

Theoretical Analysis

1 THEOREM PROOFS

1.1 Fitting of arbitrary relationships

In this section, we aim to prove the following theorem:

THEOREM 1.1. *If embeddings $\{\theta_j\}_{j=1}^m \in \mathbb{R}^{m \times d}$ are regularized such that $\|\theta_j\|_2 \leq 1, \forall j \in \{1, \dots, m\}$, the rotation-based RFM can model the inner-product based interaction pattern $\Delta_R = \mathbf{e}_1^{\alpha_{j,1}} \odot \mathbf{e}_2^{\alpha_{j,2}} \odot \dots \odot \mathbf{e}_m^{\alpha_{j,m}}$, with the mean prediction deviation $\mathbb{E}(\bar{R}) = \mathbb{E}(|\bar{\Delta}_{RFM} - \bar{\Delta}_R|) < O(s/\sqrt{d})$. Here $s = \sum_{k=1}^m \alpha_{j,k}$ is the interaction order, $\mathbf{e}_j \in \mathbb{R}^d, j \in \{1, \dots, m\}$, $\alpha_{j,k} = f(\mathbf{e}_j, \mathbf{e}_k)$, and $f: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \{0, 1\}$ is any given feature dependency function.*

For the j -th feature field, each feature is represented as a one-hot vector $\mathbf{x}_j \in \{0, 1\}^{n_j}$, where n_j is the feature number in the j -th field, and $N = \sum_{j=1}^m n_j$ is the total number of features. Afterwards, the embedding look-up operation $E(\cdot)$ is employed to map the one-hot vector \mathbf{x}_j to a low-dimensional embedding \mathbf{e}_j , i.e., $\mathbf{e}_j = E(\mathbf{x}_j)$. Formally, the one-hot encoded vector \mathbf{x}_j of the l -th feature in the j -th field is defined as $\mathbf{x}_j[k] = \begin{cases} 1, & k = l \\ 0, & k \neq l \end{cases}$, and we define the index of the l -th feature in the j -th field as $\text{ID}(\mathbf{x}_j) = \sum_{k=1}^{j-1} n_k + l$, its inverse function as $\text{OneHot}(\sum_{k=1}^{j-1} n_k + l) = \mathbf{x}_j$, and the function $\mathcal{G}(s) = \arg\min_j \sum_{k=1}^j n_k \geq s$ returns the field index of the global index s . We place all features along an axis, and the truth table \mathcal{T} of the feature dependency function f is denoted as a matrix $\mathcal{T} \in \{0, 1\}^{N \times N}$, where:

$$\mathcal{T}(s, t) = \begin{cases} f\left(E\left(\text{OneHot}(s)\right), E\left(\text{OneHot}(t)\right)\right), & \mathcal{G}(s) \neq \mathcal{G}(t) \\ 0, & \mathcal{G}(s) = \mathcal{G}(t) \end{cases}$$

Given the input feature embeddings $\{\theta_j\}_{j=1}^m$, and their one-hot representations $\{\mathbf{x}_j\}_{j=1}^m$, we can obtain the relation vector \mathbf{r}_j of the feature \mathbf{x}_j :

$$\mathbf{r}_j[k] = \begin{cases} \mathcal{T}(k, \text{ID}(\mathbf{x}_j)), & \mathcal{G}(k) \neq j \\ \frac{1}{2} + \frac{1}{2} \cdot \mathbf{x}_j[k - \sum_{l=1}^{j-1} n_l], & \mathcal{G}(k) = j \end{cases}$$

where $k \in \{1, \dots, N\}$. The vector \mathbf{r}_j measures the relation between the feature \mathbf{e}_j and the features from other fields. Meanwhile, the feature dimensions of the same field are naturally masked with $\frac{1}{2}$, except for itself, which has a value of 1. We use the vector \mathbf{r}_j as the auxiliary dimensions for the input features, $\tilde{\theta}_j = [\theta_j, \epsilon \cdot \mathbf{r}_j]$, where ϵ is a sufficiently small number. We construct the matrix \mathbf{M} as follows:

$$\mathbf{M} = \begin{bmatrix} \mathbb{O}_{N \times d} & \frac{\pi}{\epsilon} \cdot \mathbf{I}_{N \times N} \end{bmatrix} \in \mathbb{R}^{N \times (d+N)}.$$

Here \mathbb{O} is an all-zero matrix. We have $\mathbf{M}\tilde{\theta}_j = \pi \cdot \mathbf{r}_j$. Here we set all query matrices and key matrices as $\mathbf{W}_j^Q = \mathbf{W}_j^K = \mathbf{M}, \forall j = \{1, 2, \dots, m\}$, and set all the value matrices \mathbf{W}_j^V as the identity matrix \mathbf{I} . We construct m attention heads, each measuring the relationship between the features in the j -th field ($j \in [1, m]$) and all the other

features. Formally, the projection matrices $\mathbf{H}_j^Q, \mathbf{H}_j^K, \mathbf{H}_j^V$ are defined as:

$$\mathbf{H}_j^Q = \mathbf{H}_j^K = \begin{bmatrix} \mathbb{O}_{n_j \times n_1} & \dots & \mathbf{I}_{n_j \times n_j} & \dots & \mathbb{O}_{n_j \times n_m} \end{bmatrix}^\top, \\ \mathbf{H}_j^V = \begin{bmatrix} \mathbf{I}_{d \times d} & \mathbb{O}_{d \times n_1} & \dots & \mathbb{O}_{d \times n_j} & \dots & \mathbb{O}_{d \times n_m} \end{bmatrix}^\top.$$

In this way, the values are projected to the original features $\mathbf{V}^{(j)} = \{\theta_k\}_{k=1}^m$. The queries and keys of the j -th head are projected to the following: $\mathbf{Q}^{(j)} = \mathbf{K}^{(j)} = \{\pi \cdot \mathbf{r}_k^{(j)}\}_{k=1}^m$. Assume that the one-hot vector $\mathbf{x}_j = [0, 0, \dots, \underbrace{1}_{l\text{-th element}}, 0, 0, \dots]^\top$, the projected

vector $\mathbf{r}_k^{(j)}$ takes the following form:

$$\mathbf{r}_k^{(j)} = \begin{cases} [\frac{1}{2}, \frac{1}{2}, \dots, \underbrace{1}_{l\text{-th element}}, \frac{1}{2}, \frac{1}{2}, \dots]^\top, & k = j \\ [0, 1, \dots, \underbrace{1}_{l\text{-th element}}, 0, 1, \dots]^\top, & k \neq j \\ \mathcal{T}(\text{ID}(\mathbf{x}_j), \text{ID}(\mathbf{x}_k)) & \end{cases}$$

We set the weight vector $\mathbf{w} = [S, \dots, S]^\top$, and $S > 0$ is a sufficiently large number. Note that $\cos(\pm \frac{\pi}{2}) = 0$. Considering the attention score from j -th query in j -th attention head, we have:

$$\begin{aligned} \alpha_{j,l}^{RFM} &= \text{Sigmoid}\left(\mathbf{w}^\top \cos(\pi \cdot \mathbf{r}_j^{(j)} - \pi \cdot \mathbf{r}_l^{(j)})\right) \\ &= \text{Sigmoid}\left(S \cdot \cos\left(\pi - \pi \cdot \mathcal{T}(\text{ID}(\mathbf{x}_j), \text{ID}(\mathbf{x}_l))\right)\right) \\ &= \mathcal{T}(\text{ID}(\mathbf{x}_j), \text{ID}(\mathbf{x}_l)) = f(\mathbf{e}_j, \mathbf{e}_l). \end{aligned}$$

We have:

$$\hat{\theta}_j = \sum_{l=1}^m \alpha_{j,l}^{RFM} \theta_l^V = \sum_{l=1}^m f(\mathbf{e}_j, \mathbf{e}_l) \theta_l.$$

In this scheme, we only consider the j -th query in the j -th attention head ($j \in [1, m]$), and set the weight \mathbf{u} as the identity matrix \mathbf{I} . Omitting the activation function yields the following expression for the output of RFM:

$$\begin{aligned} \hat{\mathbf{y}} &= \mathbf{u}^\top (\cos(\hat{\theta}_j) + \sin(\hat{\theta}_j)) \\ &= \cos\left(\sum_{l=1}^m f(\mathbf{e}_j, \mathbf{e}_l) \theta_l\right) + \sin\left(\sum_{l=1}^m f(\mathbf{e}_j, \mathbf{e}_l) \theta_l\right) \\ &= \mathcal{H}\left(\cos\left(\sum_{l=1}^m f(\mathbf{e}_j, \mathbf{e}_l) \theta_l\right) + i \sin\left(\sum_{l=1}^m f(\mathbf{e}_j, \mathbf{e}_l) \theta_l\right)\right) \\ &= \mathcal{H}\left(\exp\left(i \sum_{l=1}^m f(\mathbf{e}_j, \mathbf{e}_l) \theta_l\right)\right) \\ &= \mathcal{H}(e^{if(\mathbf{e}_j, \mathbf{e}_1)\theta_1} \odot e^{if(\mathbf{e}_j, \mathbf{e}_2)\theta_2} \odot \dots \odot e^{if(\mathbf{e}_j, \mathbf{e}_m)\theta_m}). \end{aligned}$$

Since the construction of the order is independent of the input features $\{\mathbf{e}_j = \theta_j\}_{j=1}^m$, the theorem 1.1 is proved by combining lemma A.1 in our paper.

1.2 Gradient Analysis

In this section, we analyze and compare the gradient properties of RFM with traditional feature interaction methods. Specifically, we investigate the effects of increasing the interaction order. For this purpose, we primarily increase the order of interaction by increasing the number of feature fields. (*i.e.*, denoted as m). In the subsequent theoretical analysis, we will prove that our method exhibits, at most, linear growth in the gradient with respect to the field number. Conversely, in traditional feature interaction approaches, the gradient exhibits exponential growth with respect to the field number. For ease of analysis, we consider a structure with only one layer of self-attentive rotation mechanism:

$$\begin{aligned} G &= \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \\ &= \sigma(\text{Re}[(\exp(i\mathbf{Q}) \text{diag}(\mathbf{w})) \exp(-i\mathbf{K})^\top]) \cdot \mathbf{V} \\ y &= f(\cos(G)) + f(\sin(G)) \end{aligned}$$

For ease of mathematical illustration, we use the notation \mathbf{X}_j to denote the original input feature embedding \mathbf{e}_j . We first calculate the gradient of our approach, *i.e.*, $\frac{\partial y}{\partial \mathbf{X}}$.

$$\text{Let } \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 & & & \\ & \mathbf{X}_2 & & \\ & & \ddots & \\ & & & \mathbf{X}_m \end{bmatrix} \in \mathbb{R}^{dm \times m},$$

$$[\mathbf{W}_1^Q, \dots, \mathbf{W}_m^Q] = \tilde{\mathbf{B}} \in \mathbb{R}^{d \times dm}$$

$$[\mathbf{W}_1^K, \dots, \mathbf{W}_m^K] = \tilde{\mathbf{C}} \in \mathbb{R}^{d \times dm}.$$

$$[\mathbf{W}_1^V, \dots, \mathbf{W}_m^V] = \tilde{\mathbf{D}} \in \mathbb{R}^{d \times dm}$$

$$\text{So, } \mathbf{Q} = (\tilde{\mathbf{B}}\mathbf{X})^\top \quad \mathbf{K} = (\tilde{\mathbf{C}}\mathbf{X})^\top \quad \mathbf{V} = (\tilde{\mathbf{D}}\mathbf{X})^\top.$$

Remark $\text{Re}[\exp(\mathbf{Q})\text{diag}(\mathbf{w})\exp(-i\mathbf{K})^\top]$ as $\textcircled{1}$.

According to Euler's formula, we have:

$$\begin{aligned} \textcircled{1} &= \text{Re}\{(\cos \mathbf{Q} + i \sin \mathbf{Q})\text{diag}(\mathbf{w})[\cos(-\mathbf{K}^\top) + i \sin(-\mathbf{K}^\top)]\} \\ &= \text{Re}\{(\cos \mathbf{Q} + i \sin \mathbf{Q})\text{diag}(\mathbf{w})[\cos(\mathbf{K}^\top) - i \sin(\mathbf{K}^\top)]\} \\ &= \cos \mathbf{Q} \text{diag}(\mathbf{w}) \cos(\mathbf{K}^\top) + \sin \mathbf{Q} \text{diag}(\mathbf{w}) \sin(\mathbf{K}^\top) \\ &= \cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w}) \cos(\tilde{\mathbf{C}}\mathbf{X}) + \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w}) \sin(\tilde{\mathbf{C}}\mathbf{X}) \end{aligned}$$

$$dG = d[\sigma(\textcircled{1}) \cdot \mathbf{V}] = [d\sigma(\textcircled{1})] \mathbf{V} + \sigma(\textcircled{1}) \cdot d\mathbf{V}$$

$$d\mathbf{V} = d[(\tilde{\mathbf{D}}\mathbf{X})^\top] = d(\mathbf{X}^\top \tilde{\mathbf{D}}^\top) = (d\mathbf{X})^\top \tilde{\mathbf{D}}^\top$$

$$d\sigma(\textcircled{1}) = \sigma'(\textcircled{1}) \odot d\textcircled{1}$$

$$\begin{aligned} d\textcircled{1} &= d[\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w}) \cos(\tilde{\mathbf{C}}\mathbf{X}) + \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w}) \sin(\tilde{\mathbf{C}}\mathbf{X})] \\ &= d\left\{\left[\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w})\right] \cdot \cos(\tilde{\mathbf{C}}\mathbf{X}) + d\left\{\left[\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w})\right] \cdot \sin(\tilde{\mathbf{C}}\mathbf{X})\right\}\right\} \\ &= \left\{d\left[\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w})\right] \cdot \cos(\tilde{\mathbf{C}}\mathbf{X}) + \cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w}) \cdot d\cos(\tilde{\mathbf{C}}\mathbf{X})\right. \\ &\quad \left.+ d\left[\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w})\right] \cdot \sin(\tilde{\mathbf{C}}\mathbf{X}) + \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w}) \cdot d\sin(\tilde{\mathbf{C}}\mathbf{X})\right\} \\ &= \left[d\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w}) \cdot \cos(\tilde{\mathbf{C}}\mathbf{X}) + \cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w}) \cdot [d\cos(\tilde{\mathbf{C}}\mathbf{X})]\right. \\ &\quad \left.+ [d\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)] \cdot \text{diag}(\mathbf{w}) \cdot \sin(\tilde{\mathbf{C}}\mathbf{X}) + \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w}) \cdot [d\sin(\tilde{\mathbf{C}}\mathbf{X})]\right] \\ &= \left[-\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w}) \cdot \cos(\tilde{\mathbf{C}}\mathbf{X}) + \cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w}) \cdot \right. \\ &\quad \left. \cdot [-\sin(\tilde{\mathbf{C}}\mathbf{X}) \odot d(\tilde{\mathbf{C}}\mathbf{X})] + [\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)] \cdot \text{diag}(\mathbf{w}) \cdot \sin(\tilde{\mathbf{C}}\mathbf{X})\right. \\ &\quad \left. + \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w}) \cdot [\cos(\tilde{\mathbf{C}}\mathbf{X}) \odot d(\tilde{\mathbf{C}}\mathbf{X})]\right] \end{aligned}$$

$$\begin{aligned} dy &= \text{tr}\left[\left(\frac{\partial y}{\partial G}\right)^\top dG\right] \\ &= \text{tr}\left[\left(\frac{\partial y}{\partial G}\right)^\top \cdot \left\{\left[\sigma'(\textcircled{1}) \odot d\textcircled{1}\right] \mathbf{V} + \sigma(\textcircled{1}) \cdot [(d\mathbf{X})^\top \tilde{\mathbf{D}}^\top]\right\}\right] \\ &= \text{tr}\left\{\left(\frac{\partial y}{\partial G}\right)^\top \cdot \left[\sigma'(\textcircled{1}) \odot d\textcircled{1}\right] \cdot \mathbf{V}\right\} + \text{tr}\left\{\left(\frac{\partial y}{\partial G}\right)^\top \cdot \sigma(\textcircled{1}) \cdot (d\mathbf{X})^\top \tilde{\mathbf{D}}^\top\right\} \\ \text{tr}\left[\left(\frac{\partial y}{\partial G}\right)^\top \cdot \sigma(\textcircled{1}) \cdot (d\mathbf{X})^\top \tilde{\mathbf{D}}^\top\right] &= \text{tr}\left[\tilde{\mathbf{D}}^\top \left(\frac{\partial y}{\partial G}\right)^\top \cdot \sigma(\textcircled{1}) (d\mathbf{X})^\top\right] \\ \text{For } \mathbf{A}, \text{tr}(\mathbf{A}) &= \text{tr}(\mathbf{A}^\top) \implies = \text{tr}\left\{(d\mathbf{X}) \left[\tilde{\mathbf{D}}^\top \left(\frac{\partial y}{\partial G}\right)^\top \sigma(\textcircled{1})\right]^\top\right\} \end{aligned}$$

$$\text{For } \mathbf{A}, \mathbf{B}, \text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA}) \implies = \text{tr}\left\{\left[\tilde{\mathbf{D}}^\top \left(\frac{\partial y}{\partial G}\right)^\top \sigma(\textcircled{1})\right]^\top d\mathbf{X}\right\}$$

Remark,

$$\text{tr}\left[\left(\frac{\partial y}{\partial G}\right)^\top \cdot \left[\sigma'(\textcircled{1}) \odot d\textcircled{1}\right] \cdot \mathbf{V}\right] = \text{part1} + \text{part2} + \text{part3} + \text{part4}$$

part1

$$\begin{aligned} &= \text{tr}\left[\left(\frac{\partial y}{\partial G}\right)^\top \cdot \left\{\sigma'(\textcircled{1}) \odot \left[-\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right] \text{diag}(\mathbf{w}) \cos(\tilde{\mathbf{C}}\mathbf{X})\right\} \mathbf{V}\right] \\ &= \text{tr}\left[\mathbf{V} \left(\frac{\partial y}{\partial G}\right)^\top \cdot \left\{\sigma'(\textcircled{1}) \odot \left[-\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right] \text{diag}(\mathbf{w}) \cos(\tilde{\mathbf{C}}\mathbf{X})\right\}\right] \\ &= \text{tr}\left[\left(\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right)^\top \cdot \left\{-\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right\} \cdot \text{diag}(\mathbf{w}) \cdot \cos(\tilde{\mathbf{C}}\mathbf{X})\right] \\ &= \text{tr}(\text{diag}(\mathbf{w}) \cdot \cos(\tilde{\mathbf{C}}\mathbf{X}) \cdot \left\{\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right\}^\top \cdot \left[-\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right]) \\ &= \text{tr}\left(\left\{\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right\} [\text{diag}(\mathbf{w}) \cdot \cos(\tilde{\mathbf{C}}\mathbf{X})]^\top \odot \left[-\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right]^\top\right) \end{aligned}$$

$$= -\text{tr}\left(\left\{\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right\} \cdot [\text{diag}(\mathbf{w}) \cdot \cos(\tilde{\mathbf{C}}\mathbf{X})]^\top \odot \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right)^\top \cdot (d\mathbf{X})^\top \tilde{\mathbf{B}}^\top$$

Remark

$$F = \left\{\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1}) \cdot [\text{diag}(\mathbf{w}) \cdot \cos(\tilde{\mathbf{C}}\mathbf{X})]^\top \odot \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right\}^\top$$

So,

$$\begin{aligned} \text{part1} &= -\text{tr}(F(d\mathbf{X})^\top \tilde{\mathbf{B}}^\top) = -\text{tr}(\tilde{\mathbf{B}}^\top F(d\mathbf{X})^\top) = -\text{tr}(d\mathbf{X} \cdot (\tilde{\mathbf{B}}^\top F)^\top) \\ &= -\text{tr}\left((\tilde{\mathbf{B}}^\top F)^\top d\mathbf{X}\right) = \text{tr}\left(-(\tilde{\mathbf{B}}^\top F)^\top d\mathbf{X}\right) \end{aligned}$$

$$\begin{aligned} \text{part2} &= \text{tr}\left(\left(\frac{\partial y}{\partial G}\right)^\top \cdot \left\{\sigma'(\textcircled{1}) \odot \left[\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w}) \cdot [-\sin(\tilde{\mathbf{C}}\mathbf{X}) \odot d(\tilde{\mathbf{C}}\mathbf{X})]\right]\right\} \cdot \mathbf{V}\right) \\ &= \text{tr}\left(\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right)^\top \cdot \left\{\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w}) \cdot [-\sin(\tilde{\mathbf{C}}\mathbf{X}) \odot d(\tilde{\mathbf{C}}\mathbf{X})]\right\} \\ &= \text{tr}\left(\left\{\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right\}^\top \cdot \cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w})\right)^\top \cdot [-\sin(\tilde{\mathbf{C}}\mathbf{X})]^\top \cdot d(\tilde{\mathbf{C}}\mathbf{X}) \end{aligned}$$

Remark

$$N = \left[\left\{\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right\}^\top \cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w})\right]^\top \odot [-\sin(\tilde{\mathbf{C}}\mathbf{X})]^\top$$

So,

$$\text{part2} = \text{tr}(N \tilde{\mathbf{C}} d\mathbf{X}) = \text{tr}\left([\mathbf{N} \tilde{\mathbf{C}}]^\top d\mathbf{X}\right) = \text{tr}\left([\tilde{\mathbf{C}}^\top N^\top]^\top d\mathbf{X}\right)$$

part3

$$\begin{aligned} &= \text{tr}\left(\mathbf{V} \left(\frac{\partial y}{\partial G}\right)^\top \cdot \left\{\sigma'(\textcircled{1}) \odot \left[\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right] \text{diag}(\mathbf{w}) \sin(\tilde{\mathbf{C}}\mathbf{X})\right\}\right) \\ &= \text{tr}\left(\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right)^\top \cdot \left[\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right] \text{diag}(\mathbf{w}) \sin(\tilde{\mathbf{C}}\mathbf{X}) \\ &= \text{tr}\left(\text{diag}(\mathbf{w}) \sin(\tilde{\mathbf{C}}\mathbf{X}) \left[\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right]^\top \cdot \left[\cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \odot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right]^\top\right) \\ &= \text{tr}\left(\left[\text{diag}(\mathbf{w}) \cdot \sin(\tilde{\mathbf{C}}\mathbf{X}) \cdot \left[\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right]^\top\right]^\top \odot \cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)\right)^\top \cdot d(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \end{aligned}$$

Remark

$$R = \left[\text{diag}(\mathbf{w}) \cdot \sin(\tilde{\mathbf{C}}\mathbf{X}) \cdot \left[\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right]^\top\right]^\top \odot \cos(\mathbf{X}^\top \tilde{\mathbf{B}}^\top)^\top$$

$$\text{So part3} = \text{tr}(\tilde{\mathbf{B}}^\top R(d\mathbf{X})^\top) = \text{tr}(d\mathbf{X} (\tilde{\mathbf{B}}^\top R)^\top) = \text{tr}((\tilde{\mathbf{B}}^\top R)^\top d\mathbf{X})$$

$$\begin{aligned} \text{part4} &= \text{tr}\left(\mathbf{V} \left(\frac{\partial y}{\partial G}\right)^\top \cdot \left\{\sigma'(\textcircled{1}) \odot \left[\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w}) \cdot [\cos(\tilde{\mathbf{C}}\mathbf{X}) \odot d(\tilde{\mathbf{C}}\mathbf{X})]\right]\right\}\right) \\ &= \text{tr}\left(\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right)^\top \cdot \left\{\sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \text{diag}(\mathbf{w}) \cdot [\cos(\tilde{\mathbf{C}}\mathbf{X}) \odot d(\tilde{\mathbf{C}}\mathbf{X})]\right\} \\ &= \text{tr}\left(\left\{\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right\}^\top \cdot \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w})\right]^\top \odot \cos(\tilde{\mathbf{C}}\mathbf{X})\right)^\top \cdot \tilde{\mathbf{C}} d\mathbf{X} \end{aligned}$$

Remark

$$J = \left\{\left[\left(\frac{\partial y}{\partial G}\right)^\top \mathbf{V}^\top\right] \odot \sigma'(\textcircled{1})\right\}^\top \cdot \sin(\mathbf{X}^\top \tilde{\mathbf{B}}^\top) \cdot \text{diag}(\mathbf{w})\right]^\top \odot \cos(\tilde{\mathbf{C}}\mathbf{X})\right\}^\top$$

$$\text{So part4} = \text{tr}\left((\tilde{\mathbf{C}}^\top J^\top)^\top d\mathbf{X}\right)$$

$$\begin{aligned}
\left(\frac{\partial y}{\partial X}\right)^\top &= \left[\tilde{D}^\top \left(\frac{\partial y}{\partial G}\right)^\top \sigma(\mathbb{1})\right] + (-\tilde{B}^\top F) + \tilde{C}^\top N^\top + \tilde{B}^\top R + \tilde{C}^\top J^\top \\
F &= \left\{ \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \cdot [\text{diag}(\mathbf{w}) \cdot \cos(\tilde{C}X)]^\top \right] \odot \sin(X^\top \tilde{B}^\top) \right\}^\top \\
&= \left\{ \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \cdot [\text{diag}(\mathbf{w}) \cdot \cos(\tilde{C}X)]^\top \right\}^\top \odot \sin(\tilde{B}X) \\
&= \left\{ \text{diag}(\mathbf{w}) \cdot \cos(\tilde{C}X) \cdot \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \right\}^\top \odot \sin(Q^\top) \\
&= \left\{ \text{diag}(\mathbf{w}) \cdot \cos(K^\top) \cdot \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \right\}^\top \odot \sin(Q^\top) \\
&= \{\text{diag}(\mathbf{w}) \text{Re}(\exp(iK^\top))\} \\
&\quad \cdot \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\text{Re}[\exp(-iK)\text{diag}(\mathbf{w})\exp(iQ^\top)]) \right] \odot \text{Im}[\exp(iQ^\top)] \\
N^\top &= \left\{ \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \cdot \cos(X^\top \tilde{B}^\top) \cdot \text{diag}(\mathbf{w}) \right\}^\top \odot [-\sin(\tilde{C}X)] \\
&= \left\{ \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}^\top) \cdot \cos(X^\top \tilde{B}^\top) \cdot \text{diag}(\mathbf{w}) \right\}^\top \odot [-\sin(K^\top)] \\
&= \left\{ \text{diag}(\mathbf{w}) \cdot \cos(\tilde{B}X) \cdot \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \right\}^\top \odot [-\sin(K^\top)] \\
&= \{\text{diag}(\mathbf{w}) \text{Re}[\exp(iQ^\top)]\} \\
&\quad \cdot \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\text{Re}[\exp(iQ)\text{diag}(\mathbf{w})\exp(-iK)^\top]) \right] \\
&\quad \odot [-\text{Im}[\exp(iK^\top)]] \\
R &= \left\{ \text{diag}(\mathbf{w}) \sin(\tilde{C}X) \cdot \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \right]^\top \odot \cos(X^\top \tilde{B}^\top) \right\}^\top \\
&= \left\{ \text{diag}(\mathbf{w}) \sin(K^\top) \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}^\top) \right]^\top \odot \cos(X^\top \tilde{B}^\top) \right\}^\top \\
&= \left[\text{diag}(\mathbf{w}) \sin(K^\top) \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}^\top) \right] \right] \odot \cos(\tilde{B}X) \\
&= \{\text{diag}(\mathbf{w}) \text{Im}[\exp(iK^\top)]\} \\
&\quad \cdot \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\text{Re}[\exp(-iK)\text{diag}(\mathbf{w})\exp(iQ^\top)]) \right] \odot \text{Re}[\exp(iQ^\top)] \\
J^\top &= \left\{ \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \cdot \sin(X^\top \tilde{B}^\top) \cdot \text{diag}(\mathbf{w}) \right\}^\top \odot \cos(\tilde{C}X) \\
&= \left\{ \text{diag}(\mathbf{w}) \sin(\tilde{B}X) \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\mathbb{1}) \right] \right\}^\top \odot \cos(K^\top) \\
&= \{\text{diag}(\mathbf{w}) \text{Im}[\exp(iQ^\top)]\} \\
&\quad \cdot \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\text{Re}[\exp(iQ)\text{diag}(\mathbf{w})\exp(-iK)^\top]) \right] \odot \text{Re}[\exp(iK^\top)] \\
\frac{\partial y}{\partial X} &= \tilde{D}^\top \left(\frac{\partial y}{\partial G} \right)^\top \sigma(\text{Re}[\exp(iQ)\text{diag}(\mathbf{w})\exp(-iK)^\top]) \\
&\quad - \tilde{B}^\top (\{\text{diag}(\mathbf{w}) \cdot \text{Re}[\exp(iK^\top)]\} \\
&\quad \cdot \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\text{Re}[\exp(-iK)\text{diag}(\mathbf{w})\exp(iQ^\top)]) \right] \odot \text{Im}[\exp(iQ^\top)]) \\
&\quad - \tilde{C}^\top (\{\text{diag}(\mathbf{w}) \cdot \text{Re}[\exp(iQ^\top)]\} \\
&\quad \cdot \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\text{Re}[\exp(iQ)\text{diag}(\mathbf{w})\exp(-iK)^\top]) \right] \odot \text{Im}[\exp(iK^\top)]) \\
&\quad + \tilde{B}^\top (\{\text{diag}(\mathbf{w}) \cdot \text{Im}[\exp(iK^\top)]\} \\
&\quad \cdot \left[\left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\text{Re}[\exp(-iK)\text{diag}(\mathbf{w})\exp(iQ^\top)]) \right] \odot \text{Re}[\exp(iQ^\top)]) \\
&\quad + \tilde{C}^\top (\{\text{diag}(\mathbf{w}) \cdot \text{Im}[\exp(iQ^\top)]\} \\
&\quad \cdot \left[\left(\frac{\partial y}{\partial G} \right)^\top V^\top \right] \odot \sigma'(\text{Re}[\exp(iQ)\text{diag}(\mathbf{w})\exp(-iK)^\top]) \right] \odot \text{Re}[\exp(iK^\top)])
\end{aligned}$$

Correspondingly, we remark the equation as: $\frac{\partial y}{\partial X} = \text{PART I} - \text{PART II} - \text{PART III} + \text{PART IV} + \text{PART V}$.

For a matrix A , we define the infinite norms of matrices as: $\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$.

Since we can learn \tilde{B} , \tilde{C} , \tilde{D} , $\text{diag}(\mathbf{w})$, X and $\frac{\partial y}{\partial G}$, we assume:
 $\|\tilde{B}^\top\|_\infty < \alpha$, $\|\tilde{C}^\top\|_\infty < \beta$, $\|\text{diag}(\mathbf{w})\|_\infty < \zeta$,
 $\|\tilde{D}^\top\|_\infty < \gamma_1$, $\|\tilde{B}\|_\infty < \gamma_2$, let $\gamma = \max_i \{\gamma_1, \gamma_2\}$, then $\|\tilde{D}^\top\|_\infty < \gamma$, $\|\tilde{D}\|_\infty < \gamma$,

Similarly, $\|(\frac{\partial y}{\partial G})^\top\|_\infty < \theta$, $\|(\frac{\partial y}{\partial G})\|_\infty < \theta$.

Note that each row of X can have at most one non-zero element due to the inherent sparsity of X . We suppose the absolute value of every element in X

is smaller than η , so we have $\|X\|_\infty < \eta$ and $\|X^\top\|_\infty < \eta$. According to the compatibility of this norm:

$$\begin{aligned}
\text{PART I} &= \left\| \tilde{D}^\top \left(\frac{\partial y}{\partial G} \right)^\top \sigma(\text{Re}[\exp(iQ)\text{diag}(\mathbf{w})\exp(-iK)^\top]) \right\|_\infty \\
&\leq \|\tilde{D}^\top\|_\infty \cdot \left\| \left(\frac{\partial y}{\partial G} \right)^\top \right\|_\infty \cdot \|\sigma(\text{Re}[\exp(iQ)\text{diag}(\mathbf{w})\exp(-iK)^\top])\|_\infty \\
&\leq \gamma \cdot \theta \cdot m
\end{aligned}$$

$$\begin{aligned}
\text{PART II} &\leq \|\tilde{B}^\top\|_\infty \|\text{diag}(\mathbf{w})\|_\infty \|\text{Re}[\exp(iK^\top)]\|_\infty \|V\|_\infty \left\| \left(\frac{\partial y}{\partial G} \right)^\top \right\|_\infty \frac{1}{4} \\
&\leq \alpha \cdot \zeta \cdot m \cdot d \cdot \eta \cdot \gamma \cdot \theta \cdot \frac{1}{4} = \frac{\alpha \gamma \theta \zeta d m}{4}
\end{aligned}$$

$$\text{Similarly, PART III} \leq \frac{\beta \zeta \gamma \eta \theta d m}{4}, \text{PART IV} \leq \frac{\alpha \zeta \eta \gamma \theta d m}{4}, \text{PART V} \leq \frac{\beta \zeta \theta \gamma \eta m}{4}.$$

$$\left\| \frac{\partial y}{\partial X} \right\|_\infty \leq \gamma \theta m + \frac{\alpha \zeta \eta \gamma \theta d m}{4} + \frac{\beta \zeta \gamma \eta \theta d m}{4} + \frac{\alpha \zeta \eta \gamma \theta d m}{4} + \frac{\beta \zeta \gamma \eta \theta m}{4}$$

$$= \gamma \theta m + \frac{\alpha \zeta \eta \gamma \theta d m}{2} + \frac{\beta \zeta \gamma \eta \theta m}{2}$$

Remark $\gamma \theta + \frac{\alpha \zeta \eta \gamma \theta d}{2} + \frac{\beta \zeta \gamma \eta \theta}{2} = C_1$, thus, $\left\| \frac{\partial y}{\partial X} \right\|_\infty \leq C_1 m$.

Note that C_1 is independent of the field number m . It can be seen that under certain regularity conditions, the gradient terms grow at most linearly with m .

For the traditional feature interaction algorithms, their gradients can be formulated as:

$$\begin{aligned}
\mathbf{g} &= X_1^{\alpha_1} \odot X_2^{\alpha_2} \odot \dots \odot X_m^{\alpha_m} \\
\frac{\partial \mathbf{g}}{\partial X_i} &= X_1^{\alpha_1} \odot X_2^{\alpha_2} \odot \dots \odot \alpha_i X_i^{\alpha_i-1} \odot X_{i+1}^{\alpha_{i+1}} \odot \dots \odot X_m^{\alpha_m}, \quad \mathbf{y} = f(\mathbf{g}).
\end{aligned}$$

We suppose there exists j , for all i we suppose $|X_{ij}| \geq M - \epsilon$, here M is strictly bigger than zero, thus:

$$\left| \frac{\partial \mathbf{g}}{\partial X_{ij}} \right| \geq \alpha_i \cdot (M - \epsilon)^{\sum_{i=1}^m \alpha_i - 1}$$

$$\left| \frac{\partial \mathbf{y}}{\partial X_{ij}} \right| \geq \alpha_i \cdot (M - \epsilon)^{\sum_{i=1}^m \alpha_i - 1} \left\| \frac{\partial \mathbf{y}}{\partial \mathbf{g}} \right\|_\infty$$

Let $t = \min_i \{\alpha_i\}$, thus we have:

$$\left| \frac{\partial \mathbf{g}}{\partial X_{ij}} \right| \geq \alpha_i \cdot (M - \epsilon)^{mt-1} \left\| \frac{\partial \mathbf{y}}{\partial \mathbf{g}} \right\|_\infty$$

We can clearly see that the gradient terms of traditional feature interaction algorithms exponentially grow with the field number m .