

IBM Research Report

Large-Scale Matrix Factorization with Distributed Stochastic Gradient Descent

Rainer Gemulla

Max-Planck-Institut für Informatik
Saarbrücken, Germany

Peter J. Haas, Erik Nijkamp, Yannis Sismanis

IBM Research Division
Almaden Research Center
650 Harry Road
San Jose, CA 95120-6099
USA



Research Division

Almaden - Austin - Beijing - Cambridge - Haifa - India - T. J. Watson - Tokyo - Zurich

Large-Scale Matrix Factorization with Distributed Stochastic Gradient Descent

Rainer Gemulla¹ Peter J. Haas² Erik Nijkamp² Yannis Sismanis²

¹Max-Planck-Institut für Informatik
Saarbrücken, Germany
rgemulla@mpi-inf.mpg.de

²IBM Almaden Research Center
San Jose, CA, USA
{phaas, enijkam, syannis}@us.ibm.com

March 1, 2011

We provide a novel algorithm to approximately factor large matrices with millions of rows, millions of columns, and billions of nonzero elements. Our approach rests on stochastic gradient descent (SGD), an iterative stochastic optimization algorithm. We first develop a novel “stratified” SGD variant (SSGD) that applies to general loss-minimization problems in which the loss function can be expressed as a weighted sum of “stratum losses.” We establish sufficient conditions for convergence of SSGD using results from stochastic approximation theory and regenerative process theory. We then specialize SSGD to obtain a new matrix-factorization algorithm, called DSGD, that can be fully distributed and run on web-scale datasets using MapReduce. DSGD has good speed-up behavior and handles a wide variety of matrix factorizations. We describe the practical techniques used to optimize performance in our DSGD implementation. Experiments suggest that DSGD converges significantly faster and has better scalability properties than alternative algorithms.

Contents

1. Introduction	3
2. Example and Prior Work	4
3. Stochastic Gradient Descent	7
3.1. Preliminaries	7
3.2. SGD for Matrix Factorization	7
4. Stratified SGD	9
4.1. The SSGD Algorithm	9
4.2. Convergence of SSGD	10
4.3. Conditions for Stratum Selection	12
5. The DSGD Algorithm	16
5.1. Interchangeability	16
5.2. A Simple Case	17
5.3. The General Case	19
6. DSGD Implementation	21
7. Experiments	23
7.1. Setup	23
7.2. Relative Performance	27
7.3. Scalability	28
7.4. Selection Schemes	30
7.5. Other Loss Functions	30
8. Conclusions	31
A. MapReduce Algorithms for Matrix Factorization	33
A.1. Specialized Algorithms	33
A.2. Generic Algorithms	39
B. Parallelization Techniques for Stochastic Approximation	41
C. Example Loss Functions and Derivatives	43

1. Introduction

Low-rank matrix factorizations are effective tools for analysis of “dyadic data,” which aims at discovering and capturing the interactions between two entities [12, 15, 18, 20, 22]. Successful applications include topic detection and keyword search (where the corresponding entities are documents and terms), news personalization (users and stories), and recommendation systems (users and items). In large applications, these problems can involve matrices with millions of rows (e.g., distinct customers), millions of columns (e.g., distinct items), and billions of entries (e.g., transactions between customers and items). At such scales, distributed algorithms for matrix factorization are essential to achieving reasonable performance [12, 13, 22, 27].

In practice, exact factorization is generally neither possible nor desired, so virtually all “matrix factorization” algorithms actually produce low-rank approximations, attempting to minimize a “loss function” that measures the discrepancy between the original input matrix and product of the factors returned by the algorithm; we use the term “matrix factorization” throughout to refer to such an approximation.

With the recent advent of the MapReduce parallel processing framework, web-scale matrix factorizations have become practicable and are of increasing interest to web companies, as well as other companies and enterprises that deal with massive data. Indeed, MapReduce can be used not only to factor an input matrix, but also to efficiently construct the input matrix from massive, detailed raw data, such as customer transactions. To facilitate distributed processing, prior approaches would pick an embarrassingly parallel matrix factorization algorithm and implement it on a MapReduce cluster; the choice of algorithm was driven by the ease with which it could be distributed. In this paper, we take a different approach and start with an algorithm that is known to have good performance in non-parallel environments. Specifically, we start with stochastic gradient descent (SGD), an iterative optimization algorithm which has been shown, in a sequential setting, to be very effective for matrix factorization [18]. Although the generic SGD algorithm is not embarrassingly parallel, we can exploit the special structure of the factorization problem to obtain a version of SGD that is fully distributed and scales to extremely large matrices.

The key idea is to first develop a “stratified” variant of SGD, called SSGD, that is applicable to general loss-minimization problems in which the loss function $L(\theta)$ can be expressed as a weighted sum of “stratum losses,” so that $L(\theta) = w_1 L_1(\theta) + \dots + w_q L_q(\theta)$. At each iteration, the algorithm takes a downhill step with respect to one of the stratum losses L_s , i.e., approximately in the direction of the negative gradient $-L'_s(\theta)$. Although each such direction is “wrong” with respect to minimization of the overall loss L , we prove that, under appropriate regularity conditions, SSGD will converge to a good solution for L if the sequence of strata is chosen carefully.

We then specialize SSGD to obtain a novel distributed matrix-factorization algorithm, called DSGD. Specifically, we express the input matrix as a union of (possibly overlapping) pieces, called “strata.” For each stratum, the stratum loss is defined as the loss computed over only the data points in the stratum (and appropriately scaled). The strata are carefully chosen so that each stratum has “ d -monomial” structure, which allows SGD to be run on the stratum in a distributed manner. For example, a stratum corresponding to the nonzero entries in a block-diagonal matrix with k blocks is d -monomial for all $d \leq k$. The DSGD algorithm repeatedly selects a stratum according to the

general SSGD procedure and processes the stratum in a distributed fashion. Stratification is a technique commonly used to reduce the variance of noisy estimates [4, Sec. V.7], such as gradient estimates in SGD; here we re-purpose the stratification technique to derive a distributed factorization algorithm with provable convergence guarantees.

Our contributions are as follows:

1. We present SSGD, a novel stratified version of SGD, that is applicable to any optimization problem in which the loss function can be represented as a weighted sum of stratum losses.
2. We formally establish sufficient conditions for the convergence of SSGD using results from stochastic approximation theory and regenerative process theory.
3. We specialize SSGD to obtain DSGD, a novel distributed algorithm for low-rank matrix factorization. Both data and factors are fully distributed. DSGD has low memory requirements and scales to matrices with millions of rows, millions of columns, and billions of nonzero elements.
4. We describe practical techniques for implementing DSGD and optimizing its performance.
5. We show that DSGD is amenable to MapReduce, a popular framework for distributed processing.
6. We compare DSGD to state-of-the-art distributed algorithms for matrix factorization. Our experiments suggest that DSGD converges significantly faster, and has better scalability.

Unlike many prior algorithms, DSGD is a generic algorithm in that it can be used for a variety of different loss functions. In this paper, we focus primarily on the class of factorizations that minimize a “nonzero loss.” This class of loss functions is important for applications in which a zero represents missing data and hence should be ignored when computing loss. A typical motivation for factorization in this setting is to estimate the missing values, e.g., the rating that a customer would likely give to a previously unseen movie.

The rest of the paper is organized as follows. In Sec. 2, we introduce the factorization problem by means of an example, and discuss prior approaches to its solution. Sec. 3 describes the basic (non-parallel) SGD algorithm and its application to the matrix factorization problem. In Sec. 4 we develop the stratified variant of SGD and establish sufficient conditions for convergence. We specialize SSGD in Sec. 5 to obtain our DSGD matrix factorization algorithm. We first discuss the special “interchangeability” structure that we exploit to permit distributed execution of SGD within a stratum, and then show how to exploit this structure by means of a simple example that corresponds to processing of a single stratum by DSGD. We then give the general algorithm, which combines distributed processing within strata and careful selection of a stratum sequence. Practical implementation considerations are discussed in Sec. 6 and our empirical study of DSGD is described in Sec. 7. We conclude in Sec. 8.

2. Example and Prior Work

To gain understanding about applications of matrix factorizations, consider the “Netflix problem” [5] of recommending movies to customers. Netflix is a company that offers tens of thousands of movies for rental. The company has more than 15M customers, each of whom can provide feedback about

their personal taste by rating movies with 1 to 5 stars. The feedback can be represented in a feedback matrix such as

$$\begin{array}{c} \\ Alice \\ Bob \\ Charlie \end{array} \begin{pmatrix} Avatar & The Matrix & Up \\ ? & 4 & 2 \\ 3 & 2 & ? \\ 5 & ? & 3 \end{pmatrix}.$$

Each entry may contain additional data, e.g., the date of rating or other forms of feedback such as click history. The goal of the factorization is to predict missing entries (?); entries with a high predicted rating are then recommended to users for viewing. This is an instance of a recommender system based on matrix factorization, and has been successfully applied in practice. See [18] for an excellent discussion of the intuition behind this approach.

The traditional matrix factorization problem can be stated as follows. Given an $m \times n$ matrix \mathbf{V} and a rank r , find an $m \times r$ matrix \mathbf{W} and an $r \times n$ matrix \mathbf{H} such that $\mathbf{V} = \mathbf{W}\mathbf{H}$. As discussed previously, our actual goal is to obtain a low-rank approximation $\mathbf{V} \approx \mathbf{W}\mathbf{H}$, where the quality of the approximation is described by an application-dependent loss function L . We seek to find

$$\operatorname{argmin}_{\mathbf{W}, \mathbf{H}} L(\mathbf{V}, \mathbf{W}, \mathbf{H}),$$

i.e., the choice of \mathbf{W} and \mathbf{H} that give rise to the smallest loss. For example, assuming that missing ratings are coded with the value 0, loss functions for recommender systems are often based on the *nonzero squared loss*

$$L_{\text{NZSL}} = \sum_{i,j: \mathbf{V}_{ij} \neq 0} (\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})^2 \quad (1)$$

and usually incorporate regularization terms, user and movie biases, time drifts, and implicit feedback.

In the following, we restrict attention to loss functions that, like L_{NZSL} , can be decomposed into a sum of *local losses* over (a subset of) the entries in \mathbf{V}_{ij} . I.e., we require that the loss can be written as

$$L = \sum_{(i,j) \in Z} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j}) \quad (2)$$

for some *training set* $Z \subseteq \{1, 2, \dots, m\} \times \{1, 2, \dots, n\}$ and *local loss function* l , where \mathbf{A}_{i*} and \mathbf{A}_{*j} denote row i and column j of matrix \mathbf{A} , respectively. Many loss functions used in practice—such as squared loss, generalized Kullback-Leibler divergence (GKL), and L_p regularization—can be decomposed in such a manner [25]; see Appendix C for more examples. Note that a given loss function L can potentially be decomposed in multiple ways. In this paper, we focus primarily on the class of *nonzero decompositions*, in which $Z = \{(i, j) : \mathbf{V}_{ij} \neq 0\}$ refers to the nonzero entries in \mathbf{V} . As mentioned above, such decompositions naturally arise when zeros represent missing data. Our algorithms can handle other decompositions as well; see our preliminary results for GKL in Sec. 7. To avoid trivialities, we assume throughout that there is at least one training point in every

row and in every column of \mathbf{V} ; e.g., every customer has rated at least one movie and every movie has been rated at least once.¹

To compute \mathbf{W} and \mathbf{H} on MapReduce, all known algorithms start with some initial factors \mathbf{W}_0 and \mathbf{H}_0 and iteratively improve them. The $m \times n$ input matrix \mathbf{V} is partitioned into $d_1 \times d_2$ blocks, which are distributed in the MapReduce cluster. Both row and column factors are blocked conformingly:

$$\begin{matrix} & \mathbf{H}^1 & \mathbf{H}^2 & \dots & \mathbf{H}^{d_2} \\ \mathbf{W}^1 & \left(\mathbf{V}^{11} & \mathbf{V}^{12} & \dots & \mathbf{V}^{1d_2} \right) \\ \mathbf{W}^2 & \left(\mathbf{V}^{21} & \mathbf{V}^{22} & \dots & \mathbf{V}^{2d_2} \right) \\ \vdots & \left(\vdots & \vdots & \ddots & \vdots \right) \\ \mathbf{W}^{d_1} & \left(\mathbf{V}^{d_11} & \mathbf{V}^{d_12} & \dots & \mathbf{V}^{d_1d_2} \right) \end{matrix},$$

where we use superscripts to refer to individual blocks. The algorithms are designed such that each block \mathbf{V}^{ij} can be processed independently in the map phase, taking only the corresponding blocks of factors \mathbf{W}^i and \mathbf{H}^j as input. Some algorithms directly update the factors in the map phase (then either $d_1 = m$ or $d_2 = n$ to avoid overlap), whereas others aggregate the factor updates in a reduce phase.

Existing algorithms can be classified into *specialized algorithms*, which are designed for a particular loss, and *generic algorithms*, which work for a wide variety of loss functions. Specialized algorithms currently exist for only a small class of loss functions. For GKL loss, Das et al. [12] provide an EM-based algorithm, and Liu et al. [22] provide a multiplicative-update method. In [22], the latter MULT approach is also applied to squared loss and nonnegative matrix factorization with an “exponential” loss function (exponential NMF). Each of these algorithms in essence takes an embarrassingly parallel matrix factorization algorithm developed previously—in [15, 16] for the EM algorithm and in [20, 21] for the MULT methods—and directly distributes it across the MapReduce cluster. Zhou et al. [27] show how to distribute the well-known alternating least squares (ALS) algorithm to handle factorization problems with a nonzero squared loss function and an optional weighted L_2 regularization term. Their approach requires a double-partitioning of \mathbf{V} : once by row and once by column. Moreover, ALS requires that each of the factor matrices \mathbf{W} and \mathbf{H} can (alternately) fit in main memory. More details on each of the foregoing algorithms can be found in Appendix A.

Generic algorithms are able to handle all differentiable loss functions that decompose into summation form. A simple approach is distributed gradient descent (DGD [13, 14, 23]), which distributes gradient computation across a compute cluster, and then performs centralized parameter updates using, for example, quasi-Newton methods such as L-BFGS-B [8]. Partitioned SGD approaches make use of a similar idea: SGD is run independently and in parallel on partitions of the dataset, and parameters are averaged after each pass over the data (PSGD [14, 24]) or once at the end (ISGD [23, 24, 28]). These approaches have not been applied to matrix factorization before. Similarly to L-BFGS-B, they exhibit slow convergence in practice (see Sec. 7) and need to

¹Clearly, recommendation is impossible for a customer who has never rated a movie or a movie that has never been rated; mathematically, the \mathbf{W} factors for an empty row or the \mathbf{H} factors for an empty column can be set to any arbitrary value without affecting the loss, so the factorization problem is not well posed in this case.

store the full factor matrices in memory. This latter limitation can be a serious drawback: for large factorization problems, it is crucial that both matrix and factors be distributed. Our present work on DSGD is a first step towards such a *fully distributed* generic algorithm with good convergence properties.

3. Stochastic Gradient Descent

In this section, we discuss how to factorize a given matrix via standard (non-parallel) SGD.

3.1. Preliminaries

The goal of SGD is to find the value $\theta^* \in \mathfrak{R}^k$ ($k \geq 1$) that minimizes a given loss $L(\theta)$. The algorithm makes use of noisy observations $\hat{L}'(\theta)$ of $L'(\theta)$, the function’s gradient with respect to θ . Starting with some initial value θ_0 , SGD refines the parameter value by iterating the stochastic difference equation

$$\theta_{n+1} = \theta_n - \epsilon_n \hat{L}'(\theta_n), \quad (3)$$

where n denotes the step number and $\{\epsilon_n\}$ is a sequence of decreasing step sizes. (We assume throughout that each ϵ_n is nonnegative and finite.) Since $-L'(\theta_n)$ is the direction of steepest descent, (3) constitutes a noisy version of gradient descent. Figure 1 illustrates this process with an example in which θ is 2-dimensional.

Stochastic approximation theory can be used to show that, under certain regularity conditions [19], the noise in the gradient estimates “averages out” and SGD converges to the set of stationary points satisfying $L'(\theta) = 0$. Of course, these stationary points can be minima, maxima, or saddle points. One may argue that convergence to a maximum or saddle point is unlikely because the noise in the gradient estimates reduces the likelihood of getting stuck at such a point. Thus $\{\theta_n\}$ typically converges to a (local) minimum of L . A variety of methods can be used to increase the likelihood of finding a global minimum, e.g., running SGD multiple times, starting from a set of randomly chosen initial solutions.

In practice, one often makes use of an additional projection Π_H that keeps the iterate in a given constraint set H . For example, there is considerable interest in nonnegative matrix factorizations [20], which corresponds to setting $H = \{\theta : \theta \geq 0\}$. The projected algorithm takes form

$$\theta_{n+1} = \Pi_H[\theta_n - \epsilon_n \hat{L}'(\theta_n)]. \quad (4)$$

In addition to the set of stationary points, the projected process may converge to the set of “chain recurrent” points [19], which are influenced by the boundary of the constraint set H .

3.2. SGD for Matrix Factorization

To apply SGD to matrix factorization, we take θ to be (\mathbf{W}, \mathbf{H}) and decompose the loss L as in (2) for an appropriate training set Z and local loss function l . For brevity, we suppress the constant matrix \mathbf{V} in our notation. Denote by $L_z(\mathbf{W}, \mathbf{H}) = L_{ij}(\mathbf{W}, \mathbf{H}) = l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$ the

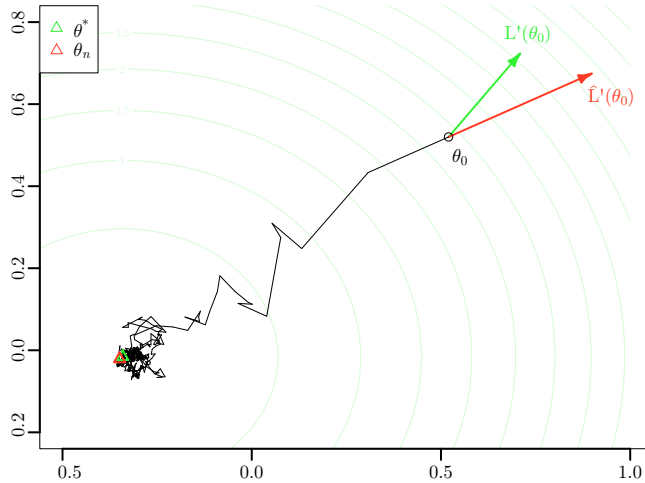


Figure 1: Example of SGD

local loss at position $z = (i, j)$. Then $L(\mathbf{W}, \mathbf{H}) = \sum_{z \in Z} L_z(\mathbf{W}, \mathbf{H})$ and hence $L'(\mathbf{W}, \mathbf{H}) = \sum_{z \in Z} L'_z(\mathbf{W}, \mathbf{H})$ by the sum rule for differentiation. DGD methods exploit the summation form of L' : they compute the local gradients L'_z in parallel and sum up. In contrast, SGD obtains noisy gradient estimates by scaling up *just one* of the local gradients, i.e.,

$$\hat{L}'(\mathbf{W}, \mathbf{H}) = N L'_z(\mathbf{W}, \mathbf{H}),$$

where $N = |Z|$ and the training point z is chosen randomly from the training set. Algorithm 1 uses SGD to perform matrix factorization.

Algorithm 1 SGD for Matrix Factorization

Require: A training set Z , initial values \mathbf{W}_0 and \mathbf{H}_0

```

while not converged do /* step */
  Select a training point  $(i, j) \in Z$  uniformly at random.
   $\mathbf{W}'_{i*} \leftarrow \mathbf{W}_{i*} - \epsilon_n N \frac{\partial}{\partial \mathbf{W}_{i*}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$ 
   $\mathbf{H}'_{*j} \leftarrow \mathbf{H}_{*j} - \epsilon_n N \frac{\partial}{\partial \mathbf{H}_{*j}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$ 
   $\mathbf{W}_{i*} \leftarrow \mathbf{W}'_{i*}$ 
end while

```

Note that, after generating a random training point $(i, j) \in Z$, we need to update only \mathbf{W}_{i*} and \mathbf{H}_{*j} , and do not need to update factors of the form $\mathbf{W}_{i'*}$ for $i' \neq i$ or $\mathbf{H}_{*j'}$ for $j' \neq j$. This computational savings follows from our representation of the global loss as a sum of local losses.

Specifically, we have used the fact that

$$\frac{\partial}{\partial \mathbf{W}_{i'k}} L_{ij}(\mathbf{W}, \mathbf{H}) = \begin{cases} 0 & \text{if } i \neq i' \\ \frac{\partial}{\partial \mathbf{W}_{ik}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j}) & \text{otherwise} \end{cases} \quad (5)$$

and

$$\frac{\partial}{\partial \mathbf{H}_{kj'}} L_{ij}(\mathbf{W}, \mathbf{H}) = \begin{cases} 0 & \text{if } j \neq j' \\ \frac{\partial}{\partial \mathbf{H}_{kj}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j}) & \text{otherwise} \end{cases} \quad (6)$$

for $1 \leq k \leq r$. SGD is sometimes referred to as *online learning* or *sequential gradient descent* [6]. Batched versions, in which multiple local losses are averaged, are also feasible but often have inferior performance in practice.

One might wonder why replacing exact gradients (GD) by noisy estimates (SGD) can be beneficial. The main reason is that exact gradient computation is costly, whereas noisy estimates are quick and easy to obtain. In a given amount of time, we can perform many quick-and-dirty SGD updates instead of a few, carefully planned GD steps. The noisy process also helps in escaping local minima (especially those with a small basin of attraction and more so in the beginning, when the step sizes are large). Moreover, SGD is able to exploit repetition within the data. Parameter updates based on data from a certain row or column will also decrease the loss in similar rows and columns. Thus the more similarity there is, the better SGD performs. Ultimately, the hope is that the increased number of steps leads to faster convergence. This behavior can be proven for some problems [7], and it has been observed in the case of large-scale matrix factorization [18].

4. Stratified SGD

In this section we develop a general stratified stochastic gradient descent (SSGD) algorithm, and give sufficient conditions for convergence. In Sec. 5 we specialize SSGD to obtain an efficient distributed algorithm (DSGD) for matrix factorization.

4.1. The SSGD Algorithm

In SSGD, the loss function $L(\theta)$ is decomposed into a weighted sum of loss functions $L_s(\theta)$ as follows:

$$L(\theta) = w_1 L_1(\theta) + w_2 L_2(\theta) + \dots + w_q L_q(\theta), \quad (7)$$

where we assume without loss of generality that $0 < w_s \leq 1$ and $\sum w_s = 1$. We refer to index s as a *stratum*, L_s as the *stratum loss* for stratum s , and w_s as the weight of stratum s . In practice, a stratum often corresponds to a part or partition of some underlying dataset. In this case, one can think of L_s as the loss incurred on the respective partition; the overall loss is obtained by summing up the per-partition losses. In general, however, the decomposition of L can be arbitrary; there may or may not be an underlying data partitioning. Also note that there is some freedom in the choice of

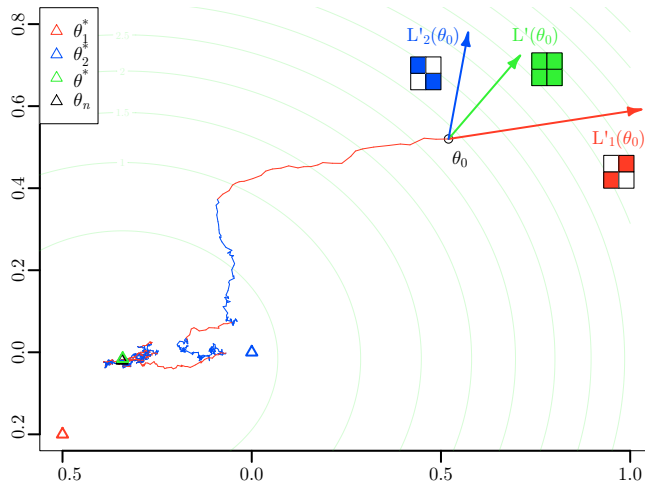


Figure 2: Example of stratified SGD

the w_s ; they may be altered to arbitrary values (subject to the constraints above) by appropriately modifying the stratum loss functions. This freedom gives room for optimization.

SSGD runs standard stochastic gradient descent on a single stratum at a time, but switches strata in a way that guarantees correctness. The algorithm can be described as follows. Suppose that there is a (potentially random) *stratum sequence* $\{\gamma_n\}$, where each γ_n takes values in $\{1, \dots, q\}$ and determines the stratum to use in the n th iteration. Using a noisy observation \hat{L}'_{γ_n} of the gradient L'_{γ_n} , we obtain the update rule

$$\theta_{n+1} = \Pi_H[\theta_n - \epsilon_n \hat{L}'_{\gamma_n}(\theta_n)]. \quad (8)$$

The sequence $\{\gamma_n\}$ has to be chosen carefully to establish convergence to the stationary (or chain-recurrent) points of L . Indeed, because each step of the algorithm proceeds approximately in the “wrong” direction, i.e., $-L'_{\gamma_n}(\theta_n)$ rather than $-L'(\theta_n)$, it is not obvious that the algorithm will converge at all. We show in Sec. 4.2 and 4.3, however, that SSGD will indeed converge under appropriate regularity conditions provided that, in essence, the “time” spent on each stratum is proportional to its weight.

Figure 2 shows an example of SSGD, where the loss (shown in green) has been decomposed into two strata (blue and red). When $\gamma_n = \text{red}$, the iterate is “pulled” towards the red optimum; when $\gamma_n = \text{blue}$, it moves towards the blue optimum. In the long run, the process converges to the overall optimum.

4.2. Convergence of SSGD

Appropriate sufficient conditions for the convergence of SSGD can be obtained from general results on stochastic approximation in Kushner and Yin [19, Sec. 5.6]. These conditions are satisfied in most

matrix factorization problems. We distinguish step-size conditions, loss conditions, stratification conditions, and stratum-sequence conditions.

The two step-size conditions involve the sequence $\{\epsilon_n\}$.

Condition 1. *The step sizes slowly approach zero in that $\epsilon_n \rightarrow 0$ and $\sum \epsilon_n = \infty$.*

Condition 2. *The step sizes decrease “quickly enough” in that $\sum \epsilon_i^2 < \infty$*

Clearly, the step sizes must decrease to 0 in order for the algorithm to converge (as specified in Condition 1). However, this convergence must occur at the correct rate. Condition 1 ensures that the SGD algorithm can move across arbitrarily large distances and thus cannot get stuck halfway to a stationary point; Condition 2 (“square summability”) ensures that the step sizes decrease to 0 fast enough so that convergence occurs. The simplest valid choice is $\epsilon_n = 1/n$.

The next pair of conditions involve the loss function.

Condition 3. *The constraint set H on which L is defined is a hyperrectangle.*

Condition 4. *$L'(\theta)$ is continuous on H .*

Note that non-differentiable points may arise in matrix factorization, e.g., when L_1 regularization is used. SSGD is powerful enough to deal with such points under appropriate regularity conditions, using “subgradients”; see [19] for details.

With respect to stratification, we require that estimates of the gradient defined with respect to a given stratum are unbiased, have bounded second moment for $\theta \in H$, and do not depend on the past. DSGD satisfies these conditions by design. More precisely, the conditions are as given below. Denote by $\hat{L}'_{\gamma_n, n}(\theta_n)$ the gradient estimate used in the n th step.

Condition 5. *The gradient estimates have bounded second moment, i.e.,*

$$\sup_n \mathbb{E} [|\hat{L}'_{\gamma_n, n}(\theta_n)|^2] < \infty$$

for all $\theta \in H$.

This condition ensures that the noise in the gradient estimates is small enough to eventually be averaged out, and is often fulfilled by choosing the constraint set H appropriately (so that $L'_s(\theta)$ is bounded for all $s \in \{1, \dots, q\}$ and $\theta \in H$).

Condition 6. *The noise in the gradient estimates is a martingale difference, i.e.,*

$$\mathbb{E} [\hat{L}'_{\gamma_n, n}(\theta_n) \mid \mathcal{F}_n] = L'_{\gamma_n}(\theta_n),$$

where $\hat{L}'_{\gamma_n, n}$ is the gradient estimate in the n th step and $\mathcal{F}_n = \sigma(\{\epsilon_i, \theta_i, \gamma_i, i \leq n\})$ is the σ -field that represents what is known at step n .

Thus, accounting for the entire history, the gradient estimate is required to be unbiased for the gradient.

Finally, we give a sufficient condition on the stratum sequence.

Condition 7. The step sizes satisfy $(\epsilon_n - \epsilon_{n+1})/\epsilon_n = O(\epsilon_n)$ and the γ_n are chosen such that the directions “average out correctly” in the sense that, for any $\theta \in H$,

$$\lim_{n \rightarrow \infty} \epsilon_n \sum_{i=0}^{n-1} [L'_{\gamma_i}(\theta) - L'(\theta)] = 0$$

almost surely.

For example, if ϵ_n were equal to $1/n$, then the n th term would represent the empirical average deviation from the true gradient over the first n steps.

We can now state our correctness result, which asserts that, under the foregoing conditions, the θ_n sequence converges almost surely to the set of limit points of an ODE that is a smoothed version of the basic SGD recursion. As shown in [19], these limit points comprise the set of stationary points of L in H , as well as a set of chain-recurrent points on the boundary of H . In our setting, the limit point to which SSGD converges is typically a good local minimum.

Theorem 1. Suppose that Conditions 1–7 hold. Then the sequence $\{\theta_n\}$ converges almost surely to the set of limit points of the projected ODE

$$\dot{\theta} = -L'(\theta) + z$$

in H , taken over all initial conditions. Here, z is the “minimum force” to keep the solution in H [19, Sec. 4.3].

The conditions used in Theorem 1 can be weakened considerably, but suffice for our purposes. The theorem follows directly from results in [19]. Indeed, as a special case of [19, Th. 6.1.1], the desired result follows from Conditions 1, 3, 4, and 5, and two “asymptotic rate of change” (ARC) assumptions: one on the gradient-estimation noise, given by (6.1.4) in [19], and one on the differences $\{L'_{\gamma_n}(\theta) - L'(\theta)\}_{n=0}^{\infty}$, given by (A6.1.3) in [19]. As discussed on pp. 137–138 of [19], the first ARC assumption is implied by Conditions 2, 5, and 6. The second ARC assumption is implied by Condition 7; see [19, p. 171].

All but the last of the sufficient conditions for convergence hold by design. Therefore, the crux of showing that SSGD converges is showing that Condition 7 holds. We address this issue next.

4.3. Conditions for Stratum Selection

The following result gives sufficient conditions on $L(\theta)$, the step sizes $\{\epsilon_n\}$, and the stratum sequence $\{\gamma_n\}$ such that Condition 7 holds. Our key assumption is that the sequence $\{\gamma_n\}$ is *regenerative* [3, Ch. VI], in that there exists an increasing sequence of almost-surely finite random indices $0 = \beta(0) < \beta(1) < \beta(2) < \dots$ that serves to decompose $\{\gamma_n\}$ into consecutive, independent and identically distributed (i.i.d.) cycles $\{C_k\}$,² with $C_k = \{\gamma_{\beta(k-1)}, \gamma_{\beta(k-1)+1}, \dots, \gamma_{\beta(k)-1}\}$.

²The cycles need not directly correspond to strata. Indeed, we make use of strategies in which a cycle comprises multiple strata.

$\dots, \gamma_{\beta(k)-1}$ } for $k \geq 1$. I.e., at each $\beta(i)$, the stratum is selected according to a probability distribution that is independent of past selections, and the future sequence of selections after step $\beta(i)$ looks probabilistically identical to the sequence of selections after step $\beta(0)$. The length τ_k of the k th cycle is given by $\tau_k = \beta(k) - \beta(k-1)$. Letting $I_{\gamma_n=s}$ be the indicator variable for the event that stratum s is chosen in the n th step, set

$$X_k(s) = \sum_{n=\beta(k-1)}^{\beta(k)-1} (I_{\gamma_n=s} - w_s)$$

for $1 \leq s \leq q$. It follows from the regenerative property that the pairs $\{(X_k(s), \tau_k)\}$ are i.i.d. for each s . The following theorem asserts that, under regularity conditions, we may pick any regenerative sequence γ_n such that $E[X_1(s)] = 0$ for all strata.

Theorem 2. *Suppose that $L(\theta)$ is differentiable on H and $\sup_{\theta \in H} |L'_s(\theta)| < \infty$ for $1 \leq s \leq q$ and $\theta \in H$. Also suppose that $\epsilon_n = O(n^{-\alpha})$ for some $\alpha \in (0.5, 1]$ and that $(\epsilon_n - \epsilon_{n+1})/\epsilon_n = O(\epsilon_n)$. Finally, suppose that $\{\gamma_n\}$ is regenerative with $E[\tau_1^{1/\alpha}] < \infty$ and $E[X_1(s)] = 0$ for $1 \leq s \leq q$. Then Condition 7 holds.*

The condition $E[X_1(s)] = 0$ essentially requires that, for each stratum s , the expected fraction of visits to s in a cycle equals w_s . By the strong law of large numbers for regenerative processes [3, Sec. VI.3], this condition—in the presence of the finite-moment condition on τ_1 —also implies that the long-term fraction of visits to s equals w_s . The finite-moment condition is typically satisfied whenever the number of successive steps taken within a stratum is bounded with probability 1.

Proof. Fix $\theta \in H$ and observe that

$$\begin{aligned} \epsilon_n \sum_{i=0}^{n-1} (L'_{\gamma_i}(\theta) - L'(\theta)) &= \epsilon_n \sum_{i=0}^{n-1} \sum_{s=1}^q (L'_s(\theta) I_{\gamma_i=s} - L'_s(\theta) w_s) \\ &= \sum_{s=1}^q L'_s(\theta) \epsilon_n \sum_{i=0}^{n-1} (I_{\gamma_i=s} - w_s). \end{aligned}$$

Since $|L'_s(\theta)| < \infty$ for each s , it suffices to show that $n^{-\alpha} \sum_{i=0}^{n-1} (I_{\gamma_i=s} - w_s) \xrightarrow{\text{a.s.}} 0$ for $1 \leq s \leq q$. To this end, fix s and denote by $c(n)$ the (random) number of complete cycles up to step n . We have

$$\sum_{i=0}^n (I_{\gamma_i=s} - w_s) = \sum_{k=1}^{c(n)} X_k(s) + R_{1,n},$$

where $R_{1,n} = \sum_{i=\beta(c(n))}^n (I_{\gamma_i=s} - w_s)$. I.e., the sum can be broken up into sums over complete cycles plus a remainder term corresponding to a sum over a partially completed cycle. Similar

calculations let us write $n = \sum_{k=1}^{c(n)} \tau_k + R_{2,n}$, where $R_{2,n} = n - \beta(c(n)) + 1$. Thus

$$\begin{aligned}
\frac{\sum_{i=0}^n (I_{\gamma_i=s} - w_s)}{n^\alpha} &= \frac{\sum_{k=1}^{c(n)} X_k(s) + R_{1,n}}{\left(\sum_{k=1}^{c(n)} \tau_k + R_{2,n}\right)^\alpha} \\
&= \frac{\sum_{k=1}^{c(n)} X_k(s)}{c(n)^\alpha} \left(\frac{\sum_{k=1}^{c(n)} \tau_k}{c(n)} + \frac{R_{2,n}}{c(n)}\right)^{-\alpha} \\
&\quad + \frac{R_{1,n}/c(n)^\alpha}{\left(\sum_{k=1}^{c(n)} \tau_k/c(n) + R_{2,n}/c(n)\right)^\alpha}.
\end{aligned} \tag{9}$$

By assumption, the random variables $\{X_k(s)\}$ are i.i.d. with common mean 0. Moreover, $|X_k(s)| \leq (1 + w_s)\tau_k$, which implies that $E[|X_1(s)|^{1/\alpha}] \leq (1 + w_s)^{1/\alpha} E[\tau_1^{1/\alpha}] < \infty$. It then follows from the Marcinkiewicz-Zygmund strong law [9, Th. 5.2.2] that $n^{-\alpha} \sum_{k=1}^n X_k(s) \xrightarrow{\text{a.s.}} 0$. Because each regeneration point, and hence each cycle length, is assumed to be almost surely finite, it follows that $c(n) \xrightarrow{\text{a.s.}} \infty$, so that $\sum_{k=1}^{c(n)} X_k(s)/c(n)^\alpha \xrightarrow{\text{a.s.}} 0$ as $n \rightarrow \infty$. Similarly, an application of the ordinary strong law of large numbers shows that $\sum_{k=1}^{c(n)} \tau_k/c(n) \xrightarrow{\text{a.s.}} E[\tau_1] > 0$. Next, note that $|R_{1,n}| \leq (1 + w_s)\tau_{c(n)+1}$, so that $R_{1,n}/c(n)^\alpha \xrightarrow{\text{a.s.}} 0$ provided that $\tau_k/k^\alpha \xrightarrow{\text{a.s.}} 0$. To establish this latter limit result, observe that for any $\epsilon > 0$, the assumed finiteness of $E[(\tau_k/\epsilon)^{1/\alpha}]$ implies [10, Th. 3.2.1] that $\sum_{k=1}^\infty \Pr[(\tau_k/\epsilon)^{1/\alpha} \geq k] < \infty$, and hence $\sum_{k=1}^\infty \Pr[\tau_k/k^\alpha \geq \epsilon] < \infty$. It then follows from the first Borel-Cantelli Lemma (see [10, Th. 4.2.1]) that $\Pr[\tau_k/k^\alpha \geq \epsilon \text{ infinitely often}] = 0$, which in turn implies [10, Th. 4.2.2] that $\tau_k/k^\alpha \xrightarrow{\text{a.s.}} 0$. A similar argument shows that $R_{2,n}/c(n) \xrightarrow{\text{a.s.}} 0$, and the desired result follows after letting $n \rightarrow \infty$ in the rightmost expression in (9). \square

The conditions on $\{\epsilon_n\}$ in Theorem 2 are often satisfied in practice, e.g., when $\epsilon_n = 1/n$ or when $\epsilon_n = 1/\lceil n/k \rceil$ for some $k > 1$ with $\lceil x \rceil$ denoting the smallest integer greater than or equal to x (so that the step size remains constant for some fixed number of steps, as in Algorithm 2 below). See Sec. 6 for further discussion.

Similarly, a wide variety of strata-selection schemes satisfy the conditions of Theorem 2. Some simple examples include (1) running precisely cw_s steps on stratum s in every ‘‘chunk’’ of c steps, and (2) repeatedly picking a stratum according to some fixed distribution $\{p_s > 0\}$ and running cw_s/p_s steps on the selected stratum s . (E.g., we can set $p_s = w_s$ so that strata are chosen proportional to their weight and a constant number of steps is run on the selected stratum, or $p_s = 1/q$ so that strata are chosen uniformly and at random but the number of steps run on the selected stratum is proportional to its weight.) For example (1), assume initially that the order of stratum visits within any two chunks is the same. Then, in the notation of Theorem 2, we have $\tau_k = c$ and $X_k(s) = 0$ with probability 1 for all k , so the conditions of the theorem hold trivially, with the chunks playing the role of regenerative cycles. In fact, the order within each chunk is irrelevant, since for any ordering the pairs $\{(Y_k(s), \tau_k)\}$ are trivially i.i.d., so that the proof of the theorem goes through essentially unchanged. For example (2), the steps at which a stratum is randomly selected clearly

form a sequence of regeneration points for $\{\gamma_n\}$. We have

$$\mathbb{E}[\tau_1] = \sum_{s=1}^q p_s \frac{cw_s}{p_s} = c \sum_{s=1}^q w_s = c$$

The sum of the random variables $I_{\gamma_n=s}$ over a cycle is cw_s/p_s if stratum s is selected and 0 otherwise, so that the expected sum is $p_s(cw_s/p_s) + (1-p_s)0 = cw_s$ and hence

$$\begin{aligned} \mathbb{E}[X_1(s)] &= \mathbb{E}\left[\sum_{\gamma_n \in \text{cycle } 1} (I_{\gamma_n=s} - w_s)\right] \\ &= \mathbb{E}\left[\sum_{\gamma_n \in \text{cycle } 1} I_{\gamma_n=s}\right] - \mathbb{E}[w_s \tau_k] = cw_s - cw_s = 0. \end{aligned}$$

Moreover, τ_1 is bounded above by $\max_s cw_s/p_s < \infty$, and hence has finite moments of all orders.

To give a better idea of the scope of stratum-selection schemes covered by Theorem 2, we conclude by discussing a randomized procedure in which, after visiting a stratum s , we visit the same stratum at the next step with probability $p_s = 1 - (cw_s)^{-1}$, and with probability $1 - p_s$ we select a new stratum randomly and uniformly from the q strata. (Thus the new stratum may correspond to the old stratum.) Here c is a constant that is large enough to ensure that $cw_s > 1$ for each s . Denote by $s_0 \in \{1, 2, \dots, q\}$ the initial stratum to be visited; we fix s_0 a priori. Then $\{\gamma_n\}$ is regenerative, with the regeneration points corresponding to the successive steps at which a new stratum is selected and the new stratum is s_0 . We can write $\tau_1 = \sum_{i=1}^N V_i$, where N is the total number of new-stratum selections, and V_i is the number of successive visits to the i th selected stratum. The random variable N has a geometric distribution with mean q and, given that stratum s is selected at the i th selection epoch, V_i has a geometric distribution with mean cw_s . Moreover, the V_i 's are mutually independent and independent of N , and V_2, \dots, V_N are i.i.d., specifically, each V_i ($i > 1$) is distributed as an average of $q - 1$ independent geometric random variables with means $\{w_s : s \neq s_0\}$. It follows that

$$\mathbb{E}[\tau_1] = \mathbb{E}[V_1] + \mathbb{E}\left[\sum_{i=2}^N V_i\right] = cw_{s_0} + \mathbb{E}[N - 1] \left((q - 1)^{-1} \sum_{s \neq s_0} cw_s \right) = \sum_s cw_s = c.$$

Let N_s denote the number of new-stratum selections equal to s in the first cycle. It is not hard to see that $N_{s_0} = 1$ with probability 1 and that $\mathbb{E}[N_s] = 1$ for $s \neq s_0$. Thus, for each s ,

$$\mathbb{E}\left[\sum_{\gamma_n \in \text{cycle } 1} I_{\gamma_n=s}\right] = cw_s \mathbb{E}[N_s] = cw_s,$$

so that $\mathbb{E}[X_1(s)] = 0$. Finally, we show that τ_1 has finite moments of all orders. Using the simple bound $(\sum_{i=1}^n x_i)^\beta \leq (n \max_{1 \leq i \leq n} x_i)^\beta \leq n^\beta \sum_{i=1}^n x_i^\beta$ for $x_1, x_2, \dots, x_n \geq 0$ and $\beta \geq 0$, we have $\mathbb{E}[\tau_1^\beta] \leq \mathbb{E}[N^\beta] m(\beta, s_0) + \mathbb{E}[N^\beta(N - 1)] \sum_{s \neq s_0} m(\beta, s)/(q - 1)$. Here $m(\beta, s)$ denotes

the (finite) β th moment of a geometric random variable with mean cw_g . The desired result then follows from the fact that N also has a geometric distribution, and hence has finite moments of all orders.

The above examples are primarily of theoretical interest. In Sec. 6, we focus on some schemes that are particularly suitable for practical implementation in the context of DSGD.

5. The DSGD Algorithm

We can exploit the structure of the matrix factorization problem to derive a distributed algorithm for rank- r matrix factorization via SGD. The idea is to specialize the SSGD algorithm, choosing the strata such that SGD can be run on each stratum in a distributed manner. We first discuss the “interchangeability” structure that we will exploit for distributed processing within a stratum.

5.1. Interchangeability

In general, distributing SGD is hard because the individual steps depend on each other: from (4), we see that θ_n has to be known before θ_{n+1} can be computed. However, in the case of matrix factorization, the SGD process has some structure that we can exploit.

We focus on loss-minimization problems of the form $\text{minimize}_{\theta \in H} L(\theta)$ where the loss function L has summation form: $L(\theta) = \sum_{z \in Z} L_z(\theta)$.

Definition 1. Two training points $z_1, z_2 \in Z$ are interchangeable if for all loss functions L having summation form, all $\theta \in H$, and $\epsilon > 0$,

$$\begin{aligned} L'_{z_1}(\theta) &= L'_{z_1}(\theta - \epsilon L'_{z_2}(\theta)) \\ \text{and} \quad L'_{z_2}(\theta) &= L'_{z_2}(\theta - \epsilon L'_{z_1}(\theta)). \end{aligned} \tag{10}$$

Two disjoint sets of training points $Z_1, Z_2 \subset Z$ are interchangeable if z_1 and z_2 are interchangeable for every $z_1 \in Z_1$ and $z_2 \in Z_2$.

As described in Sec. 5.2 below, we can swap the order of consecutive SGD steps that involve interchangeable training points without affecting the final outcome.

Now we return to the setting of matrix factorization, where the loss function has the form $L(\mathbf{W}, \mathbf{H}) = \sum_{(i,j) \in Z} L_{ij}(\mathbf{W}, \mathbf{H})$ with $L_{ij}(\mathbf{W}, \mathbf{H}) = l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$. The following theorem gives a simple criterion for interchangeability.

Theorem 3. Two training points $z_1 = (i_1, j_1) \in Z$ and $z_2 = (i_2, j_2) \in Z$ are interchangeable if they share neither row nor column, i.e., $i_1 \neq i_2$ and $j_1 \neq j_2$.

Proof. The result is a direct consequence of the decomposition of the global loss into a sum of local losses. Specifically, it follows from (5) and (5) that the partial derivatives of L_{ij} (1) depend only on \mathbf{V}_{ij} , \mathbf{W}_{i*} , and \mathbf{H}_{*j} , and (2) are nonzero only with respect to $\mathbf{W}_{i_1}, \dots, \mathbf{W}_{i_r}, \mathbf{H}_{1j}, \dots, \mathbf{H}_{rj}$. When $i_1 \neq i_2$ and $j_1 \neq j_2$, both (\mathbf{W}, \mathbf{H}) and $(\mathbf{W}, \mathbf{H}) - \epsilon L'_{z_1}(\mathbf{W}, \mathbf{H})$ agree on the values of \mathbf{W}_{i_2*} and \mathbf{H}_{*j_2} for any choice of (\mathbf{W}, \mathbf{H}) , which establishes the second part of (10). Analogous arguments hold for the first part. \square

It follows that if two blocks of \mathbf{V} share neither rows or columns, then the sets of training points contained in these blocks are interchangeable.

5.2. A Simple Case

We introduce the DSGD algorithm by considering a simple case that essentially corresponds to running DSGD using a single “ d -monomial” stratum (see Sec. 5.3). The goal is to highlight the technique by which DSGD runs the SGD algorithm in a distributed manner within a stratum. For a given training set Z , denote by \mathbf{Z} the corresponding *training matrix*, which is obtained by zeroing out the elements in \mathbf{V} that are not in Z ; these elements usually represent missing data or held-out data for validation. In our simple scenario, Z corresponds to our single stratum of interest, and the corresponding training matrix \mathbf{Z} is block-diagonal:

$$\begin{matrix} & \mathbf{H}^1 & \mathbf{H}^2 & \dots & \mathbf{H}^d \\ \mathbf{W}^1 & \left(\begin{array}{cccc} \mathbf{Z}^1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}^2 & \dots & \vdots \\ \vdots & \vdots & \ddots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & \mathbf{Z}^d \end{array} \right) & & & \end{matrix}, \quad (11)$$

where \mathbf{W} and \mathbf{H} are blocked conformingly. Denote by Z^b the set of training points in block \mathbf{Z}^b . We exploit the key property that, by Theorem 3, sets Z^i and Z^j are interchangeable for $i \neq j$. For some $T \in [1, \infty)$, suppose that we run T steps of SGD on \mathbf{Z} , starting from some initial point $\theta_0 = (\mathbf{W}_0, \mathbf{H}_0)$ and using a fixed step size ϵ . We can describe an instance of the SGD process by a *training sequence* $\omega = (z_0, z_1, \dots, z_{T-1})$ of T training points. Figure 1 shows an example of such a training sequence. Define $\theta_0(\omega) = \theta_0$ and

$$\theta_{n+1}(\omega) = \theta_n(\omega) + \epsilon Y_n(\omega),$$

where the update term $Y_n(\omega) = -NL'_{\omega_n}(\theta_n(\omega))$ is the scaled negative gradient estimate as in standard SGD. We can write

$$\theta_T(\omega) = \theta_0 + \epsilon \sum_{n=0}^{T-1} Y_n(\omega). \quad (12)$$

To see how to exploit the interchangeability structure, consider the subsequence $\sigma_b(\omega) = \omega \cap Z^b$ of training points from block \mathbf{Z}^b ; the subsequence has length $T_b(\omega) = |\sigma_b(\omega)|$. The following theorem asserts that we can run SGD on each block independently, and then sum up the results.

Theorem 4. *Using the definitions above,*

$$\theta_T(\omega) = \theta_0 + \epsilon \sum_{b=1}^d \sum_{k=0}^{T_b(\omega)-1} Y_k(\sigma_b(\omega)). \quad (13)$$

Proof. We establish a one-to-one correspondence between the update terms $Y_n(\omega)$ in (12) and $Y_k(\sigma_b(\omega))$ in (13). Denote by $z_{b,k}$ the $(k+1)$ st element in $\sigma_b(\omega)$, i.e., the $(k+1)$ st element from block \mathbf{Z}^b in ω . Denote by $\pi(z_{b,k})$ the 0-based position of this element in ω . We have $\omega_{\pi(z_{b,k})} = z_{b,k}$. Now consider the first element $z_{b,0}$ from block b . We have $z_n \notin \mathbf{Z}^b$ for all previous elements $n < \pi(z_{b,0})$. Since the training matrix is block-diagonal, blocks have pairwise disjoint rows and pairwise disjoint columns. Thus by Theorem 3, $z_{b,0}$ is interchangeable with each of the z_n for $n < \pi(z_{b,0})$. We can therefore eliminate the z_n one by one:

$$\begin{aligned} Y_{\pi(z_{b,0})}(\omega) &= -NL'_{z_{b,0}}(\theta_{\pi(z_{b,0})}(\omega)) = -NL'_{z_{b,0}}(\theta_{\pi(z_{b,0})-1}(\omega)) \\ &= \dots = -NL'_{z_{b,0}}(\theta_0) = Y_0(\sigma_b(\omega)). \end{aligned}$$

Using the same argument, we can safely remove *all* update terms from elements not in block \mathbf{Z}^b . By induction on k , we obtain

$$\begin{aligned} Y_{\pi(z_{b,k})}(\omega) &= -NL'_{z_{b,k}} \left(\theta_0 + \epsilon \sum_{n=0}^{\pi(z_{b,k})-1} Y_n(\omega) \right) \\ &= -NL'_{z_{b,k}} \left(\theta_0 + \epsilon \sum_{l=0}^{k-1} Y_{\pi(z_{b,l})}(\sigma_b(\omega)) \right) \\ &= Y_k(\sigma_b(\omega)). \end{aligned} \tag{14}$$

The assertion of the theorem now follows from

$$\begin{aligned} \theta_T(\omega) &= \theta_0 + \epsilon \sum_{n=0}^{T-1} Y_n(\omega) = \theta_0 + \epsilon \sum_{b=1}^d \sum_{k=0}^{T_b(\omega)-1} Y_{\pi(z_{b,k})}(\omega) \\ &= \theta_0 + \epsilon \sum_{b=1}^d \sum_{k=0}^{T_b(\omega)-1} Y_k(\sigma_b(\omega)), \end{aligned}$$

where we first reordered the update terms and then used (14). \square

We used the fact that \mathbf{Z} is block-diagonal only to establish interchangeability between blocks. This means that Theorem 4 also applies when the matrix is not block-diagonal, but can be divided into a set of interchangeable submatrices in some other way.

We now describe how to exploit Theorem 4 for distributed processing on MapReduce. We block \mathbf{W} and \mathbf{H} conformingly to \mathbf{Z} —as in (11)—and divide processing into d independent map tasks $\Gamma_1, \dots, \Gamma_d$ as follows. Task Γ_b is responsible for subsequence $\sigma_b(\omega)$: It takes \mathbf{Z}^b , \mathbf{W}^b , and \mathbf{H}^b as input, performs the block-local updates $\sigma_b(\omega)$, and outputs updated factor matrices $\mathbf{W}_{\text{new}}^b$ and

$\mathbf{H}_{\text{new}}^b$.³ By Theorem 4, we have

$$\mathbf{W}' = \begin{pmatrix} \mathbf{W}_{\text{new}}^1 \\ \vdots \\ \mathbf{W}_{\text{new}}^d \end{pmatrix} \quad \text{and} \quad \mathbf{H}' = (\mathbf{H}_{\text{new}}^1 \quad \cdots \quad \mathbf{H}_{\text{new}}^d),$$

where \mathbf{W}' and \mathbf{H}' are the matrices that one would obtain by running sequential SGD on ω . Since each task accesses different parts of both training data and factor matrices, the data can be distributed across multiple nodes and the tasks can run simultaneously. In Sec. 6, we describe how to efficiently implement the above idea.

5.3. The General Case

We now present the complete DSGD matrix-factorization algorithm. The key idea is to stratify the training set Z into a set $S = \{Z_1, \dots, Z_q\}$ of q strata so that each individual stratum $Z_s \subseteq Z$ can be processed in a distributed fashion. We do this by ensuring that each stratum is “ d -monomial” as defined below. The d -monomial property generalizes the block-diagonal structure of the example in Sec. 5.2, while still permitting the techniques of that section to be applied. The strata must *cover* the training set in that $\bigcup_{s=1}^q Z_s = Z$, but overlapping strata are allowed. The parallelism parameter d is chosen to be greater than or equal to the number of available processing tasks.

Definition 2. A stratum Z_s is d -monomial if it can be partitioned into d nonempty subsets $Z_s^1, Z_s^2, \dots, Z_s^d$ such that $i \neq i'$ and $j \neq j'$ whenever $(i, j) \in Z_s^{b_1}$ and $(i', j') \in Z_s^{b_2}$ with $b_1 \neq b_2$. A training matrix \mathbf{Z}_s is d -monomial if it is constructed from a d -monomial stratum Z_s .

There are many ways to stratify the training set according to Def. 2. In our current work, we perform *data-independent blocking*; more advanced strategies may improve the speed of convergence further. We first randomly permute the rows and column of \mathbf{Z} , and then create $d \times d$ blocks of size $(m/d) \times (n/d)$ each; the factor matrices \mathbf{W} and \mathbf{H} are blocked conformingly. This procedure ensures that the expected number of training points in each of the blocks is the same, namely, N/d^2 . Then, for a permutation j_1, j_2, \dots, j_d of $1, 2, \dots, d$, we can define a stratum as $Z_s = Z^{1j_1} \cup Z^{2j_2} \cup \dots \cup Z^{dj_d}$, where the *substratum* Z^{ij} denotes the set of training points that fall within block \mathbf{Z}^{ij} . We can represent a stratum Z_s by a template $\tilde{\mathbf{Z}}_s$ that displays each block \mathbf{Z}^{ij} corresponding to a substratum Z^{ij} of Z_s , with all other blocks represented by zero matrices. When $d = 2$, for example, we obtain two strata represented by the templates

$$\tilde{\mathbf{Z}}_1 = \begin{pmatrix} \mathbf{Z}^{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}^{22} \end{pmatrix} \quad \text{and} \quad \tilde{\mathbf{Z}}_2 = \begin{pmatrix} \mathbf{0} & \mathbf{Z}^{12} \\ \mathbf{Z}^{21} & \mathbf{0} \end{pmatrix}.$$

The set S of possible strata contains $d!$ elements, one for each possible permutation of $1, 2, \dots, d$. Note that different strata may overlap when $d > 2$. Also note that there is no need to materialize these strata: They are constructed on-the-fly by processing only the respective blocks of \mathbf{Z} .

³Since training data is sparse, a block \mathbf{Z}^b may contain no training points; in this case we cannot execute SGD on the block, so the corresponding factors simply remain at their initial values.

Given a set of strata and associated weights $\{w_s\}$, we decompose the loss into a weighted sum of per-stratum losses as in (7): $L(\mathbf{W}, \mathbf{H}) = \sum_{s=1}^q w_s L_s(\mathbf{W}, \mathbf{H})$. (As in Sec. 3.2, we suppress the fixed matrix \mathbf{V} in our notation for loss functions.) We use per-stratum losses of form

$$L_s(\mathbf{W}, \mathbf{H}) = c_s \sum_{(i,j) \in Z_s} L_{ij}(\mathbf{W}, \mathbf{H}), \quad (15)$$

where c_s is a stratum-specific constant; see the discussion below. When running SGD on a stratum, we use the gradient estimate

$$\hat{L}'_s(\mathbf{W}, \mathbf{H}) = N_s c_s L'_{ij}(\mathbf{W}, \mathbf{H}) \quad (16)$$

of $L'_s(\mathbf{W}, \mathbf{H})$ in each step, i.e., we scale up the local loss of an individual training point by the size $N_s = |Z_s|$ of the stratum. For example, from the $d!$ strata described previously, we can select d disjoint strata Z_1, Z_2, \dots, Z_d such that they cover Z . Then any given loss function L of the form (2) can be represented as a weighted sum over these strata by choosing w_s and c_s subject to $w_s c_s = 1$. Recall that w_s can be interpreted as the “time” spent on each stratum in the long run. A natural choice is to set $w_s = N_s/N$, i.e., proportional to the stratum size. This particular choice leads to $c_s = N/N_s$ and we obtain the standard SGD gradient estimator $\hat{L}'_s(\mathbf{W}, \mathbf{H}) = N L'_{ij}(\mathbf{W}, \mathbf{H})$. As another example, we can represent L as a weighted sum in terms of all $d!$ strata; in light of the fact that each substratum Z^{ij} lies in exactly $(d-1)!$ of these strata, we choose $w_s = N_s/((d-1)N)$ and use the value of $c_s = N/N_s$ as before.

The individual steps in DSGD are grouped into *subepochs*, each of which amounts to processing one of the strata. In more detail, DSGD makes use of a sequence $\{(\xi_k, T_k)\}$, where ξ_k denotes the *stratum selector* used in the k th subepoch, and T_k the number of steps to run on the selected stratum. Note that this sequence of pairs uniquely determines an SSGD stratum sequence as in Sec. 4.1: $\gamma_1 = \dots = \gamma_{T_1} = \xi_1$, $\gamma_{T_1+1} = \dots = \gamma_{T_1+T_2} = \xi_2$, and so forth. The $\{(\xi_k, T_k)\}$ sequence is chosen such that the underlying SSGD algorithm, and hence the DSGD factorization algorithm, is guaranteed to converge; see Sec. 4.3. Once a stratum ξ_k has been selected, we perform T_k SGD steps on Z_{ξ_k} ; this is done in a parallel and distributed way using the technique of Sec. 5.2. DSGD is shown as Algorithm 2, where we define an *epoch* as a sequence of d subepochs. As will become evident in Sec. 6 below, an epoch roughly corresponds to processing the entire training set once.

When executing Algorithm 2 on d nodes in a shared-nothing environment such as MapReduce, the input matrix need only be distributed once. Then the only data that are transmitted between nodes during subsequent processing are (small) blocks of factor matrices. Indeed, if node i stores blocks $\mathbf{W}^i, \mathbf{Z}^{i1}, \mathbf{Z}^{i2}, \dots, \mathbf{Z}^{id}$ for $1 \leq i \leq d$, then only matrices $\mathbf{H}^1, \mathbf{H}^2, \dots, \mathbf{H}^d$ need be transmitted.

Since, by construction, parallel processing within the k th selected stratum leads to the same update terms as for the corresponding sequential SGD algorithm on Z_{ξ_k} , we have established the connection between DSGD and SSGD. Thus the convergence of DSGD is implied by the convergence of the underlying SSGD algorithm; see Sec. 4.2.

Algorithm 2 DSGD for Matrix Factorization

Require: Z, W_0, H_0 , cluster size d
 $W \leftarrow W_0$
 $H \leftarrow H_0$
Block $Z / W / H$ into $d \times d / d \times 1 / 1 \times d$ blocks
while not converged **do** /* *epoch* */
 Pick step size ϵ
 for $s = 1, \dots, d$ **do** /* *subepoch* */
 Pick d blocks $\{Z^{1j_1}, \dots, Z^{dj_a}\}$ to form a stratum
 for $b = 1, \dots, d$ **do** /* *in parallel* */
 Run SGD on the training points in Z^{bj_b} (step size = ϵ)
 end for
 end for
end while

6. DSGD Implementation

In this section, we discuss some practical issues around DSGD, including initialization considerations and practical methods for choosing the training sequence for the parallel SGD step, selecting strata, and picking the step size ϵ . As above, a “subepoch” corresponds to processing a stratum and an “epoch”—roughly equivalent to a complete pass through the training data—corresponds to processing a sequence of d strata.

Initialization. Some care must be taken when choosing initial factor values W_0 and H_0 . In the case of nonzero squared loss L_{NZSL} , for example, choosing $W_0 = \mathbf{0}$ and $H_0 = \mathbf{0}$ results in the factors remaining equal to zero at all future DSGD iterations. For GKL loss, we cannot have $W_{i*} = \mathbf{0}$ for any i or H_{*j} for any j , since then the loss function is ill defined. In our implementation, we generate initial factor values using a pseudorandom number generator, which ensures that all initial values are nonzero.

Training sequence. When processing a subepoch (i.e., a stratum), we do not generate a global training sequence and then distribute it among blocks. Instead, each task generates a *local* training sequence directly for its corresponding block. This reduces communication cost and avoids the bottleneck of centralized computation. Practical experience suggests that good results are achieved when (1) the local training sequence covers a large part of the local block, and (2) the training sequence is randomized. We consider the following strategies for processing block Z^{ij} :

- *Sequential selection (SEQ)*. Scan Z^{ij} in the order it is stored.
- *With-replacement selection (WR)*. Randomly select training points from Z^{ij} ; each point may be selected multiple times.
- *Without-replacement selection (WOR)*. Randomly select training points from Z^{ij} ; each point is selected precisely once.

The first two strategies are extremes: SEQ satisfies (1) but not (2), whereas WR satisfies (2) but not (1). A compromise—which worked best in our experiments—is WOR; it ensures that many different

training points are selected while at the same time maximizing randomness. Note that Theorem 2 implicitly assumes WR but can be extended to cover SEQ and WOR as well. (In brief, redefine a stratum to consist of a single training point and redefine the stratum weights w_s accordingly.)

Update terms. When processing a training point (i, j) during an SGD step on stratum s , we use the gradient estimate $\hat{L}'_s(\theta) = NL'_{ij}(\theta)$ as in standard SGD; this corresponds to a choice of $c_s = N/N_s$ in (16). For (i, j) picked uniformly and at random from Z_s , the estimate is unbiased for the gradient of the stratum loss $L_s(\theta)$ given in (15).

Stratum selection. Recall that the stratum sequence (ξ_k, T_k) determines which of the strata is chosen in each subepoch and how many steps are run on that stratum. We choose training sequences such that $T_k = N_{\xi_k} = |Z_{\xi_k}|$; this ensures that we can make use of all the training points in the stratum. For the data-independent blocking scheme given in Sec. 5.3, each block Z^{ij} occurs in $(d-1)!$ of the $d!$ strata. Thus we do not need to process all strata to cover the entire training set. As above, we want to process a large part of the training set in each epoch, while at the same time maximizing randomization. To select a set of d strata to visit during an epoch, we use strategies similar to those for intra-block training point selection:

- *Sequential selection (SEQ).* Pick a sequence of d strata that jointly cover the entire training matrix. Then cycle through this sequence and ignore all other strata.
- *With replacement selection (WR).* Repeatedly pick a stratum uniformly and at random from the set of all strata until d strata have been processed.
- *Without replacement selection (WOR).* Pick a sequence of d strata such that the d strata jointly cover the entire training set; the sequence is picked uniformly and at random from all such sequences of d strata.⁴

Taking the scaling constant c_s in (15) as N/N_s , we can see that all three strategies are covered by Theorem 2, where each epoch corresponds to a regenerative cycle. We argue informally as follows. Recall that if Theorem 2 is to apply, then w_s must correspond to the long-term fraction of steps run on stratum Z_s . For SEQ, this means that all but d of the weights are zero, and the remaining weights satisfy $w_s = N_s/N$. For WR and WOR, we have $w_s = N_s/((d-1)!N)$, since we select each stratum s equally often in the long run, and always perform N_s steps on stratum s . The question is then whether these choices of w_s lead to a legitimate representation of L as in (7). One can show that $\{w_s\}$ satisfies (7) for all Z and L of form (2) if and only if

$$\sum_{s: Z_s \supseteq Z^{ij}} w_s c_s = 1 \quad (17)$$

for each substratum Z^{ij} . Direct verification shows that (17) holds for the above choices of w_s when $c_s = N/N_s$.

⁴This can be performed efficiently by randomly permuting the rows and columns of a matrix of form

$$\begin{pmatrix} 1 & 2 & \cdots & d \\ 2 & 3 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ d & 1 & \cdots & d-1 \end{pmatrix}.$$

The (k, i) -entry then contains the column index of the block to pick from row i in the k th subepoch.

Step sizes. The stochastic approximation literature often works with step size sequences roughly of form $\epsilon_n = 1/n^\alpha$ with $\alpha \in (0.5, 1]$; and Theorem 2 guarantees asymptotic convergence for such choices. In practice, one may want to deviate from these choices to achieve faster convergence over the finite number of steps that are actually executed. We use an adaptive method for choosing the step size sequence. We exploit the fact that—in contrast to SGD in general—we can determine the current loss after every epoch. Thus we can check whether an epoch decreased or increased the loss. With this observation in mind, we employ a heuristic called *bold driver*, which is often used for gradient descent. Starting from an initial step size ϵ_0 , we (1) increase the step size by a small percentage (say, 5%) whenever we see a decrease of loss, and (2) we drastically decrease the step size (say, by 50%) if we observe an increase of loss. Within each epoch, the step size remains fixed. Given a reasonable choice of ϵ_0 , the bold driver method worked extremely well in our experiments. To pick ϵ_0 , we leverage the fact that we have many compute nodes available. We replicate a small sample of Z (say, 0.1%) to each node. We then try different step sizes in parallel. Initially, we make a pass over the sample for step sizes $1, 1/2, 1/4, \dots, 1/2^{d-1}$; this is done in parallel at all d nodes. The step size that gives the best result is selected as ϵ_0 . As long as our loss decreases, we repeat a variation of this process after every epoch, where we try step sizes within a factor of $[1/2, 2]$ of the current step size. Eventually, the so-chosen step size will become too large and the value of the loss will increase. Intuitively, this happens when the iterate has moved closer to the global solution than to the local solution of the sample. As soon as we observe an increase of loss, we switch to the bold driver method for the rest of the process.

7. Experiments

We compared various matrix factorization algorithms with respect to their convergence properties, runtime efficiency, and scalability. We found that the convergence speed of DSGD is on par or better than alternative methods, even when these methods are specialized to the loss function. In terms of overall performance, we found that DSGD is significantly faster, produces more stable results, and has better scalability properties.

7.1. Setup

We implemented our new DSGD method on top of MapReduce, along with the PSGD, ISGD, DGD, ALS, and MULT methods discussed in Sec. 2. The DGD algorithm uses the L-BFGS quasi-Newton method as in [13]. DSGD, PSGD, and L-BFGS are generic methods that work with a wide variety of loss functions, whereas ALS and MULT are restricted to quadratic loss functions and GKL, respectively. We used two different implementations and compute clusters; one for in-memory experiments and one for large scale-out experiments on very large datasets.

The in-memory implementation is based on R and C, and uses R’s `snowfall` package to implement MapReduce. It targets datasets that are small enough to fit in aggregate memory, i.e., with up to a few billion nonzero entries. We block and distribute the input matrix across the cluster before running each experiment. The factor matrices are communicated via Samba mount points.

The R cluster consists of 16 nodes, each running two Intel Xeon E5530 processors with 8 cores at 2.4GHz. Every node has 48GB of memory.

The second implementation is based on Hadoop [2], an open-source MapReduce implementation. The Hadoop cluster is equipped with 40 nodes, each with two Intel Xeon E5440 processors and 4 cores at 2.8GHz and 32 GB of memory. A straightforward DSGD implementation on Hadoop achieves scalability and fault-tolerance, but suffers from some of Hadoop’s internal overheads. We expect that some of these overheads will diminish as Hadoop matures. For the time being, we employ the following optimizations. First, we store the data in primitive arrays rather than Java objects; this avoids the performance bottleneck caused by “immutable record decoders” [17]. Second, we use block-wise I/O to read entire matrix blocks at once; this avoids expensive per-record processing costs. Finally, we submit a single MapReduce job per epoch (complete data-matrix scan) rather than per subepoch (stratum scan); this reduces Hadoop’s high overhead in spawning jobs. To do this, we make use of a custom locking scheme that preserves fault-tolerance.

For our experiments with all SGD-based approaches, we used adaptive step size computation based on a sample of roughly 1M data points, switching to the bold driver as soon as an increase in loss was observed. The time for step size selection is included in all of our performance plots. For ISGD, we performed (parallel) step size computation using the parameters of only the first partition; this step size worked well for all partitions. Unless stated otherwise, we used WOR selection for both training sequences (all approaches) and stratum sequences (DSGD only).

We used the Netflix competition dataset [5] for our experiments on real data. The dataset contains a small subset of movie ratings given by Netflix users, specifically, 100M anonymized, time-stamped ratings from roughly 480k customers on roughly 18k movies. For larger-scale experiments on the in-memory implementation, we used a synthetic dataset with 10M rows, 1M columns, and 1B nonzero entries. We first generated matrices \mathbf{W}^* and \mathbf{H}^* by repeatedly sampling values from the Gaussian(0,10) distribution. We then sampled 1B entries from the product $\mathbf{W}^* \mathbf{H}^*$, and added Gaussian(0,1) noise to each sample; this procedure ensured the existence of a reasonable low-rank factorization. For all experiments, we centered the input matrix around its mean. The starting points \mathbf{W}_0 and \mathbf{H}_0 were chosen by sampling entries uniformly and at random from $[-0.5, 0.5]$; we used the same starting point for each algorithm to ensure fair comparison. Unless stated otherwise, we used rank $r = 50$.

We used four well-known loss functions in our experiments: plain nonzero squared loss (L_{NZSL}), nonzero squared loss with an L_2 regularization term (L_{L2}), nonzero squared loss with a nonzero-weighted L_2 term (L_{NZL2}), and generalized KL divergence (L_{GKL}):

$$\begin{aligned}
 L_{\text{NZSL}} &= \sum_{(i,j) \in Z} (\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})^2 \\
 L_{\text{L2}} &= L_{\text{NZSL}} + \lambda(\|\mathbf{W}\|_{\text{F}}^2 + \|\mathbf{H}\|_{\text{F}}^2) \\
 L_{\text{NZL2}} &= L_{\text{NZSL}} + \lambda(\|\mathbf{N}_1 \mathbf{W}\|_{\text{F}}^2 + \|\mathbf{H} \mathbf{N}_2\|_{\text{F}}^2) \\
 L_{\text{GKL}} &= \sum_{(i,j) \in Z} (\mathbf{V}_{ij} \log \mathbf{V}_{ij} / [\mathbf{W}\mathbf{H}]_{ij} - \mathbf{V}_{ij}) + \sum_{i,j} [\mathbf{W}\mathbf{H}]_{ij},
 \end{aligned} \tag{18}$$

where \mathbf{N}_1 (\mathbf{N}_2) is a diagonal matrix that rescales each row (column) of \mathbf{W} (\mathbf{H}) by the number of

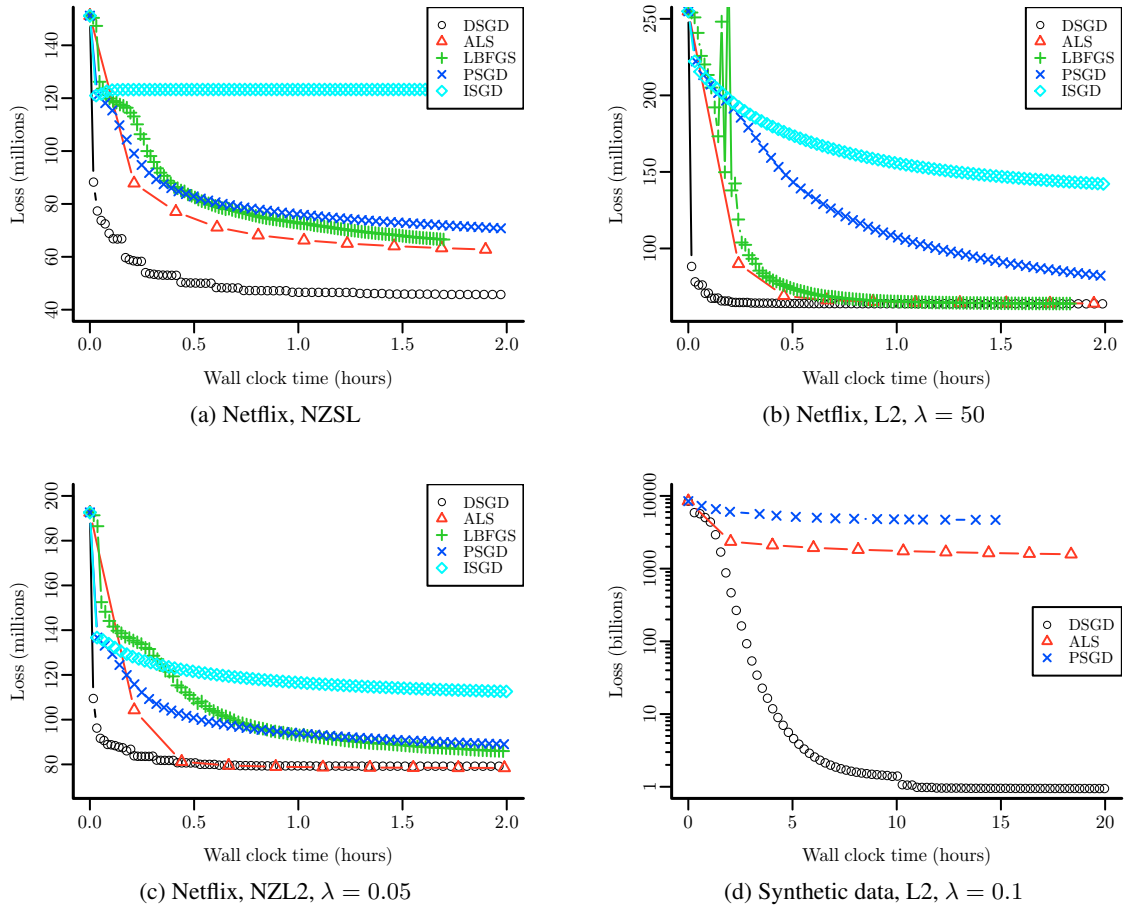
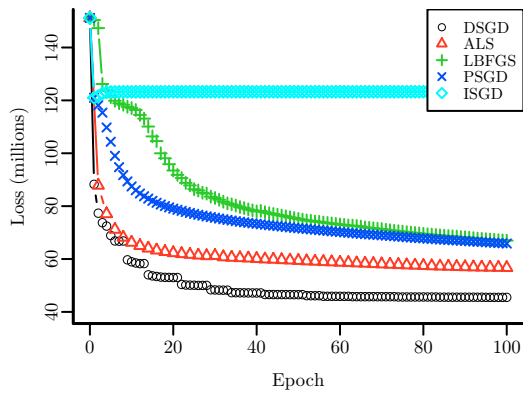
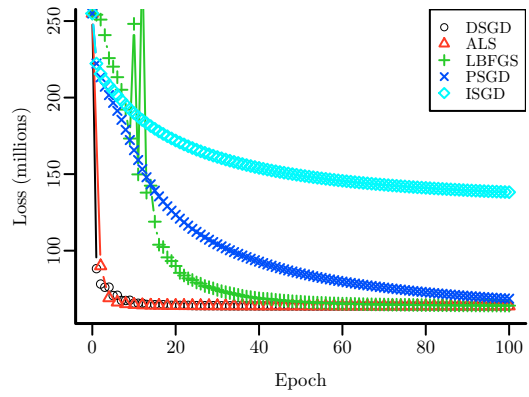


Figure 3: Performance in terms of wall-clock time

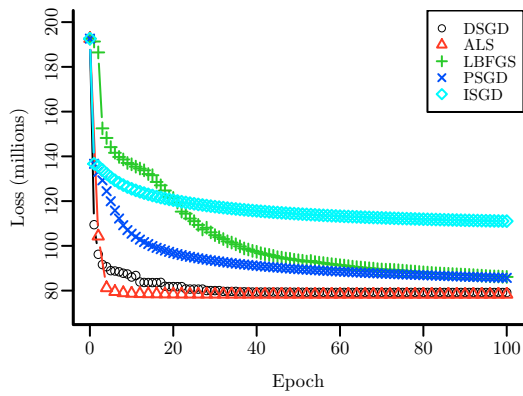
nonzero entries of \mathbf{Z} in that row (column), and $\|\mathbf{A}\|_F$ denotes the Frobenius norm of a matrix \mathbf{A} (see App. C). L_{NZL2} has been used successfully on the Netflix data [18], and L_{GKL} has applications in text indexing [15]. We used “principled” values of λ throughout. E.g., we used values of $\lambda = 50$ and $\lambda = 0.05$ for L_{L2} and L_{NZL2} on Netflix data, respectively, and $\lambda = 0.1$ for L_{L2} on synthetic data, the former values because they yielded the best movie recommendations on held-out data, and the latter because it is “natural” in that the resulting minimum-loss factors correspond to the “maximum a posteriori” Bayesian estimator of \mathbf{W} and \mathbf{H} under the Gaussian-based procedure used to generate the synthetic data.



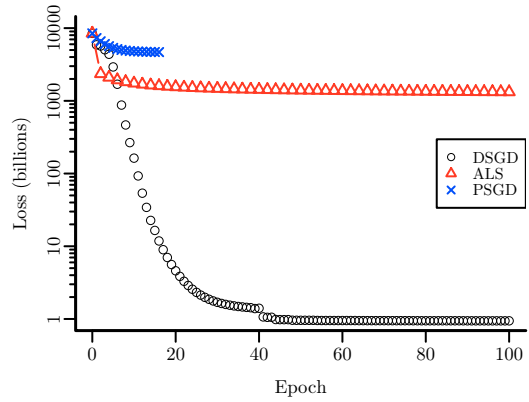
(a) Netflix, NZSL



(b) Netflix, L2, $\lambda = 50$



(c) Netflix, NZL2, $\lambda = 0.05$



(d) Synthetic data, L2, $\lambda = 0.1$

Figure 4: Performance in terms of epochs

7.2. Relative Performance

We first evaluated the relative performance of the matrix factorization algorithms. For various loss functions and datasets, we ran 100 epochs—i.e, scans of the data matrix—with each algorithm and measured the elapsed wall-clock time, as well as the value of the loss after every epoch. We used 64-way distributed processing on 8 nodes (with 8 concurrent map tasks per node).⁵

Representative results are given in Figure 3, which displays the achieved loss as a function of wall clock time, and in Figure 4, which displays the loss as a function of the number of epochs.

As can be seen from Figure 3, DSGD converges about as fast as—or faster than—alternative methods, with the DSGD performing markedly better for L_{NZSL} loss over the Netflix data and for L_{L2} loss over both the Netflix dataset and the large synthetic dataset. Comparing the various SGD-based approaches, we observe that ISGD and PSGD exhibit consistently inferior performance, and offer the following explanation. The matrix-factorization problem is “non-identifiable” in that the loss function has many global minima that correspond to widely different values of (\mathbf{W}, \mathbf{H}) . Averages of (\mathbf{W}, \mathbf{H}) values from different partitions, as computed by ISGD and PSGD, do not correspond to good overall solutions, and the algorithms may not converge to a local minimum of the global loss, or may converge very slowly. E.g., in the example of Figure 2, ISGD would converge to a point on the line between the red and blue minima. Of the ISGD and PSGD algorithms, it is not surprising that ISGD has the worst convergence behavior; recall that ISGD computes a local optimum of each partition of the dataset via SGD, and averages only once after all 100 epochs.⁶ PSGD improves on ISGD by averaging parameters after every epoch, but it is still outperformed by most other approaches. DSGD performs best; its usage of stratification instead of averaging significantly improves convergence speed. For the remainder of our discussion, we focus on DSGD as the best SGD-based algorithm and compare it with L-BFGS and ALS. L-BFGS is clearly inferior to the other two algorithms. Indeed, we do not give results for L-BFGS in Figures 3d or 4d because its centralized parameter-update step ran out of memory when faced with very large data. In the other three experiments, L-BFGS is able to execute more epochs per unit time than the other algorithms—e.g., for the L_{NZSL} experiment, DSGD ran 43 epochs, ALS ran 10 epochs, and L-BFGS ran 61 epochs in the first hour—but the per-epoch decrease in loss is relatively small. In general, the foregoing differences in runtime are explained by different computational costs (highest for ALS, which has to solve $m + n$ least-squares problems per epoch) and synchronization costs (highest for PSGD, which has to average all parameters in each epoch).

ALS achieves performance roughly comparable to DSGD for L_{NZL2} loss over Netflix data, taking about 15 more minutes to get within the vicinity of the minimal loss but ultimately yielding a slightly lower loss (about 1% less than that of DSGD) after two hours. ALS is clearly inferior to DSGD, however, in the other three experiments. The differences are more noticeable in Figure 3 than in Figure 4, since they reflect the larger execution time per epoch for ALS as indicated above; see

⁵Note that for all approaches but ISGD and ALS, 64-way distributed processing is excessive for the Netflix data; the execution time is dominated by latencies. We nonetheless used 64-way processing to get a consistent view over datasets of various sizes.

⁶The intermediate points shown in Figure 3 have been obtained by pausing ISGD after every epoch in order to average parameters and compute the loss. The time to do so is not included in the wall-clock time of ISGD.

also Figure 5a. For the first experiment— L_{NZSL} loss over Netflix data—the lack of a regularization term makes the factorization difficult for ALS because the search space is relatively large and there are many equivalent solutions. Specifically, we observed that ALS conducted large moves through the parameter space; the factors grew without bound. On the synthetic data with L_{L2} , ALS is very effective in the first epoch, but then converges slowly. We have observed similar behavior on very small matrices, with ALS getting stuck when moving along “valleys” of small loss in which both W and H change simultaneously. DSGD does not suffer from these problems and has superior convergence properties.

In summary, the overall performance of DSGD was consistently more stable than that of alternative algorithms. The speed of convergence was comparable or faster.

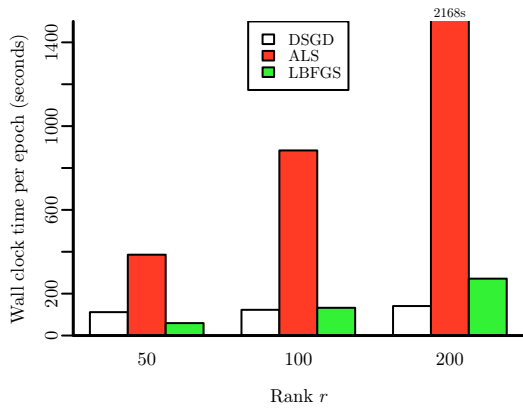
7.3. Scalability

Next, we studied various scalability issues in our shared-nothing MapReduce environment. We examined the effect of scaling up the approximation rank r , and then explored the scalability of DSGD on a Hadoop cluster by scaling up the dataset size, the number of cores, and then both the dataset size and number of cores. Overall, the gradient descent methods scale better with increasing rank than ALS, and DSGD has good scalability properties on Hadoop, provided that the amount of data processed per core does not become so small that system overheads start to dominate.

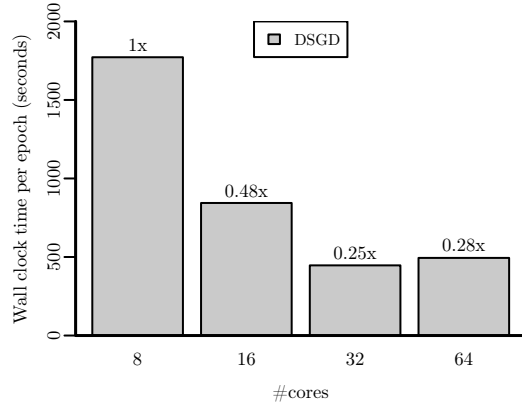
To explore the effect of increasing the approximation rank, we used the Netflix data with L_{NZL2} loss and the same R-cluster setup as before (64 cores, 8 nodes); note that the value $r = 50$ corresponds to our relative-performance experiments. The results are displayed in Figure 5a. As observed previously, ALS is significantly slower than DSGD and L-BFGS. ALS spends most of its time on constructing and solving (in the least-squares sense) systems of linear equations. Since the number of both equations and variables increases with rank—construction is $O(Nr^2)$, solving is $O((m+n)r^3)$ —the performance degrades significantly as r increases. L-BFGS performs centralized updates of the factors; these centralized updates become a bottleneck as the rank (and thus factor size) increases. The impact of increased rank on DSGD appears rather mild, mainly because factors are fully distributed. As the rank increases further and gradient estimation becomes the major bottleneck, we expect to see a more pronounced increase in runtime.

Our remaining experiments focus on the performance of DSGD on the Hadoop cluster. Figure 5b depicts scale-up results as we repeatedly double the number of cores while keeping the data size constant at 6.4B entries. Figure 5c plots the runtime per epoch as we repeatedly quadruple the data (while keeping the number of cores constant at 32), and Figure 5d shows scale-out results, in which we scale both data and cores simultaneously.

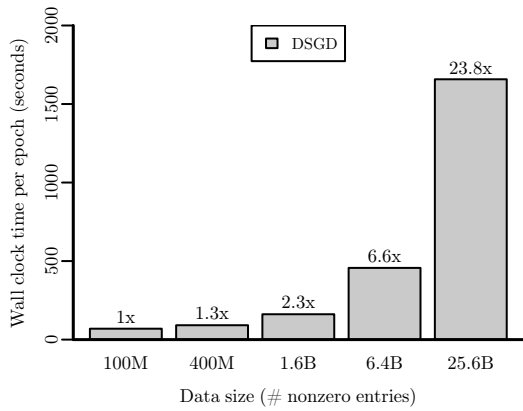
As can be seen in Figure 5b, DSGD initially achieves roughly linear speed-up as the number of cores is repeatedly doubled, up to 32 cores. After this point, speed-up performance starts to degrade. The reason for this behavior is that, when the number of cores becomes large, the amount of data processed per core becomes small—e.g., 64-way DSGD requires 64^2 blocks, so that the amount of data per block is only $\approx 25\text{MB}$. The actual time to execute DSGD on the data becomes negligible, and the overall processing times become dominated by Hadoop overheads, especially the time required to spawn a task. (Hadoop is designed for tasks that run at least on the order of



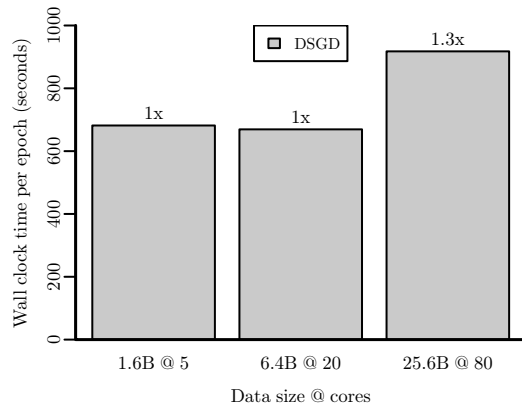
(a) Effect of r (R@64)



(b) Increasing cores (Hadoop, 6.4B entries)



(c) Increasing data (Hadoop @ 32)



(d) Increasing data and cores (Hadoop)

Figure 5: Scalability

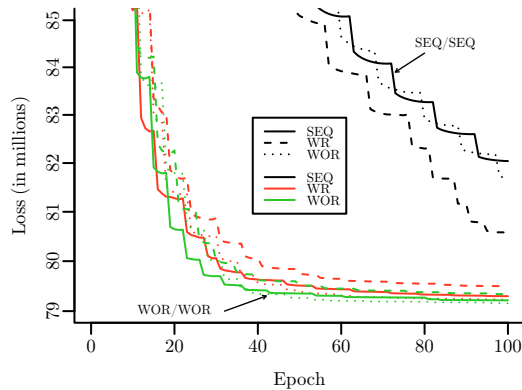


Figure 6: Effect of stratum selection (line color) and training sequence (line type)

minutes.) A similar phenomenon can be seen in Figure 5c, where the elapsed time is sublinear in the data size for small datasets ($\leq 1.6\text{B}$ entries); for larger datasets, overheads no longer mask performance, and the runtime increases linearly with dataset size. In our scale-out experiments (Figure 5d), the impact of Hadoop overheads is more muted, since scaling up the data size offsets the overhead effect caused by increasing the number of cores. E.g., the processing time initially remains constant as the dataset size and number of cores are each scaled up by a factor of 4, with the overall runtime increasing by a modest 30% as we scale to very large datasets on large clusters. The foregoing overhead effects can potentially be ameliorated by improving scheduling in Hadoop or using an alternative parallel runtime system such as Spark [26].

7.4. Selection Schemes

Finally, we evaluated the impact of different strategies for selecting both strata and training sequences. The runtime cost for the various alternatives are comparable, but the speed of convergence differs significantly. Figure 6 shows exemplary results for 64-way SGD on the Netflix data with L_{NZL2} , where we plot all combinations of options for training-point and stratum selection. Sequential stratum selection performed worst: the curves corresponding to this scheme cluster at the upper right of the plot, whereas the curves for the randomized strategies cluster at the lower left; i.e., any form of randomized stratum selection helped significantly. For randomized selection schemes, the WOR strategy for selecting training points yielded the best results. Overall, WOR selection for both training points and strata ensures a good balance between randomization and processing many different data points.

7.5. Other Loss Functions

To show that DSGD can be applied to a variety of loss functions, we implemented L_{GKL} —a loss function that is used for nonnegative matrix factorization and that does *not* ignore zeros in the \mathbf{V} matrix—and ran the resulting factorization algorithm on the Netflix data, along with the loss-specific

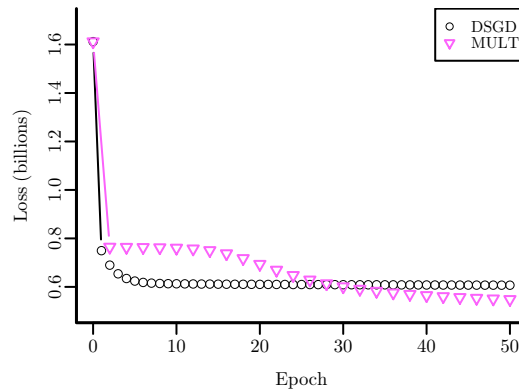


Figure 7: GKL on Netflix data (R@64)

MULT algorithm of Das et al. [12]. DSGD performs respectably compared to MULT (Figure 7), reaching the vicinity of the minimum loss more rapidly—roughly 7 epochs for DSGD versus 27 epochs for MULT—and achieving an ultimate loss that is only modestly greater than that of MULT. Our implementation is a first cut; we are currently refining it further.

8. Conclusions

We introduced DSGD, a novel algorithm for large-scale matrix factorization. DSGD is fully distributed and can handle matrices with millions of rows, millions of columns, and billions of nonzero entries. In contrast to most alternative algorithms, DSGD is generic in that it supports a wide variety of loss functions that arise in practice. Our experiments indicate that DSGD is on par or faster than specialized algorithms in terms of runtime, convergence properties, and memory requirements.

References

- [1] R. Albright, J. Cox, D. Duling, A. Langville, and C. Meyer. Algorithms, initializations, and convergence for the nonnegative matrix factorization. Technical Report Math 81706, NCSU, 2006.
- [2] Apache Hadoop. <https://hadoop.apache.org>.
- [3] S. Asmussen. *Applied Probability and Queues*. Springer, 2nd edition, 2003.
- [4] S. Asmussen and P. W. Glynn. *Stochastic Simulation: Algorithms and Analysis*. Springer, 2007.
- [5] J. Bennett and S. Lanning. The Netflix prize. In *KDD Cup and Workshop*, 2007.

- [6] C. M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer, 2007.
- [7] L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In *NIPS*, volume 20, pages 161–168, 2008.
- [8] R. H. Byrd, P. Lu, J. Nocedal, and C. Zhu. A limited memory algorithm for bound constrained optimization. *SIAM J. Sci. Comput.*, 16(5):1190–1208, 1995.
- [9] Y. S. Chow and H. Teicher. *Probability Theory: Independence, Interchangeability, Martingales*. Springer, 2nd edition, 1988.
- [10] K. L. Chung. *A Course in Probability Theory*. Elsevier, third edition, 2001.
- [11] A. Cichocki and R. Zdunek. Regularized alternating least squares algorithms for non-negative matrix/tensor factorization. In *ISNN '07: Proc. of the 4th international symposium on Neural Networks*, pages 793–802, 2007.
- [12] A. S. Das, M. Datar, A. Garg, and S. Rajaram. Google news personalization: scalable online collaborative filtering. In *WWW*, pages 271–280, 2007.
- [13] S. Das, Y. Sismanis, K. S. Beyer, R. Gemulla, P. J. Haas, and J. McPherson. Ricardo: Integrating R and Hadoop. In *SIGMOD*, pages 987–998, 2010.
- [14] K. B. Hall, S. Gilpin, and G. Mann. MapReduce/Bigtable for distributed optimization. In *NIPS LCCC Workshop*, 2010.
- [15] T. Hofmann. Probabilistic latent semantic indexing. In *SIGIR*, pages 50–57, 1999.
- [16] T. Hofmann. Latent semantic models for collaborative filtering. *ACM Trans. Inf. Syst.*, 22(1):89–115, 2004.
- [17] D. Jiang, B. C. Ooi, L. Shi, and S. Wu. The performance of MapReduce: An in-depth study. *PVLDB*, 3(1):472–483, 2010.
- [18] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *IEEE Computer*, 42(8):30–37, 2009.
- [19] H. J. Kushner and G. Yin. *Stochastic Approximation and Recursive Algorithms and Applications*. Springer, 2nd edition, 2003.
- [20] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- [21] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *NIPS*, pages 556–562, 2000.

- [22] C. Liu, H.-c. Yang, J. Fan, L.-W. He, and Y.-M. Wang. Distributed nonnegative matrix factorization for web-scale dyadic data analysis on mapreduce. In *WWW*, pages 681–690, 2010.
- [23] G. Mann, R. McDonald, M. Mohri, N. Silberman, and D. Walker. Efficient large-scale distributed training of conditional maximum entropy models. In *NIPC*, pages 1231–1239, 2009.
- [24] R. McDonald, K. Hall, and G. Mann. Distributed training strategies for the structured perceptron. In *HLT*, pages 456–464, 2010.
- [25] A. P. Singh and G. J. Gordon. A unified view of matrix factorization models. In *ECML PKDD*, pages 358–373, 2008.
- [26] M. Zaharia, N. M. M. Chowdhury, M. Franklin, S. Shenker, and I. Stoica. Spark: Cluster computing with working sets. Technical Report UCB/EECS-2010-53, EECS Department, University of California, Berkeley, May 2010.
- [27] Y. Zhou, D. Wilkinson, R. Schreiber, and R. Pan. Large-scale parallel collaborative filtering for the Netflix Prize. In *AAIM*, pages 337–348, 2008.
- [28] M. A. Zinkevich, M. Weimer, A. J. Smola, and L. Li. Parallelized stochastic gradient descent. In *NIPS*, pages 2595–2603, 2010.

A. MapReduce Algorithms for Matrix Factorization

We review some algorithms for finding an $m \times r$ matrix \mathbf{W} and an $r \times n$ matrix \mathbf{H} such that $\mathbf{V} \approx \mathbf{WH}$ for a given $m \times n$ input matrix \mathbf{V} , in the sense of minimizing a specified loss function $L(\mathbf{V}, \mathbf{WH})$ computed over N training points. We focus on algorithms that have been shown to work in a distributed setting—in which d tasks can perform computations in parallel—on very large matrices. All algorithms are of an iterative nature. They start with some initial factors \mathbf{W}_0 and \mathbf{H}_0 , which are repeatedly updated until some convergence criteria is met. Interesting properties of such algorithms include supported loss functions, convergence properties, whether they support non-negativity or box constraints, time and space complexity, and how distribution is achieved. An overview of the algorithms is given in Table 1.

A.1. Specialized Algorithms

The specialized algorithms below are each designed for a certain class of loss functions. In all cases, distributed processing is achieved by partitioning the input matrix and splitting up linear algebra computations across nodes.

Table 1: MapReduce Algorithms for Matrix Factorization*

Algorithm	Loss	Regularizer	Partitioning	[Scans	M	R	Time] per Iteration
ALS [27]	SL, NZSL	L2, NZL2	V & W by rows, V & H by columns (as above)	2	2	0	$O(d^{-1}[Nr^2 + (m+n)r^3])$
EM [12]	SL	L2	V blocked rect., W & H conformingly	2	2	0	$O(d^{-1}[Nr + (m+n)r^2] + r^3)$
MULT [22]	GKL	-	V blocked rect., W & H conformingly	1	1	1	$O(d^{-1}Nr)$
	GKL	-	V blocked rect., W & H conformingly	2	2	2	$O(d^{-1}Nr)$
	SL	-	(as above)	2	2	2	$O(d^{-1}[Nr + (m+n)r^2])$
GD [13]	diff.	diff.	V blocked rect., W & H conformingly	1	1	1	$O(d^{-1}N)T_{L'} + T_{\text{UPD}}$
PSGD	diff.	diff.	V partitioned arbitrary, W & H replicated	1	1	1	$O(d^{-1}N)T_{L'} + O((m+n)r)$
DSGD	diff.	diff.	V blocked square, W & H conformingly	1	d^{**}	0	$O(d^{-1}N)T_{L'}$
	NZSL	L2, NZL2	(as above)	1	d^{**}	0	$O(d^{-1}Nr)$

*The Scans, M, and R columns refer to the number of data scans, map jobs, and reduce jobs per iteration of the algorithm, respectively. "diff." means that the algorithm can handle arbitrary differentiable loss functions. $T_{L'}$ and T_{UPD} denote the time required by gradient-based methods to compute a local gradient L'_{ij} and the time to update all parameters once gradients have been computed, respectively. Some algorithms have multiple entries; these entries may refer to different losses or regularizers (MULT), to specific loss or regularizers on which improved runtime properties are achieved (ALS), or to an example (DSGD).

**One Map job per subepoch, each processing d blocks. In our Hadoop implementation, we use customized locking strategies to run just one Map job per iteration.

Alternating Least Squares. In its standard form, the method of alternating least squares optimizes

$$L_{\text{SL}} = \sum_{i,j} (V_{ij} - [\mathbf{W}\mathbf{H}]_{ij})^2.$$

The method alternates between finding the best value for \mathbf{W} given \mathbf{H} , and finding the best value of \mathbf{H} given \mathbf{W} . This amounts to computing the least squares solutions to the following systems of linear equations

$$\begin{aligned} \mathbf{V}_{i*} - \underline{\mathbf{W}}_{i*} \mathbf{H}_n &= \mathbf{0}, \\ \mathbf{V}_{*j} - \mathbf{W}_{n+1} \underline{\mathbf{H}}_{*j} &= \mathbf{0}, \end{aligned}$$

where the unknown variable is underlined. This specific form suggests that each row of \mathbf{W} can be updated by accessing only the corresponding row in the data matrix \mathbf{V} , while each column in \mathbf{H} can be updated by accessing the corresponding column in \mathbf{V} . This facilitates distributed processing; see below. The equations can be solved using a method of choice. We obtain

$$\begin{aligned} \mathbf{W}_{n+1}^T &\leftarrow (\mathbf{H}_n \mathbf{H}_n^T)^{-1} \mathbf{H}_n \mathbf{V}^T, \\ \mathbf{H}_{n+1} &\leftarrow (\mathbf{W}_{n+1}^T \mathbf{W}_{n+1})^{-1} \mathbf{W}_{n+1}^T \mathbf{V}. \end{aligned}$$

for the unregularized loss shown above. When an additional L_2 regularization term of form $\lambda(\|\mathbf{W}\|_{\text{F}}^2 + \|\mathbf{H}\|_{\text{F}}^2)$ is added, we obtain

$$\mathbf{W}_{n+1}^T \leftarrow (\mathbf{H}_n \mathbf{H}_n^T + \lambda \mathbf{I})^{-1} \mathbf{H}_n \mathbf{V}^T, \quad (19a)$$

$$\mathbf{H}_{n+1} \leftarrow (\mathbf{W}_{n+1}^T \mathbf{W}_{n+1} + \lambda \mathbf{I})^{-1} \mathbf{W}_{n+1}^T \mathbf{V}. \quad (19b)$$

Since the update term of \mathbf{H}_{n+1} depends on \mathbf{W}_{n+1} , the input matrix has to be processed twice to update both factor matrices.

In contrast to SVD, ALS does not produce an orthogonal factorization and it might get stuck in local minima. However, ALS can handle a wide range of variations for which SVD is not applicable, but which are important in practice. Examples include non-negativity constraints [1], sparsity constraints [1, 11], weights [11], regularization [11, 27], or the restriction to nonzero entries [27]. In general, ALS is applicable when the loss function is quadratic in both \mathbf{W} and \mathbf{H} .

Zhou et al. [27] proposed a distributed version of ALS for L_{NZSL} , see Eq. (1). We first describe a variation for L_{SL} , and then outline how to optimize L_{NZSL} . In both cases, the algorithm runs two Map-only jobs per iteration, one to compute \mathbf{W}_{n+1} and one to compute \mathbf{H}_{n+1} . Each of the jobs uses a different partitioning of \mathbf{V} , either by rows or by columns. For example, we use the row partitioning to compute \mathbf{W}_{n+1} : Each mapper reads a set of rows of \mathbf{V} , the corresponding rows of \mathbf{W}_n , and the entire matrix \mathbf{H}_n . The i th mapper then solves the part of equation (19a) that concerns its rows:

$$\mathbf{W}_{n+1,i*}^T \leftarrow \underbrace{(\mathbf{H}_n \mathbf{H}_n^T + \lambda \mathbf{I})^{-1}}_{\text{coefficient matrix}} \underbrace{\mathbf{H}_n \mathbf{V}_{i*}^T}_{\text{right-hand side}},$$

where \mathbf{W}_{n+1,i^*} denotes the rows of \mathbf{W}_{n+1} read by the i th mapper, and similarly for \mathbf{V}_{i^*} . For L_{SL} , the coefficient matrix is shared across rows, which allows us to reuse computation (e.g., a QR factorization of the coefficient matrix). The matrix \mathbf{H}_{n+1} is computed analogously. The overall time complexity is $O(d^{-1}[Nr + (m+n)r^2] + r^3)$ time.⁷ For L_{NZSL} , we modify the algorithm so that it uses a different coefficient matrix for each system of equations, i.e., for each mapper. Intuitively, we remove equations that correspond to zero entries of \mathbf{V} from the least-squares problem. This is achieved by using

$$\mathbf{W}_{n+1,i^*}^{\text{T}} \leftarrow \underbrace{(\mathbf{H}_n^{(i)}[\mathbf{H}_n^{(i)}]^{\text{T}} + \lambda\mathbf{I})^{-1}}_{\text{coefficient matrix}} \underbrace{\mathbf{H}_n \mathbf{V}_{i^*}^{\text{T}}}_{\text{right-hand side}},$$

where $\mathbf{H}_n^{(i)}$ consists of just the columns of \mathbf{H}_n that have non-zero entries in the respective column of \mathbf{V}_{i^*} . The computation of \mathbf{H}_{n+1} is analogous. The time complexity increases to $O(d^{-1}[Nr^2 + (m+n)r^3])$.⁸

Expectation Maximization (EM). Hofmann et al. [15, 16] proposed an EM algorithm to minimize the KL divergence L_{KL} in the context of “probabilistic latent semantic analysis” (pLSA). As we discuss below, the algorithm can be seen as a matrix factorization algorithm. Let \mathbf{V} be non-negative and $\sum_{i,j} \mathbf{V}_{ij} = 1$.⁹ Then, \mathbf{V} corresponds to a probability distribution over pairs (i, j) . pLSA factors this probability distribution as follows

$$\Pr[i, j] \approx \sum_z \Pr[z] \Pr[i | z] \Pr[j | z], \quad (20)$$

where $z \in \{z_1, \dots, z_r\}$ is a latent variable and follows a multinomial distribution over r topics. If we identify $\Pr[i, j] = \mathbf{V}_{ij}$, $\Pr[z] = \mathbf{Z}_{zz}$ (where \mathbf{Z} is a diagonal $r \times r$ matrix), $\Pr[i | z] = \mathbf{W}'_{iz}$, and $\Pr[j | z] = \mathbf{H}'_{zj}$, we obtain the following equivalent matrix factorization

$$\mathbf{V} \approx \mathbf{W}' \mathbf{Z} \mathbf{H}'.$$

Note that \mathbf{W}' (\mathbf{H}') describes a conditional probability distribution; thus each of its columns (rows) sums to 1.

Model fitting is performed as follows. In the E-step, we compute the probability $\Pr[z | i, j]$ that entry (i, j) is explained by topic z :

$$\Pr[z | i, j] = \frac{\Pr[z] \Pr[i | z] \Pr[j | z]}{\sum_{z'} \Pr[z'] \Pr[i | z'] \Pr[j | z']}. \quad (21)$$

⁷The right-hand sides can be constructed in $O(Nr)$ time. The coefficient matrix can be constructed in $O((m+n)r^2)$ time, and reduced to an upper-triangular form in $O(r^3)$ time (not parallel). For each of the $m+n$ equation systems, back-substitution takes $O(r^2)$ time.

⁸We can construct all coefficient matrices in $O(Nr^2)$ time. Solving each system takes $O(r^3)$ time.

⁹If $R = \sum_{i,j} \mathbf{V}_{ij} \neq 1$, we normalize \mathbf{V} by dividing each element by R .

In the M-step, the parameters are updated. Using normalization constants $K_z = \sum_{i,j} \mathbf{V}_{ij} \Pr[z | i, j]$, we set

$$\Pr[j | z] = \frac{1}{K_z} \sum_i \mathbf{V}_{ij} \Pr[z | i, j] \quad (22a)$$

$$\Pr[i | z] = \frac{1}{K_z} \sum_j \mathbf{V}_{ij} \Pr[z | i, j] \quad (22b)$$

$$\Pr[z] = K_z \quad (22c)$$

It can be shown that the KL divergence between the distribution \mathbf{V} and the fitted distribution (20) is non-increasing in every EM iteration. The EM algorithm converges to a stationary point of L_{KL} .

We now transform the EM algorithm into the language of linear algebra. This will allow us to uncover similarities between the EM algorithm and the multiplicative update rules described later, and also facilitates exposition of distributed EM. In what follows, we set $\mathbf{W} = \mathbf{W}'$ and $\mathbf{H} = \mathbf{Z}\mathbf{H}'$, i.e., we factor \mathbf{Z} into the parameter matrix \mathbf{H} . \mathbf{Z} can be readily factored out after convergence, if desired. The E-step (21) becomes

$$\Pr[z | i, j] = \left[\frac{\mathbf{W}_{*z} \mathbf{H}_{z*}}{\mathbf{W}\mathbf{H}} \right]_{ij}, \quad (23)$$

where division is performed element-wise. Let $\mathbf{K}_n = \text{diag}(K_{n,1}, \dots, K_{n,r})$ be the matrix of the normalization constants used in the $(n+1)$ st M-step. Inserting (23) directly into the equations (22), we obtain the following update rules:

$$\mathbf{W}'_{n+1} \leftarrow [\mathbf{W}_n \circ (\mathbf{V}/\mathbf{W}_n \mathbf{H}_n) \mathbf{H}_n^{\text{T}}] \mathbf{K}_n^{-1} \quad (24a)$$

$$\mathbf{H}'_{n+1} \leftarrow \mathbf{K}_n^{-1} [\mathbf{W}_n^{\text{T}} (\mathbf{V}/\mathbf{W}_n \mathbf{H}_n) \circ \mathbf{H}_n] \quad (24b)$$

$$\mathbf{Z}_{n+1} = \mathbf{K}_n, \quad (24c)$$

where \circ denotes element-wise multiplication. Note that \mathbf{K}_n^{-1} is easy to compute as \mathbf{K}_n is a diagonal matrix. Since all resulting matrices describe (conditional) probability distributions, they are normalized appropriately. By construction, we have

$$\begin{aligned} \mathbf{K}_n &= \text{diag}(\text{colSums} [\mathbf{W}_n \circ (\mathbf{V}/\mathbf{W}_n \mathbf{H}_n) \mathbf{H}_n^{\text{T}}]) \\ &= \text{diag}(\text{rowSums} [\mathbf{W}_n^{\text{T}} (\mathbf{V}/\mathbf{W}_n \mathbf{H}_n) \circ \mathbf{H}_n]), \end{aligned}$$

where $\text{colSums}[\mathbf{A}]_j = \sum_i \mathbf{A}_{ij}$ and $\text{rowSums}[\mathbf{A}]_i = \sum_j \mathbf{A}_{ij}$ for a matrix \mathbf{A} . If we compute \mathbf{W}_{n+1} and \mathbf{H}_{n+1} directly, we arrive at the following final update rules

$$\tilde{\mathbf{W}}_{n+1} \leftarrow \mathbf{W}_n \circ (\mathbf{V}/\mathbf{W}_n \mathbf{H}_n) \mathbf{H}_n^{\text{T}} \quad (25a)$$

$$\mathbf{W}_{n+1} \leftarrow \tilde{\mathbf{W}}_{n+1} \text{diag}(1/\text{colSums}[\tilde{\mathbf{W}}_{n+1}]) \quad (25b)$$

$$\mathbf{H}_{n+1} \leftarrow \mathbf{W}_n^{\text{T}} (\mathbf{V}/\mathbf{W}_n \mathbf{H}_n) \circ \mathbf{H}_n, \quad (25c)$$

where $\tilde{\mathbf{W}}_{n+1}$ is an intermediate variable that is normalized to obtain \mathbf{W}_{n+1} .

Das et al. [12] show how to distribute the EM algorithm using MapReduce. The idea is to partition matrix \mathbf{V} into $d_1 \times d_2$ blocks, \mathbf{W} into d_1 conforming blocks, and \mathbf{H} into d_2 conforming blocks. Only one MapReduce job is needed to perform one EM iteration. To make this work, we first push the normalization of $\tilde{\mathbf{W}}_{n+1}$ in equations (25) into the $(n+2)$ nd iteration:

$$\tilde{\mathbf{W}}_{n+1} \leftarrow \tilde{\mathbf{W}}_n \mathbf{K}_n^{-1} \circ (\mathbf{V} / \tilde{\mathbf{W}}_n \mathbf{K}_n^{-1} \mathbf{H}_n) \mathbf{H}_n^T \quad (26a)$$

$$\mathbf{K}_{n+1} \leftarrow \text{diag}(\text{colSums}[\tilde{\mathbf{W}}_{n+1}]) \quad (26b)$$

$$\mathbf{H}_{n+1} \leftarrow [\tilde{\mathbf{W}}_n \mathbf{K}_n^{-1}]^T (\mathbf{V} / \tilde{\mathbf{W}}_n \mathbf{K}_n^{-1} \mathbf{H}_n) \circ \mathbf{H}_n. \quad (26c)$$

where $\tilde{\mathbf{W}}_0 = \mathbf{W}_0$ and $\mathbf{K}_0 = \text{diag}(\text{colSums} \mathbf{W}_0)$. Each mapper reads a block of \mathbf{V} , the corresponding blocks of $\tilde{\mathbf{W}}_n$ and \mathbf{H}_n , and the entire matrix \mathbf{K}_n . For each nonzero element of \mathbf{V} , the mapper computes the quantity

$$p_{z|ij} = \Pr[z | i, j] = \mathbf{W}_{n,iz} \mathbf{K}_{n,zz}^{-1} (\mathbf{V}_{ij} / \mathbf{W}_{n,i*} \mathbf{K}_n^{-1} \mathbf{H}_{n,*j}) \mathbf{H}_{n,zj}^T,$$

and outputs pairs $(i, p_{z|ij})$, $(j, p_{z|ij})$, $(z, p_{z|ij})$. Thus there is one *group* per row, one per column, and one per topic. The reducer for row i computes $\tilde{\mathbf{W}}_{n+1,i*}$ using the transformation:

$$\begin{aligned} \tilde{\mathbf{W}}_{n+1,iz} &\leftarrow \tilde{\mathbf{W}}_{n,iz} \mathbf{K}_{n,zz}^{-1} \circ (\mathbf{V}_{i*} / \tilde{\mathbf{W}}_{n,i*} \mathbf{K}_n^{-1} \mathbf{H}_n) \mathbf{H}_{n,z*}^T \\ &= \sum_j \tilde{\mathbf{W}}_{n,iz} \mathbf{K}_{n,zz}^{-1} \circ (\mathbf{V}_{ij} / \tilde{\mathbf{W}}_{n,i*} \mathbf{K}_n^{-1} \mathbf{H}_{n,*j}) \mathbf{H}_{n,zj}^T \\ &= \sum_j p_{z|ij}. \end{aligned}$$

Similarly,

$$\begin{aligned} \mathbf{H}_{n+1,zj} &\leftarrow \sum_i p_{z|ij} \\ \mathbf{K}_{n+1,zz} &\leftarrow \sum_{i,j} p_{z|ij} \end{aligned}$$

Thus reducers simply sum up the entries within each group. Since summation is distributive, preaggregation can be performed in a combine step at each mapper. To output the result in blocked form, groups are assigned to reducers according to the blocking—e.g., a single reducer processes all rows (groups) of the first block of $\tilde{\mathbf{W}}_{n+1}$. Assuming $m, n = O(N)$, the overall time complexity per iteration is $O(d^{-1}Nr)$.

Multiplicative updates (MULT). Lee et.al [20] proposed multiplicative update rules for non-negative matrix factorization under L_{GKL} . Later [21], they refined the L_{GKL} rules and developed rules for L_{SL} . The refined update rules can be seen as a rescaled version of gradient descent, where the step sizes are computed individually for each parameter. In all cases, the rules are multiplicative in that each factor gets multiplied by some update term, which varies from method to method. Each

iteration is non-increasing in the loss, and the factor matrices converge to a stationary point of the loss function.

We first consider the rules for GKL given in [20]:

$$\tilde{\mathbf{W}}_{n+1} \leftarrow \mathbf{W}_n \circ (\mathbf{V}/\mathbf{W}_n \mathbf{H}_n) \mathbf{H}_n^T \quad (27a)$$

$$\mathbf{W}_{n+1} \leftarrow \tilde{\mathbf{W}}_{n+1} \text{diag}(1/\text{colSums}[\tilde{\mathbf{W}}_{n+1}]) \quad (27b)$$

$$\mathbf{H}_{n+1} \leftarrow \mathbf{W}_{n+1}^T (\mathbf{V}/\mathbf{W}_{n+1} \mathbf{H}_n) \circ \mathbf{H}_n, \quad (27c)$$

These are almost the update rules (25) of EM, but \mathbf{W}_{n+1} is used instead of \mathbf{W}_n when computing \mathbf{H}_{n+1} . As a consequence, updates of \mathbf{W} and \mathbf{H} cannot be performed simultaneously anymore, and two scans of \mathbf{V} are required. The refined rules for GKL are

$$\begin{aligned} \mathbf{W}_{n+1} &\leftarrow \mathbf{W}_n \circ (\mathbf{V}/\mathbf{W}_n \mathbf{H}_n) \mathbf{H}_n^T \text{diag}(1/\text{rowSums}[\mathbf{H}_n]), \\ \mathbf{H}_{n+1} &\leftarrow \mathbf{H}_n \circ \text{diag}(1/\text{colSums}[\mathbf{W}_{n+1}]) \mathbf{W}_{n+1}^T (\mathbf{V}/\mathbf{W}_{n+1} \mathbf{H}_n). \end{aligned}$$

These update rules can be seen as symmetric versions of (27); they work better in practice. Thus, when we refer to MULT for L_{GKL} , we refer to the refined rules above. The time complexity remains $O(Nr)$. The rules for L_{SL} are given by:

$$\begin{aligned} \mathbf{H}_{n+1} &\leftarrow \mathbf{H}_n \circ \frac{\mathbf{W}_n^T \mathbf{V}}{\mathbf{W}_n^T \mathbf{W}_n \mathbf{H}_n}, \\ \mathbf{W}_{n+1} &\leftarrow \mathbf{W}_n \circ \frac{\mathbf{V} \mathbf{H}_{n+1}^T}{\mathbf{W}_n \mathbf{H}_{n+1} \mathbf{H}_{n+1}^T}. \end{aligned}$$

Sparsity of \mathbf{V} is readily exploited: $\mathbf{W}^T \mathbf{V}$ and $\mathbf{V} \mathbf{H}_{n+1}^T$ each can be computed in a single scan of the nonzero elements of \mathbf{V} . The overall time complexity is $O(Nr + (m+n)r^2)$.

Liu et al. [22] give MapReduce versions of MULT; the underlying ideas are the same as for distributed EM [12]. Minor modifications are needed to match the refined rules or the L_{SL} rules.

A.2. Generic Algorithms

Generic algorithms can be used to solve our generalized problem statement. We first discuss gradient descent, which has been applied successfully to large-scale matrix factorization. We then summarize some recent work in the area of distributed SGD.

Gradient descent (GD). Gradient descent is a well-known optimization technique. The idea is to compute the direction of steepest descent at the current value of the parameters, and then take a small step in this direction. We describe the algorithm in terms of the local loss functions L_{ij} . The gradients are given by:

$$\frac{\partial}{\partial \mathbf{W}_{i^*}} L(\mathbf{W}, \mathbf{H}) = \frac{\partial}{\partial \mathbf{W}_{i^*}} \sum_{(i',j) \in Z} L_{i'j}(\mathbf{W}_{i'^*}, \mathbf{H}_{*j}) = \sum_{j \in Z_{i^*}} \frac{\partial}{\partial \mathbf{W}_{i^*}} L_{ij}(\mathbf{W}_{i^*}, \mathbf{H}_{*j}),$$

where $Z_{i^*} = \{j: (i, j) \in Z\}$. This means that the gradient w.r.t. \mathbf{W}_{i^*} depends on \mathbf{H} and row i of both the loss matrix L and \mathbf{W} . Similarly, we have

$$\frac{\partial}{\partial \mathbf{H}_{*j}} L(\mathbf{W}, \mathbf{H}) = \sum_{i \in Z_{*j}} \frac{\partial}{\partial \mathbf{W}_{*j}} L_{ij}(\mathbf{W}_{i^*}, \mathbf{H}_{*j}),$$

where $Z_{*j} = \{i: (i, j) \in Z\}$. The gradient w.r.t. to \mathbf{W} is given by the matrix of first-order partial derivatives

$$\mathbf{W}_n^\nabla = \frac{\partial}{\partial \mathbf{W}} L(\mathbf{W}_n, \mathbf{H}_n) = \begin{pmatrix} \frac{\partial}{\partial \mathbf{W}_{1^*}} L(\mathbf{W}_n, \mathbf{H}_n) \\ \frac{\partial}{\partial \mathbf{W}_{2^*}} L(\mathbf{W}_n, \mathbf{H}_n) \\ \vdots \\ \frac{\partial}{\partial \mathbf{W}_{m^*}} L(\mathbf{W}_n, \mathbf{H}_n) \end{pmatrix}.$$

Similarly,

$$\mathbf{H}_n^\nabla = \frac{\partial}{\partial \mathbf{H}} L(\mathbf{W}_n, \mathbf{H}_n) = \left(\frac{\partial}{\partial \mathbf{H}_{*1}} L(\mathbf{W}_n, \mathbf{H}_n) \quad \cdots \quad \frac{\partial}{\partial \mathbf{H}_{*n}} L(\mathbf{W}_n, \mathbf{H}_n) \right).$$

Then, GD performs the following iteration

$$\begin{aligned} \mathbf{W}_{n+1} &= \mathbf{W}_n - \epsilon_n \mathbf{W}_n^\nabla \\ \mathbf{H}_{n+1} &= \mathbf{H}_n - \epsilon_n \mathbf{H}_n^\nabla. \end{aligned}$$

If the conditions on $\{\epsilon_n\}$ and L given in Sec. 4.2 hold, and a projection term is added to the above equations, the algorithm converges asymptotically to the set of limit points of the projected ODE $\dot{\theta} = -L'(\theta) + z$, where $\theta = (\mathbf{W}, \mathbf{H})$ and, as before, z is the minimum force to keep the solution in the constraint region.

GD converges very slowly and many iterations may be required to approach a stationary point. To get reasonable convergence speed, Newton or quasi-Newton methods replace the step sizes by (estimates of) the inverse Hessian of L at the current parameter estimate. A scalable and memory-efficient method is L-BFGS-B [8]. If L is non-convex, GD may get stuck in a local optimum. Especially in its standard form, GD is thus not well suited for loss functions that have many “small bumps”, i.e., many bad local minima.

Both first-order and second-order GD methods can be distributed using MapReduce [23, 13]. The key idea of this approach to arbitrarily partition the loss matrix across a cluster of nodes. Each node sums up the partial gradients of its part of the data, the partial gradients are then summed up at the reducers. The authors argue that this can be done conveniently by using a query language on top of MapReduce. Once the gradient has been computed, a master node will perform the update of the parameter values using, for example, the L-BFGS-B method.

Partitioned stochastic gradient descent (PSGD, ISGD). PSGD and ISGD both are recent approaches to distribute SGD without using stratification; see Algorithm 3. The idea is to partition the data randomly into d partitions, and run SGD independently and in parallel on each partition.

Algorithm 3 PSGD/ISGD for Matrix Factorization

Require: Z , W_0 , H_0 , degree of parallelism d
 $W \leftarrow W_0$
 $H \leftarrow H_0$
Randomly divide Z into d partitions Z_1, \dots, Z_d
while not converged **do** /* *epoch* */
 Pick step size ϵ
 Distribute W and H /* *ISGD: only in first iteration* */
 for $b = 1, \dots, d$ **do** /* *in parallel* */
 Run SGD on the training points in Z^b (step size = ϵ)
 end for
 Collect W and H from all nodes and average /* *ISGD: only in last iteration* */
end while

Results are averaged in a *parameter mixing* step after either each epoch (PSGD [14, 24]) or once after convergence on each partition (ISGD [23, 24, 28]); observe that PSGD requires periodic synchronization between the partition-processing tasks whereas, for ISGD, processing on the different partitions can proceed in a mutually independent fashion. Both approaches can be implemented naturally on MapReduce. Compared to DSGD, ISGD is slightly more efficient (since there is no synchronization) whereas PSGD is slightly less efficient (since there are additional averaging steps). In the setting of matrix factorization, it is possible to reorder a set of rows of W and corresponding columns of H without affecting the value of the loss function. It follows that the loss function has many global minima that correspond to many different values of the factor matrices. In ISGD, the SGD processes on different partitions tend to converge to different solutions; the average of these local solutions is usually not a global loss minimizer. PSGD performs better and does converge, but it is outperformed by DSGD and, in some cases, even L-BFGS; see Sec. 7.

B. Parallelization Techniques for Stochastic Approximation

For completeness, we summarize standard approaches for parallel stochastic approximation [19, Ch. 12]. These approaches do not map naturally to MapReduce, and are designed for the case in which the computation of the update term is expensive (e.g., requires a simulation) and communication is rather cheap (e.g., few parameters). Neither assumption holds for matrix factorization, but our approaches are inspired by the techniques below.

Pipelined stochastic gradient descent. The *pipelined computation model* is based on a “delayed update” scheme in which the gradient estimate used in the n th step is based on the parameter value from the $(n - d)$ th step, where d is the number of available processors, e.g.,

$$\theta_{n+1} = \theta_n - \epsilon Y_{n-d},$$

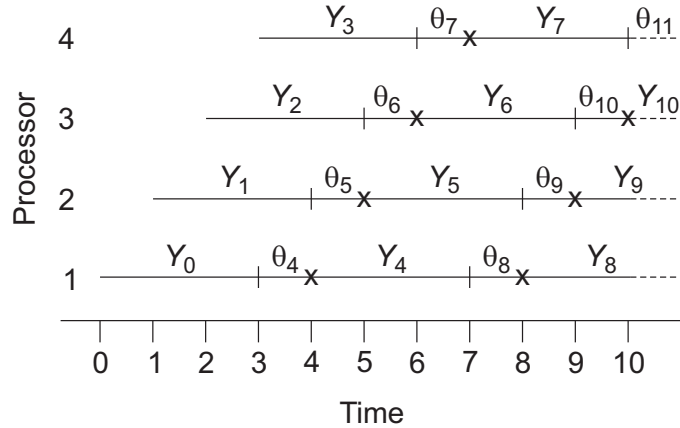


Figure 8: Pipelined SGD ($d = 3, t_g = 3, t_u = 1$)

where $Y_k = \hat{L}'(\theta_k)$ for $k \geq 0$. (We fix the step size ϵ for ease of notation.) This permits a scheme in which $d + 1$ processors compute gradient estimates and parameter values in a pipelined fashion. Figure 8 illustrates the technique for the case $d = 3$, assuming that it takes $t_g = 3$ time units to compute a gradient and $t_u = 1$ time unit to update a parameter value. The scheme is initiated by choosing θ_0 and setting $\theta_3 = \theta_2 = \theta_1 = \theta_0$. At time $n = 0$, processor 1 begins to compute the update term Y_0 . This computation completes at time $n = 3$ (marked with a “|”), at which point processor 1 begins to compute $\theta_4 = \theta_3 - \epsilon Y_0$, finishing at time $n = 4$ (marked with an “x”). Similarly, processor 2 begins to compute the update term Y_1 at time $n = 1$. At time $n = 4$, this computation completes, and processor 2 uses the value of θ_4 that has just been computed by processor 1 to compute $\theta_5 = \theta_4 - \epsilon Y_1$, finishing this computation at time $n = 5$. The other processors behave similarly. This scheme can be extended to handle variable updating delays.

Decentralized stochastic gradient descent In the *distributed and decentralized network model*, both parameter updates and computation of update terms take place in a distributed and decentralized fashion. This powerful model is a generalization of the pipelined computation model. Both the components of the parameter vector and the loss function are distributed across the set of d processors. Processor p is responsible for component θ^p and partial loss $L_p(\theta)$ such that:

$$\theta = \begin{pmatrix} \theta^1 \\ \theta^2 \\ \vdots \\ \theta^d \end{pmatrix} \quad \text{and} \quad L(\theta) = \sum_{p=1}^d L_p(\theta).$$

The processors operate in parallel and (potentially) at different speeds; communication is asynchronous. Thus at any given point in “real time” (wall clock time), each processor resides at a different time point in “iterate time” (step number). More specifically, the $(n + 1)$ st iteration at processor p involves the following steps:

1. *Estimate $\hat{\theta}_{n,p}$.* Obtain an estimate $\hat{\theta}_{n,p}$ of the entire parameter vector by using θ_n^p for the p th component (the current value of the component managed by p) and the most recently received value for all other components (from step 5).
2. *Compute update terms.* For each processor q (including p), compute (or estimate) update term

$$Y_{n,p}^q = -\frac{\partial L_p(\hat{\theta}_{n,p})}{\partial \theta^q}.$$

3. *Communicate update terms.* For each processor q , send $Y_{n,p}^q$ to q .
4. *Compute θ_{n+1}^p .* Add all unprocessed update terms Y_*^p to θ_n^p . This includes update terms received from other nodes as well as the update term $Y_{n,p}^p$ computed in the previous step. Do not wait for any “missing” update terms, and if multiple update terms have been received from a single processor q , process them all.

$$\theta_{n+1}^p = \theta_n^p + \epsilon \sum \{ \text{unprocessed update terms } Y_*^p \}.$$

5. *Broadcast θ_{n+1}^p .* Broadcast the new parameter component to all other nodes.

Weak convergence results for this process model are discussed in [19, Ch. 12]. The proofs proceed by a careful treatment of “iterate time”, and then use appropriate time scale changes to obtain “real time” results.

C. Example Loss Functions and Derivatives

Table 2 displays the definitions of several commonly used loss functions as mentioned in Section 7: nonzero squared loss (L_{NZSL}), non-zero squared loss with L_2 regularization (L_{L2}), nonzero squared loss with a nonzero-weighted L_2 term (L_{NZL2}), KL divergence (L_{KL}), and generalized KL divergence (L_{GKL}). In the table, $\|\mathbf{A}\|_{\text{F}}$ denotes the Frobenius norm of a matrix \mathbf{A} : $\|\mathbf{A}\|_{\text{F}} = (\sum_i \sum_j a_{ij}^2)^{1/2}$.

Our methods for obtaining a rank- r approximate factorization $\mathbf{V} \approx \mathbf{W}\mathbf{H}$ of an $m \times n$ input matrix \mathbf{V} require that we represent each loss function L as a sum of local losses over points in the training set Z , i.e., $L(\mathbf{W}, \mathbf{H}) = \sum_{(i,j) \in Z} L_{ij}(\mathbf{W}, \mathbf{H})$, where $L_{ij}(\mathbf{W}, \mathbf{H}) = l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$ for an appropriate function l , so that the gradient of L can be decomposed as a sum of local-loss gradients: $L'(\mathbf{W}, \mathbf{H}) = \sum_{(i,j) \in Z} L'_{ij}(\mathbf{W}, \mathbf{H})$. For each loss function L considered, Table 2 gives formulas for the components of the local-loss gradient L'_{ij} . In these formulas, N_{i*} and N_{*j} denote the number of nonzero elements in row i and column j of the matrix \mathbf{V} . Moreover, $J = \{1, \dots, m\} \times \{1, \dots, n\}$ and $B = \{1, \dots, d\} \times \{1, \dots, d\}$. Finally, for a $u \times v$ matrix \mathbf{A} , the quantities $\text{rowSums}(\mathbf{A})$ and $\text{colSums}(\mathbf{A})$ denote the $u \times 1$ column vector containing the row sums of \mathbf{A} and the $1 \times v$ row vector containing the column sums of \mathbf{A} ; thus, for example, $\text{colSums}(\mathbf{W})$ and $\text{rowSums}(\mathbf{H})$ are each vectors of length r .

As can be seen, special care has to be taken with regularization terms when representing L as a sum of local losses. In the case of L_{L2} , for example, we proceed as follows. Recall from Sec. 2 our

running assumption that there is at least one training point in every row and in every column of \mathbf{V} . Since $Z = \{(i, j) : \mathbf{V}_{ij} > 0\}$, this means that $N_{i*} > 0$ and $N_{*j} > 0$ for all i and j . Then we have

$$\begin{aligned} \|\mathbf{W}\|_{\mathbb{F}}^2 &= \sum_{i=1}^m \|\mathbf{W}_{i*}\|_{\mathbb{F}}^2 = \sum_{i=1}^m N_{i*} \frac{\|\mathbf{W}_{i*}\|_{\mathbb{F}}^2}{N_{i*}} \\ &= \sum_{i=1}^m \sum_{j=1}^n I[(i, j) \in Z] \frac{\|\mathbf{W}_{i*}\|_{\mathbb{F}}^2}{N_{i*}} = \sum_{(i, j) \in Z} \frac{\|\mathbf{W}_{i*}\|_{\mathbb{F}}^2}{N_{i*}} \end{aligned} \quad (28)$$

and, similarly, $\|\mathbf{H}\|_{\mathbb{F}}^2 = \sum_{(i, j) \in Z} \|\mathbf{H}_{*j}\|_{\mathbb{F}}^2 / N_{*j}$. (Here $I[A]$ denotes the indicator function of A .) These results lead directly to the representation of L_{L2} in Table 2. A similar algebraic manipulation is used to represent L_{NZL2} .

The case of L_{GKL} merits additional discussion. First recall that, after randomly permuting the rows and column of the training matrix \mathbf{Z} , we partition \mathbf{Z} into d^2 blocks, and partition each of the factor matrices conformingly into d blocks as $\mathbf{W} = (\mathbf{W}^1, \dots, \mathbf{W}^d)^{\text{T}}$ and $\mathbf{H} = (\mathbf{H}^1, \dots, \mathbf{H}^d)$. As before, we denote by Z^b the set of training points that lie in block $b \in B$. A stratum comprises d blocks, selected such that each pair of blocks has no row or column indices in common. We originally required that the local loss at a point $(i, j) \in Z$ be of the form $l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$ mentioned above, which ensures that SGD can be run independently within each block in the stratum. Observe, however, that we need only require that the local losses be of the form $l(\mathbf{V}_{ij}, \mathbf{W}^{b(i*)}, \mathbf{H}^{b(*j)})$, where $b(i*)$ denotes the block of \mathbf{W} that contains \mathbf{W}_{i*} and $b(*j)$ denotes the block of \mathbf{H} that contains \mathbf{H}_{*j} . This looser definition preserves the interchangeability structure, so that updates to parameter values for a given block will not affect parameter values corresponding to other blocks, and SGD can still be executed in a distributed manner within the stratum. We represent L_{GKL} as a sum of losses in this looser sense. In addition to allowing a more general representation of each local loss, we also use two different decompositions of L_{GKL} . The first (resp., second) representation is more amenable to calculating derivatives with respect to the \mathbf{W}_{ik} (resp., \mathbf{H}_{kj}) factors.

Denote by N_{i*}^b (resp., N_{*j}^b) the number of training points in substratum Z^b that appear in row i (resp., column j) of the training matrix \mathbf{Z} :

$$N_{i*}^b = |\{j : (i, j) \in Z^b\}| \quad \text{and} \quad N_{*j}^b = |\{i : (i, j) \in Z^b\}|.$$

We assume that $N_{i*}^b > 0$ for each row i that intersects block Z^b and that $N_{*j}^b > 0$ for each column j that intersects Z^b ; otherwise, we can always add additional (zero-valued) training points to the training set Z .¹⁰ For $b \in B$, denote by $J_1^b = \{i : (i, j) \in Z^b\}$ and $J_2^b = \{j : (i, j) \in Z^b\}$ the sets of first indices and second indices of points—both zero and nonzero—in block b . Set $Q^b = \sum_{i \in J_1^b} \sum_{j \in J_2^b} [\mathbf{W}\mathbf{H}]_{ij}$ and note that

$$Q^b = \sum_{i \in J_1^b} \sum_{j \in J_2^b} \mathbf{W}_{i*} \mathbf{H}_{*j} = \sum_{i \in J_1^b} \mathbf{W}_{i*} \cdot \text{rowSums}(\mathbf{H}^b) = \sum_{(i, j) \in Z^b} \frac{\mathbf{W}_{i*}}{N_{i*}^b} \cdot \text{rowSums}(\mathbf{H}^b),$$

¹⁰For each such added point (i, j) , we define the quantity $\mathbf{V}_{ij} \log(\mathbf{V}_{ij} / [\mathbf{W}\mathbf{H}]_{ij})$ that appears in the definitions of L_{ij}^{W} and L_{ij}^{H} below—as well as the quantity $\mathbf{V}_{ij} / [\mathbf{W}\mathbf{H}]_{ij}$ that appears in the definitions of $\partial L_{ij}^{\text{W}} / \partial \mathbf{W}_{ik}$ and $\partial L_{ij}^{\text{H}} / \partial \mathbf{H}_{kj}$ —to be equal to 0 for all \mathbf{W} and \mathbf{H} .

where the final equality follows from a manipulation as in (28). We then have

$$\begin{aligned}
L_{\text{GKL}} &= \sum_{(i,j) \in Z} (\mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} - \mathbf{V}_{ij}) + \sum_{(i,j) \in J} [\mathbf{W}\mathbf{H}]_{ij} \\
&= \sum_{b \in B} \left(\sum_{(i,j) \in Z^b} (\mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} - \mathbf{V}_{ij}) + Q^b \right) \\
&= \sum_{b \in B} \sum_{(i,j) \in Z^b} \left(\mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} - \mathbf{V}_{ij} + \frac{\mathbf{W}_{i^*}}{N_{i^*}^b} \cdot \text{rowSums}(\mathbf{H}^b) \right).
\end{aligned}$$

Thus we obtain the representation $L_{\text{GKL}} = \sum_{(i,j) \in Z} L_{ij}^W$, where

$$L_{ij}^W = \mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} - \mathbf{V}_{ij} + \frac{\mathbf{W}_{i^*}}{N_{i^*}^b} \cdot \text{rowSums}(\mathbf{H}^{b(i^*)}).$$

This is a local loss of the “loose” form discussed above. This representation of L_{GKL} is convenient for computing derivatives with respect to the \mathbf{W} factors. Observe that for each $(i, j) \in Z$, we have $\partial L_{ij}^W / \partial \mathbf{W}_{i'k} = 0$ for $i' \neq i$, so that—if we are running SGD on a block Z^b and are estimating the gradient based on a sampled training point $(i, j) \in Z^b$ —only the elements of \mathbf{W}_{i^*} need to be updated, and we retain computational efficiency as in Algorithm 1. In a similar manner, we can derive an alternate representation $L_{\text{GKL}} = \sum_{(i,j) \in Z} L_{ij}^H$, where

$$L_{ij}^H = \mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} - \mathbf{V}_{ij} + \frac{\mathbf{H}_{*j}}{N_{*j}^b} \cdot \text{colSums}(\mathbf{W}^{b(i^*)}).$$

This decomposition is useful for computing derivatives with respect to the \mathbf{H} factors. With an abuse of notation, we write $\partial L_{ij} / \partial \mathbf{W}_{ik}$ for $\partial L_{ij}^W / \partial \mathbf{W}_{ik}$ and $\partial L_{ij} / \partial \mathbf{H}_{kj}$ for $\partial L_{ij}^H / \partial \mathbf{H}_{kj}$. Although the “gradient” L'_{ij} that we have just defined is not actually the gradient of some well defined local loss function L_{ij} , it nonetheless holds that $\sum_{(i,j) \in Z^b} L'_{ij}$ is equal to the gradient of L_{GKL} on Z^b , which is all that we need for DSGD. The final derivative formulas are given in Table 2.

Table 2: Examples of loss functions and derivatives

Loss Function	Definition and Derivatives
L_{NZSL}	$L_{\text{NZSL}} = \sum_{(i,j) \in Z} (\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})^2$ $\frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij}) \mathbf{H}_{kj}$ $\frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij}) \mathbf{W}_{ik}$

Table 2: Examples of loss functions and derivatives (continued)

Loss Function	Definition and Derivatives
L_{L2}	$L_{L2} = L_{NZSL} + \lambda (\ \mathbf{W}\ _F^2 + \ \mathbf{H}\ _F^2)$ $= \sum_{(i,j) \in Z} \left[(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})^2 + \lambda \left(\frac{\ \mathbf{W}_{i*}\ _F^2}{N_{i*}} + \frac{\ \mathbf{H}_{*j}\ _F^2}{N_{*j}} \right) \right]$ $\frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij}) \mathbf{H}_{kj} + 2\lambda \frac{\mathbf{W}_{ik}}{N_{i*}}$ $\frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij}) \mathbf{W}_{ik} + 2\lambda \frac{\mathbf{H}_{kj}}{N_{*j}}$
L_{NZL2}	$L_{NZL2} = L_{NZSL} + \lambda (\ \mathbf{N}_1 \mathbf{W}\ _F^2 + \ \mathbf{H} \mathbf{N}_2\ _F^2)$ $= \sum_{(i,j) \in Z} \left[(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})^2 + \lambda (\ \mathbf{W}_{i*}\ _F^2 + \ \mathbf{H}_{*j}\ _F^2) \right]$ $\frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij}) \mathbf{H}_{kj} + 2\lambda \mathbf{W}_{ik}$ $\frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij}) \mathbf{W}_{ik} + 2\lambda \mathbf{H}_{kj}$
L_{KL}	$L_{KL} = \sum_{(i,j) \in Z} (\mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}})$ $\frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij} = -\mathbf{H}_{kj} \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}}$ $\frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij} = -\mathbf{W}_{ik} \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}}$

Table 2: Examples of loss functions and derivatives (continued)

Loss Function	Definition and Derivatives
L_{GKL}	$L_{\text{GKL}} = \sum_{(i,j) \in Z} (\mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} - \mathbf{V}_{ij}) + \sum_{(i,j) \in J} [\mathbf{W}\mathbf{H}]_{ij}$ $= \sum_{b \in B} \sum_{(i,j) \in Z^b} \left(\mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} - \mathbf{V}_{ij} + \frac{W_{i*}^b}{N_{i*}^b} \cdot \text{rowSums}(\mathbf{H}^b) \right)$ $= \sum_{b \in B} \sum_{(i,j) \in Z^b} \left(\mathbf{V}_{ij} \log \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} - \mathbf{V}_{ij} + \frac{H_{*j}^b}{N_{*j}^b} \cdot \text{colSums}(\mathbf{W}^b) \right)$ $\frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij} \stackrel{\text{def}}{=} \frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij}^{\mathbf{W}} = -\mathbf{H}_{kj} \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} + \frac{\text{rowSums}(\mathbf{H}^{b(i*)})_k}{N_{i*}^b}$ $\frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij} \stackrel{\text{def}}{=} \frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij}^{\mathbf{H}} = -\mathbf{W}_{ik} \frac{\mathbf{V}_{ij}}{[\mathbf{W}\mathbf{H}]_{ij}} + \frac{\text{colSums}(\mathbf{W}^{b(i*)})_k}{N_{*j}^b}$