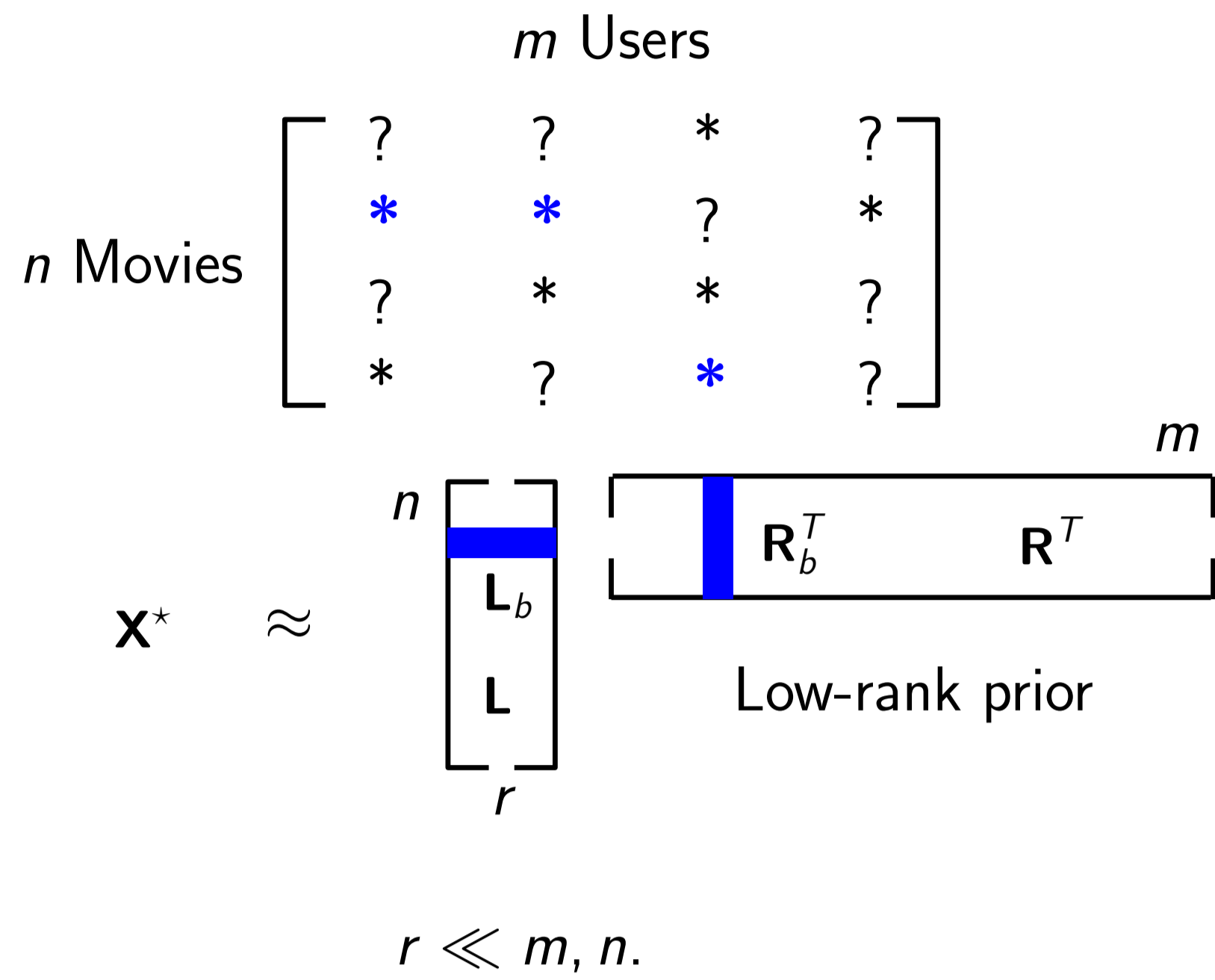


Bamdev Mishra<sup>\*</sup> and Rodolphe Sepulchre<sup>†</sup><sup>\*</sup>Amazon Development Centre India, Bangalore, Karnataka 560055, India (Bamdevm@amazon.com)<sup>†</sup>University of Cambridge, Department of Engineering, Cambridge CB2 1PZ, UK (R.Sepulchre@eng.cam.ac.uk)

## Motivation



## Formulation

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times m}} \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{X}) - \mathcal{P}_\Omega(\mathbf{X}^*)\|_F^2$$

subject to  $\text{rank}(\mathbf{X}) = r.$

- The rank constrained is parameterized as  $\mathbf{X} = \mathbf{L}\mathbf{R}^T.$

- Equivalently, the problem that is tackled by stochastic gradient descent (SGD) is

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times m}} \frac{1}{2} \sum_{ij \in \Omega} (\mathbf{L}_i \mathbf{R}_j^T - \mathbf{X}_{ij}^*)^2.$$

- $\Omega$  is the set of known entries.

## Contributions

- We propose a **scaled** variant of SGD that accelerates the standard SGD algorithm.
- We show that the proposed updates respect the **scale invariance** of the factorization model.
- The proposed updates are **scalable**.  
The computational complexity per epoch is  $O(|\Omega|(r^3/b + b_L r^2/b + b_R r^2/b + r + \log b))$ , where  $b$  is the batch size.
- Code: [www.bamdevmishra.com/codes/scaledSGD](http://www.bamdevmishra.com/codes/scaledSGD).

## Proposed scaled SGD algorithm

- Pick  $b$  known entries with their indices.
- Set up the completion subproblem by finding the indices corresponding to the submatrices  $\mathbf{L}_b$  and  $\mathbf{R}_b$ , which need to be modified. Consequently, find the subset  $\Omega_b$  of indices out of the total  $b_L b_R$  indices.
- Compute the residual  $\mathbf{S}_b = \mathcal{P}_{\Omega_b}(\mathbf{L}_b \mathbf{R}_b^T - \mathbf{X}_b^*).$
- Given a step-size  $t$ , update  $\mathbf{L}_b$  and  $\mathbf{R}_b$  as

$$\mathbf{L}_{b+} = \mathbf{L}_b - t \mathbf{S}_b \mathbf{R}_b \left( \frac{b\mu}{\max(m,n)} (\mathbf{R}^T \mathbf{R}) + (1-\mu) (\mathbf{R}_b^T \mathbf{R}_b) \right)^{-1}$$

$$\mathbf{R}_{b+} = \mathbf{R}_b - t \mathbf{S}_b^T \mathbf{L}_b \left( \frac{b\mu}{\max(m,n)} (\mathbf{L}^T \mathbf{L}) + (1-\mu) (\mathbf{L}_b^T \mathbf{L}_b) \right)^{-1}.$$

- Update  $\mathbf{L}^T \mathbf{L}$  and  $\mathbf{R}^T \mathbf{R}.$
- Repeat.

Choice of  $\mu$ 

- $\mu$  combines complementary **local** with **global** curvature information.
- $\mu = 0$  implies we use only local curvature information.
- $\mu = 1$  implies we use only global curvature information.
- In problem instances where a large number of entries are already known, i.e.,  $|\Omega|$  is large, the influence of  $\mu < 1$  is minimal. However, for ill-conditioned data, making use of local information is more critical, and a smaller value of  $\mu$  is more appropriate, e.g.,  $\mu = 0.5.$

Choice of  $b$ 

For  $b = |\Omega|$ , the proposed algorithm behaves like a batch gradient descent algorithm. The choice of  $b = 1$  is more appropriate for a fully online setting.

Scale invariance  $b$ 

The proposed updates also resolve the issue of scale invariance arising from non-uniqueness of matrix factorization. The issue refers to the behavior of algorithms which should behave *equivalently* when initialized, say, either with  $(\mathbf{L}_0, \mathbf{R}_0)$  or with  $(\mathbf{L}_0 \mathbf{M}^{-1}, \mathbf{R}_0 \mathbf{M}^T)$  for all non-singular matrices  $\mathbf{M}.$  **This is achieved by our updates.**

## Numerical comparisons

