I \otimes U)$
- $\tilde{V} := \tilde{W}(V \otimes I)$
- $Z := (W \odot W \odot (M - UV^T)) \otimes I$

7. Numerics

A. Numerical issues

- Mainly, computation of the Hessian (or its approximation) and its inversion.
- Use of the QR factorization ($\tilde{U} = \tilde{U}_Q \tilde{U}_R$) improves accuracy and performance.
- Cholesky decomposition is used for computing ΔU .

B. Profiler-guided optimization

- Standard Matlab tricks for code speedup.
- Mex files for some Kronecker product computations.

Algorithm	Successes on Din (Out of 20)		Runtime on Din (ms / iter)	
	Original	Modified	Original	Modified
CH_LM_S	4	4	380	140
CH_LM_M	1	6	369	143
CH_LM_M_GN	15	19	205	109

C. Removal of redundant computations

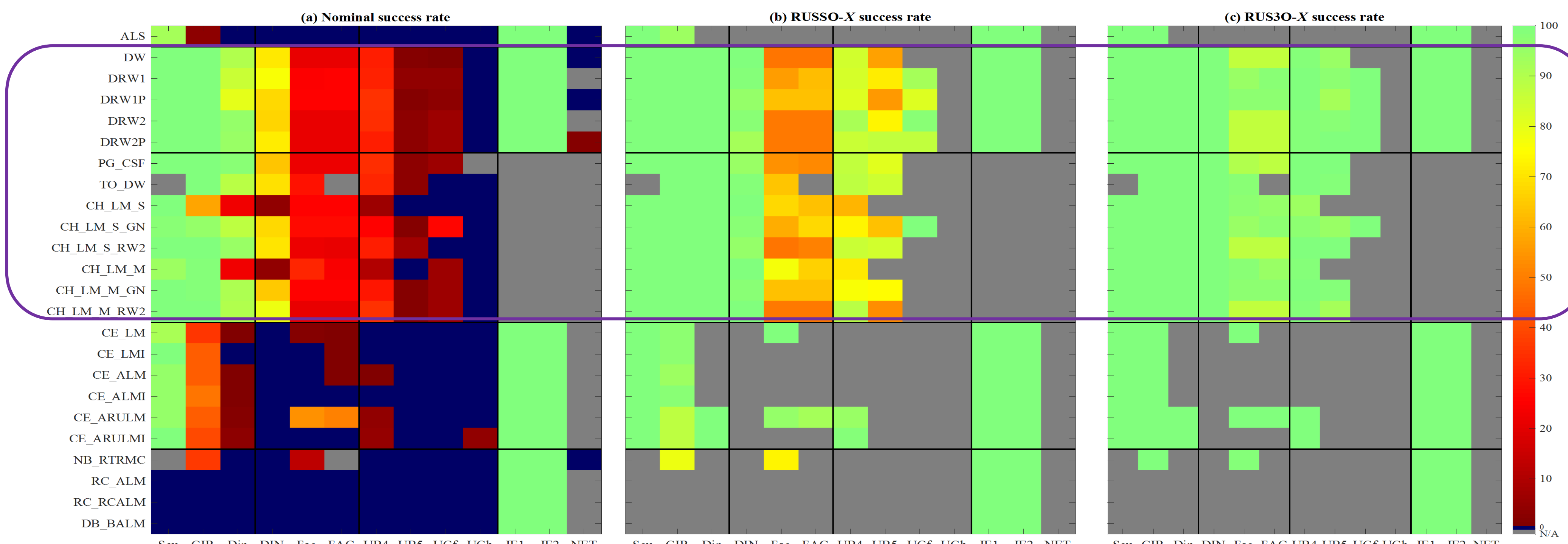
- VarPro JTJ is of the form $A \otimes B$ where A and B are symmetric. i.e. internal symmetry in addition to the ordinary JTJ symmetry. (max. 4 fold speedup.)
- \tilde{U}_Q is block-diagonal. Each sub-block \tilde{U}_{jQ} is a function of the respective column of W; $\tilde{U}_{jQ} = qf(\tilde{W}_j U)$ where $\tilde{W}_j = \text{diag}(\text{vec}(W(:, j)))$. Hence, if $W(:, i) = W(:, j)$, $\tilde{U}_{iQ} = \tilde{U}_{jQ}$. i.e. redundant QR computations can be reduced depending on the dataset. (Up to 5 times faster for DIN)

8. A list of datasets used

ID	Nature	m	n	r	Fill (%)
Scu	Photometric stereo	46	16,301	3	66.5
GIR	Non-rigid SfM	166	240	6	69.8
Din	Rigid SfM	72	319	4	23.1
DIN	Rigid SfM	72	4,983	4	9.2
Fac	Photometric stereo	20	2,596	4	64.9
FAC	Photometric stereo	20	2,944	4	58.3
UB4	Non-rigid SfM	110	1,760	4	14.4
UB5	Non-rigid SfM	110	1,760	5	14.4
UGf	Non-rigid SfM	380	4,885	6	9.1
UGb	Non-rigid SfM	380	6,310	4	6.9
JE1	Recommender	100	24,983	7	72.5
JE2	Recommender	100	23,500	7	72.7
NET	Recommender	2,000	50,000	4	2.7

9. Comparison of success rates

When run from 20 to 100 different random starting points, the VarPro-based algorithms (row 2-14) perform best in terms of success rates. RUSSO-X and RUS3O-X improves it further. (RUSkO-X stands for “run algorithm-X until seen k^{th} optimum”).

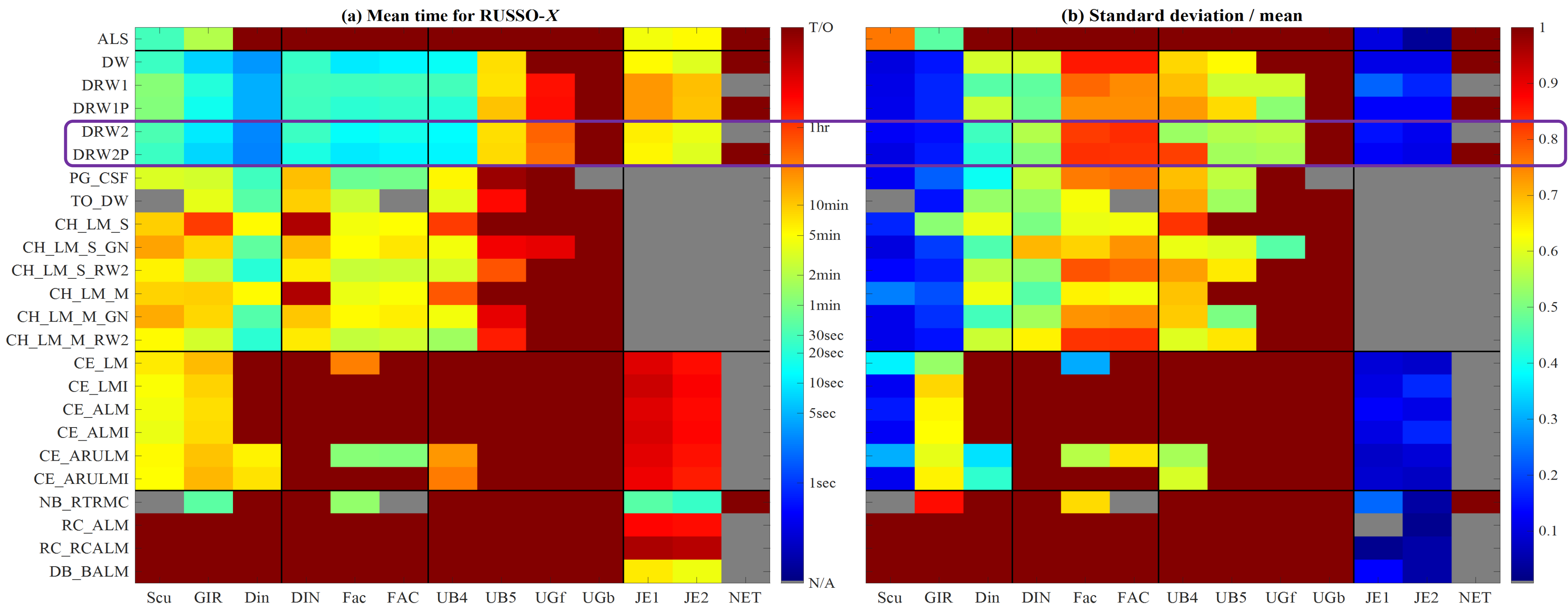


Conclusions

- Approximations:** Recent top matrix factorization algorithms, which have been derived in a multitude of different ways, can be re-derived in a unified way such that similarities and differences are more easily observed.
- Numerics:** Exploiting matrix structures and using block-QR decomposition improves the variable projection (VarPro) algorithms both in terms of speed and convergence.
- Manifold optimization:** The Riemannian manifold optimization framework can be incorporated to give a slight more improvement.
- Random restarts:** A simple meta-algorithm, RUSSO-X can be wrapped around any existing matrix factorization algorithm and improve its success rate.
- Most importantly, VarPro is key to success** for matrix factorization, everything else is detail.

10. Mean time and standard deviation for RUSSO-X

All are run in single-threaded mode. Overall, the best performing family is the “DRW” set. DRW2 and DRW2P are the winners.



References

- [1] P.-A. Absil, R. Mahony, and R. Sepulchre. Optimization Algorithms on Matrix Manifolds. Princeton University Press, 2008.
- [2] N. Boumal and P.-A. Absil. RTRMC: A Riemannian trust-region method for low-rank matrix completion. NIPS 2011.
- [3] A. M. Buchanan and A. W. Fitzgibbon. Damped Newton Algorithms for Matrix Factorization with Missing Data. CVPR, 2005.
- [4] R. Cabral, F. De la Torre, J. P. Costeira, and A. Bernardino. Unifying Nuclear Norm and Bilinear Factorization Approaches for Low-rank Matrix Decomposition. ICCV, 2013.
- [5] P. Chen. Optimization Algorithms on Subspaces: Revisiting Missing Data Problem in Low-rank Matrix. IJCV, 2008.
- [6] G. H. Golub and V. Pereyra. The Differentiation of Pseudo-Inverses and Nonlinear Least Squares Problems Whose Variables Separate. SINUM, 1973.
- [7] P. F. Gotardo and A. M. Martinez. Computing Smooth Time Trajectories for Camera and Deformable Shape in Structure from Motion with Occlusion. PAMI, 2011.
- [8] R. Vidal and R. Hartley. Motion Segmentation with Missing Data Using PowerFactorization and GPFA.
- [9] T. Okatani, T. Yoshida, and K. Deguchi. Efficient Algorithm for Low-rank Matrix Factorization with Missing Components and Performance Comparison of Latest Algorithms. ICCV, 2011.
- [10] A. Ruhe and P.A. Wedin. Algorithms for Separable Nonlinear Least Squares Problems. SIREV, 1980.
- [11] T. Wiberg. Computation of principal components when data are missing. 1976.