
Supplementary Materials for “RaFM: Rank-Aware Factorization Machines”

1. Proof of Theorem 6

Proof. Recall the estimated gradient:

$$\widehat{grad} = \frac{\partial L(\mathcal{B}_{1,m}, y)}{\partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} + \frac{1}{N} \sum_{\mathbf{x}'} \frac{\partial L(\mathcal{B}'_{1,m}, y)}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p}} \mathbf{G}^{-1} \frac{\partial^2 L(\mathcal{B}_{1,p-1}, \mathcal{B}_{1,p})}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p} \partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}}$$

According to the chain rule and the fact that $\mathcal{B}_{1,p-1}$ does not contain $\mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}$, we have

$$\begin{aligned} & \frac{\partial^2 L(\mathcal{B}_{1,p-1}, \mathcal{B}_{1,p})}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p} \partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} \\ &= L_{12}(\mathcal{B}_{1,p-1}, \mathcal{B}_{1,p}) \left(\frac{\partial \mathcal{B}_{1,p-1}}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p}} \right)^\top \frac{\partial \mathcal{B}_{1,p}}{\partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} \\ & \quad + L_{22}(\mathcal{B}_{1,p-1}, \mathcal{B}_{1,p}) \left(\frac{\partial \mathcal{B}_{1,p}}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p}} \right)^\top \frac{\partial \mathcal{B}_{1,p}}{\partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} \\ & \quad + L_2(\mathcal{B}_{1,p-1}, \mathcal{B}_{1,p}) \frac{\partial^2 \mathcal{B}_{1,p}}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p} \partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} \end{aligned} \tag{S1}$$

where L_2 means the partial derivative of L with regard to its first value, while L_{12} and L_{22} are the second-order partial derivatives of L . Note that the last line of (S1) equals 0 because that each pairwise interaction in the RaFM will not contain vectors from different FMs. Therefore,

$$\frac{\partial^2 L(\mathcal{B}_{1,p-1}, \mathcal{B}_{1,p})}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p} \partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} = \mathbf{H} \frac{\partial \mathcal{B}_{1,p}}{\mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}}$$

where \mathbf{H} is a $(|\mathcal{F}_p| D_{p-1}) \times 1$ matrix:

$$\mathbf{H} = L_{12}(\mathcal{B}_{1,p-1}, \mathcal{B}_{1,p}) \left(\frac{\partial \mathcal{B}_{1,p-1}}{\mathbf{v}^{(p-1)}|_{\mathcal{F}_p}} \right)^\top + L_{22}(\mathcal{B}_{1,p-1}, \mathcal{B}_{1,p}) \left(\frac{\partial \mathcal{B}_{1,p}}{\mathbf{v}^{(p-1)}|_{\mathcal{F}_p}} \right)^\top$$

Then we have

$$\begin{aligned} \widehat{grad} &= L_1(\mathcal{B}_{1,m}, y) \frac{\partial \mathcal{B}_{1,m}}{\partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} \\ & \quad + \frac{1}{N} \sum_{\mathbf{x}'} L_1(\mathcal{B}'_{1,m}, y) \frac{\partial \mathcal{B}'_{1,m}}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p}} \mathbf{G}^{-1} \mathbf{H} \frac{\partial \mathcal{B}_{1,p}}{\mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} \\ &= L_1(\mathcal{B}_{1,m}, y) \frac{\partial \mathcal{B}_{1,m}}{\partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} + \lambda \frac{\partial \mathcal{B}_{1,p}}{\mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} \end{aligned}$$

where

$$\lambda = \frac{1}{N} \sum_{\mathbf{x}'} L_1(\mathcal{B}'_{1,m}, y) \frac{\partial \mathcal{B}'_{1,m}}{\partial \mathbf{v}^{(p-1)}|_{\mathcal{F}_p}} \mathbf{G}^{-1} \mathbf{H}$$

Note that λ is the multiplication of a $1 \times |\mathcal{F}_p| D_{p-1}$ matrix, a $|\mathcal{F}_p| D_{p-1} \times |\mathcal{F}_p| D_{p-1}$ matrix, and a $|\mathcal{F}_p| D_{p-1} \times 1$ matrix, and thus is a scalar. Moreover, the derivative of $\mathcal{B}_{1,m}$ and $\mathcal{B}_{1,p}$ with respect to $\mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}$ is the same. Therefore, the direction of \widehat{grad} is parallel to that of $L_1(\mathcal{B}_{1,m}, y) \left(\frac{\partial \mathcal{B}_{1,m}}{\partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}} \right)$, i.e. $\partial L(\mathcal{B}_{1,m}, y) / \partial \mathbf{v}^{(p)}|_{\mathcal{F}_p - \mathcal{F}_{p+1}}$. \square

2. Performance Bound of the Learning Algorithm

Theorem S1. Assume there exist two nonnegative functions $d(\cdot)$ and $\Delta(\cdot, \cdot)$ such that d is monotonically increasing, and for all $f_1(\cdot), f_2(\cdot)$ we have

$$\begin{aligned} d\left(\frac{1}{N}\sum_{\mathbf{x}}L(f_1(\mathbf{x}), y)\right) &\leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(f_2(\mathbf{x}), y)\right) \\ &\quad + d\left(\frac{1}{N}\sum_{\mathbf{x}}\Delta(f_1(\mathbf{x}), f_2(\mathbf{x}))\right) \end{aligned} \quad (\text{S2})$$

then regarding the training error of the k -th FM model, i.e. $\mathcal{B}_{k,k}$, we have

$$\begin{aligned} d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{k,k}^*, y)\right) &\leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{1,m}^*, y)\right) \\ &\quad + \sum_{p=k}^{m-1} d\left(\frac{1}{N}\sum_{\mathbf{x}}\Delta(\mathcal{B}_{1,p}^*, \mathcal{B}_{1,p+1}^*)\right) \end{aligned} \quad (\text{S3})$$

where $\mathcal{B}_{1,m}^*$ is the optimal $\mathcal{B}_{1,m}$, and $\mathcal{B}_{1,p}^*$ is defined in the same way.

Remark: In practice, Δ represents the error of expressing f_2 by f_1 . Eq. (S2) is an extension to the triangle inequality, and can be applied to both regression tasks and classification tasks. In regression tasks, L is the square loss, then we can set d as the square root function and $\Delta = L$. In classification tasks, L is the logarithm loss, then we can let d be an identity function, and define Δ as

$$\Delta(f_1(\mathbf{x}), f_2(\mathbf{x})) = C_{\theta, \delta} \mathcal{D}_{KL}[f_1(\mathbf{x}) \| f_2(\mathbf{x})] + \log \delta \quad (\text{S4})$$

where $\delta > 1$, $C_{\theta, \delta} = \frac{\log \delta}{\theta \log \delta + (1-\theta) \log \frac{1-\theta}{1-\theta/\delta}}$, and $\theta = \min_{\mathbf{x}} [yf_2(\mathbf{x}) + (1-y)(1-f_2(\mathbf{x}))]$. \mathcal{D}_{KL} is the KL divergence of two binomial variables. The readers can refer to Section 3 in the supplementary material for the proof of Eq. (S2) for logarithm loss.

In order to prove Theorem S1, we first provide the following lemma:

Lemma S2. The following inequalities hold

$$d\left(\frac{1}{N}\sum_{\mathbf{x}}l(\mathcal{B}_{l,k}^*, y)\right) \leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{l-1,k}^*, y)\right) \quad (\text{S5})$$

Proof. Note that we have $\mathcal{B}_{l,k} = \mathcal{B}_{l-1,k}$ provided that

$$\mathbf{v}_i^{(l)} = \begin{bmatrix} \mathbf{0}_{D_l - D_{l-1}} \\ \mathbf{v}_i^{(l-1)} \end{bmatrix}, \forall i \in \mathcal{F}_l$$

And such solution also satisfies the constraint (12) in the main body of the paper. Therefore $\mathcal{B}_{l-1,k}$ is a submodel of $\mathcal{B}_{l,k}$, and thus the optimal training error of $\mathcal{B}_{l,k}$ is smaller than $\mathcal{B}_{l-1,k}$, and (S5) follows. \square

Proof of Theorem S1. According to (S5) we have

$$\begin{aligned} d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{k,k}^*, y)\right) &\leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{k-1,k}^*, y)\right) \leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{k-2,k}^*, y)\right) \\ &\quad \dots \\ &\leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{1,k}^*, y)\right) \end{aligned}$$

Moreover, according to (S2) we have

$$\begin{aligned}
d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{1,k}^*,y)\right) &\leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{1,k+1}^*,y)\right) + d\left(\frac{1}{N}\sum_{\mathbf{x}}\Delta(\mathcal{B}_{1,k}^*,\mathcal{B}_{1,k+1}^*)\right) \\
&\leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{1,k+2}^*,y)\right) + \sum_{p=k}^{k+1} d\left(\frac{1}{N}\sum_{\mathbf{x}}\Delta(\mathcal{B}_{1,p}^*,\mathcal{B}_{1,p+1}^*)\right) \\
&\dots \\
&\leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{1,m}^*,y)\right) + \sum_{p=k}^{m-1} d\left(\frac{1}{N}\sum_{\mathbf{x}}\Delta(\mathcal{B}_{1,p}^*,\mathcal{B}_{1,p+1}^*)\right)
\end{aligned}$$

Therefore we have

$$d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{k,k}^*,y)\right) \leq d\left(\frac{1}{N}\sum_{\mathbf{x}}L(\mathcal{B}_{1,m}^*,y)\right) + \sum_{p=k}^{m-1} d\left(\frac{1}{N}\sum_{\mathbf{x}}\Delta(\mathcal{B}_{1,p}^*,\mathcal{B}_{1,p+1}^*)\right)$$

□

3. Quasi-Triangle Inequality for Logarithmic Loss

The following proposition is an extension of the triangle inequality for log loss.

Proposition S3. Suppose $y \in \{0, 1\}$, $0 < \hat{y}_1, \hat{y}_2 < 1$, and define the log loss function $L(\hat{y}_i, y)$ and the KL divergence $\mathcal{D}_{KL}(\hat{y}_1 \parallel \hat{y}_2)$ as

$$\begin{aligned}
L(\hat{y}_i, y) &= -y \log \hat{y}_i - (1-y) \log(1-\hat{y}_i) \\
\mathcal{D}_{KL}(\hat{y}_1 \parallel \hat{y}_2) &= \hat{y}_1 \log \frac{\hat{y}_1}{\hat{y}_2} + (1-\hat{y}_1) \log \frac{1-\hat{y}_1}{1-\hat{y}_2}
\end{aligned}$$

then $\forall \delta > 1, 0 < \theta \leq y\hat{y}_1 + (1-y)(1-\hat{y}_1)$, we have

$$L(\hat{y}_2, y) \leq L(\hat{y}_1, y) + C_{\theta, \delta} \mathcal{D}_{KL}(\hat{y}_1 \parallel \hat{y}_2) + \log \delta \quad (\text{S6})$$

where

$$C_{\theta, \delta} = \frac{\log \delta}{\theta \log \delta + (1-\theta) \log \frac{1-\theta}{1-\theta/\delta}} \quad (\text{S7})$$

Before proving Proposition S3, we first provide some lemmas.

Lemma S4. $\forall \theta > 0, \delta > 1$, we have $C_{\theta, \delta} > 0$, and $C_{\theta, \delta}$ monotonically decreases when θ increases.

Proof. Consider $g(\theta) = 1/C_{\theta, \delta} = \theta + \left[(1-\theta) \log \frac{1-\theta}{1-\theta/\delta} / \log \delta\right]$, then we have

$$\begin{aligned}
g'(\theta) \log \delta &= \log \delta - \log \frac{1-\theta}{1-\theta/\delta} - 1 + \frac{1}{\delta} \frac{1-\theta}{1-\theta/\delta} \\
&= -\log \left(\frac{1-\theta}{\delta} \frac{1-\theta}{1-\theta/\delta}\right) + \left(\frac{1}{\delta} \frac{1-\theta}{1-\theta/\delta} - 1\right) > 0
\end{aligned}$$

here we use the fact that $\log x < (x-1)$ unless $x=1$. Due to that $\log \delta > 0$, we have $g'(\theta) \geq 0$, thus $g(\theta)$ is an increasing function. Moreover we have

$$g(\theta) > g(0) = 0$$

Therefore, $C_{\theta, \delta} = 1/g(\theta)$ is a decreasing function with respect to θ , and $C_{\theta, \delta} > 0$. □

Lemma S5. For $0 < \hat{y}_1, \hat{y}_2 < 1$, we have

$$\log \frac{\hat{y}_1}{\hat{y}_2} \leq C_{\hat{y}_1, \delta} \mathcal{D}_{KL}(\hat{y}_1 \| \hat{y}_2) + \log \delta \quad (\text{S8})$$

Proof. When $\hat{y}_1 < \delta \hat{y}_2$, we have $\log(\hat{y}_1/\hat{y}_2) \leq \log \delta$, thus (S8) holds due to the nonnegativity of the KL divergence. Now we discuss the case when $\hat{y}_1 \geq \delta \hat{y}_2$. Consider the ratio between $\mathcal{D}_{KL}(\hat{y}_1 \| \hat{y}_2)$ and $\log(\hat{y}_1/\hat{y}_2)$:

$$\frac{\mathcal{D}_{KL}(\hat{y}_1 \| \hat{y}_2)}{\log(\hat{y}_1/\hat{y}_2)} = \hat{y}_1 + (1 - \hat{y}_1) \frac{\log(1 - \hat{y}_1) - \log(1 - \hat{y}_2)}{\log \hat{y}_1 - \log \hat{y}_2} = \hat{y}_1 + (1 - \hat{y}_2) h(\hat{y}_1, \hat{y}_2)$$

where $h(\hat{y}_1, \hat{y}_2) = \log \frac{1 - \hat{y}_1}{1 - \hat{y}_2} / \log \frac{\hat{y}_1}{\hat{y}_2}$. It is easy to show that

$$\frac{\partial h}{\partial \hat{y}_2} = - \frac{\mathcal{D}_{KL}(\hat{y}_2 \| \hat{y}_1)}{(\log \hat{y}_1 - \log \hat{y}_2)^2 \hat{y}_2 (1 - \hat{y}_2)} \leq 0$$

Therefore according to $\hat{y}_1 \geq \delta \hat{y}_2$, we have $h(\hat{y}_1, \hat{y}_2) \geq h(\hat{y}_1, \hat{y}_1/\delta)$, and

$$\frac{\mathcal{D}_{KL}(\hat{y}_1 \| \hat{y}_2)}{\log(\hat{y}_1/\hat{y}_2)} \geq \hat{y}_1 + (1 - \hat{y}_1) \frac{\log(1 - \hat{y}_1) - \log(1 - \hat{y}_1/\delta)}{\log \hat{y}_1 - \log(\hat{y}_1/\delta)} = \frac{1}{C_{\hat{y}_1, \delta}}$$

Therefore (S8) also holds when $\hat{y}_1 \geq \delta \hat{y}_2$. □

Proof of Proposition S3. We first prove the case when $y = 1$. In this case, we have $\theta \leq \hat{y}_1$, and

$$\begin{aligned} L(\hat{y}_2, y) - L(\hat{y}_1, y) &= \log \frac{\hat{y}_1}{\hat{y}_2} \\ &\leq C_{\hat{y}_1, \delta} \mathcal{D}_{KL}(\hat{y}_1 \| \hat{y}_2) + \log \delta \\ &\leq C_{\theta, \delta} \mathcal{D}_{KL}(\hat{y}_1 \| \hat{y}_2) + \log \delta \end{aligned}$$

where the first and second inequalities are according to Lemmas S5 and S4. For $y = 0$, we can let $y' = 1 - y$, $\hat{y}'_1 = 1 - \hat{y}_1$, and $\hat{y}'_2 = 1 - \hat{y}_2$, and use the same discussion for y' , \hat{y}'_1 and \hat{y}'_2 . □