

LOW RANK FIELD-WEIGHTED FACTORIZATION MACHINES FOR LOW LATENCY ITEM RECOMMENDATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Factorization machine (FM) variants are widely used in item recommendation systems that operate under strict throughput and latency requirements, such as online advertising systems. Their main strength is their ability to model pairwise feature interactions while being resilient to data sparsity by learning factorized representation, and having computational graphs that allow fast inference. Moreover, when items are ranked as a part of a query for each incoming user, these graphs facilitate computing the portion stemming from the user and context fields only once per query, and the computational cost for each ranked item is proportional only to the number of fields that vary among the ranked items. Consequently, in terms of inference cost, the number of user or context fields is practically unlimited.

More advanced variants of FMs, such as field-aware and field-weighted FMs, provide better accuracy by learning a representation of field-wise interactions, but require computing all pairwise interaction terms explicitly. In particular, the computational cost during inference is proportional to the square of the number of all fields, including user, context, and item. This is prohibitive in many systems when the number of fields is large, and imposes a limit on the number of user and context fields. To mitigate this caveat, heuristic pruning of low intensity field interactions is commonly used to accelerate inference.

In this work we propose an alternative to the pruning heuristic in field-weighted FMs using a diagonal plus symmetric low-rank decomposition, that reduces the computational cost of inference, by allowing it to be proportional to the number of item fields only. Using a set of numerical experiments, we show that aggressive rank reduction outperforms similarly aggressive pruning, both in terms of accuracy and item recommendation speed. Beyond computational complexity analysis, we corroborate our claim of faster inference experimentally having deployed our solution to a major online advertising system, where we observed significant ranking latency improvements.

1 INTRODUCTION

Recommendation systems driven by machine-learned predictive models are widely used throughout the industry for a large variety of applications, from movie recommendation, to bidding in ad auctions. In some applications, recommendation quality is the main objective, where sophisticated and computationally complex deep learning techniques are used to capture the affinity between the users and the recommended items. But other applications require striking an intricate balance between the accuracy of the predictive models, their training and inference speed. For example, real-time bidding systems in programmatic advertising are required to compute a ranking score for a large number of ads in a matter of a few milliseconds, and to train quickly in order to adapt to the ever-changing ad marketplace conditions. Such systems often deploy variants of the celebrated factorization machine (FM) models (Rendle, 2010a) in order to overcome data sparsity issues, while excelling at achieving a good balance between prediction accuracy and speed.

Rendle (2010a) also showed that while FMs model pairwise feature interactions, whose number is quadratic in the number of features, there is an equivalent formulation of FMs whose computational complexity is *linear* in the number features. This already facilitates fast training. Moreover, when ranking items for a given user in a given context, the user and context features are the same for

all items. The equivalent formulation allows caching the user and context computation results once, meaning that the computational cost per item is linear in the number of item features only. Therefore, factorization machines are also extremely efficient for ranking.

Many real world applications involve large-scale multi-field data (e.g. gender, item category). However, FMs have a limited predictive accuracy, since they fail to capture the fact that the same feature can behave differently when interacting with features from different fields. To resolve this issue, many variants that incorporate field information have been proposed in recent years, including field-aware (Juan et al., 2016a), field-weighted (Pan et al., 2018a), and field-embedded (Pande, 2021; Sun et al., 2021a) factorization machines. Unfortunately, none of these variants admit an equivalent formulation whose computational complexity is linear in the number of features. This poses a challenge when building cost-effective, large-scale real-time item recommendation systems, as the additional gain from incorporating field information comes at the cost of additional computing power.

The focus of this work is the *field-weighted factorization machine* (FwFM) variant, under the large-scale and low inference latency regime. It is similar to a regular factorization machine, with an additional symmetric matrix of parameters modeling the strength of pairwise field interactions. It is attractive in practice due to the fact that in terms of memory consumption and the number of parameters, it is on par with a regular FM. In addition, it is less prone to over-fitting, as pointed out by Juan et al. (2016a) and Pan et al. (2018a), and admits a simple heuristic for reducing the computational complexity at the inference stage by pruning low-magnitude field interactions. However, pruning has some cost in terms of accuracy. In this work, we devise a different way to significantly reduce the computational cost of field-weighted factorization machines by decomposing the matrix of pairwise field interactions. Motivated by the visualization of the field interaction matrices in Pan et al. (2018a), that resembles a block-like structure due to field groups exhibiting similar interaction behavior, we employ the well-known diagonal plus low-rank (DPLR) matrix decomposition, to produce an FwFM variant we call DPLR-FwFM. We show that our idea facilitates utilizing the additional accuracy provided by incorporating field information, while benefiting from a per-item computational cost that is linear in the number of item fields, albeit slower than regular FM by a factor proportional to the rank of the low-rank part of the decomposition.

We evaluate our technique on public datasets, and on proprietary data from an large-scale online advertising system of a major company. Our results demonstrate that, in practice, DPLR-FwFM with extremely low ranks outperform aggressively pruned FwFM models both in terms of latency and accuracy. Finally, we deploy our solution in an online advertising system of a major company, and show that the latency incurred by a prediction module based on our DPLR-FwFM model is better than that of the module currently used in production and is based on a pruned FwFM.

To summarize, the main contributions of our paper are:

1. Reformulate FwFM models to make their computational cost on par, or higher by a small constant factor of our choice, compared to regular FMs, and significantly cheaper than FwFMs. This, while benefiting from the higher accuracy provided by incorporating field information.
2. We show that the accuracy of the obtained models is on par or higher than pruned FwFM models on public and proprietary data-sets.
3. We demonstrate that, in practice, in a real-world online advertising system, our approach can significantly reduce the computational costs, thus achieving lower latency with the same compute power, or reducing the compute power required to achieve a given latency.

2 RELATED WORK

Recommendation technologies such as *Collaborative Filtering* (CF) (Goldberg et al., 1992), help users discover new items based on past preferences. *Matrix Factorization* (MF) is a leading approach in CF, addressing data expansion and sparsity by using a latent factor model (Koren et al., 2009). Such systems find applications in various domains, including movie recommendation (Bell & Koren, 2007), music recommendation (Aizenberg et al., 2012), ad matching (Aharon et al., 2013) and more. *Factorization machine* (FM) variants excel in benchmark tabular classification tasks, particularly when dealing with a large number of interactions and requiring fast predictions. This

family includes several members as the original FM (Rendle, 2010b), *field-aware factorization machine* (FFM) (Juan et al., 2016b), *field weighted factorization machine* (FwFM) (Pan et al., 2018b) and the *field-matrixed factorization machine* (FmFM) (Sun et al., 2021b). For handling a wider range of features with non-linear interactions, Cheng et al. (2016) presented *Wide & Deep* learning approach, which trains a network that combines a linear model and a deep neural network. Their work was followed by a wave of deep learning techniques aimed at improving prediction accuracy, while also being sensitive to low-latency efficiency. An essential application that motivated such works, and also the focus of this work, is CTR prediction for online advertising. It has been the topic of many techniques as DeepFM (Guo et al., 2017), xDeepFM (Lian et al., 2018), DCN (Wang et al., 2017), AutoInt (Song et al., 2019), DeepLight (Deng et al., 2021) and more. While these methods offer improved predictions, their high latency and long training time cause many applications to still prefer FM-based approaches.

Low-rank factorizations have emerged as a versatile and powerful tool in *machine learning* (ML) and *artificial intelligence* (AI), finding applications across a spectrum of domains. For instance, low-rank matrix-factorization in the context of recommendation systems (Koren et al., 2009; Rendle, 2010b), in natural language processing for topic modeling, and dimensionality reduction (Blei et al., 2003), and more recently, for large-scale pre-training and fine-tuning of models on diverse natural language understanding tasks (Hu et al., 2021). In cases where the involved matrices are full-ranked and low-rank factorizations are less effective, *diagonal plus low-rank factorization* (DPLR) may be considered (see the work of Saunderson et al. (2012) and references therein). DPLR involves approximating a given matrix, often a large, and high-dimensional one, by decomposing it into the sum of two matrices: a diagonal matrix and a low-rank matrix. In statistics, DPLR is used to approximate high-dimensional covariance matrices of multivariate normal distributions (Liutkus & Yoshii, 2017; Ong et al., 2018). In ML, DPLR is used for enhancing the efficiency of models (see the work of Bingham et al. (2019) for an ML library that uses DPLR covariance matrices). For instance, DPLR-based methods have been employed in deep learning architectures for more efficient and accurate models (Zhao et al., 2016; Tomczak et al., 2020; Mishkin et al., 2018). Inspired by an observation of Lin et al. (2018), that the regular FM field-interaction matrix equals a rank-one matrix minus the identity matrix, we apply DPLR to FwFM and improve its training and inference efficiency.

One FM like approach that resembles our work is of Almagor & Hoshen (2022), where the authors apply a low-rank factorization to a generalized FM family representation (that includes FM, FFM, FwFM and FmFM as special cases). At first glance, their method seems to have much in common with our work, as both works use low-rank factorization for the fields interaction matrix. However, a closer examination reveals that while Almagor & Hoshen (2022) method may reduce the cost of computing the interaction of a *single pair* of fields, our method *avoids computing pairwise interactions altogether*. Thus, the methods differ significantly.

3 BACKGROUND AND PROBLEM FORMULATION

In this section, we provide a detailed account of the computational efficiency properties of factorization machines in recommending items and the efficiency drop introduced by switching to field-aware variants. We then formulate the problem of reducing the gap between the two.

3.1 EFFICIENCY OF FACTORIZATION MACHINES

Given a feature vector $\mathbf{x} \in \mathbb{R}^n$, the FM model proposed by Rendle (2010a) computes

$$\Phi_{\text{FM}}(\mathbf{x}; b_0, \mathbf{b}, \mathbf{w}_1, \dots, \mathbf{w}_n) = b_0 + \langle \mathbf{b}, \mathbf{x} \rangle + \sum_{i=1}^n \sum_{j=i+1}^n \langle x_i \mathbf{w}_i, x_j \mathbf{w}_j \rangle,$$

where $b_0 \in \mathbb{R}$, $\mathbf{b} \in \mathbb{R}^n$, and $\mathbf{w}_1, \dots, \mathbf{w}_n \in \mathbb{R}^k$ are trainable parameters. Overall, the model has $1 + n + nk$ trainable parameters, and the direct formula above for Φ_{FM} can be computed in $O(n^2k)$ time. But as Rendle (2010a) pointed out, the pairwise interaction terms can be re-written as

$$\sum_{i=1}^n \sum_{j=i+1}^n \langle x_i \mathbf{w}_i, x_j \mathbf{w}_j \rangle = \frac{1}{2} \left(\left\| \sum_{i=1}^n x_i \mathbf{w}_i \right\|^2 - \sum_{i=1}^n \left\| x_i \mathbf{w}_i \right\|^2 \right) \quad (1)$$

This dramatically reduces the time complexity to $O(nk)$.

In practice, the feature vector encodes a row in a tabular data-set of past interactions between users and items, whose columns, often named *fields*, contain categorical features. Some columns describe the context (including the user), whereas others describe the item. Thus, \mathbf{x} formally consists of field-wise one-hot encodings, for example:

$$\mathbf{x} = \underbrace{\left(\underbrace{0, 1, 0, 0, 0}_{\text{field 1}}, \underbrace{1, 0, 0, 0}_{\text{field 2}}, \dots, \underbrace{0, 0, 1, 0, 0}_{\text{field } m_c}, \underbrace{0, 0, 1, 0}_{\text{field } m_c + 1}, \dots, \underbrace{0, 0, 0, 0, 1}_{\text{field } m} \right)^T}_{\text{context fields}} \underbrace{\hspace{10em}}_{\text{item fields}}$$

Therefore, when computing Φ_{FM} we can only use the m nonzero entries of \mathbf{x} corresponding to the m fields. We note that there may be fields having multiple values, such as a list of movie genres, but we defer their treatment to the discussion of FwFM models in the sequel.

Formally, for a given feature vector \mathbf{x} , let ℓ_1, \dots, ℓ_m denote the nonzero indices, and define $\mathbf{v}_i \equiv \mathbf{w}_{\ell_i}$. Hence, the formula in equation 1 can be reformulated by summing over the fields, rather than the features. Moreover, we can split the summation over all fields to a summation over the context fields $\mathcal{C} = \{1, \dots, m_c\}$ and the item fields $\mathcal{I} = \{m_c + 1, \dots, m\}$, and obtain:

$$\sum_{i=1}^n \sum_{j=i+1}^n \langle x_i \mathbf{w}_i, x_j \mathbf{w}_j \rangle = \sum_{i=1}^m \sum_{j=i+1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle \quad (2a)$$

$$= \frac{1}{2} \left(\left\| \sum_{i=1}^m \mathbf{v}_i \right\|^2 - \sum_{i=1}^m \|\mathbf{v}_i\|^2 \right) \quad (2b)$$

$$= \frac{1}{2} \left(\left\| \sum_{i \in \mathcal{C}} \mathbf{v}_i + \sum_{i \in \mathcal{I}} \mathbf{v}_i \right\|^2 - \sum_{i \in \mathcal{C}} \|\mathbf{v}_i\|^2 - \sum_{i \in \mathcal{I}} \|\mathbf{v}_i\|^2 \right) \quad (2c)$$

By equation 2c, we see that when ranking a large number of items for a given context, the sums over \mathcal{C} can be computed only *once*. For each item, the computational complexity is $O(|\mathcal{I}|k)$, namely, and it depends largely on the number of *item fields*. Thus, in practice, we can use a model with a large number of user or context features to achieve a high degree of personalization, without incurring a significant computational cost during item ranking.

3.2 INEFFICIENCY OF FIELD-WEIGHTED FACTORIZATION MACHINES

Field-weighted factorization machines (FwFM) (Pan et al., 2018a) model the varying behavior of a feature belonging to some field when interacting with features from different fields in the form of a trainable symmetric field interaction matrix $\mathbf{R} \in \mathbb{R}^{m \times m}$. Its components $R_{i,j}$ model the *intensity* of the interaction between field i and field j . The model’s output is

$$\Phi_{\text{FwFM}}(\mathbf{x}; b_0, \mathbf{b}, \mathbf{w}_1, \dots, \mathbf{w}_n, \mathbf{R}) = b_0 + \langle \mathbf{b}, \mathbf{x} \rangle + \sum_{i=1}^n \sum_{j=i+1}^n \langle x_i \mathbf{w}_i, x_j \mathbf{w}_j \rangle R_{f_i, f_j},$$

where f_i is the field corresponding to feature i . Pan et al. (2018a) demonstrate that this model family achieves a significantly higher accuracy compared to a regular FM, and is comparable with other field-aware variants. One attractive property of this variant is its number of parameters and memory requirements - it has only $\frac{m(m-1)}{2}$ additional parameters compares to a regular FM comprising the above-diagonal entries of the field interaction matrix \mathbf{R} .

Similarly to equation 2a, under the one-hot encoding assumption, given an input \mathbf{x} , the pairwise interaction term can be written in terms of the field vectors $\mathbf{v}_1, \dots, \mathbf{v}_m$ that are a subset of the embedding vectors $\mathbf{w}_1, \dots, \mathbf{w}_n$ that correspond to the m nonzero entries of \mathbf{x} :

$$\sum_{i=1}^n \sum_{j=i+1}^n \langle x_i \mathbf{w}_i, x_j \mathbf{w}_j \rangle R_{f_i, f_j} = \sum_{i=1}^m \sum_{j=i+1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle R_{i,j} \quad (3)$$

However, the time complexity of computing its output Φ_{FwFM} is dominated by the $O(m^2k)$ complexity of computing the pairwise term in equation 3, which is quadratic in the total number of

fields. Compared to the $O(|\mathcal{I}|k)$ per item complexity of a regular FM, which is linear in the number of item fields, FwFMs pose a serious challenge for applications where inference speed is critical. While the quadratic complexity property is shared by many other field-aware variants, such as Juan et al. (2016a); Pande (2021); Sun et al. (2021a), in this work, we focus on addressing this challenge for FwFM models.

The one-hot encoding assumption does not cover the case when a field has multiple values, such as a list of movie genres. In this case, we assume that multiple components of a field’s encoding, that correspond to multiple values, may be nonzero. We handle this case by assuming that a field does not interact with itself, i.e. $R_{f,f} = 0$ for any field f . A direct consequence is that the pairwise interaction term can be written in terms of field vectors, as in equation 3, but each vector v_i may be a weighted sum of the corresponding feature embedding vectors. For example, a movie with 3 genres may be encoded by placing $\frac{1}{3}$ in the corresponding components of \mathbf{x} , which will result in an average of the genre embedding vectors.

3.3 PROBLEM FORMULATION

A typical approach used in practice to speed up inference is pruning the matrix \mathbf{R} by zeroing-out entries whose magnitude is below a threshold, as suggested by Pan et al. (2018a) and Sun et al. (2021a). However, as we show in section 5, aggressive pruning reduces the accuracy of the model. In this work, we aim to devise a method to achieve inference speed that is proportional to the number of *item fields*, and is only ρ times slower than a regular FM, where $0 < \rho \ll m$ is a configurable factor. This, while retaining an accuracy that is comparable with that of a regular FwFM model.

4 LOW RANK FIELD-WEIGHTED FACTORIZATION MACHINES

In this section we reformulate the pairwise field interaction term in equation 3 in an equivalent, but significantly more efficient manner. Throughout this section, we assume that the field vectors $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^k$ are embedded into the *rows* of the matrix $\mathbf{V} \in \mathbb{R}^{m \times k}$:

$$\mathbf{V} = \begin{bmatrix} -\mathbf{v}_1 - \\ -\mathbf{v}_2 - \\ \vdots \\ -\mathbf{v}_m - \end{bmatrix} \quad (4)$$

Moreover, since equation 3 uses only the upper triangular part of \mathbf{R} , we may choose the remaining entries arbitrarily to our convenience, and throughout this section we assume that \mathbf{R} is symmetric with a zero diagonal. Under this assumption, equation 3 can be re-written as

$$\sum_{i=1}^m \sum_{j=i+1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle R_{i,j} = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle R_{i,j}. \quad (5)$$

In the sequel, we describe our method to efficiently compute the double sum in the right-hand side of equation 5. We first review the mathematical background, and then the complete method.

4.1 MATHEMATICAL FOUNDATION

In this section we present a technical result that provides the motivation for factorizing the matrix \mathbf{R} to obtain an efficient algorithm. Then, we present an example explaining why the widely used *low-rank factorization* lacks expressive power, and use it to motivate a *diagonal plus low-rank factorization*.

The following identity let us reformulate the pairwise interaction formula as the one in equation 5 in a matrix form, that is inspired by the reformulation in Lin et al. (2018) for regular FMs.

Identity 1. Let $\mathbf{A} \in \mathbb{R}^{m \times k}$. Then, for any matrix $\mathbf{Q} \in \mathbb{R}^{m \times m}$ we have

$$\sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{A}_{i,:}, \mathbf{A}_{j,:} \rangle Q_{i,j} = \text{Tr}(\mathbf{A}^T \mathbf{Q} \mathbf{A}) \quad (6)$$

The proof can be found in Appendix A.

Since \mathbf{R} is a real symmetrical square matrix its eigenvalue decomposition is composed of real eigenvalues. If the eigenvalues decay quickly, we can use an approximation obtained by using the ρ largest magnitude eigenvalues, for some $0 < \rho \ll m$, in the form $\mathbf{R} = \mathbf{U}^T \text{diag}(\mathbf{e})\mathbf{U}$, where $\mathbf{U} \in \mathbb{R}^{\rho \times m}$ and $\mathbf{e} \in \mathbb{R}^\rho$.

However, a low rank decomposition lacks the expressive power we need. To see why, consider a regular FM, which is equivalent to an FwFM where all field interactions are 1, meaning that the field-interaction matrix is:

$$\mathbf{R}_{\text{FM}} = \begin{pmatrix} 0 & 1 & \dots & 1 \\ 1 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 0 \end{pmatrix}$$

We have $\text{rank}(\mathbf{R}_{\text{FM}}) = m$, and its eigenvalues are $\{m - 1, -1, \dots, -1\}$, and hence a low-rank approximation within a reasonable accuracy is impossible. Since a low-rank decomposition does not even have enough expressive to model a regular FM, an FwFM is clearly out of reach. The obstacle stems from the zero diagonal of \mathbf{R} , which is an ‘‘anomaly’’ in the matrix. But as Lin et al. (2018) pointed out, this anomaly can be easily handled using a diagonal matrix:

$$\mathbf{R}_{\text{FM}} = \mathbf{1}\mathbf{1}^T - \mathbf{I}, \quad (7)$$

where $\mathbf{1}$ is a column vector whose components are all 1. In other words, the matrix \mathbf{R}_{FM} can be decomposed into a sum of a diagonal matrix and a low-rank matrix. Motivated by the above, we use a diagonal plus low rank (DPLR) decomposition to facilitate fast inference, as is shown in the proposition below. The low-rank part resembles the eigenvalue decomposition, to be able to model arbitrary symmetric matrices, including indefinite ones.

Proposition 1. *Let \mathbf{V} be as defined in equation 4, let $\mathbf{R} \in \mathbb{R}^{m \times m}$ be a symmetric matrix, and suppose that $\mathbf{U} \in \mathbb{R}^{\rho \times m}$, $\mathbf{e} \in \mathbb{R}^\rho$, and $\mathbf{d} \in \mathbb{R}^m$ are such that*

$$\mathbf{R} = \mathbf{U}^T \text{diag}(\mathbf{e})\mathbf{U} + \text{diag}(\mathbf{d}). \quad (8)$$

Define $\mathbf{P} = \mathbf{U}\mathbf{V}$. Then,

$$\sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle R_{i,j} = \sum_{i=1}^m d_i \|\mathbf{v}_i\|^2 + \sum_{i=1}^{\rho} e_i \|\mathbf{P}_{i,:}\|^2 \quad (9)$$

The proof can be found in Appendix A. Now we are ready to present our solution based on proposition 1.

4.2 THE MODIFIED MODEL AND FAST INFERENCE ALGORITHM

Our solution consists of a modification of the FwFM model, and an algorithm that achieves inference during ranking in $O(\rho|\mathcal{I}|k)$ per item.

4.2.1 THE DIAGONAL PLUS LOW RANK FwFM MODEL

Instead of learning the field interaction matrix \mathbf{R} directly, we learn its decomposition in the diagonal plus low-rank form. Since $\text{diag}(\mathbf{R}) = \mathbf{0}$, the diagonal component is fully determined by the low-rank decomposition. Formally, we replace the learned parameter \mathbf{R} , with the learned parameters $\mathbf{U} \in \mathbb{R}^{\rho \times m}$ and $\mathbf{e} \in \mathbb{R}^\rho$, where ρ is a hyper-parameter. The matrix \mathbf{R} is formally defined as

$$\mathbf{R} = \mathbf{U}^T \text{diag}(\mathbf{e})\mathbf{U} + \text{diag}(\mathbf{d}), \quad \text{where } \mathbf{d} \equiv -\text{diag. of}(\mathbf{U}^T \text{diag}(\mathbf{e})\mathbf{U}). \quad (10)$$

In other words, we ensure that \mathbf{R} is defined to be a symmetric matrix with a zero diagonal by forming a symmetric eigenvalue-like decomposition, and subtracting the diagonal of the resulting matrix. The formal definition of \mathbf{R} allows us to apply proposition 1, but we do not actually need to compute the matrix \mathbf{R} itself at any stage. Instead, we use the proposition, and replace the FwFM pairwise interaction term in equation 3, with the outcome of the proposition. Namely, we compute $\mathbf{P} = \mathbf{U}\mathbf{V}$ in $O(\rho mk)$ time, and then compute $\frac{1}{2} (\sum_{i=1}^m d_i \|\mathbf{v}_i\|^2 + \sum_{i=1}^{\rho} e_i \|\mathbf{P}_{i,:}\|^2)$ in $O(\rho k + mk)$ time.

We call the resulting model a diagonal plus low-rank field weighted factorization machine, or DPLR-FwFM. Although this is not our main objective, a nice byproduct is that training epochs also become slightly faster, since the above two steps cost $O(\rho mk)$ time instead the naïve $O(m^2k)$. This is important for some online advertising systems that require deploying a 'fresh' model as soon as possible, to quickly adapt to the quickly evolving marketplace conditions (Sculley et al., 2015; Lommatzsch et al., 2017; Zhang et al., 2020). However, for inference during *ranking* we need to take the result of the proposition one step further.

4.2.2 FAST INFERENCE FOR ITEM RANKING

Observe that the field vectors embedded into the rows of \mathbf{V} can be decomposed into context and item vectors:

$$\mathbf{V} = \begin{bmatrix} \mathbf{V}_C \\ \mathbf{V}_I \end{bmatrix}$$

The corresponding columns of \mathbf{U} can be decomposed similarly:

$$\mathbf{U} = [\mathbf{U}_C \quad \mathbf{U}_I]$$

Consequently, matrix \mathbf{P} can be written as:

$$\mathbf{P} = \mathbf{U}_C \mathbf{V}_C + \mathbf{U}_I \mathbf{V}_I,$$

and the sum $\sum_{i=1}^m d_i \|\mathbf{v}_i\|^2$ can also be decomposed as

$$\sum_{i=1}^m d_i \|\mathbf{v}_i\|^2 = \sum_{i \in \mathcal{C}} d_i \|\mathbf{v}_i\|^2 + \sum_{i \in \mathcal{I}} d_i \|\mathbf{v}_i\|^2$$

To summarize, during ranking, the pairwise interaction term of each item during ranking using DPLR-FwFM can be computed by the following algorithm:

Algorithm 1 - DPLR FwFM item pairwise interaction for inference with cached context

Input:

- The field matrix \mathbf{V} , partitioned into $\mathbf{V}_C, \mathbf{V}_I$
- $\mathbf{U} \in \mathbb{R}^{\rho \times m}$, $\mathbf{d} \in \mathbb{R}^m$, $\mathbf{e} \in \mathbb{R}^\rho$, that form the decomposition $\mathbf{R} = \mathbf{U}^T \text{diag}(\mathbf{e})\mathbf{U} + \text{diag}(\mathbf{d})$, with \mathbf{U} partitioned into $\mathbf{U}_C, \mathbf{U}_I$

Output: The pairwise interactions $\sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle R_{i,j}$

Steps:

1. Once per context, compute:
 - (a) Compute $\mathbf{P}_C = \mathbf{V}_C \mathbf{U}_C$
 - (b) Compute $s_C = \sum_{i \in \mathcal{C}} d_i \|\mathbf{v}_i\|^2$
2. Compute $\mathbf{P} = \mathbf{P}_C + \mathbf{U}_I \mathbf{V}_I$
3. Return $s_C + \sum_{i \in \mathcal{I}} d_i \|\mathbf{v}_i\|^2 + \sum_{i=1}^\rho e_i \|\mathbf{P}_{i,:}\|^2$

The first step is computed only once per context in $O(\rho|\mathcal{C}|k)$ time, and its results are cached. The last two steps are computed for every item, and their complexity is $O(\rho|\mathcal{I}|k)$ per item, as we desire.

5 EXPERIMENTS

In this section we present experimental results demonstrating that our approach significantly reduces the cost of item recommendation serving without compromising accuracy, in comparison to the pruning approach. Beyond the results presented here, in appendix C we demonstrate that finding a DPLR factorization of the field interaction matrix of a regular FwFM, instead of training a DPLR representation, may not be a good approach in practice.

5.1 RESULTS FOR PUBLIC DATA-SETS

We compare our approach to an FwFM, an FM, and a pruned FwFM using the Criteo Display Advertising Challenge¹ data-set, comprising of 39 fields and about 30M samples, and the MovieLens 1M (Harper & Konstan, 2015) data-set, with 8 fields and 1M samples. Out of the available MovieLens data-sets, we chose the one with the most number of informative context and item fields. The timestamp is converted into year, month, day of week, and hour of day, creating a 11 field data-set. The "genres" field in the MovieLens data-set has multiple values, and as described in section 3.2, we average the genre embedding vectors to produce one genre vector.

Both data-sets were randomly split into 80% training, 10% validation, and 10% test sets. The validation sets were used for tuning the learning rate using Optuna (Akiba et al., 2019). Features with less than 10 occurrences in the training set, and those appearing in the test and validation set that did not appear in the training set, are replaced with a special 'rare feature'. Numerical features are binned similarly to the strategy used by the Criteo Challenge winners, via the $x \rightarrow \lfloor \ln^2(x) \rfloor$ function. All models use 8-dimensional latent vectors. Code will be made available upon acceptance.

We compare DPLR with pruning by matching the number of entries we leave in the pruned field interaction matrix to the number of parameters in the DPLR model. For a rank ρ , we use $\rho(m+1)$ parameters, and the corresponding pruned model is left with the largest magnitude $\rho(m+1)$ field interaction coefficients. Equivalently, we leave $100 \times \frac{2\rho(m+1)}{m(m-1)}$ percent of the field interactions.

The results are summarized in table 1. We see that aggressive pruning, leaving less than 20% of the interactions, indeed degrades performance, and our DPLR solution with a rank of at least two outperforms the pruned models. Surprisingly, DPLR models with a rank of 1 perform worse than an FM, although in theory their expressive power is stronger.

Table 1: Results for public data-sets. A lower LogLoss or MSE, and a higher AUC are better. The last column shows the improvement percentage of the DPLR over the equivalently pruned model.

	Rank	% left	FM	FwFM	DPLR	Pruned	DPLR vs Pruned (%)
LogLoss (Criteo)	1	5.40%	0.4467	0.4428	0.4470	0.4516	1.02%
	2	10.80%	0.4467	0.4428	0.4458	0.4471	0.30%
	3	16.19%	0.4467	0.4428	0.4456	0.4458	0.05%
	4	21.59%	0.4467	0.4428	0.4450	0.4448	-0.03%
	5	26.99%	0.4467	0.4428	0.4446	0.4439	-0.17%
AUC (Criteo)	1	5.40%	0.8051	0.8059	0.8023	0.7984	0.49%
	2	10.80%	0.8051	0.8059	0.8065	0.7661	5.27%
	3	16.19%	0.8051	0.8059	0.8057	0.8005	0.65%
	4	21.59%	0.8051	0.8059	0.8059	0.8073	-0.17%
	5	26.99%	0.8051	0.8059	0.8071	0.8079	-0.10%
MSE (MovieLens)	1	21.82%	0.7301	0.7289	0.7353	0.7469	1.55%
	2	43.64%	0.7301	0.7289	0.7262	0.7308	0.63%

5.2 PROPRIETARY ONLINE ADVERTISING SYSTEM

We validate our approach on an online advertising system of a large commercial company, by comparing the trained model's accuracy, and the ad ranking latency. Our tests are performed on FwFM models for *click-through rate* (CTR) having 82 fields. Before deployment, the FwFM's field interaction matrix is pruned to keep only 10% of the largest magnitude entries to conform to the strict latency requirements.

¹<https://www.kaggle.com/c/criteo-display-ad-challenge>

5.2.1 ACCURACY EXPERIMENT

The model trains periodically in a 'sliding window' mode, meaning that in each period associated with time T , the model is trained on logged data from the month before T , and is evaluated on logged data from the day before T . After sub-sampling to allow training within a reasonable amount of time, the training set comprises tens of millions logged user-ad interactions. We train and evaluate the model on 7 consecutive intervals using the LogLoss and AUC metrics, and average the results by weighting them according to the number of evaluation samples. The performance of the models, and the improvement of DPLR over pruned models, are summarized in Table 2. Surprisingly, the DPLR models perform *better* on the evaluation set than a regular FwFM models. We believe it means that a low-rank field interaction matrix is a good regularization prior for this specific domain. Nevertheless, as our results on public data-sets show, this is not the case in general.

Table 2: AUC and LogLoss improvements versus an FwFM on proprietary data. Higher is better.

Rank	1	2	3	4	5	6
LogLoss lift (%)	-0.478%	0.228%	0.363%	0.409%	0.445%	0.404%
AUC lift (%)	-0.132%	0.071%	0.098%	0.112%	0.122%	0.115%

5.2.2 LATENCY EXPERIMENT

We deploy the model with rank 3, to a test environment that serves a small portion of the traffic, a few thousand ad ranking queries per minute. The rank was chosen to correspond to pruning 90% of the field interaction entries in terms of the number of parameters. Out of the model's 63 fields, 38 are item fields, and therefore we can theoretically expect approximately 40% latency inference speed improvement in the best case.

We measure the latency incurred by the inference for ad CTR prediction, and the total latency incurred by ranking all eligible ads for a given query. Table 3 summarizes the results, and appendix B presents full time-series plots. Evidently, the inference latency is improved by a few dozen percents, whereas the query latency by $\sim 5\%$, since CTR prediction is only one component of the ad query serving algorithm.

Table 3: Aggregated latency of a single inference operation (left) and a single ranking operation (right). Our system measures, in each minute, the average, 95th percentile, and the 99th percentile latency for inference, and the 95th percentile latency for ranking. The measurements were aggregated over a 10 hour period. The lifts row presents latency improvement - a higher lift is better.

	Inference per ad			Ranking
	Average	95 th percentile	99 th percentile	95 th percentile
Lift (%)	34.27%	29.11%	25.57%	5.45%

6 CONCLUSIONS AND FUTURE WORK

In this work we proposed learning a diagonal plus low-rank decomposition of the field interaction matrix of FwFM models as an alternative to the commonly used pruning heuristic in large scale low-latency recommendation systems. We demonstrated that our approach has the potential to outperform pruned models in both item recommendation speed and model accuracy.

An immediate future direction is attempting to apply a more general decomposition to the larger family of Sun et al. (2021a) models, that includes field-aware and field-weighted FMs via, for example, a higher-order tensor decomposition. Moreover, looking at equation 9, we observe that the columns of U have a notion of 'field importance' - if all the entries of $U_{:,i}$ are negligible, the effect of the i -th field is negligible, and it can be discarded. Hence, another future direction is mathematically or experimentally quantifying this notion of importance, as an alternative to the approach in Kaplan et al. (2021). This is especially important for real-time systems, where just reducing the number of fields may have a tremendous effect on latency and recommendation costs.

REFERENCES

- Michal Aharon, Natalie Aizenberg, Edward Bortnikov, Ronny Lempel, Roi Adadi, Tomer Benyamini, Liron Levin, Ran Roth, and Ohad Serfaty. OFF-set: one-pass factorization of feature sets for online recommendation in persistent cold start settings. In *RecSys'2013*, pp. 375–378, 2013.
- Natalie Aizenberg, Yehuda Koren, and Oren Somekh. Build your own music recommender by modeling internet radio streams. In *Proceedings of the 21st international conference on World Wide Web*, pp. 1–10. ACM, 2012.
- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.
- Chen Almagor and Yedid Hoshen. You say factorization machine, i say neural network-it’s all in the activation. In *Proceedings of the 16th ACM Conference on Recommender Systems*, pp. 389–398, 2022.
- Robert M Bell and Yehuda Koren. Lessons from the netflix prize challenge. *Acm Sigkdd Explorations Newsletter*, 9(2):75–79, 2007.
- Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul A. Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep universal probabilistic programming. *J. Mach. Learn. Res.*, 20:28:1–28:6, 2019. URL <http://jmlr.org/papers/v20/18-403.html>.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishu Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*, pp. 7–10, 2016.
- Wei Deng, Junwei Pan, Tian Zhou, Deguang Kong, Aaron Flores, and Guang Lin. Deeplight: Deep lightweight feature interactions for accelerating ctr predictions in ad serving. In *Proceedings of the 14th ACM international conference on Web search and data mining*, pp. 922–930, 2021.
- David Goldberg, David Nichols, Brian M Oki, and Douglas Terry. Using collaborative filtering to weave an information tapestry. *Communications of the ACM*, 35(12):61–70, 1992.
- Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. Deepfm: a factorization-machine based neural network for ctr prediction. *arXiv preprint arXiv:1703.04247*, 2017.
- F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Yuchin Juan, Yong Zhuang, Wei-Sheng Chin, and Chih-Jen Lin. Field-aware factorization machines for ctr prediction. In *Proceedings of the 10th ACM Conference on Recommender Systems*, RecSys ’16, pp. 43–50, New York, NY, USA, 2016a. Association for Computing Machinery. ISBN 9781450340359. doi: 10.1145/2959100.2959134. URL <https://doi.org/10.1145/2959100.2959134>.
- Yuchin Juan, Yong Zhuang, Wei-Sheng Chin, and Chih-Jen Lin. Field-aware factorization machines for ctr prediction. In *Proceedings of the 10th ACM conference on recommender systems*, pp. 43–50, 2016b.
- Yohay Kaplan, Yair Koren, Rina Leibovits, and Oren Somekh. Dynamic length factorization machines for ctr prediction. In *2021 IEEE International Conference on Big Data (Big Data)*, pp. 1950–1959. IEEE, 2021.

- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- Jianxun Lian, Xiaohuan Zhou, Fuzheng Zhang, Zhongxia Chen, Xing Xie, and Guangzhong Sun. xdeepfm: Combining explicit and implicit feature interactions for recommender systems. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1754–1763, 2018.
- Xiao Lin, Wenpeng Zhang, Min Zhang, Wenwu Zhu, Jian Pei, Peilin Zhao, and Junzhou Huang. Online compact convexified factorization machine. In *Proceedings of the 2018 World Wide Web Conference, WWW '18*, pp. 1633–1642, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee. ISBN 9781450356398. doi: 10.1145/3178876.3186075. URL <https://doi.org/10.1145/3178876.3186075>.
- Antoine Liutkus and Kazuyoshi Yoshii. A diagonal plus low-rank covariance model for computationally efficient source separation. In *2017 IEEE 27th International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1–6. IEEE, 2017.
- Andreas Lommatzsch, Benjamin Kille, and Sahin Albayrak. Incorporating context and trends in news recommender systems. In *Proceedings of the international conference on web intelligence*, pp. 1062–1068, 2017.
- Aaron Mishkin, Frederik Kunstner, Didrik Nielsen, Mark Schmidt, and Mohammad Emtiyaz Khan. Slang: Fast structured covariance approximations for bayesian deep learning with natural gradient. *Advances in Neural Information Processing Systems*, 31, 2018.
- Victor M-H Ong, David J Nott, and Michael S Smith. Gaussian variational approximation with a factor covariance structure. *Journal of Computational and Graphical Statistics*, 27(3):465–478, 2018.
- Junwei Pan, Jian Xu, Alfonso Lobos Ruiz, Wenliang Zhao, Shengjun Pan, Yu Sun, and Quan Lu. Field-weighted factorization machines for click-through rate prediction in display advertising. In *Proceedings of the 2018 World Wide Web Conference, WWW '18*, pp. 1349–1357, New York, NY, USA, 2018a. Association for Computing Machinery. ISBN 9781450356398. doi: 10.1145/3178876.3186040. URL <https://doi.org/10.1145/3178876.3186040>.
- Junwei Pan, Jian Xu, Alfonso Lobos Ruiz, Wenliang Zhao, Shengjun Pan, Yu Sun, and Quan Lu. Field-weighted factorization machines for click-through rate prediction in display advertising. In *Proceedings of the 2018 World Wide Web Conference*, pp. 1349–1357, 2018b.
- Harshit Pande. Field-embedded factorization machines for click-through rate prediction, 2021.
- Kaare Brandt Petersen and Michael Syskind Pedersen. The matrix cookbook, 2012.
- Steffen Rendle. Factorization machines. In *2010 IEEE International Conference on Data Mining*, pp. 995–1000, 2010a. doi: 10.1109/ICDM.2010.127.
- Steffen Rendle. Factorization machines. In *2010 IEEE International conference on data mining*, pp. 995–1000. IEEE, 2010b.
- James Saunderson, Venkat Chandrasekaran, Pablo A Parrilo, and Alan S Willsky. Diagonal and low-rank matrix decompositions, correlation matrices, and ellipsoid fitting. *SIAM Journal on Matrix Analysis and Applications*, 33(4):1395–1416, 2012.
- David Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-Francois Crespo, and Dan Dennison. Hidden technical debt in machine learning systems. *Advances in neural information processing systems*, 28:2503–2511, 2015.
- Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. Autoint: Automatic feature interaction learning via self-attentive neural networks. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pp. 1161–1170, 2019.

- Yang Sun, Junwei Pan, Alex Zhang, and Aaron Flores. Fm2: Field-matrixed factorization machines for recommender systems. In *Proceedings of the Web Conference 2021*, WWW '21, pp. 2828–2837, New York, NY, USA, 2021a. Association for Computing Machinery. ISBN 9781450383127. doi: 10.1145/3442381.3449930. URL <https://doi.org/10.1145/3442381.3449930>.
- Yang Sun, Junwei Pan, Alex Zhang, and Aaron Flores. Fm2: Field-matrixed factorization machines for recommender systems. In *Proceedings of the Web Conference 2021*, pp. 2828–2837, 2021b.
- Marcin Tomczak, Siddharth Swaroop, and Richard Turner. Efficient low rank gaussian variational inference for neural networks. *Advances in Neural Information Processing Systems*, 33:4610–4622, 2020.
- Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. Deep & cross network for ad click predictions. In *Proceedings of the ADKDD'17*, pp. 1–7. 2017.
- Yang Zhang, Fuli Feng, Chenxu Wang, Xiangnan He, Meng Wang, Yan Li, and Yongdong Zhang. How to retrain recommender system? a sequential meta-learning method. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 1479–1488, 2020.
- Yong Zhao, Jinyu Li, and Yifan Gong. Low-rank plus diagonal adaptation for deep neural networks. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5005–5009. IEEE, 2016.

A PROPOSITION PROOFS

A.1 PROOF OF IDENTITY 1

Proof. Recall the definition of the Frobenius inner product of matrices:

$$\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i=1}^m \sum_{j=1}^n A_{i,j} B_{i,j} = \text{Tr}(\mathbf{A}^T \mathbf{B}),$$

and the circular shift invariance property of the trace operator (Petersen & Pedersen, 2012, Section 1.1):

$$\text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{BCA})$$

For any $1 \leq i, j \leq m$ it holds that $\langle \mathbf{A}_{i,:}, \mathbf{A}_{j,:} \rangle = (\mathbf{AA}^T)_{i,j}$, and thus

$$\sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{A}_{i,:}, \mathbf{A}_{j,:} \rangle Q_{i,j} = \sum_{i=1}^m \sum_{j=1}^m (\mathbf{AA}^T)_{i,j} Q_{i,j} = \langle \mathbf{AA}^T, \mathbf{Q} \rangle = \text{Tr}(\mathbf{AA}^T \mathbf{Q}),$$

where the last two equalities follow from the definition of the Frobenius inner product, and the fact that \mathbf{AA}^T is symmetric. Finally, the circular shift invariance property implies that $\text{Tr}(\mathbf{AA}^T \mathbf{Q}) = \text{Tr}(\mathbf{A}^T \mathbf{QA})$, completing the proof. \square

A.2 PROOF OF PROPOSITION 1

Proof. Invoking identity 1 with $\mathbf{A} = \mathbf{V}$ and $\mathbf{Q} = \mathbf{R}$, and then using the decomposition in equation 8 we have

$$\begin{aligned} \sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle R_{i,j} &= \text{Tr}(\mathbf{V}^T \mathbf{RV}) \\ &= \text{Tr}(\mathbf{V}^T (\text{diag}(\mathbf{d}) + \mathbf{U}^T \cdot \text{diag}(\mathbf{e}) \cdot \mathbf{U}) \mathbf{V}) \\ &= \text{Tr}(\mathbf{V}^T \text{diag}(\mathbf{d}) \mathbf{V}) + \text{Tr}(\underbrace{(\mathbf{UV})^T}_{\mathbf{P}^T} \text{diag}(\mathbf{e}) \underbrace{(\mathbf{UV})}_{\mathbf{P}}). \end{aligned}$$

Invoking identity 1 with $\mathbf{A} = \mathbf{V}$ and $\mathbf{Q} = \text{diag}(\mathbf{d})$, we obtain

$$\text{Tr}(\mathbf{V}^T \text{diag}(\mathbf{d}) \mathbf{V}) = \sum_{i=1}^m d_i \|\mathbf{v}_i\|^2,$$

and again with $\mathbf{A} = \mathbf{P}$ and $\mathbf{Q} = \text{diag}(\mathbf{e})$ we obtain

$$\text{Tr}(\mathbf{P}^T \text{diag}(\mathbf{e}) \mathbf{P}) = \sum_{i=1}^{\rho} e_i \|\mathbf{P}_{i,:}\|^2.$$

Therefore,

$$\sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle R_{i,j} = \sum_{i=1}^m d_i \|\mathbf{v}_i\|^2 + \sum_{i=1}^{\rho} e_i \|\mathbf{P}_{i,:}\|^2$$

\square

B LATENCY CHARTS

In figure 1 we observe that the improvement in the latency of the low-rank solution is consistent over time.

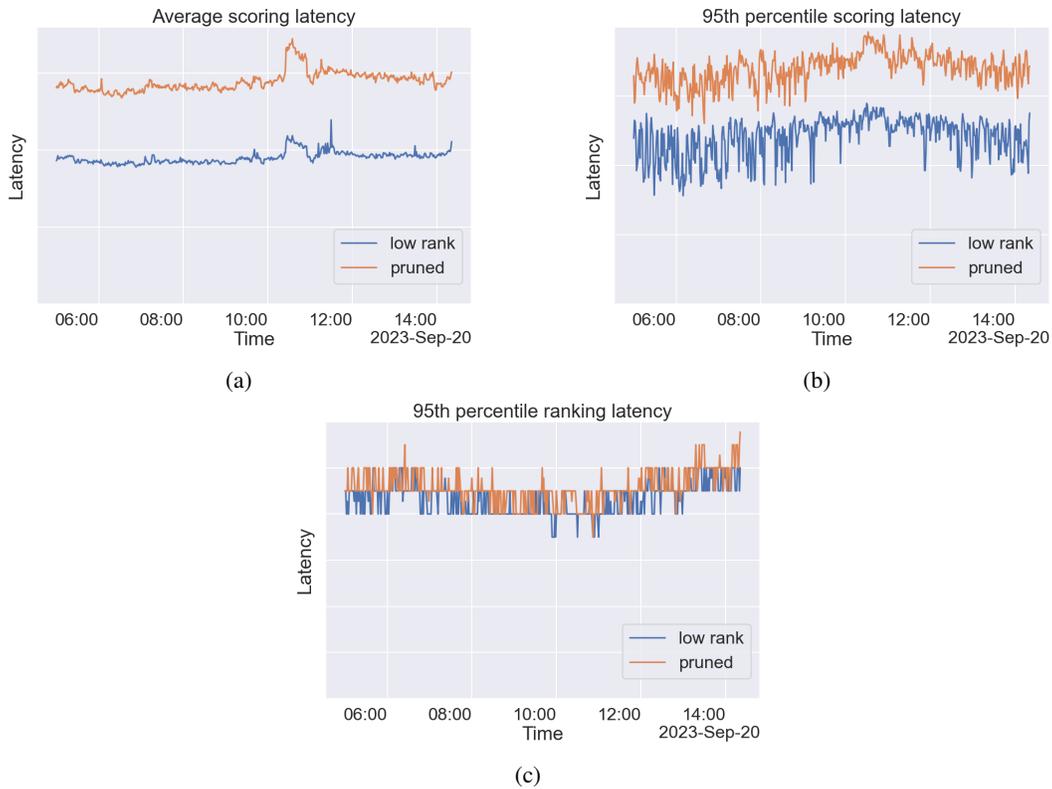


Figure 1: Latency graphs over a 10 hour period.

C POST-HOC FACTORIZATION OF THE FIELD INTERACTION MATRIX

An alternative approach to training a DPLR representation of the field interaction matrix can be training a regular FwFM model, and computing the best DPLR representation we can afterwards. We call this the “post-hoc” approach.

Unfortunately, this may not be a good idea. Suppose we obtained an approximation $\tilde{\mathbf{R}}$ of the model’s field interaction matrix \mathbf{R} . Denoting $\tilde{\mathbf{R}} = \mathbf{R} + \mathbf{E}$, where \mathbf{E} is the approximation error, by identity 1 we obtain

$$\begin{aligned} \sum_{i=1}^m \sum_{j=1}^m \langle \mathbf{v}_i, \mathbf{v}_j \rangle \tilde{R}_{i,j} &= \text{Tr}(\mathbf{V}^T \tilde{\mathbf{R}} \mathbf{V}) \\ &= \text{Tr}(\mathbf{V}^T \mathbf{R} \mathbf{V}) + \text{Tr}(\mathbf{V}^T \mathbf{E} \mathbf{V}) \\ &\leq \text{Tr}(\mathbf{V}^T \mathbf{R} \mathbf{V}) + \sum_{i=1}^m \lambda_i(\mathbf{V} \mathbf{V}^T) \sigma_i(\mathbf{E}), \end{aligned}$$

where the last inequality stems from the Von Neumann trace inequality. Thus, the approximation error boils down to the singular value spectrum of the error matrix. We took the field interaction weights obtained from an FwFM trained on the Criteo data-set, and computed the singular value spectrum of the approximation errors obtained from two approximations: (a) a DPLR approximation of rank 5, computed by minimizing the nuclear norm of the error, and (b) a pruned field interaction matrix where the top 200 entries (same number of parameters as the rank-5 DPLR approximation) were left. The error singular spectrum of both approximations is plotted in Figure 2. We observe that the large eigenvalues of the post-hoc DPLR approximation error are *much* larger than the ones of the pruned approximation error.

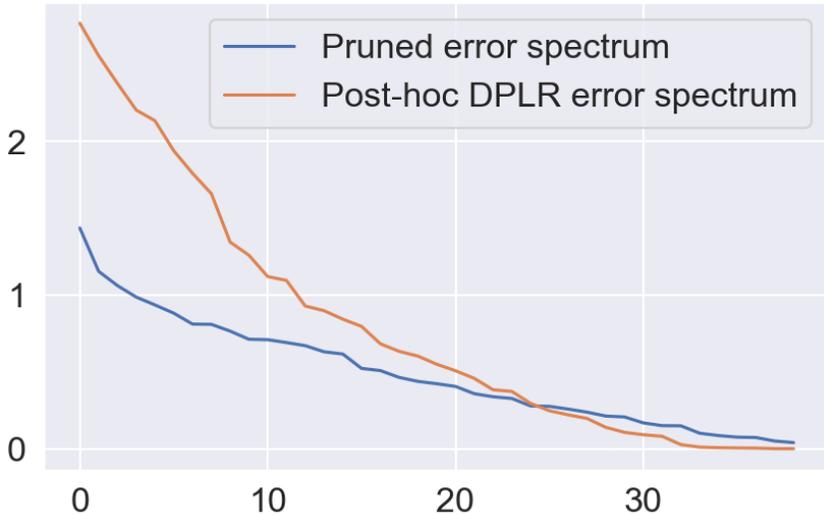


Figure 2: Singular value spectrum of two approximations of the true field interaction matrix.