

Predicting the Learning Rate of Gradient Descent for Accelerating Matrix Factorization

C. S. B. Nóbrega, L. B. Marinho

Federal University of Campina Grande, Brazil
caionobrega@copin.ufcg.edu.br, lbmarinho@dsc.ufcg.edu.br

Abstract. Matrix Factorization (MF) has become the predominant technique in recommender systems. The model parameters are usually learned by means of numerical methods, such as gradient descent. The learning rate of gradient descent is typically set to lower values in order to ensure that the algorithm will not miss a local optimum. As a consequence, the algorithm may take several iterations to converge. Ideally, one wants to find the learning rate that will lead to a local optimum in the first iterations, but that is very difficult to achieve given the high complexity of the search space. Starting with an exploratory analysis on several recommender systems datasets, we observed that there is an overall linear relationship between the learning rate and the number of iterations needed until convergence. Another key observation is that this relationship holds across the different recommender datasets, with only minor variations. From this, we propose to use simple linear regression models for predicting, for an unknown dataset, a good learning rate to start with. The idea is to estimate a learning rate that will get us as close as possible to a local optimal in the first iteration, without overshooting it. We evaluate our approach on 8 real-world recommender datasets and compare it against the standard algorithm and adaptive learning rate strategies from the literature. We show that, for some datasets, we can reduce the number of iterations up to 40% when compared to the standard approach.

Categories and Subject Descriptors: I.2.6 [Artificial Intelligence]: Learning

Keywords: Gradient Descent, Learning Rate, Matrix Factorization, Recommender Systems

1. INTRODUCTION

Matrix factorization is one of the most successful techniques in recommender systems nowadays [Koren et al. 2009]. MF features high accuracy, scalability, robustness against sparsity, and implementation ease, which explains its success and popularity. Moreover, the winning method of the Netflix challenge¹, a very important recommender systems challenge that provided major advances in the area, used MF as a key component. By now, there are many variations of MF [Paterek 2007] and [Gantner et al. 2010], each one handling a different aspect of recommender systems. In this paper, we investigate MF for the classic problem of rating prediction, i.e., how to predict the ratings that users would give to items they still not accessed (e.g. movies still not watched).

The model parameters of MF are usually learned by means of numerical methods such as gradient descent, given that the loss function is non-convex. The learning rate of gradient descent is typically set to lower values (usually 0.01 or 0.02 [DeCoste 2006]) in order to ensure that the algorithm does not miss a local optimum and diverges. As a consequence, the algorithm may take several iterations to converge, which ends up increasing the computational costs of the training phase. [Rendle and Schmidt-Thieme 2008], for example, report 200 iterations until convergence in the Netflix dataset with the learning rate set to 0.01.

¹<http://www.netflixprize.com/>

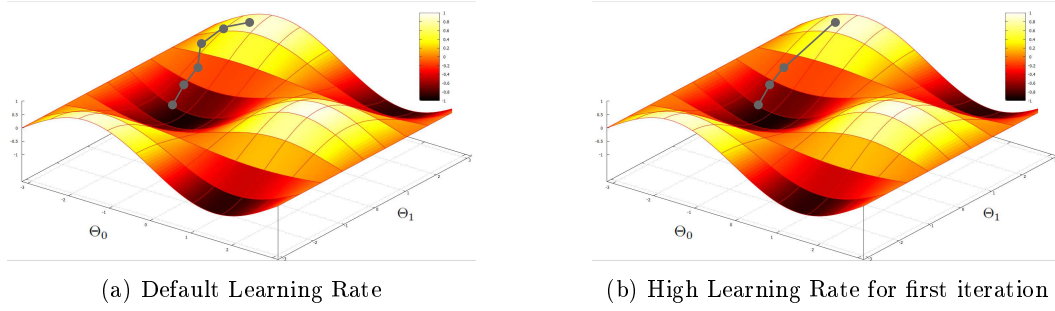


Fig. 1: Illustration of gradient descent for learning hypothetical parameters Θ_0 e Θ_1 based on the loss function $J(\Theta_0, \Theta_1)$: (a) using default learning rate over all iterations. (b) using a better value for the first iteration

Ideally, one wants to find the learning rate that leads to a local optimum in the first iterations, but that is very difficult to achieve given the high complexity of the search space. However, if in one shot one gets close enough to the local optimum, one will need only a few iterations until convergence, hence saving iterations and speeding up the learning process. Figure 1 illustrates an hypothetical execution of gradient descent for learning the parameters Θ_0 and Θ_1 under the loss function $J(\Theta_0, \Theta_1)$ and the idea of how to reduce the number of iterations (represented by points) with a good initial value for the learning rate.

In this paper, we investigate the intrinsic relationship between the learning rate and the number of iterations of MF applied to recommender systems datasets. Our problem is to predict a good estimation of the initial value of the learning rate in order to minimize the number of iterations, leading to rapid convergence on yet unseen recommendation datasets. We observed that, for most of the recommender systems datasets we used (7 of 8), there is a positive linear relationship between the learning rate and the number of iterations needed until convergence. Surprisingly, this relationship holds across the different datasets, with only minor variations. We assume that MF presents similar learning behavior when applied to similar datasets, i.e., datasets that share the same entities (users and items), domain (recommender systems), and rating scale (5 star rating scale). From this observation, we propose to use simple linear regression models for predicting, for an unknown dataset, which learning rate would lead to the minimum number of iterations.

We conduct experiments in three real world datasets and show that we can reduce the number of iterations up to 40% when compared to the standard approach which uses a fixed learning rate value. Furthermore, we compare our approach to three learning rate adaptation techniques from the literature.

This paper is organized as follows. In Section 2, we present the related work. In Section 3, we formalize the research problem this paper investigates. In Section 4, we investigate the relationship between the learning rate and the number of iterations needed until convergence in the most popular recommender systems datasets available to the public. In Section 5 we present the results and conclusions of the experiments conducted. Finally, in Section 6, we conclude the paper and discusses opportunities for future work.

2. RELATED WORK

In general, the value of the learning rate is determined by a search among candidate values. Using a validation set, a MF model for each candidate is learned and evaluated. The candidate that leads to the best model is chosen. This approach is expensive since it trains several models, especially when the search space is large. Moreover, this is a very common approach to find other hyperparameters, such as the ones related to regularization. In fact, [Rendle 2012] propose an adaptive approach for speeding up the learning of the regularization hyperparameters. In this paper we are only interested

in the learning rate of gradient descent.

Gradient descent is a generic approach used as the subroutine of many learning and optimization algorithms, e.g., artificial neural networks (ANN) [Moreira and Fiesler 1995]. In fact, some of the most recent advances on gradient descent algorithms come from different research areas other than recommender systems. A common approach is to adapt the learning rate over the iterations of the gradient descent, i.e., algorithms dynamically incorporate knowledge of the geometry of the data observed in earlier iterations to perform more informative updates [Bartlett et al. 2007] and [Auer et al. 2000]. The method introduced in [Duchi et al. 2011], for example, starts with a default learning rate that is divided by the norm of the previous gradients in subsequent iterations. MF models, on the other hand, have a large number of parameters which makes it difficult to grasp the geometry of the search space. Indeed, as learning rate adaptation has proven to be an effective solution to the problem of tuning the learning rate in the ANN field, [Luo et al. 2013] propose different versions of gradient descent to learn the parameters of the MF. They adapted 3 techniques, namely, Deterministic Step Size Adaptation (DSSA), Incremental Delta Bar Delta (IDBD) and Stochastic Meta Descent (SMD) and compare to his own method, Gradient Cosine Adaptation (GCA). In all these methods there is a learning rate for each parameter of MF, which are updated based on the gradient of the previous iterations. For instance, in DSSA, if two successive updates of gradient are made in the same direction, then the learning rate is increased to make the learning process faster, and vice versa. One drawback of the aforementioned techniques, is that they consider a learning rate for each parameter, which leads to a memory overhead since this information needs to be cached during the learning phase.

Except [Luo et al. 2013], none of the related works deal with MF or recommender systems. Moreover, the aforementioned works follow a bottom-up approach trying to adapt the learning rate based on previous values, while we follow a top-down approach, trying to predict, from the beginning, the learning rate that will lead to the minimum number of iterations.

Another line of research worth mentioning concerns distributed versions of gradient descent, as presented in [Zinkevich et al. 2010] for the domain of artificial neural networks.

3. PROBLEM SETTING

In this paper, we tackle the problem of accelerating matrix factorization, estimating a good value to initialize the learning rate of gradient descent. The idea is to use a learning rate that gets as close as possible to a local optimum. As pointed out in Section 2, the learning rate is a low constant value determined by cross validation (usually 0.01 or 0.02 [DeCoste 2006]). Before formalizing our problem, we briefly recall MF and gradient descent.

Let U denote the set of users and I the set of items, and S the sparse matrix, where rows represent users, columns represent items, and values represent the ratings that users gave to items. As the name suggests, MF attempts to decompose $S \in \mathbb{R}^{|U| \times |I|}$ into the product of the two lower rank matrices $P \in \mathbb{R}^{|U| \times k}$ and $Q \in \mathbb{R}^{|I| \times k}$, such that $S \approx PQ^T$. Here k denotes the number of latent factors used, i.e., the latent factors that explain the interactions between users and items. After the factorization, users and items are represented by vectors lying in the same space of k dimensions. Predictions are now done by simple dot products between user and items vectors, i.e., rows of P and columns of Q^T resp. MF is cast as an optimization problem (cf. equation 1). For the rating prediction problem the loss function *err* is the well known squared error. Since matrix S is sparse, the error is computed only w.r.t. to its observed values.

$$\underset{P, Q}{\operatorname{argmin}} \operatorname{err}(P \times Q^T; S) \quad (1)$$

As mentioned before, the model parameters, here represented by the vector $\Theta = (P, Q)$, are usually learned by gradient descent. Gradient descent uses the partial derivative of the loss function w.r.t.

to each parameter to guide the search for the local optimum. In each iteration gradient descent gives another step towards it, where the size of the step is given by the learning rate α . The equation 2 depicts, in high level, the update step of gradient descent.

$$\Theta = \Theta - \alpha \left(\frac{\partial}{\partial \Theta} \text{err}(P \times Q^T; S) \right) \quad (2)$$

In the following, we define the overall problem we investigate in this paper. Let β denote the learning rate used in the first iteration of equation 2 and α the learning rate of the remaining iterations. Let $\text{iter}(\beta, \alpha, P, Q)$ be a function that returns the number of iterations needed to learn P and Q (by means of equation 1) with a given configuration of β and α . In this paper, we assume that α is fixed and focus on finding the β that minimizes the number of iterations until convergence. The basic idea is to start with a β that gets as close as possible to a local optimum and then use a conservative α for the remaining iterations until convergence. In all, our problem is formalized as follows:

$$\underset{\beta}{\operatorname{argmin}} \text{iter}(\beta, \alpha, P, Q) \quad (3)$$

4. LEARNING RATE PREDICTION

Ideally, we would like to find a learning rate capable of reaching a local optimum in the first iterations of the gradient descent. As we do not know what value yields this result, we do a search. In this search, we start with a given learning rate to be used only in the first iteration, here called β , and then, after the very first iteration, we set it back to a conservative (or default) value called α . What we want to do is to get as close as possible to a local optimum, without overshooting it. The rest of the way we do it as before, i.e., using small step sizes. Doing that we expect to decrease the overall number of iterations needed until convergence.

To investigate the impact of different β values in the learning phase of MF, we first define a metric that indicates the percentage of iterations' reduction in comparison to the standard procedure using a fixed and small α . For example, given a dataset, assume that MF needs 20 iterations to achieve convergence in the validation set with the default learning rate α , and 12 iterations with a given β as initial learning rate. Now, the reduction is of 0.4 or 40%, if compared to the original number of iterations, and is calculated as $1 - (12/20)$.

For further investigating the relationship between β and the percentage of iterations' reduction, we varied the β values from 0.01 to 0.1 with increments of 0.005 and used the default $\alpha = 0.01$ suggested by the literature to proceed until convergence. We have chosen this range, because we noticed that, for most of the datasets, when $\beta > 0.1$ the gradient descent diverges, which might indicate that we are overshooting a local optimum. In figure 2 we plot the β values used (x axis) versus the percentage of iteration reduction (y axis) for all datasets.

Table I: Dataset characteristics: number of users, items, and ratings (instances) in the training and validation sets.

<i>Dataset</i>	$ U $	$ I $	# rating instances	Domain
Movielens-100k	943	1,682	100,000	Movies
Movielens-1M	6,040	3,706	1,000,209	Movies
Movielens-10M	69,878	10,677	10,000,054	Movies
Netflix	480,189	17,770	100,480,507	Movies
Yahoo-Movies	7,642	11,916	221,364	Movies
Epinions	40,163	139,738	664,824	General Products
Amazon	2,146,276	1,231,018	5,744,088	General Products
Dating	135,359	168,791	17,359,346	Dating

Fig. 2: Relationship between β and the percentage of iterations' reduction. The red lines show a linear model fitted to this data.

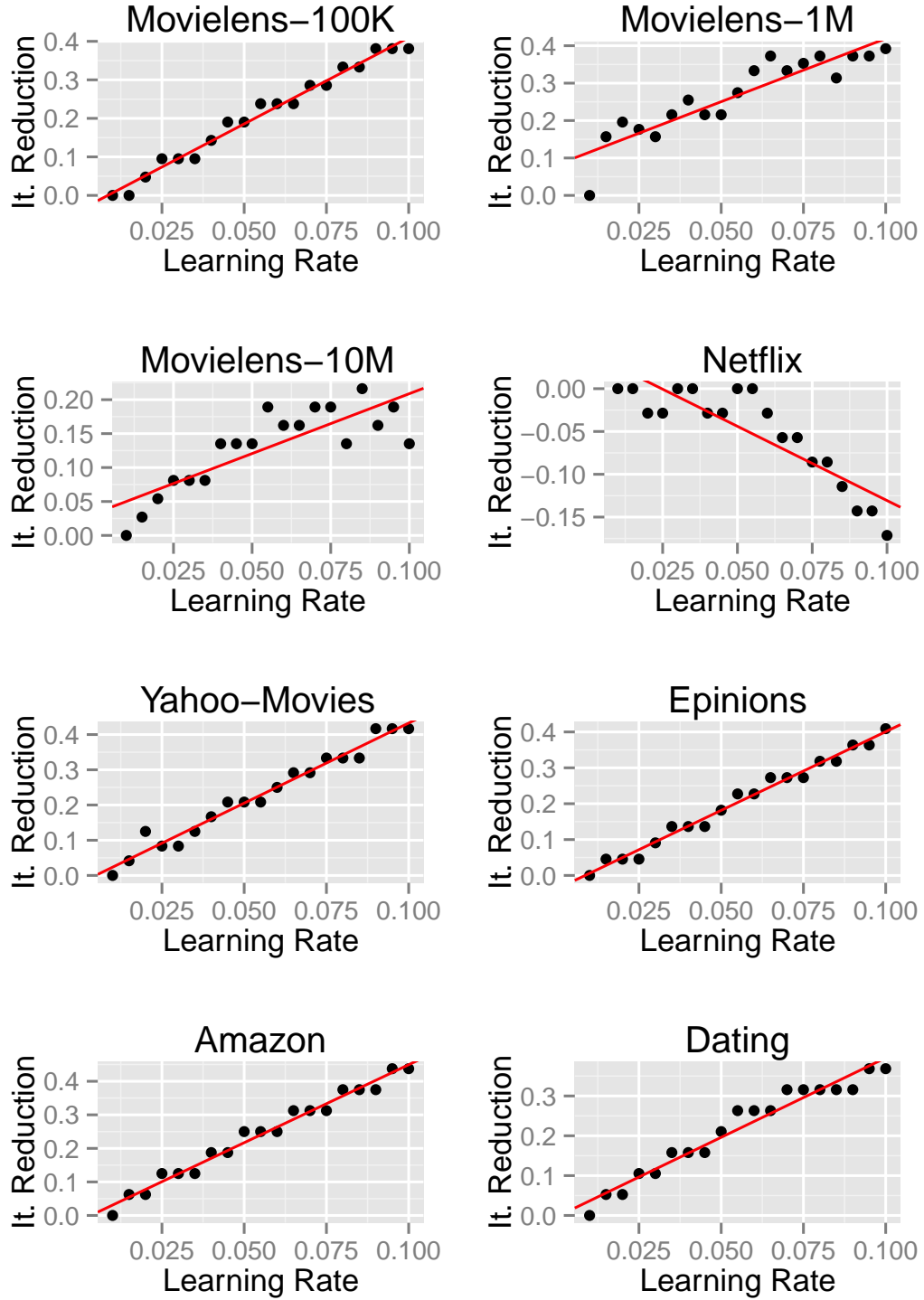


Table II: Evaluation of linear models based on statistic R^2 . Max. reduction indicates the greater percentage of iterations reduced when compared to the #iterations using default α (learning rate)

<i>Dataset</i>	R^2	% max. reduction	#iter. with default $\alpha = 0.01$
Movielens-100K	0.9801	38.1%	21
Movielens-1M	0.8184	39.2%	51
Movielens-10M	0.662	21.6%	37
Netflix	0.7525	0%	35
Yahoo-Movies	0.9694	41.6.5%	24
Epinions	0.9838	40.9%	22
Amazon	0.9779	43.7%	16
Dating	0.9514	36.8%	19

In fact, for all datasets we have used in our experiments, we observed a positive linear relationship between the β and the percentage of iterations' reduction, i.e., as higher the β , the higher the iterations reduction. To confirm this relationship, we fitted a simple linear regression model to each dataset and observed the goodness of fit (R^2 values). The R^2 varies between 0 and 1, indicating, in percentage terms, how the model can explain the observed values. The higher the R^2 , the better it fits the data. The R^2 values are close to 1, which indicates strong linear correlation (cf. Table II). In the last column of table II we show the number of iterations needed by the standard MF using the default $\alpha = 0.01$.

The exception is the Netflix dataset that do not follow the same pattern of the other datasets. We do not have a concrete explanation for that yet, but one guess is that 0.01 already represents a good initial learning rate, and thus does not leave much room for improving upon.

5. EXPERIMENTS

In this section we describe how to use the prediction model for accelerating the learning phase of MF and compare this approach with other learning rate adaptation strategies introduced in [Luo et al. 2013].

5.1 Experimental Reproducibility

For the sake of experimental reproducibility, we have chosen well-known recommender systems datasets, all of them public and available online². We employ 5-fold cross-validation and use the MF default setup values of the recommender system library MyMediaLite [Gantner et al. 2011]. The number of latent factors k , regularization term and stop criterion are 10, 0.015 and 0.001 respectively. We also implemented, on top of MyMediaLite, the learning rate adaptation strategies described in [Luo et al. 2013] with the following meta-parameter settings (and using the same nomenclature): $\alpha = 0.0005$ and $\beta = 0.0005$ for DSSA, $\theta = 0.002$ for IDBD, $K = 0.9$ and $\theta = 0.005$ for SMD. For making the results comparable, we mapped all datasets into the [1, 5] rating scale. Table I presents the data characteristics of all the datasets used.

5.2 Learning Rate Prediction

As discussed in the previous section, MF presents very similar learning behavior across different recommender datasets concerning the choices of β for the range [0.01, 0.1] (cf. figure 3). For most of the datasets, this behavior is highly explainable by a simple linear regression model. This suggests that we can use the linear model fitted to one dataset for making predictions for others, yet unseen,

²The Dating dataset is available at <http://www.occamlab.com/petricek/data>; the Yahoo-Movies dataset is available at <http://webscope.sandbox.yahoo.com/catalog.php>; the other datasets used are available at <http://konect.uni-koblenz.de>

datasets. Given a new recommender systems dataset, the question we want to answer is this: “What learning rate should I use if I want to reduce the number of iterations in $x\%$ in relation to the standard algorithm?”

For evaluating our approach we devised the following experimental protocol. We have held one dataset for testing, and used the other datasets, individually, for prediction. We then compare the number of iterations needed by MF using the standard gradient descent with the number of iterations achieved with the β predicted by each of the other datasets. We excluded the Netflix since it presents a negative linear correlation, as pointed out in Section 4.

We used the datasets Amazon, Movielens-10M and Dating as test datasets. The choice is supported by the fact that, although both are recommender systems datasets, they represent different domains (cf. table I). For the test datasets, we try to predict the β that achieves 40% of reduction in the number of iterations, which is close to the empirical upper bound. Notice that the empirical upper bound of reduction for the Amazon is 47.3%, for MovieLens-10M is 21.6% and 36.8% for Dating.

Table III shows the results. For the Amazon dataset, the best predicting dataset was Movielens-10M. Surprisingly, this prediction led to a reduction of 62.5%, which is more than the empirical upper bound. It was an exception, since values above 0.1 led to divergence. For the Movielens-10M, the best predicting dataset was the Movielens-1M. In this case, the reduction was 21.6% which is far from 40%, but this is expected since the empirical upper bound suggests that 21.6% is already the best we can get. For the Dating dataset, the best predicting datasets were Movielens-100k and Epinions, again returning the empirical upper bound.

Figure 3 illustrates the reduction of iterations in Amazon when using the β predicted by Movielens-10M for a 40% input value. MF with the standard gradient descent reached convergence in 13 iterations, while it drops to 6 with our approach. It is interesting to notice that different local minimum are achieved, where our approach reaches an even better local minimum.

We can see that all the results are very similar. In fact, the confidence interval of the percentage of iteration's reduction for the Amazon dataset is [36.2, 51.5], [14.1, 19.1] for MovieLens, and [26.9, 36.1] for the Dating dataset, with significance level of 0.05, which proves that all models provided similar good prediction. We can speculate that this relationship stems from features shared by the datasets

Table III: Prediction learning rate for Amazon, MovieLens-10M and Dating datasets with linear models extracted from others datasets.

<i>Dataset Test</i>	<i>Dataset Train</i>	predicted β	% iter. reduction	#iter. with predicted β
Amazon	MI-100K	0.097	43.7%	9
	MI-1M	0.087	37.5%	10
	MI-10M	0.159	62.5%	6
	Yahoo-Movies	0.091	37.5%	10
	Epinions	0.099	43.7%	9
	Dating	0.098	38.4%	6
Movielens-10M	MI-100k	0.097	18.9%	30
	MI-1M	0.087	21.6%	29
	Yahoo-Movies	0.091	16.2%	31
	Epinions	0.099	13.5%	32
	Amazon	0.088	16.2%	31
	Dating	0.098	13.5%	32
Dating	MI-100k	0.097	36.8%	12
	MI-1M	0.087	31.5%	13
	MI-10M	0.159	21%	15
	Yahoo-Movies	0.091	31.5%	13
	Epinions	0.099	36.8%	12
	Amazon	0.088	31.5%	13

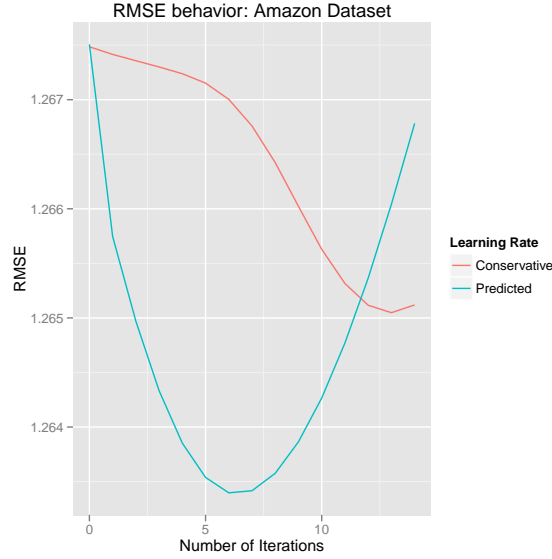


Fig. 3: RMSE (root mean squared error) relative to each iteration of gradient descent using the standard gradient descent versus our approach

e.g. all of these datasets have the same range of ratings [1, 5]. However, despite similar predictions, these datasets differ in the domain. This finding is interesting because this approach tends to happen regardless the domain considered.

5.3 Comparing to Learning Rate Adaptation Strategies

We compared our approach, Learning Rate Estimation (LRE), against three learning rate adaptation strategies from the literature, namely, DSSA, IDBD and SMD.

For the DSSA strategy, the update of learning rate is based on the signs of two successive gradients updates. The idea is to penalize the learning rate if the successive directions (represented by the signs) are distinct, multiplying it by a value slightly lower than 1, which decreases the step size. On the other hand, if successive updates of gradients are made in the same direction, then the learning rate is multiplied by a value slightly larger than 1, thus increasing the step size.

In the IDBD strategy, all learning rates are in the exponential form $\alpha = e^\sigma$. The exponential relationship between the learning rate α , and σ parameter that is actually modified has two advantages. First, it assures that learning rate will always be positive. Second, it is a mechanism for making geometric steps in σ , which is desirable because some learning rates must become very small while others remain large. The central idea of IDBD is to model the learning rate by considering the effect of all previous learning rates, represented by σ , on the current parameters.

In the SMD strategy, all learning rates are adjusted in log-space and are optimized over an exponentially decaying trace of gradients, similarly to the IDBD strategy. However, differently from IDBD where each learning rate takes into account the previous values, SMD takes into consideration the effects of one specified learning rate on the other related parameters.

In the evaluation, we used our best predicted learning rate, estimated in the last section, and set the adaptation strategies to the default learning rate (0.01), in order to make the results comparable with traditional approach and our conservative step used after first iteration. The table IV show the results. For all datasets, our approach needed fewer iterations than the methods compared, keeping similar RMSE values, i.e., our approach reached convergence faster and keep equivalent accuracy.

Table IV: Compare our approach with learning rate adaptation strategies.

<i>Dataset Test</i>	<i>Approach</i>	<i>RMSE</i>	<i>#it</i>
Amazon	MF-LRE	1.2620	6
	MF-DSSA	1.2680	16
	MF-IDBD	1.2686	16
	MF-SMD	1.2679	16
Movielens-10M	MF-LRE	0.8025	29
	MF-DSSA	0.7980	36
	MF-IDBD	0.7960	36
	MF-SMD	0.7981	41
Dating	MF-LRE	0.8731	12
	MF-DSSA	0.873	19
	MF-IDBD	0.8728	19
	MF-SMD	0.8723	19

It is worth mentioning that the compared methods are more complex than ours in terms of space since they have to store additional matrices to keep the learning rates and the gradients, one for each parameter. While our approach has space complexity of $O((|U|+|I|) \times k + |T|)$, where $|T|$ is the number of ratings, the adaptation approaches have space complexity in the order of $O((|U| + |I|) \times 3k + |T|)$.

6. CONCLUSIONS AND FUTURE WORK

In this work, we introduced an approach for predicting the learning rate of matrix factorization's gradient descent. The main idea is to use this prediction to get as close as possible to a local optimum in the first iteration.

Starting from an exploratory investigation on different recommender datasets, we observed that there is an overall linear relationship between the learning rate and the number of iterations needed until convergence. From this, we propose to use simple linear for predicting, for an unknown dataset, which learning rate leads to the minimum number of iterations. We have tested this hypothesis on 8 real-world datasets, which share features as the domain e.g. Movielens and Netflix are movie datasets and Epinions and Amazon are products datasets. We show that, for some datasets, we can reduce the number of iterations up to 40% when compared to the standard approach.

Moreover, we compared our approach against three learning rate adaptation strategies from the literature and showed that our method outperforms all of them in all the evaluated datasets.

As future work, we plan to combine our approach with adaptive algorithms in order to further reduce the number of iterations needed until convergence. So, in addition to the initial value, we would have a dynamic learning rate capable to adapt itself during the model learning. Finally, we plan to investigate this behavior in another types of datasets, such as implicit feedback recommendation datasets.

REFERENCES

- AUER, P., CESA-BIANCHI, N., AND GENTILE, C. Adaptive and self-confident on-line learning algorithms, 2000.
- BARTLETT, P., HAZAN, E., AND RAKHLIN, A. Adaptive online gradient descent. Tech. Rep. UCB/EECS-2007-82, 2007.
- DECOSTE, D. Collaborative prediction using ensembles of maximum margin matrix factorizations. In *Proceedings of the 23rd international conference on Machine learning*. ICML '06. ACM, New York, NY, USA, pp. 249–256, 2006.
- DUCHI, J., HAZAN, E., AND SINGER, Y. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.* vol. 12, pp. 2121–2159, July, 2011.
- GANTNER, Z., RENDLE, S., FREUDENTHALER, C., AND SCHMIDT-THIEME, L. MyMediaLite: A free recommender system library. In *5th ACM International Conference on Recommender Systems (RecSys 2011)*. Chicago, USA, 2011.

- GANTNER, Z., RENDLE, S., AND SCHMIDT-THIEME, L. Factorization models for context-/time-aware movie recommendations. In *Proceedings of the Workshop on Context-Aware Movie Recommendation*. CAMRa '10. ACM, New York, NY, USA, pp. 14–19, 2010.
- KOREN, Y., BELL, R., AND VOLINSKY, C. Matrix factorization techniques for recommender systems. *Computer* 42 (8): 30–37, 2009.
- LUO, X., XIA, Y., AND ZHU, Q. Applying the learning rate adaptation to the matrix factorization based collaborative filtering. *Know.-Based Syst.* vol. 37, pp. 154–164, Jan., 2013.
- MOREIRA, M. AND FIESLER, E. Neural networks with adaptive learning rate and momentum terms. Tech. rep., 1995.
- PATEREK, A. Improving regularized singular value decomposition for collaborative filtering. In *Proceedings of the KDD Cup Workshop at the 13th ACM International Conference on Knowledge Discovery and Data Mining*. SIGKDD '07. ACM, New York, NY, USA, pp. 39–42, 2007.
- RENDLE, S. Learning recommender systems with adaptive regularization. In *Proceedings of the fifth ACM international conference on Web search and data mining*. New York, NY, USA, pp. 133–142, 2012.
- RENDLE, S. AND SCHMIDT-THIEME, L. Online-updating regularized kernel matrix factorization models for large-scale recommender systems. In *Proceedings of the 2008 ACM conference on Recommender systems*. RecSys '08. ACM, New York, NY, USA, pp. 251–258, 2008.
- ZINKEVICH, M., WEIMER, M., SMOLA, A., AND LI, L. Parallelized stochastic gradient descent. In *Advances in Neural Information Processing Systems 23*. pp. 2595–2603, 2010.