

# Low-rank Approximations for Incomplete Matrices

Liliya Ageeva  
Sergey Makarychev  
Aleksandr Rozhnov  
Anton Zhevnerchuk

Mentor: Maksim Rakhuba

There are some challenging industrial problems in which only incomplete data is available and the goal is to “complete” the data.

There are some challenging industrial problems in which only incomplete data is available and the goal is to “complete” the data.

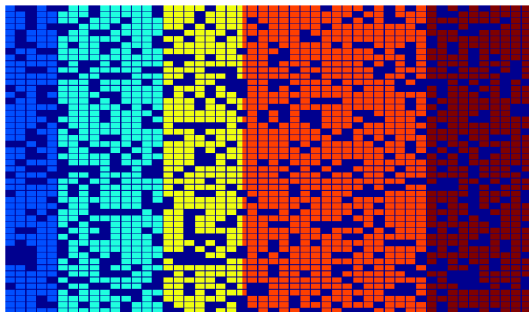
- Recommender systems

There are some challenging industrial problems in which only incomplete data is available and the goal is to “complete” the data.

- Recommender systems
- Repairing damaged files

There are some challenging industrial problems in which only incomplete data is available and the goal is to “complete” the data.

- Recommender systems
- Repairing damaged files



## General Formulation

### Input:

- Data dimensions ( $M, N$ )
- Set of known entries  $\Omega$  and values at them (sparse matrix  $X$ )

### Output:

- Low-rank approximation  $Z$

## General Formulation

### Input:

- Data dimensions  $(M, N)$
- Set of known entries  $\Omega$  and values at them (sparse matrix  $X$ )

### Output:

- Low-rank approximation  $Z$

## Approach I

Fix rank, minimize error at the known entries.

## Approach II

Set maximum acceptable error at the known entries, minimize rank.

- Alternating Least Squares (Approach I)
- Riemannian optimization [Vandereycken, 2012] (Approach I)
- Soft-Input [Mazumder, Hastie, Tibshirani, 2010] (Approach II)



$$\underset{Z}{\text{minimize}} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - Z_{ij})^2 \quad \text{s.t.} \quad \text{rank}(Z) = K$$

## Alternating Least Squares (ALS)

- Find  $Z$  in the form  $Z = U^T V$ ,  $U \in \mathbb{R}^{K \times M}$ ,  $V \in \mathbb{R}^{K \times N}$
- Update  $U$  and  $V$  independently until convergence
- At each step optimal  $U$  and  $V$  can be found analytically

$$\underset{Z}{\text{minimize}} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - Z_{ij})^2 \quad \text{s.t.} \quad \text{rank}(Z) = K$$

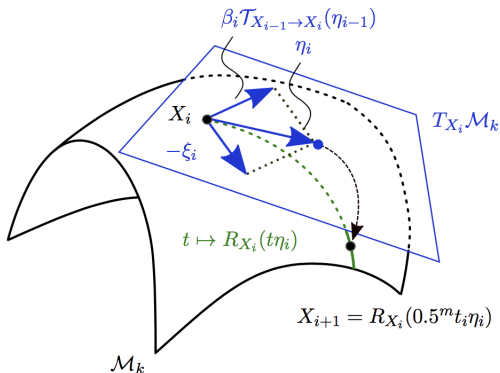
## Alternating Least Squares (ALS)

- Find  $Z$  in the form  $Z = U^T V$ ,  $U \in \mathbb{R}^{K \times M}$ ,  $V \in \mathbb{R}^{K \times N}$
- Update  $U$  and  $V$  independently until convergence
- At each step optimal  $U$  and  $V$  can be found analytically

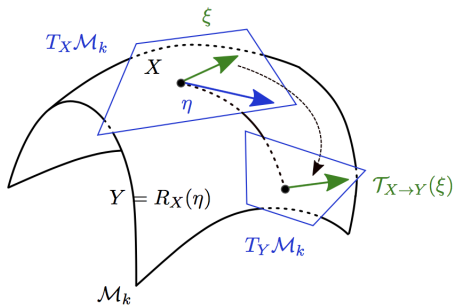
**Avoiding overfitting:** add regularization term  $\lambda(\|U\|_F^2 + \|V\|_F^2)$ , still explicit formulas for optimal  $U$  and  $V^T$  at each step.

- Matrices of fixed-rank  $k$  forms a smooth manifold of dimensionality  $(m + n - k)k$ .
- Tangent space of the same dimensionality.
- Algorithm closely resembles a typical non-linear CG algorithm with Armijo line-search for unconstrained optimization.

## Visualization of non-linear CG on a Riemannian manifold



## Vector transport on a Riemannian manifold



Initial formulation:

$$\underset{Z}{\text{minimize}} \quad \text{rank}(Z) \quad \text{s.t.} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - Z_{ij})^2 \leq \delta \quad (1)$$

Initial formulation:

$$\underset{Z}{\text{minimize}} \quad \text{rank}(Z) \quad \text{s.t.} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - Z_{ij})^2 \leq \delta \quad (1)$$

Convex relaxation of (1):

$$\underset{Z}{\text{minimize}} \quad \|Z\|_* \quad \text{s.t.} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - Z_{ij})^2 \leq \delta \quad (2)$$

Initial formulation:

$$\underset{Z}{\text{minimize}} \quad \text{rank}(Z) \quad \text{s.t.} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - Z_{ij})^2 \leq \delta \quad (1)$$

Convex relaxation of (1):

$$\underset{Z}{\text{minimize}} \quad \|Z\|_* \quad \text{s.t.} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - Z_{ij})^2 \leq \delta \quad (2)$$

Equivalent reformulation of (2):

$$\underset{Z}{\text{minimize}} \quad \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - Z_{ij})^2 + \lambda \|Z\|_* \quad (3)$$



If the full  $X$  is known and  $U\Lambda V^T$  is SVD for  $X$ , then the solution to (3) is given by

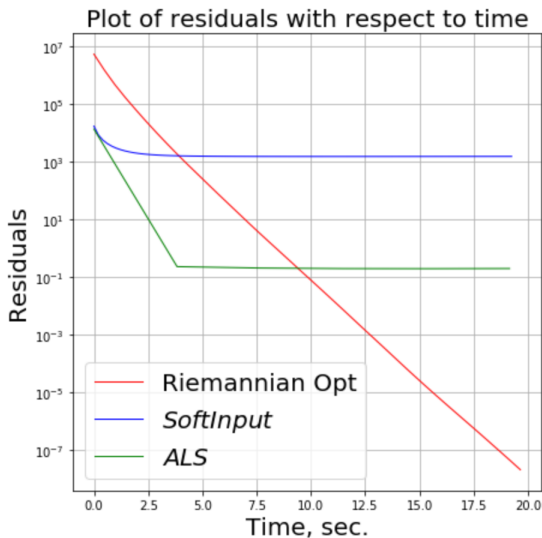
$$Z = U\Lambda_\lambda V^T, \text{ where } \Lambda_\lambda = \text{diag}\left((\sigma_1 - \lambda)_+, \dots, (\sigma_{\min\{M,N\}} - \lambda)_+\right).$$

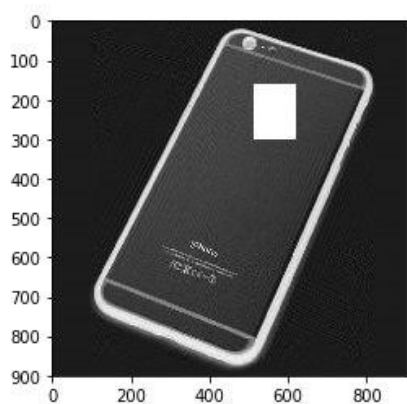
### SOFT-INPUT

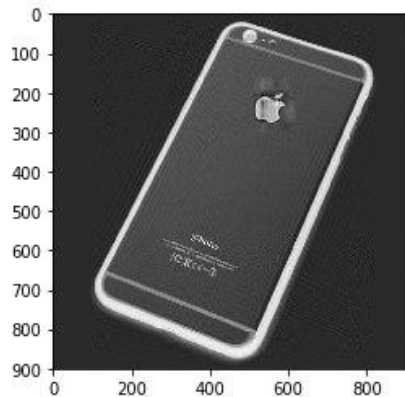
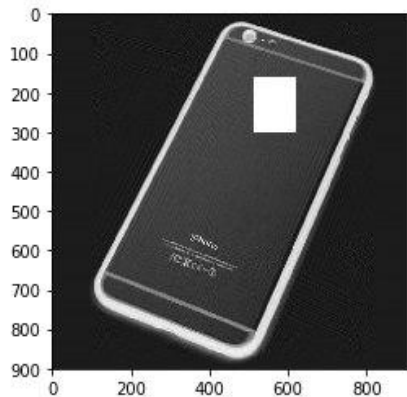
- Set  $Z_0$  to be a zero-matrix.
- At each step  $k$  approximate unknown entries of  $X$  by  $Z_{k-1}$ , set  $Z_k$  to be a solution for a problem with all values given.

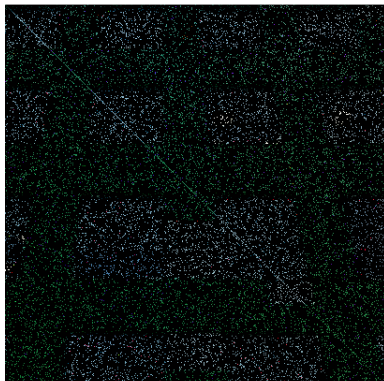
**NOTE:** Since approximated full matrix is of form  $Z_k + (X - Z_{\Omega,k})$ , MATVEC can be cheap.

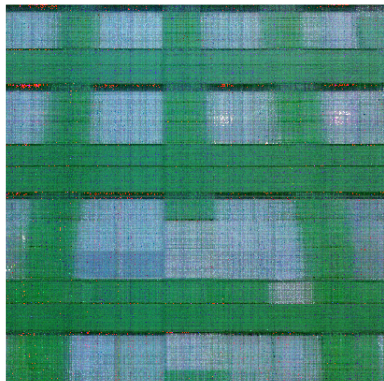
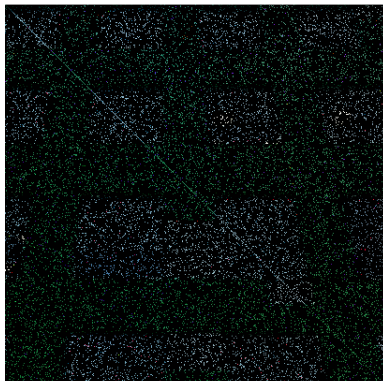
	<b>ALS</b>	<b>Riemannian</b>	<b>SOFT-INPUT</b>
Parallelizable	✓	✓	✗
Adaptive rank choice	✗	✗	✓
Flops per iteration	$O( \Omega K^2 + K^3)$	$O((M + N)K^2 +  \Omega K)$	$O((M + N +  \Omega )K^3)$

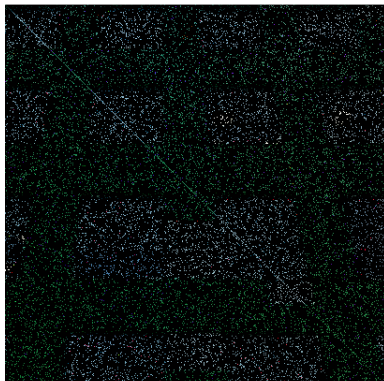




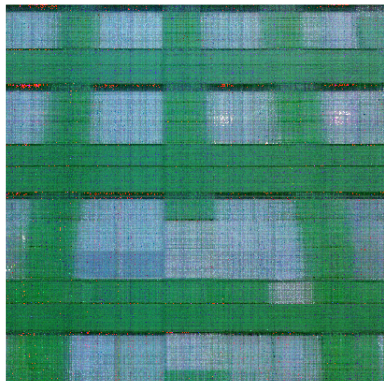








75 kB



451 kB



- Implementation of three competitive completion algorithms
- Comparative analysis of implemented algorithms
- Application to repairing damaged files