




A Comparative Study of Matrix Completion for Different Missing Data Patterns

Ningyuan Huang
Yi Li



A word cloud of data science and statistics terms arranged in a diamond shape. The terms include: SVD, MCAR, clinical-data, multiple-imputation, collaborative-filtering, matrix-completion, movie-rating, MICE, NMAR, and ZIP.

Missing data mechanism

- 
- Missing completely at random (MCAR)
 - Missing at random (MAR)
 - Not missing at random (NMAR)

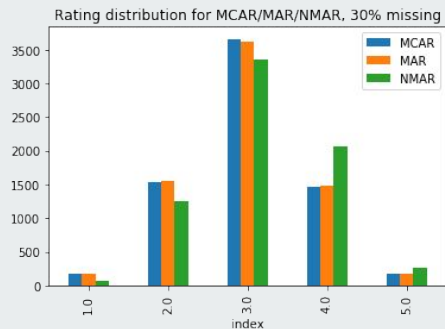
Evaluation: test error

$$\frac{\|P_{\Omega}^{\perp}(Y - M)\|_F^2}{\|P_{\Omega}^{\perp}(Y)\|^2}$$

State-of-art methods

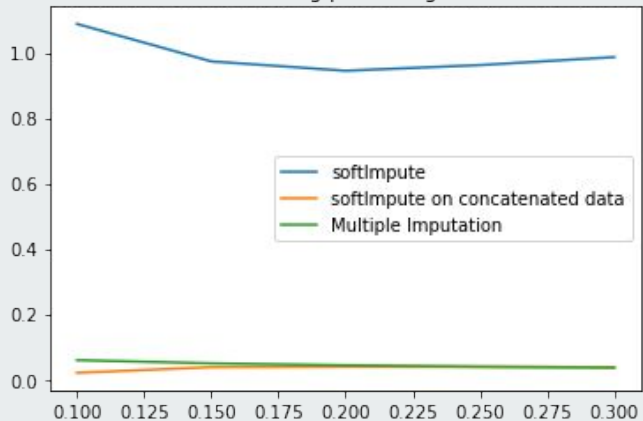
- Multiple Imputation
- Nuclear-norm minimization (SoftImpute)
 - ◆ Objective: $\min_M \frac{1}{2} \|P_{\Omega}(Y - M)\|_F^2 + \lambda \|M\|_*$
 - ◆ Algorithm: $\hat{Y} \leftarrow P_{\Omega}(Y) + P_{\Omega}^{\perp}(\hat{M})$
 $\hat{Y} = USV^T$
 $\hat{M} \leftarrow US_{\lambda}V^T$
- Implicit joint modeling (SoftImpute-concat)

Synthetic rating data

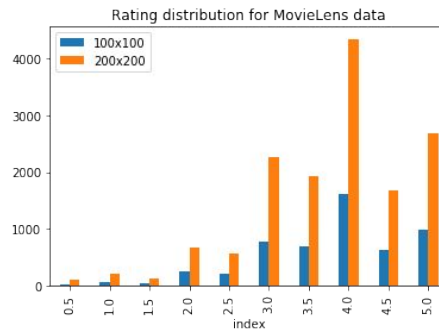


- 100 x 100
- True rank = 10
- Simulated from matrix factorization

test errors versus missing percentage, NMAR simulation

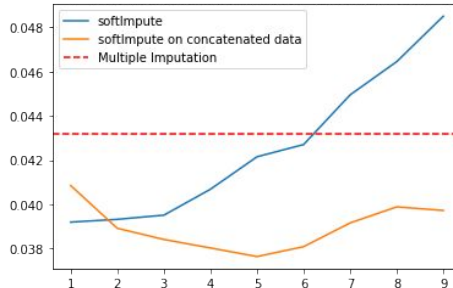


MovieLens data

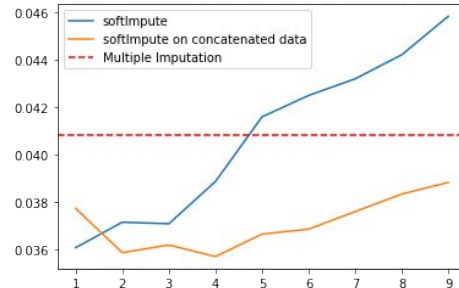


- 100 x 100, 200 x 200
- NMAR pattern (skewed to higher rating)

test errors versus rank k, MovieLens 100x100 data



test errors versus rank k, MovieLens 200x200 data



EMBARC Data

- 8-week RCT of Sertraline enrolled 287 patients with Major Depressive Disorder.
- Baseline variables were divided to 3 data sets (Clinical, EEG and fMRI) with different size and proportion of data missing.
- MCAR, MAR and NMAR generated from fully observed subsets.
- Different missing rates generated on fMRI.

TABLE I
Missing characteristics of EMBARC data

Data	Max(mean)	MCAR	MAR	NMAR
Clinical (240,32)	10%(1%)	5%	5% PLA 10% TRT	Top 5%
EEG (213,16)	21%(12%)	20%	15% PLA 25% TRT	Top 20%
fMRI (146, 208)	40%(31%)	40%	35% PLA 45% TRT	Top 40%

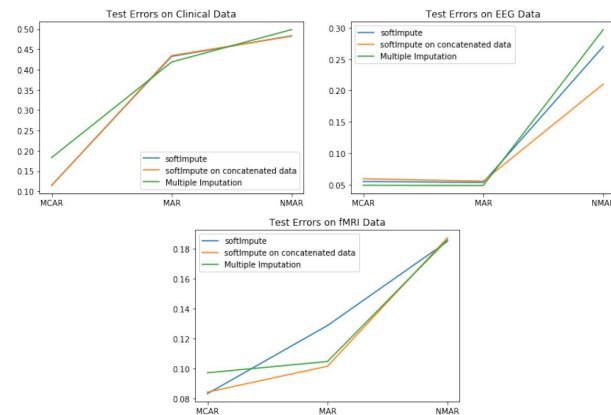


Fig. 4. Test errors of 3 methods on EMBARC data under different missing mechanisms

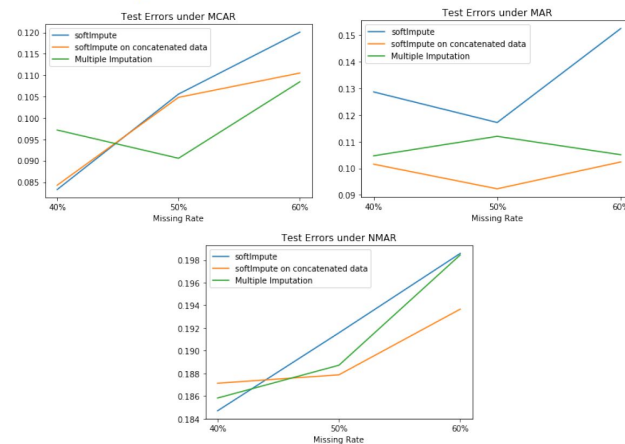


Fig. 5. Test errors of 3 methods under different missing rate generated on fMRI data.

Why SoftImpute-Concat Works?

Consider M and its SVD:

$$M = \begin{bmatrix} 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ 3 & -3 & 3 & -3 \end{bmatrix} = \begin{bmatrix} 0.5 & 0 \\ 0.5 & 0 \\ 0.5 & 0 \\ 0.5 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 8 & 0 \\ 0 & 6 \end{bmatrix} \begin{bmatrix} 0.5 & 0.5 & 0.5 & 0.5 \\ 0.5 & -0.5 & 0.5 & -0.5 \end{bmatrix}$$

Can we recover it?

$$M = \begin{bmatrix} \text{👤} & 2 & 2 & 2 \\ 2 & \text{👤} & 2 & 2 \\ 2 & 2 & \text{👤} & 2 \\ 2 & 2 & 2 & \text{👤} \\ 3 & \text{👤} & 3 & -3 \end{bmatrix}$$

M concatenated with mask

$$M_{con} = \begin{bmatrix} 2 & 2 & 2 & 2 & 1 & 0 & 0 & 0 \\ 2 & 2 & 2 & 2 & 0 & 1 & 0 & 0 \\ 2 & 2 & 2 & 2 & 0 & 0 & 1 & 0 \\ 2 & 2 & 2 & 2 & 0 & 0 & 0 & 1 \\ 3 & -3 & 3 & -3 & 0 & 1 & 0 & 0 \end{bmatrix} = UDV^T$$

$$D = \text{diag}\{8.06, 6.08, 1, 1, 0.99\}$$

$$U = \begin{bmatrix} -0.50 & 0.02 & -0.00 & 0.82 & -0.29 \\ -0.50 & -0.01 & 0.00 & 0.00 & 0.87 \\ -0.50 & 0.02 & -0.71 & -0.41 & -0.29 \\ -0.50 & 0.02 & 0.71 & -0.41 & -0.29 \\ -0.02 & -1.00 & -0.00 & 0.00 & -0.02 \end{bmatrix}$$



Conclusion

TABLE III

Matrix completion method comparison: The number of + indicates the performance, from weak (+) to strong (+++)

Methods	Low-rank	NMAR	High missing rate
<i>SoftImpute</i>	++	+	+
<i>SoftImpute-concat</i>	++	+++	++
<i>Multiple Imputation</i>	+	++	++

Future Work

- Theoretical guarantee for SoftImpute-concat
- Evaluate mixed data matrix, longitudinal data, large-sparse matrix
- Derive confidence interval for estimated missing values

