

Problem

Approximating a matrix by decomposing it into two low rank factor matrices. Used in recommendation systems.



Challenges:

- Large matrix dims (e.g: netflix dataset: $\sim 18000 \times 480000$).
- Efficient computations needed.
- Privacy and security concerns!!

Problem Formulation

- Model as optimization problem. Objective function derived by analyzing S^{upper} and S^{lower} .
- For S^{upper} , for blocks (i, j) and (i + 1, j) convergence in Ws; for the blocks (i, j) and (i, j + 1) = >convergence in Us.
- Cost of a structure as comprising of two components: *f* and *d*.
- *f* => measures how close it is to the original matrix .
- d = measures consensus between two adjacent Us (denoted as d^U) or Ws (denoted as d^W).

For a structure pivoted at (i, j):

$$f_{ij} = \left\| \mathbf{X}_{ij} - \mathbf{U}_{ij} \mathbf{W}^{\top}_{ij} \right\|_{F}^{2}, d_{ij}^{U} = \left\| \mathbf{U}_{ij} - \mathbf{U}_{ij+1} \right\|_{F}^{2}, \text{ and } d_{ij}^{W} = \left\| \mathbf{W}_{ij} - \mathbf{W}_{i+1j} \right\|_{F}^{2},$$

Consequently, the the total cost (g) for a structure turns out to be:

$$g_{ij}^{\text{upper}} = f_{ij} + f_{i+1j} + f_{ij+1} + \rho \| \mathbf{U}_{ij} \|$$

where ρ is the *weight* factor. For S^{lower} , we can derive the costs in similar fashion. For decomposition of **X** into $p \times q$ end goal: minimize the sum of costs for all S^{upper} and S^{lower} possible, i.e.,

$$\min_{\mathbf{U}_{ij},\mathbf{W}_{ij}} \sum_{i=1,j=1}^{p,q} g_{ij}^{\text{upper}} + g_{ij}^{\text{lower}} + \lambda \|\mathbf{U}_{ij}\|_F^2 + \lambda \|\mathbf{W}_{ij}\|_F^2,$$

 λ is the regularization parameter added according to [6].

References

- [1] Genevieve Gorrell. Generalized hebbian algorithm for incremental singular value decomposition in natural language processing. In 11th Conference of the European Chapter of the Association for Computational Linguistics (EACL), 2006.
- [2] Naiyang Guan, Dacheng Tao, Zhigang Luo, and Bo Yuan. Nenmf: An optimal gradient method for nonnegative matrix factorization. *IEEE Transactions on Signal Processing* 60, 6:2882–2898, 2012. [3] Cheng-Tao Chu, Sang Kyun Kim, Yi-An Lin, Yuan Yuan Yu, Gary Bradski, Andrew Ng, and Kunle Olukotun. Map-reduce for machine learning on multicore. In Neural Information Processing Systems (*NIPS*), pages 281–288, 2006.
- [4] Quoc Le, Marc'Aurelio Ranzato, Rajat Monga, Matthieu Devin, Kai Chen, Greg Corrado, Jeff Dean, and Andrew Ng. Building high-level features using large scale unsupervised learning. In International *Conference in Machine Learning*, pages 81–88, 2012. [5] István Hegedűs, Márk Jelasity, Levente Kocsis, and András A. Benczúr. Fully distributefully distributed robust singular value decomposition. In 14-th IEEE International Conference on Peer-to-Peer
- Computing (P2P), 2014. [6] Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. Technical report, arXiv preprint arXiv:1606.04838, 2016.

A two-dimensional decomposition approach for matrix completion through gossip Mukul Bhutani and Bamdev Mishra

Previous Work

- 1. Generally treated as an optimization problem; solved using gradient search [1, 2].
- 2. Parallel versions of gradient search still require a central server [3, 4].
- 3. [5] followed $\mathbf{X} = \mathbf{U}\mathbf{W}^{\top}$ such that each row of **X** and **U** is stored in different nodes. A public matrix W is exchanged between the nodes. Random walks are done to bring convergence in W. However, here too a single agent takes care of the complete row.

$$\mathbf{U}_{ij+1} \|_{F}^{2} + \rho \| \mathbf{W}_{ij} - \mathbf{W}_{i+1j} \|_{F}^{2}$$



Algorithm 1: Basic update algorithm via SGD		RMSE on some popular datasets.						
SGD								
input : Decomposed blocks for X and rank r . output: Us, Ws.			1	Number of	Some popular datasets. aber of blocks $p \times q$ $\times 3$ 4×4 5×5 10×10 ovieLens 1 million 0.99 1.04 0.99 1.13 0.99 1.03 1.00 1.22 0.99 1.03 0.99 1.34 ovieLens 20 million 0.92 0.93 0.99 1.01 0.92 0.93 1.02 1.11 0.94 0.93 1.05 1.24 Netflix 0.98 1.14 1.02 1.02			
1 Initialize all Us and Ws.		Rank	2×2	3×3	4×4	5×5	10×10	
2 while convergence is not reached do 3 $S^{\text{struct}} = \text{randomly pick a valid structure.}$ 4 $[\text{Us}, \text{Ws}] = \text{updateThroughSGD}(\text{Xs}, S^{\text{struct}}).$ 5 Check for convergence. 6 end		5 10 15	0.87 0.86 0.86	MovieLe 0.99 0.99 0.99	ns 1 milli 1.04 1.03 1.03	on 0.99 1.00 0.99	1.13 1.22 1.34	
				MovieLe	of blocks $p \times q$ 4×4 5×5 10×10 ens 1 million 1.04 0.99 1.13 1.03 1.00 1.22 1.03 0.99 1.34 ens 20 million 0.93 0.99 1.01 0.93 0.99 1.01 0.93 1.02 1.11 0.93 1.05 1.24 Netflix 1.13 1.06 1.02 1.14 1.02 1.02			
Note: The number of times a particular structure may be selected is not equal for all and hence a normalization constant should be appropriately		5 10 15	0.95 0.96 0.96	0.92 0.93 0.94	0.93 0.93 0.93	0.99 1.02 1.05	1.01 1.11 1.24	
		5 10	1.03 1.00	N 0.98 0.98	etflix 1.13 1.14	1.06 1.02	1.02 1.02	
induption. Det paper for details.		15	1.00	1.11	1.16	1.02	1.03	



Our Decomposition Pattern



Algorithm

Experiments: synthetic datasets

Empirical proof of convergence of the algorithm. Cost as function of number of iterations.

$m \times n$ (input matrix dimensions) $p \times q$ (dimensions of decomposed grid)	$Exp#1 \\ 500 \times 500 \\ 4 \times 4$	$\begin{array}{c} \text{Exp#2} \\ 500 \times 500 \\ 4 \times 5 \end{array}$	Exp#3 500 × 500 5 × 5	$\begin{array}{c} \text{Exp#4}\\ 500 \times 500\\ 6 \times 6\end{array}$	$\begin{array}{r} \text{Exp#5}\\ 5000 \times 5000\\ 5 \times 5\end{array}$
NumIterations					
0	1.45e+05	1.45e+05	1.45e+05	1.44e+05	6.42e+05
80000	6.92e-03	1.32e-01	1.45e+00	4.74e+02	1.26e+05
160000	9.62e-06	7.65e-05	1.44e-04	9.94e-01	2.83e+02
240000	convergence	1.07e-05	1.25e-05	1.04e-02	2.85e-01
260000	0	convergence	1.21e-05	4.41e-03	7.39e-02
280000		0	convergence	1.96e-03	2.09e-02
300000				9.28e-04	6.44e-03
400000				convergence	convergence

1. X is decomposed into a $p \times q$ dimensional rectangular grid of blocks.

Each \mathbf{X}_{ij} can then be factored as $\mathbf{U}_{ij} \ \mathbf{W}^{\top}_{ij}$ as

3. Each block just gossips with its neighbors and tries to reach to a consensus.

rows => consensus in U & each column => consensus for **W**.

5. All these Us and Ws combined together to from the universal U and W.

6. This communication pattern leads to groups of blocks (S^{upper} & S^{lower}) which can be thought of as gossiping.

Experimentation: real datasets