

CAN TRANSFORMERS PERFORM PCA ?

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformers demonstrate significant advantage as the building block of Large Language Models. Recent efforts are devoted to understanding the learning capacities of transformers at a fundamental level. This work attempts to understand the intrinsic capacity of transformers in performing dimension reduction from complex data. Theoretically, our results rigorously show that transformers can perform Principle Component Analysis (PCA) similar to the Power Method, given a supervised pre-training phase. Moreover, we show the generalization error of transformers decays by $n^{-1/5}$ in L_2 . Empirically, our extensive experiments on the simulated and real world high dimensional datasets justify that a pre-trained transformer can successfully perform PCA by simultaneously estimating the first k eigenvectors and eigenvalues. These findings demonstrate that transformers can efficiently extract low dimensional patterns from high dimensional data, shedding light on the potential benefits of using pre-trained LLM to perform inference on high dimensional data.

1 INTRODUCTION

Large Language Models (LLMs) have demonstrated significant success in learning and performing inference on real world high dimensional datasets. Most modern LLMs use transformers (Vaswani, 2017) as their backbones. Transformers are a class of models that demonstrate significant advantages over the previous neural network models using recurrent architectures, achieving many state-of-the-art performance in learning tasks including natural language processing (Wolf et al., 2020) and computer vision (Khan et al., 2022). However, the underlying mechanism for the success of transformers remains largely a mystery to researchers. One critical and most fundamental questions is why this model works well when adapting to large volume of data used in the training of LLMs.

It is well known in the machine learning and statistics community that high dimensional inference are subject to *the Curse of Dimensionality*. Hence, practioners use various of methods to perform dimension reduction in the covariate space before they perform subsequent inferential tasks. However, performing inference using pre-trained LLMs does not require researchers to perform dimension reduction manually (Ma et al., 2023). Instead, LLMs are able to extract the essential information from the input texts. This work attempts to understand how transformers are able to extract low dimensional patterns from large volume of data via answering the question raised in the title. We perform our study through theoretically analyzing the transformer model to show its approximation and generalization capacities in performing the PCA task.

PCA is one of the most important dimension reduction methods used in practice. In particular, it has fundamental utilities in machine learning (Bishop & Nasrabadi, 2006; Hastie et al., 2009), high dimensional statistics (Fan et al., 2020; Wainwright, 2019), and econometrics (Bai et al., 2008). For the most part, PCA is a sub-procedure for algorithms solving more complicated problems. Therefore, understanding how well transformers perform PCA is also of great importance to find potential new utilities where LLMs succeed in.

Although PCA is a standard unsupervised learning task that does not require a particular machine learning model, we show that pretraining a transformer can equip it with the capacity of performing PCA on unseen instances. The pretraining step studied in this work follows a supervised learning paradigm where the input is the covariate matrix and the label is designed to be its top k eigenvectors. Some existing works also make attempt to study the universal approximation power of transformers on classes of functions including (Pérez et al., 2021; Wei et al., 2022; Yun et al., 2019). However, the

054 mapping from matrix to eigenvectors is not easily integrated into their framework. Moreover, their
 055 work does not imply the results obtained here as we provide approximation errors for the principle
 056 eigenvectors and the corresponding generalization error bounds given by the pretraining procedure.
 057

058 **Contributions.** We summarize our major contributions as follows:
 059

- 060 1. We rigorously show that a pre-trained transformer can perform PCA and give approxima-
 061 tion error bound. The proof is constructive where we utilize the logical similarities between
 062 the forward propagation on transformers and the Power Method to bound the approxima-
 063 tion error;
- 064 2. We further provide upper bounds on the generalization error for the empirical risk mini-
 065 mizer in the pre-training task. Coupling with the approximation error and making tradeoff
 066 between the different terms, we show that transformers can generalize with L_2 error rate as
 067 fast as $n^{-1/5}$ with high probability where n is the number of pre-trained samples;
- 068 3. We systematically evaluate the performance of PCA with different parameter value com-
 069 binations. These empirical results demonstrate that transformers can perform very well in
 070 extracting principal eigenvectors and eigenvalues from data, even in regions where theoret-
 071 ical results are hard to obtain, given a proper pretraining procedure.
 072

073 1.1 RELATED WORKS

074 This work is related to a few different branches in the literature.
 075

076 **In Context Learning of Transformers.** Some recent works studied the in-context learning (ICL)
 077 capacities of Transformers (Garg et al., 2022; Bai et al., 2024). In particular, (Bai et al., 2024) con-
 078 sidered the approximation and generalization properties of transformers on the ICL tasks, including
 079 many linear regression and logistic regression setups. The problem of PCA is a standard unsuper-
 080 vised learning problem. Hence, it differs from ICL in that there is no individual label that the model
 081 needs to learn. Akyürek et al. (2022); Von Oswald et al. (2023) considered the approximation of
 082 transformers on gradient descent when performing ICL. In this work, the proof machine utilizes the
 083 Power Method. We also notice that performing gradient descent is difficult to obtain the eigenvec-
 084 tors as no explicit functional form is given. To the best of authors’ knowledge, this is the first work
 085 that provides theoretical guarantees for the transformers’ approximation of the Power method in the
 086 literature. Other related works on the more practical side of ICL can be found in Dong et al. (2022)
 087 and reference therein.
 088

089 **Other Theoretical Works on Transformers.** Many other attempts are made to theoretically un-
 090 derstand transformers. Yun et al. (2019) studied the universal approximation properties of trans-
 091 formers on sequence-to-sequence functions. Pérez et al. (2021); Bhattamishra et al. (2020); Liu
 092 et al. (2022) studied the computational power of transformers. Hron et al. (2020) studied the limit of
 093 infinite width multi/single head attentions. Yao et al. (2021) showed that transformers can process
 094 bounded hierarchical languages and demonstrate better space complexity than the recurrent neural
 095 networks.
 096

097 **Notations** In this work we follow the following notation conventions. The vector valued variable
 098 is given by boldfaced characters. We denote $[n] := \{1, \dots, n\}$ and $[i : j] := \{i, i + 1, \dots, j\}$ for
 099 $i < j$. The universal constants are given by C and is ad hoc. For a vector \mathbf{v} we denote $\|\mathbf{v}\|_2$ as
 100 its L_2 norm. For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ we denote its operator norm as $\|\mathbf{A}\|_2 := \sup_{\mathbf{v} \in \mathbb{S}^{n-1}} \|\mathbf{A}\mathbf{v}\|_2$.
 101 Given two sequences a_n and b_n , we denote $a_n \lesssim b_n$ or $a_n = O(b_n)$ if $\limsup_{n \rightarrow \infty} \frac{a_n}{b_n} < \infty$ and
 102 $a_n = o(b_n)$ if $\limsup_{n \rightarrow \infty} \frac{a_n}{b_n} = 0$.

103 **Organizations** The rest of the paper is organized as follows: Section 2 describes the learning
 104 problems, the idea of our constructive proof, and reviews standard contexts; Section 3 provides
 105 rigorous theoretical results; Section 4 provides extensive experimental details and results; Section
 106 5 discusses the limitations and potential future works. The detailed proofs and additional figures
 107 in experiments are delayed to the appendix. The supplementary materials include the code for the
 experiments.

2 APPROXIMATE PCA BY PRE-TRAINED TRANSFORMERS

This section discusses how we construct a multi-layered transformer model such that forward propagate along the it gives us the left principle eigenvectors of the input matrix. Our discussions is splitted into 3 subsections: In 2.1 we review the mathematical forms of the Transformer model in this work; In 2.2 we review the classical power method algorithm to perform PCA and connects it with the multiphase Transformer design; In 2.3 we demonstrate how we perform supervised pre-training to achieve a model in 2.2.

2.1 THE TRANSFORMERS

We consider the context learning problem on the transformer model. Under this formulation, we have the following definition regarding an attention layer. These definitions are similar to that given by Bai et al. (2024).

Definition 1 (Attention Layer). *A self-attention layer with M heads is denoted as $Attn_{\theta_1}(\cdot)$ with parameters $\theta_1 = \{(\mathbf{V}_m, \mathbf{Q}_m, \mathbf{K}_m)\}_{m \in [M]} \subset \mathbb{R}^{D \times D}$. On input sequence $\mathbf{H} \in \mathbb{R}^{D \times N}$,*

$$Attn_{\theta_1}(\mathbf{H}) = \mathbf{H} + \frac{1}{N} \sum_{m=1}^M (\mathbf{V}_m \mathbf{H}) \sigma((\mathbf{Q}_m \mathbf{H})^\top (\mathbf{K}_m \mathbf{H})),$$

where σ is the *relu* activation function.

Remark 1. *Note that, instead of the concatenated feature given by multi-head attention, we consider simple average on the multi-head output. And the activation function we considered is Relu instead of Softmax that appears in most empirical works. (Shen et al., 2023) empirically verified that Relu transformers are strong alternatives to Softmax transformers. We also omit the layer-wise normalization used to stabilize the training procedure. These adaptation are designed for the technical convinience. In the simulation section we carefully evaluate the effect of these additional features.*

The following two layers defines the classical MLP layers with residual connections.

Definition 2 (MLP Layer). *A MLP layer with hidden dimension D' is denoted as $MLP_{\theta_2}(\cdot)$ with parameter $\theta_2 \in (\mathbf{W}_1, \mathbf{W}_2) \in \mathbb{R}^{D' \times D} \times \mathbb{R}^{D \times D'}$. On any input sequence $\mathbf{H} \in \mathbb{R}^{D \times N}$, we define*

$$MLP_{\theta_2}(\mathbf{H}) := \mathbf{H} + \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{H}).$$

Then we use the above two definitions on the MLP and the Attention layers to define the Transformers.

Definition 3 (Transformer). *We define a transformer $TF_{\theta}(\cdot)$ as a composition of self-attention layers with MLP layers. Consider output dimension to be \tilde{D} , the . In particular, a L -layered Transformer is defined by*

$$TF_{\theta}(\mathbf{H}) := \tilde{\mathbf{W}}_0 \times MLP_{\theta_2^L}(Attn_{\theta_1^L}(\cdots MLP_{\theta_2^2}(Attn_{\theta_1^1}(\mathbf{H}))) \times \tilde{\mathbf{W}}_1,$$

where $\tilde{\mathbf{W}}_0 \in \mathbb{R}^{d_1 \times D}$ and $\tilde{\mathbf{W}}_1 \in \mathbb{R}^{N \times d_2}$.

We use θ to denote all the parameters in the transformer and the super-index ℓ to denote the parameter matrix corresponds to the ℓ -th layer. Under such definition, the parameter θ is given by

$$\theta = \{ \{ \{ \mathbf{Q}_m^\ell, \mathbf{K}_m^\ell, \mathbf{V}_m^\ell \}_{m \in [M]}, \mathbf{W}_1^\ell, \mathbf{W}_2^\ell \}_{\ell \in [L]}, \tilde{\mathbf{W}}_0, \tilde{\mathbf{W}}_1 \}.$$

Remark 2. *The model's Lipschitzness can be strictly governed by the following aspects: (1) The number of layers; (2) The number of heads; (3) The maximum operator norm of the parameters. These results further lead to an upper bound on the generalization error. Collecting the above three aspects, we define the following operator norm of the parameters*

$$\|\theta\|_{op} := \max_{\ell \in [L]} \left\{ \max_{m \in [M^\ell]} \{ \|\mathbf{Q}_m^\ell\|_2, \|\mathbf{K}_m^\ell\|_2 \} + \sum_{m=1}^{M^\ell} \|\mathbf{V}_m^\ell\|_2 + \|\mathbf{W}_1^\ell\|_2 + \|\mathbf{W}_2^\ell\|_2 \right\},$$

where M^ℓ is the number of heads of the ℓ -th attention layer. It is shown in (Bai et al., 2024) that such norm relates to the Lipschitz constant of transformers.

The two additional matrices $\tilde{\mathbf{W}}_0$ and $\tilde{\mathbf{W}}_1$ serve for the dimension adjustment purpose such that the output of $TF_{\theta}(\cdot)$ will be of dimension $\mathbb{R}^{d_1 \times d_2}$.

2.2 THE POWER METHOD

The power method is an efficient iterative algorithm to solve for the principle eigenvectors and eigenvalues of an asymmetric matrix (Golub & Van Loan, 2013). The formal statement is given by algorithm 1.

Algorithm 1: Power Method for the Left Singular Vectors

Data: Matrix $\mathbf{X} \in \mathbb{R}^{d \times N}$, Number of Iterations τ

Symmertize $A = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{d \times d}$;

Let the set of eigenvectors be $\mathcal{V} = \{\}$. Initialize $A_1 \leftarrow A$;

for $\ell \leftarrow 1$ **to** k **do**

 Sample a random vector $\mathbf{v}_{0,\ell} \in \mathbb{S}^{N-1}$. Initialize $\mathbf{v}_\ell^{(0)} \leftarrow \mathbf{v}_{0,\ell}$;

for $t \leftarrow 1$ **to** τ **do**

 Apply the procedure to obtain the principle eigenvector $\mathbf{v}_\ell^{(t)} = \frac{A_\ell \mathbf{v}_\ell^{(t-1)}}{\|A_\ell \mathbf{v}_\ell^{(t-1)}\|_2}$;

 Let $\mathcal{V} \leftarrow \mathcal{V} \cup \{\mathbf{v}_\ell^{(\tau)}\}$;

 Compute the eigenvalue estimate $\hat{\lambda}_\ell \leftarrow \|A_\ell \mathbf{v}_\ell^{(\tau)}\|_2$;

 Update the matrix by $A_{\ell+1} = A_\ell - \hat{\lambda}_\ell \mathbf{v}_\ell^{(\tau)} \mathbf{v}_\ell^{(\tau)\top}$;

return \mathcal{V} ;

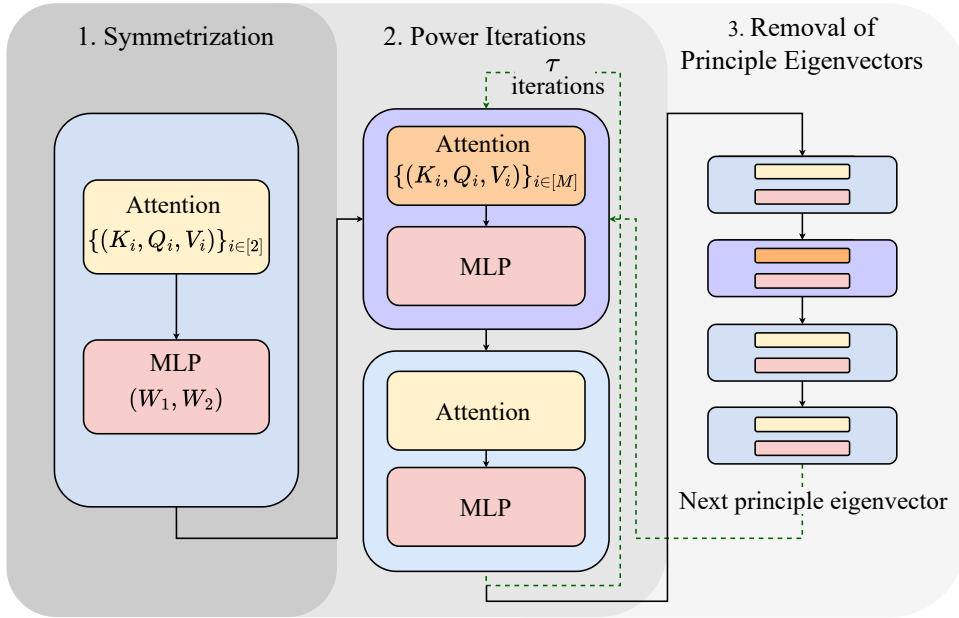


Figure 1: **The Approximation by Transformers.** The above diagram illustrates the design of our transformer model. There are three important sub-networks in the design: (1) *The Symmetrization* sub-network symmetrizes \mathbf{X} and stamp $\mathbf{X}\mathbf{X}^\top$ in the output, which corresponds to the first step in the Power Method. (2) *The Power Iterations* sub-network performs in total of τ iterations for each of the principle eigenvectors, corresponds to the iterative update step in the Power Method. (3) *The Removal of Principle Eigenvectors* subnetwork performs the estimate of $\hat{\lambda}_\ell$ and the update of matrix A_ℓ in the Power Method. Finally, we apply \tilde{W}_0 and \tilde{W}_1 to adjust the dimension of the output. The different colors in the diagram corresponds to the different type of layers: (1) *Yellow Blocks* denote the Attention layer with 2 heads. (2) *Orange Blocks* denote the multihead transformers with larger $M \gg 2$. (3) *Pink Blocks* denote the MLP layer.

The algorithm itself first generate the symmetrized covariate matrix $A = XX^\top$. Then, for each eigenvector that we hope to recover, the power method generates a random vector uniformly distributed on the unit sphere. Then we iterately take the matrix product between the symmetric matrix A and $v_{0,\ell}$ followed by normalizing it. Note that such iterative matrix product can finally converge to the top-1 principle eigenvector as the iterations go further.

The iterative structure is akin to the forward propagation of the transformer achitecture. Given by such similarities, we show that there exists a parameter setup for Transformers that the forward propogation performs principle component analysis. The challenge is in constructing parameters that approximate the complete procedure of the power method.

In figure 1, we provide an *approximate* algorithm for the power method through the lens of a forward propagation on the transformer. This algorithm dissects the approximation of power method into the propagation along multiple sub-networks, each phase corresponds to a single step in the power method. Combining them, we show in section 3 that Transformers achieve good approximation guarantees.

2.3 PRETRAINING VIA SUPERVISED LEARNING

The standard PCA problem is unsupervised where no labels given. However, the Transformers are usually used in the supervised learning setup. To make full use of Transformers in the PCA task, we need to perform supervised pre-training. In our theoretical analysis, we construct the input of the Transformer as a *context-augmented matrix* given by the following

$$\mathbf{H} = \begin{bmatrix} \mathbf{X} \\ \mathbf{P} \end{bmatrix} \in \mathbb{R}^{D \times N}, \quad \mathbf{P} = \begin{bmatrix} \tilde{\mathbf{p}}_{1,1}, \dots, \tilde{\mathbf{p}}_{1,N} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \in \mathbb{R}^{(D-d) \times N},$$

where the matrix \mathbf{P} contains contextual information, which is specified in section 3. The design also makes sure \mathbf{P} is unrelated to \mathbf{X} . In the experiments, we show that the auxillary matrix \mathbf{P} is not necessary for the pre-trained Transformer to perform PCA with high accuracy. For the output, our theoretical analysis gives the following matrix

$$TF_{\theta}(\mathbf{H}) = [\hat{\mathbf{v}}_1^\top \quad \dots \quad \hat{\mathbf{v}}_k^\top]^\top \in \mathbb{R}^{dk}$$

which corresponds to the estimated principle eigenvectors of the matrix \mathbf{X} .

The Learning Problem. Consider a set of samples $\{\mathbf{X}_i\}_{i \in [u]}$ i.i.d. sampled from some distribution $p_{\mathbf{X}}$, we construct their oracle top- k principle components as $\mathbf{V}_i = [\mathbf{v}_1^{i,\top} \quad \dots \quad \mathbf{v}_k^{i,\top}]^\top$ and the context-augmented input matrix as \mathbf{H}_i for each \mathbf{X}_i . Then, the pretraining procedure is given by minimizing the following objective for some convex loss function $L(\cdot, \cdot) : \mathbb{R}^{dk} \times \mathbb{R}^{dk} \rightarrow \mathbb{R}$,

$$\hat{\theta} = \arg \min_{\theta \in \Theta(B_{\theta}, B_M)} \sum_{i=1}^u L(TF_{\theta}(\mathbf{H}_i), \mathbf{V}_i). \quad (1)$$

Here we consider $\Theta(B_{\theta}) := \{\theta : \|\theta\| \leq B_{\theta}, \max_{\ell} M_{\ell} \leq B_M\}$ to be the space of parameters. We also consider guarantees in the L_2 norm which states that $L(\mathbf{x}_1, \mathbf{x}_2) := \|\mathbf{x}_1 - \mathbf{x}_2\|_2$ in the theoretical part. Since $\hat{\theta}$ given by minimizing the empirical risk is not obtainable in practice, our theory only gives guarantee on the empirical risk minimizer. We further show that the local minimizers obtained through stochastic gradient optimization achieve good empirical performance in section 4.

3 THEORETICAL RESULTS

This section presents our theoretical results and the idea of taking each steps in the proof. Our proof constructs a particular instance of the transformers and show that the forward propagation on our constructed instance approximates the Power Method. We also carefully design the contextual matrix \mathbf{P} , explained as follows.

The Design of Auxillary Matrix. Our design of the matrix \mathbf{P} consists of three parts:

1. *Place Holder.* For $\ell \in \{1\} \cup [4 : k + 3]$ and $i \in [N]$, we let $\tilde{\mathbf{p}}_{\ell,i} = \mathbf{0} \in \mathbb{R}^{d \times 1}$. The place holders in \mathbf{P} records the intermediate results in the forward propagation.
2. *Identity Matrix.* We let $[\tilde{\mathbf{p}}_{2,1} \ \dots \ \tilde{\mathbf{p}}_{2,N}] = [\mathbf{I}_d \ \mathbf{0}_{d \times (N-d)}]$. The identity matrix in \mathbf{P} helps us screen out all the covariates \mathbf{X} in the forward propagation.
3. *Random Samples on the Hypersphere.* We let $\tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,k}$ be the i.i.d. samples uniformly distributed on \mathbb{S}^{d-1} . The random samples on the sphere corresponds to the initial vectors $\mathbf{v}_{0,\ell}$ for $\ell \in [k]$ in algorithm 1.

Given the above construction on the auxillary matrix \mathbf{P} , we are ready to state the existence theorem in this work, given as follows.

Theorem 3.1 (Transformer Approximation of the Power Iteration). *Denote the eigenvalues of $\mathbf{X}\mathbf{X}^\top$ to be $\lambda_1 > \lambda_2 > \dots > \lambda_k > \dots$. Let $\Delta := \min_{1 \leq i < j \leq k} |\lambda_i - \lambda_j|$. Assume that the eigenvalues of \mathbf{X} satisfy $\|\mathbf{X}\|_2 \leq B_X$. Assume that the initialized vectors $\tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N}$ satisfy $\tilde{\mathbf{p}}_{3,i}^\top \mathbf{v}_i \geq \delta$ for all $i \in [k]$ and make the rest of the vectors $\mathbf{0}$. Then, there exists a transformer model with number of layers $L = 2\tau + 4k + 1$ and number of heads $M \leq \lambda_1^d \frac{C}{\epsilon^2}$ with $\tau \leq \frac{\log(1/\epsilon_0\delta)}{\epsilon_0}$ such that for all $\epsilon_0, \epsilon > 0$, the final output $\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_k$ given by the transformer model achieve*

$$\|\hat{\mathbf{v}}_{\eta+1} - \mathbf{v}_{\eta+1}\|_2 \leq C\tau\epsilon\lambda_1^2 + \frac{C\lambda_1\sqrt{\epsilon_0}}{\Delta} \prod_{i=1}^{\eta} \frac{5\lambda_{i+1}}{\Delta}.$$

Moreover, consider the accuracy of multiple vs as a whole. There exists θ such that

$$L(TF_{\theta}(\mathbf{H}), \mathbf{V}) \leq C\tau\epsilon k\lambda_1^2 + C \left(\frac{\epsilon_0\lambda_1^2}{\Delta^2} \sum_{\eta=1}^{k-1} \prod_{i=1}^{\eta} \frac{25\lambda_{i+1}^2}{\Delta^2} \right)^{1/2}.$$

Remark 3. *The approximation error consists of two terms. The frist term comes from the approximation of the Power Method iterations by transformers. The second term comes from the error caused by finite iteration τ . To acquire a more direct account of the error terms and its order of magnitude, we consider a special case where the eigenvalues $\lambda_1 \asymp \lambda_2 \asymp \dots \asymp \lambda_k \asymp \Delta$. Then our results boils down to*

$$\|TF_{\theta}(\mathbf{H}) - [\mathbf{v}_1^\top, \mathbf{v}_2^\top, \dots, \mathbf{v}_k^\top]^\top\|_2 \leq C\tau\epsilon k\lambda_1^2 + C \frac{\lambda_1}{\Delta} \sqrt{k\epsilon_0}.$$

These results hide dimension d in the universal constant. Hence the dimension significantly affects the approximation properties of transformers. Our experimental results in section 4 also indicate that learning high dimensional principle eigenvectors is challenging.

We show that the conditions on $\tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N}$ can be achieved through sampling from isotropic Gaussians, given by the following lemma.

Lemma 3.1. *Let $\mathbf{y} \in \mathbb{R}^d$ be a random vector with isotropic Gaussian as its probability density. Consider $\mathbf{x} = \frac{\mathbf{y}}{\|\mathbf{y}\|_2}$. Let \mathbf{v} be any unit length vector, then we have for all $\delta < \frac{1}{2}d^{-1}$, $\mathbb{P}(|\mathbf{v}^\top \mathbf{x}| \leq \delta) \leq \frac{1}{\sqrt{\pi}}\sqrt{\delta} + \exp(-C\delta^{-\frac{1}{2}})$. Therefore, for all $\delta < \frac{1}{2}d^{-1}$, the event in theorem 3.1 is achieved with*

$$\mathbb{P}\left(\exists i \in [k] \text{ such that } \mathbf{x}_i^\top \mathbf{v}_i \leq \frac{\delta}{\sqrt{d}}\right) \leq \frac{k\sqrt{\delta}}{\sqrt{\pi}} + k \exp(-C\delta^{-1}).$$

Given the approximation error provided by theorem 3.1, we further provide the generalization error bound for the ERM defined by equation 1. This requires us to consider the following regularity conditions on the underlying distribution of $\mathbf{X}\mathbf{X}^\top$ (which also translates to the distribution of \mathbf{X}).

Assumption 1. *The distribution of $\mathbf{X}\mathbf{X}^\top$ supports on*

$$\mathbb{X} := \left\{ A : A \in \mathbf{S}_{++}^d, B_X \geq \lambda_1(A) > \lambda_2(A) > \dots > \lambda_k(A), \inf_{1 \leq i < j \leq k} \lambda_i(A) - \lambda_j(A) \geq \Delta \right\}.$$

Remark 4. *The above assumption can be easily generalized to distribution that supports on \mathbb{X} with high probability. Examples of such distribution include the Wishart distribution under the Gaussian design. In this work, we stick to the simplest case where the maximum eigenvalue is bounded from above.*

Given the above assumption, we are ready to state the generalization bound.

Proposition 1. *With probability at least $1 - \xi$, the ERM solution $\hat{\theta}$ satisfies*

$$\mathbb{E} \left[L(TF_{\hat{\theta}}(\mathbf{H}), \mathbf{V}) \mid \hat{\theta} \right] \leq \inf_{\theta \in \Theta(B_{\theta}, B_M)} \mathbb{E} [L(TF_{\theta}(\mathbf{H}), \mathbf{V})] + C \sqrt{\frac{k^3 L B_M d^2 \log(B_{\theta} + B_X + k) + \log(1/\xi)}{n}}.$$

Together with the bound given by theorem 3.1 and lemma 3.1, which essentially give a high probability upper bound on $\inf_{\theta \in \Theta(B_{\theta}, B_M)}$ we can derive a general upper bound on the generalization error, given as follows.

Corollary 3.1.1. *Under assumption 1, with probability at least $1 - \xi - \frac{k\sqrt{\delta}}{\sqrt{\pi}} - k \exp(-C\delta^{-1/2})$ for all $\delta < d^{-1}$ we have for all $\epsilon, \epsilon_0 > 0$,*

$$\mathbb{E} \left[L(TF_{\hat{\theta}}(\mathbf{H}), \mathbf{V}) \mid \hat{\theta} \right] \leq \mathbb{E} \left[C\tau\epsilon k\lambda_1^2 + C \left(\frac{\epsilon_0 \lambda_1^2}{\Delta^2} \sum_{\eta=1}^{k-1} \prod_{i=1}^{\eta} \frac{25\lambda_{i+1}^2}{\Delta^2} \right)^{1/2} \right] + C \sqrt{\frac{k^3 \log(\delta/\epsilon_0) \lambda_1^d d^2 \log(B_{\theta} + B_X + k) + \log(1/\xi)}{n\epsilon_0\epsilon^2}}.$$

Remark 5. *If we consider optimizing the bound w.r.t. ϵ_0 and ϵ , we obtain that $\mathbb{E} \left[L(TF_{\hat{\theta}}(\mathbf{H}), \mathbf{V}) \mid \hat{\theta} \right] \lesssim n^{-1/5}$ given that the rest of the parameters are of constant scales. It is not known if the results are improvable or not and the authors believe this question worth future explorations.*

4 SIMULATIONS

In this section, we verify the theoretical result in section 3 on synthetic and real-world datasets. Our experiments include both prediction of eigenvalues and eigenvectors. For synthetic datasets, we generate samples according to normal distributions. We focus on evaluating the effects of three major parameters: (1) The Impact of D ; (2) The Impact of Number of Layers; (3) The Impact of k_{train} ¹. For real-world datasets, we perform experiments on MNIST (LeCun et al., 1998) and Fashion-MNIST (Xiao et al., 2017). *All the results presented in this section are errors on the testing set.*

Data Preparation. For synthetic data $\mathbf{X} \in \mathbb{R}^{D \times N}$, we generate each column with a randomly initialized multivariate Gaussian distribution $\mathcal{N}(\mu, \Sigma)$. We then generate the labels as the top- k eigenvalues λ and eigenvectors V of the empirical covariance matrix $\mathbf{X}^T \mathbf{X} / (N - 1)$ via `numpy.linalg.eigh`. For real world dataset, we apply SVD to reduce the dimensionality of both datasets and evaluate whether the transformer, previously trained on multivariate Gaussian data, are capable of performing PCA on those real-world datasets. If the transformer successfully learns to perform PCA, we expect comparable performance on both synthetic and real-world data. For more details on data generation and configuration, please refer to table 2 in appendix C.2.

Model. We use the GPT2-architecture transformer (Radford et al., 2019) as our backbone model. We follow most settings in Garg et al. (2022), but replace the Softmax attention with ReLU attention as constructed in definition 1. We also provide a empirical comparison between Softmax and ReLU attention in Figure ?? in the Appendix. We use a slightly different architectures to predict eigenvalues and eigenvectors. For eigenvalues prediction, we flatten the transformer output

¹We denote k_{train} as the value of k used in training

Table 1: **Cosine Similarity for Different k .** We denote k_{train} as the number of eigenvectors to predict during training. For example, for $k_{\text{train}} = 4$, the model is trained to predict 4 eigenvectors.

k-th eigenvec.	k=1	k=2	k=3	k=4
$k_{\text{train}} = 4$	0.891(0.006)	0.616(0.038)	0.282(0.047)	0.120(0.022)
$k_{\text{train}} = 3$	0.908(0.011)	0.706(0.023)	0.366(0.018)	-
$k_{\text{train}} = 2$	0.903(0.006)	0.647(0.019)	-	-
$k_{\text{train}} = 1$	0.894(0.009)	-	-	-

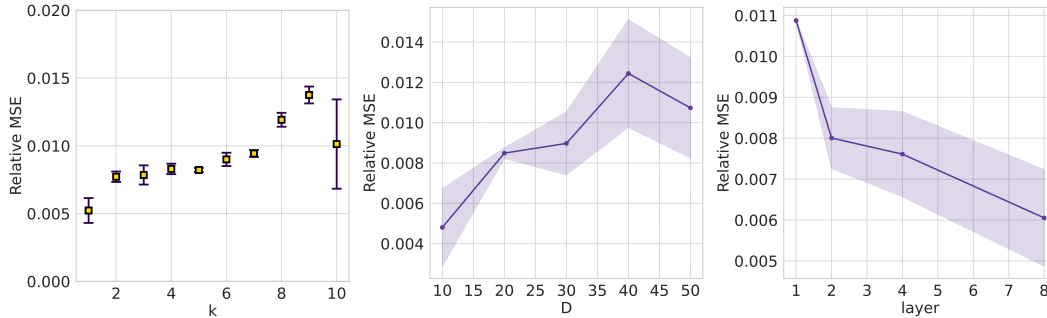


Figure 2: **Comparisons of Eigenvalue Prediction on Synthetic Data.** (1) *Left: Evidence of Transformer’s Ability to Predict Multiple Eigenvalues.* We use a small transformer (layer = 3, head = 2, embedding = 64) to predict top 10 eigenvalues with $D = 20$ and $N = 50$. All of the (Relative MSE) of 10 eigenvalues are below 2%, verifying that the transformer can predict eigenvalue very well. Additionally, the error of prediction grows slightly with k . We note that: (i) Higher-order eigenvalues require additional iterations and models with more layers. (ii) Smaller eigenvalues are more sensitive to the fluctuations in the predicted values under the relative MSE metric. (2) *Middle: Predictions of eigenvalues with different input dimension D .* We use a small transformer and use $N = 10$ in this experiment. We show that prediction errors increase significantly as dimension scales up, corroborating our theoretical remark 3. (3) *Right: Predictions of eigenvalues with different number of layers.* We use the same input as the previous multiple eigenvalues predictions experiment, and use a small transformer to predict top-3 eigenvalues. As the number of layers grow, the model performs better on eigenvalue prediction, which aligns with the result in theorem 3.1.

$TF_{\theta}(\mathbf{H}) \in \mathbb{R}^{N \times D}$ and use a linear layer $W_{\lambda} \in \mathbb{R}^{(N \cdot D) \times k}$ to readout the top k eigenvalues. As for eigenvectors, we use one more linear layer $W_v \in \mathbb{R}^{(N \cdot D) \times (k \cdot D)}$ to readout k eigenvectors concatenated in a 1-dimension vector. We use a transformer with layer = 3, head = 2, and embedding size = 64 to speed up the training process for most settings and find that it is sufficient to predict multiple eigenvalues and top-1 eigenvector well, see below sections for detailed discussion.

Metrics. For eigenvalues, we use relative mean squared error (RMSE) as loss function $\mathcal{L}_{\text{RMSE}}$ and evaluation metric. For the loss of predicting top- K eigenvalue, the loss function is defined as following

$$\mathcal{L}_{\text{RMSE}}(\lambda_i, \hat{\lambda}_i) := \frac{1}{K} \sum_{i=1}^K \frac{\lambda_i - \hat{\lambda}_i}{\lambda_i + \epsilon}.$$

For eigenvectors, we use cosine similarity as loss function and evaluation metric. For predicting k eigenvectors, the loss function is defined as

$$\mathcal{L}_{\text{cos}}(v_i, \hat{v}_i) := \frac{1}{K} \sum_{i=1}^K 1 - \frac{v_i \cdot \hat{v}_i}{\max(\|v_i\|_2 \cdot \|\hat{v}_i\|_2, \epsilon)},$$

where v_i represent the i -th eigenvector. The design of these loss functions not only matches the intuition of eigenvectors and eigenvalues, but also stabilize training by normalizing the loss values.

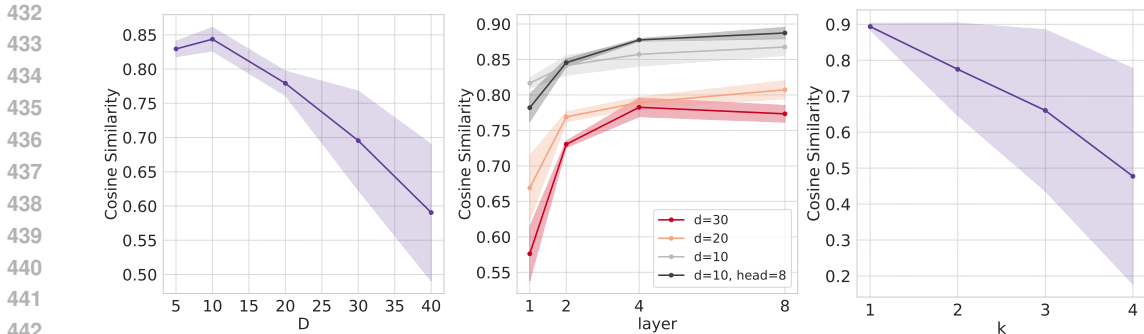


Figure 3: **Comparison of Eigenvector Prediction on Synthetic Data.** (1) Left: Prediction of top-1 eigenvector with different input dimension D . We use a small transformer and $N = 10$. As the dimension D scales up the eigenvector prediction suffers significantly. (2) Middle: Prediction of top=1 eigenvector with varying number of layers. We start from small transformer and use $N = 5$. The result demonstrates an ‘elbow effect’, where we show that the increase of L significantly boost the performance when L is small but halt to progress for larger L . We believe this can be explained by the bias-variance tradeoff. (3) Right: Predictions of eigenvectors with different numbers of k . We use $N = 10$ and $D = 10$ in this experiment. We use a larger transformer with layer= 12, heads= 8, and an embedding size= 256 in this experiment. The result demonstrates a decreasing prediction accuracy and increasing standard deviation as k increases. We list the individual cosine similarities of the predicted k -th eigenvectors in table 1. All the evaluations in the above three figures are averaged on three runs with different random seed.

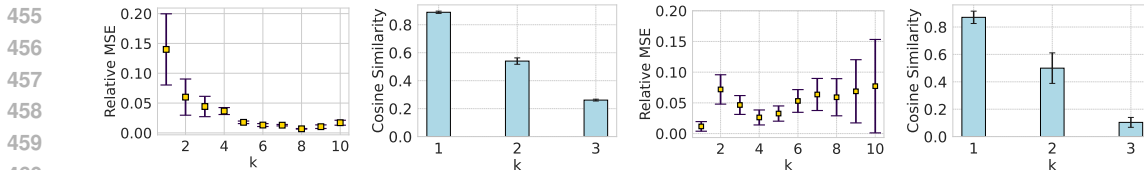


Figure 4: **Comparison of Eigenvalues and Eigenvectors Prediction on Real World Data.** (1) Left: Predicting Top-10 Eigenvalues on MNIST (2) Second Left: Predicting Top-3 Eigenvectors on MNIST (3) Second Right: Predicting Top-10 Eigenvalues on FMNIST (4) Right: Predicting Top-3 Eigenvectors on FMNIST. We show that on real world datasets, Transformers perform similarly to the synthetic datasets. The experimental setup for real world data is analogous to the ones performed for synthetic data.

4.1 SYNTHETIC DATA VS REAL WORLD DATA

Synthetic Dataset. The results on synthetic data are in figure 2. We first observe that transformers are capable of predicting top-10 eigenvalues with small error ($< 2\%$ error). The result also corresponds to theorem 3.1, indicating transformers are able to perform the power iteration method and generate eigenvalues with small error. For the impact of D , we observe the second subfigure in figure 2. In general, we discover an increasing trend of RMSE when D increases, this coincides with the theoretical findings stated in remark 3. We also observe that the error of prediction slightly increases with k , which is natural as the prediction dimension grows larger. For the impact of layers (right subfigure of figure 2), we observe that as the number of layer increases, RMSE shows significant reduction. This matches our theoretical construction as we show the iteration of power methods correspond to the number of layers, see figure 1 for the visualization of our transformer model. One thing to highlight is higher-order eigenvalues (larger k) have smaller magnitudes, which are more sensitive to fluctuations in the predicted values when using the relative MSE metric. This explains the higher variance/error of higher order eigenvalues. For eigenvectors, we also observe that transformers are capable of predicting principle eigenvectors. In particular, the cosine similarity between predicted eigenvector and ground truth is close to 1 when D is small.

Real World Dataset. The results on real world dataset are in figure 4. We observe that transformers are also capable of predicting top- k eigenvalues well on both MNIST and FMNIST. Despite the

486 difference in data distribution on training and test data, transformers are able to produce small error
 487 on predicting eigenvalues. Overall, we show that pretrained transformers learn PCA, and is able to
 488 generalize to other datasets as well. For eigenvectors, we can see that trained transformers show
 489 similar behavior on real world datasets when comparing to the synthetic ones. Indicating that our
 490 model actually learn to perform PCA instead of learn certain inductive bias.

492 4.2 PREDICTION WITH DIFFERENT PARAMETER COMBINATION.

494 **Prediction with different D .** The results are given in figure 3. In this experiment, we test the
 495 influence of increasing feature dimension D affects the ability of a Transformer model to predict
 496 the principal eigenvector of a data matrix. We use the simplest setting with a small transformer and
 497 test on $D = 5, 10, 20, 30, 40$, $N = 10$ and predict top-1 eigenvector. As the feature dimension
 498 D increases, we observe a clear trend of performance degradation of the Transformer’s ability to
 499 predict the principal eigenvector accurately. This confirms our theoretical results stated in remark 3
 500 that the feature dimension D affects the approximation properties of Transformers significantly.

502 **Prediction with different L .** The results are given in figure 3. In this experiment, we change the
 503 number of layers of transformer with head = 2 and embedding = 64, and set $N = 5$ to speed up the
 504 experiment. We observe that as the number of layers increases, the testing error also decreases, but
 505 the decreasing scale is less obvious when the number of layers becomes larger. However, the rate
 506 of improvement diminishes as the number of layers becomes larger. This suggests that increasing
 507 model depth alone is not sufficient for significantly enhancing eigenvector prediction. To verify our
 508 guess, we increase the number of heads from 2 to 8 and find that the cosine similarity increases
 509 further and with a slightly steeper incline. Note that there is a sharper decrease between layer= 1
 510 and layer= 2 across different d . This finding supports theorem 3.1 that we need at least 2 layers of
 511 the transformer to perform one iteration of the power method.

512 **Prediction with different k .** The results are given in table 1, and the right subfigure in figure 3
 513 where the cosine similarity in y-axis is averaged over k eigenvectors. We use $N = 10$ and $D = 10$
 514 in this experiment. We use a larger transformer with layer= 12, heads= 8, and an embedding
 515 size= 256 in this experiment. As shown in figure 3, the model’s ability to predict top k eigenvectors
 516 decreases as more eigenvectors are predicted, with increasing standard deviation. Table 1 lists the
 517 individual cosine similarities of the predicted k -th eigenvectors. The results show that most errors
 518 come from high-order eigenvectors. When trained to predict $k_{\text{train}} = 4$ eigenvectors, the model
 519 performs as well at predicting the top 1 eigenvector as when trained on $k_{\text{train}} = 1$. The result
 520 shows that most prediction errors come from high-order eigenvectors. This suggests that the pivotal
 521 difficulty is in the prediction of higher order eigenvectors.

523 5 DISCUSSIONS

524 This section discusses the limitations in this work and potential future working directions.

528 **Limitations.** Our limitations in the theoretical results can be summarized as follows: **(1)** From
 529 the theoretical perspective, our results guarantee the performance of ERM solutions whereas the
 530 true estimator is obtained through stochastic gradient descent method; **(2)** Our theoretical results
 531 utilize the context-augmented matrix \mathcal{P} , which is verified removable from our empirical results. It
 532 is conjectured that this is also not necessary in theory.

534 **Future Works.** Beyond resolving the limitations in this work, other future working directions
 535 from this work include: **(1)** Extend the results for relu transformers to softmax transformers. This
 536 step requires researchers to develop a new approximation bound for the softmax function; **(2)** Certify
 537 whether the rate $n^{-1/5}$ is sharp or not. The authors believe that this rate is improvable but it remains
 538 quite challenging; **(3)** The Spectral Method. Many modern high dimensional statistical questions
 539 can be resolved using the spectral method, which relies on the PCA as a sub-procedure. It is of
 general interest to see if these problems can be solved similarly by Transformers.

REFERENCES

- 540
541
542 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algo-
543 rithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*,
544 2022.
- 545
546 Jushan Bai, Serena Ng, et al. Large dimensional factor analysis. *Foundations and Trends® in*
547 *Econometrics*, 3(2):89–163, 2008.
- 548
549 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:
550 Provable in-context learning with in-context algorithm selection. *Advances in neural information*
551 *processing systems*, 36, 2024.
- 552
553 Satwik Bhattamishra, Arkil Patel, and Navin Goyal. On the computational power of transformers
554 and its implications in sequence modeling. *arXiv preprint arXiv:2006.09286*, 2020.
- 555
556 Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, vol-
557 ume 4. Springer, 2006.
- 558
559 Avrim Blum, John Hopcroft, and Ravindran Kannan. *Foundations of data science*. Cambridge
560 University Press, 2020.
- 561
562 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu,
563 and Zhifang Sui. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- 564
565 Jianqing Fan, Runze Li, Cun-Hui Zhang, and Hui Zou. *Statistical foundations of data science*.
566 Chapman and Hall/CRC, 2020.
- 567
568 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn
569 in-context? a case study of simple function classes. *Advances in Neural Information Processing*
570 *Systems*, 35:30583–30598, 2022.
- 571
572 Evarist Giné and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical mod-*
573 *els*, volume 40. Cambridge university press, 2016.
- 574
575 Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- 576
577 Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of*
578 *statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- 579
580 Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: Nngp and
581 ntk for deep attention networks. In *International Conference on Machine Learning*, pp. 4376–
582 4386. PMLR, 2020.
- 583
584 Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and
585 Mubarak Shah. Transformers in vision: A survey. *ACM computing surveys (CSUR)*, 54(10s):
586 1–41, 2022.
- 587
588 Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selec-
589 tion. *Annals of Statistics*, pp. 1302–1338, 2000.
- 590
591 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
592 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 593
594 Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers
595 learn shortcuts to automata. *arXiv preprint arXiv:2210.10749*, 2022.
- 596
597 Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large
598 language models. *Advances in neural information processing systems*, 36:21702–21720, 2023.
- 599
600 Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is turing-complete. *Journal of Ma-*
601 *chine Learning Research*, 22(75):1–35, 2021.
- 602
603 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
604 models are unsupervised multitask learners. 2019.

594 Kai Shen, Junliang Guo, Xu Tan, Siliang Tang, Rui Wang, and Jiang Bian. A study on relu and
595 softmax in transformer. *arXiv preprint arXiv:2302.06461*, 2023.
596

597 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
598

599 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordv-
600 intsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient
601 descent. In *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.

602 Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cam-
603 bridge university press, 2019.

604 Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on
605 approximating turing machines with transformers. *Advances in Neural Information Processing*
606 *Systems*, 35:12071–12083, 2022.

607

608 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
609 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art
610 natural language processing. In *Proceedings of the 2020 conference on empirical methods in*
611 *natural language processing: system demonstrations*, pp. 38–45, 2020.

612 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark-
613 ing machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

614

615 Shunyu Yao, Binghui Peng, Christos Papadimitriou, and Karthik Narasimhan. Self-attention net-
616 works can process bounded hierarchical languages. *arXiv preprint arXiv:2105.11115*, 2021.

617 Yi Yu, Tengyao Wang, and Richard J Samworth. A useful variant of the davis–kahan theorem for
618 statisticians. *Biometrika*, 102(2):315–323, 2015.

619

620 Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J Reddi, and Sanjiv Kumar.
621 Are transformers universal approximators of sequence-to-sequence functions? *arXiv preprint*
622 *arXiv:1912.10077*, 2019.
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648	A Additional Theoretical Background	13
649		
650	B Proofs	13
651	B.1 Proof of Proposition 1	13
652	B.2 Proof of Theorem 3.1	14
653	B.3 Proof of Lemma 3.1	22
654	B.4 Proof of Lemma B.1	23
655	B.5 Proof of Lemma B.2	24
656	C Experimental Details	26
657	C.1 Setup.	26
658	C.2 Data.	26
659	C.3 Hyperparameters.	26
660	C.4 Additional Experimental Results	26

A ADDITIONAL THEORETICAL BACKGROUND

Definition 4 (Sufficiently Smooth d -variate function). Denote $B_\infty^d(R) := [-R, R]^d$ as the standard ℓ_∞ ball in \mathbb{R}^d . We say a function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ is (R, C_ℓ) smooth if for $s = \lceil (d-1)/2 \rceil + 2$, g is a C^s function on $B_\infty^d(\mathbb{R})$ and

$$\sup_{\mathbf{z} \in B_\infty^d(R)} \|\nabla^d g(\mathbf{z})\|_\infty = \sup_{\mathbf{z} \in B_\infty^d(R)} \max_{j_1, \dots, j_i \in [d]} |\partial_{x_{j_1} \dots x_{j_i}} g(\mathbf{x})| \leq L_i$$

for all $i \in \{0, 1, \dots, s\}$, with $\max_{0 \leq i \leq s} L_i R^i \leq C_\ell$.

Definition 5 (Approximability by sum of Relus (Bai et al., 2024)). A function $g : \mathbb{R}^k \rightarrow \mathbb{R}$ is $(\epsilon_{approx}, R, M, C)$ -approximable by sum of Relus if there exists a function $f_{M,C}$ such that

$$f_{M,C}(\mathbf{z}) = \sum_{m=1}^M c_m \sigma(\mathbf{a}_m^\top [\mathbf{z}; 1]) \text{ with } \sum_{m=1}^M |c_m| \leq C, \max_{m \in [M]} \|\mathbf{a}_m\|_1 \leq 1, \quad \mathbf{a}_m \in \mathbb{R}^{k+1}, c_m \in \mathbb{R},$$

such that $\sup_{\mathbf{z} \in B_\infty^k(R)} |g(\mathbf{z}) - f_{M,C}(\mathbf{z})| \leq \epsilon_{approx}$.

B PROOFS

B.1 PROOF OF PROPOSITION 1

Proof. The proof follows from (Wainwright, 2019), using the fact that for all $\theta \in \Theta(B_\theta, B_M)$, we have

$$\frac{1}{n} \sum_{j=1}^n L(TF_{\hat{\theta}}(\mathbf{H}_i), \mathbf{V}_i) \leq \frac{1}{n} \sum_{j=1}^n L(TF_\theta(\mathbf{H}_i), \mathbf{V}_i),$$

it is not hard to show that

$$\mathbb{E} [L(TF_{\hat{\theta}}(\mathbf{H}), \mathbf{V})] \leq \inf_{\theta \in \Theta(B_\theta, B_M)} \mathbb{E} [L(TF_\theta(\mathbf{H}), \mathbf{V})] + 2 \sup_{\theta \in \Theta(B_\theta, B_M)} |X_\theta|,$$

where $X_\theta = \frac{1}{n} \sum_{j=1}^n L(TF_\theta(\mathbf{H}_i), \mathbf{V}_i) - \mathbb{E}[L(TF_\theta(\mathbf{H}), \mathbf{V})]$ is the empirical process indexed by θ . The tail bound for empirical process requires us to verify a few regularity conditions (Giné & Nickl, 2016) on the function L and the set Θ

1. The metric entropy of an operator norm ball $\log N(\delta, B_{\|\cdot\|_{op}}(r), \|\cdot\|_{op}) \leq CLB_M D^2 \log(1 + 2(B_\theta + B_X + k)/\delta)$.
2. $L(TF_\theta(\mathbf{H}), \mathbf{V}) \leq C\sqrt{k}$.
3. The Lipschitz condition of Transformers satisfies that for all $\theta_1, \theta_2 \in \Theta(B_\theta, B_M)$, we have $L(TF_{\theta_1}(\mathbf{H}), \mathbf{V}) - L(TF_{\theta_2}(\mathbf{H}), \mathbf{V}) \leq CLB_1^L \|\theta_1 - \theta_2\|_{op}$ where $B_1 = B_\theta^4 B_X^3$.

The first and second verifications follow immediately from J.2 in (Bai et al., 2024). The third verification is given upon noticing that as $L(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2$,

$$\sup_{\boldsymbol{\theta}, \mathbf{H}, \mathbf{V}} L(TF_{\boldsymbol{\theta}}(\mathbf{H}), \mathbf{V}) \leq C\sqrt{k}, \quad \|\nabla_{\mathbf{x}} L\| \leq C.$$

Further note that $\|\tilde{\mathbf{W}}_0\|_2 \asymp \|\tilde{\mathbf{W}}_1\|_2 \asymp 1$. Given the above result, and corollary J.1 in (Bai et al., 2024), we can show that

$$L(TF_{\boldsymbol{\theta}_1}(\mathbf{H}), \mathbf{V}) - L(TF_{\boldsymbol{\theta}_2}(\mathbf{H}), \mathbf{V}) \leq CLB_1^L \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_{op},$$

where $B_1 = B_{\theta}^4 B_X^3$. Therefore, using the uniform concentration bound given by proposition A.4 we can show that with probability at least $1 - \xi$, we have

$$\sup_{\boldsymbol{\theta} \in \Theta(B_{\theta}, B_M)} |X_{\theta}| \leq C\sqrt{k} \sqrt{\frac{LB_M D^2 \log(B_{\theta} + B_X + k) + \log(1/\delta)}{n}}.$$

Therefore, replacing D with Ckd we complete the proof. \square

B.2 PROOF OF THEOREM 3.1

Proof. Our proof can be dissected into the following steps: 1. We construct a Transformer with fixed parameters that performs (1) The computation of the symmetrized covariate matrix; (2) The approximation of the power method; (3) The removal of the principle eigenvectors; (4) Adjust the dimension of the output through multiplying the two matrices $\tilde{\mathbf{W}}_0$ and $\tilde{\mathbf{W}}_1$ on the left and right.

1. The Covariate Matrix.

To compute the covariate matrix $\mathbf{X}\mathbf{X}^{\top}$, we construct $\mathbf{H} = \begin{bmatrix} \mathbf{X}_1, \dots, \mathbf{X}_N \\ \tilde{\mathbf{p}}_{1,1}, \dots, \tilde{\mathbf{p}}_{1,N} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} = \begin{bmatrix} \mathbf{X} \\ \mathbf{P} \end{bmatrix}$ we let

$m = 2$ and

$$\mathbf{V}_1^{cov} = I_D = -\mathbf{V}_2^{cov}, \quad \mathbf{Q}_1^{cov, \top} \mathbf{K}_1^{cov} = -\mathbf{Q}_2^{\top} \mathbf{K}_2 = \begin{bmatrix} \mathbf{0}_{N+1 \times d}, I_d, \mathbf{0} \\ \mathbf{0}, \mathbf{0}, \mathbf{0} \end{bmatrix} \in \mathbb{R}^{D \times D},$$

$$\tilde{\mathbf{p}}_{1,\ell,j} = \mathbf{0}, \quad \tilde{\mathbf{p}}_{2,\ell,j} = \begin{cases} \mathbb{1}_{\ell=j} & \text{when } \ell \leq d \\ 0 & \text{when } \ell > d \end{cases} \quad (2)$$

Under the above construction, we obtain that

$$\mathbf{Q}_1^{\top} \mathbf{K}_1 \mathbf{H} = \begin{bmatrix} I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{D \times N}, \quad \mathbf{Q}_2^{\top} \mathbf{K}_2 \mathbf{H} = \begin{bmatrix} -I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{D \times N},$$

$$\sigma(\mathbf{H}^{\top} \mathbf{Q}_1^{\top} \mathbf{K}_1 \mathbf{H}) + \sigma(\mathbf{H}^{\top} \mathbf{Q}_2^{\top} \mathbf{K}_2 \mathbf{H}) = [\mathbf{X}^{\top}, \mathbf{0}] \in \mathbb{R}^{N \times N}.$$

We further obtain that

$$\frac{1}{N} \sum_{m=1}^M (\mathbf{V}_m \mathbf{H}) \times \sigma((\mathbf{Q}_m \mathbf{H})^{\top} (\mathbf{K}_m \mathbf{H})) = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{X}\mathbf{X}^{\top} \in \mathbb{R}^{d \times d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{D \times D}.$$

Therefore, the output is given by $\tilde{\mathbf{H}}^{cov} = \begin{bmatrix} \mathbf{X} \\ \mathbf{X}\mathbf{X}^{\top}, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix}$.

2. The Power Iteration.

Then we consider constructing a single attention layer that approximates the power iteration. This step involves two important operations: (1) Obtaining the vector given by $\tilde{X}\tilde{X}^\top v$. (2) Approximation of the value of the inverse norm given by $1/\|\tilde{X}\tilde{X}^\top v\|_2$. We show that one can use the multihead Relu Transformer to achieve both goals simultaneously, whose parameters are given by

$$\begin{aligned} V_1^{pow,1} &= -V_2^{pow,1} = \begin{bmatrix} \mathbf{0}_{(3d+1)\times(2d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{(d)\times(2d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ Q_1^{pow,1} &= -Q_1^{pow,1} = \begin{bmatrix} \mathbf{0}_{(d+1)\times(d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ K_1^{pow,1} &= K_2^{pow,1} = \begin{bmatrix} \mathbf{0}_{(3d+1)\times(3d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(3d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \tilde{p}_{4,j} = \mathbf{0} \text{ for all } j \in [N]. \end{aligned}$$

Given the above formulation, we are able to show that

$$\begin{aligned} Q_2^{pow,1} \tilde{H}^{cov} &= -Q_1^{pow,1} \tilde{H}^{cov} = \begin{bmatrix} \mathbf{0}_{(2d+1)\times N} \\ \tilde{X}\tilde{X}^\top, \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \\ K_2^{pow,1} \tilde{H}^{cov} &= K_1^{pow,1} \tilde{H}^{cov} = \begin{bmatrix} \mathbf{0}_{2d+1} \\ \tilde{p}_{3,1}, \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \end{aligned}$$

which implies that

$$\sum_{m \in \{1,2\}} \sigma((Q_m^{pow,1} \tilde{H}^{cov})^\top K_m^{pow,1} \tilde{H}^{cov}) = \begin{bmatrix} \mathbf{0} & \mathbf{0}_{d \times (2d+1)} \\ \tilde{X}\tilde{X}^\top \tilde{p}_{3,1} & \mathbf{0}_{d \times (N-1)} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

Then we can show that

$$\begin{aligned} \tilde{H}^{pow,1} - \tilde{H}^{cov} &= \sum_{m \in \{1,2\}} V_m^{pow,1} \tilde{H}^{cov} \times \sigma((Q_m^{pow,1} \tilde{H}^{cov})^\top K_m^{pow,1} \tilde{H}^{cov}) \\ &= \begin{bmatrix} \mathbf{0}_{3d+1} \\ \tilde{X}\tilde{X}^\top \tilde{p}_{3,1}, \mathbf{0}_{d \times (N-1)} \\ \mathbf{0} \end{bmatrix}. \end{aligned}$$

Therefore, we conclude that the output of the first power iteration layer is given by

$$\tilde{H}^{pow,1} = \begin{bmatrix} \tilde{X} \\ \tilde{y}^\top \\ \tilde{X}\tilde{X}^\top, \mathbf{0} \\ \tilde{p}_{2,1}, \dots, \tilde{p}_{2,N} \\ \tilde{p}_{3,1}, \dots, \tilde{p}_{3,N} \\ \tilde{X}\tilde{X}^\top \tilde{p}_{3,1}, \mathbf{0} \\ \tilde{p}_{5,1}, \dots, \tilde{p}_{5,N} \\ \vdots \\ \tilde{p}_{\ell,1}, \dots, \tilde{p}_{\ell,N} \end{bmatrix}.$$

Then, using lemma B.2, we design an extra attention layer that performs the normalizing procedure, with the following parameters for all $m \in [M]$,

$$\begin{aligned} V_m^{pow,2} &= \begin{bmatrix} \mathbf{0}_{d \times (4d+1)} & c_m I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad Q_m^{pow,2} = \begin{bmatrix} \mathbf{0}_{d \times (2d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ K_m^{pow,2} &= \begin{bmatrix} \mathbf{0}_{1 \times (3d+1)} & \mathbf{a}_m^\top & \mathbf{0} \\ \vdots & \vdots & \vdots \\ \mathbf{0}_{1 \times (3d+1)} & \mathbf{a}_m^\top & \mathbf{0} \\ \mathbf{0}_{(D-d) \times (3d+1)} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{aligned}$$

Under the above construction, we obtain that

$$(\mathbf{Q}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1})^\top = \begin{bmatrix} I_{d \times d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{K}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1} = \begin{bmatrix} \mathbf{a}_m^\top \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1} & \mathbf{0} \\ \vdots & \\ \mathbf{a}_m^\top \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1} & \mathbf{0} \\ \mathbf{0}_{(D-d) \times 1} & \mathbf{0} \end{bmatrix}.$$

Then, given $\mathbf{V}_m^{pow,2}$ we can show that under the condition given by lemma B.2, we have

$$\left\| \sum_{m=1}^M \mathbf{V}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1} \sigma \left((\mathbf{Q}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1})^\top (\mathbf{K}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1}) \right) - \begin{bmatrix} \mathbf{0}_{4d+1} \\ \frac{\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2} - \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}, \mathbf{0} \\ \mathbf{0} \end{bmatrix} \right\|_\infty < \epsilon,$$

Moreover, we can further achieve that

$$\left\| \sum_{m=1}^M \mathbf{V}_m^{pow,2j} \tilde{\mathbf{H}}^{pow,1} \sigma \left((\mathbf{Q}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1})^\top (\mathbf{K}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1}) \right) - \begin{bmatrix} \mathbf{0}_{4d+1} \\ \frac{\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2} - \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}, \mathbf{0} \\ \mathbf{0} \end{bmatrix} \right\|_2 < \epsilon \|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2.$$

Hence, using the fact that $\tilde{\mathbf{H}}^{pow,2} = \tilde{\mathbf{H}}^{pow,1} + \sum_{i=1}^m \mathbf{V}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1} \sigma \left((\mathbf{Q}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1})^\top (\mathbf{K}_m^{pow,2} \tilde{\mathbf{H}}^{pow,1}) \right)$, we obtain that

$$\left\| \tilde{\mathbf{H}}^{pow,2} - \begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{y}} \\ \mathbf{X} \mathbf{X}^\top, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1} \\ \frac{\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2}, \dots, \mathbf{0} \\ \vdots \end{bmatrix} \right\|_2 < \epsilon \|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2.$$

Then we construct another attention layer, which performs similar calculations as that of $pow, 1$ but switch the rows of $\tilde{\mathbf{p}}_{3,1}$ with that of $\frac{\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2}$. Our construction for the third layer is given by

$$\begin{aligned} \mathbf{V}_1^{pow,3} &= -\mathbf{V}_2^{pow,3} = \begin{bmatrix} \mathbf{0}_{(3d+1) \times (2d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d \times (2d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{Q}_1^{pow,3} &= -\mathbf{Q}_2^{pow,3} = \begin{bmatrix} \mathbf{0}_{(3d+1) \times (d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d \times (d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_1^{pow,3} &= \mathbf{K}_2^{pow,3} = \begin{bmatrix} \mathbf{0}_{(4d+1) \times (4d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d \times (4d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \tilde{\mathbf{p}}_{4,j} = \mathbf{0} \text{ for all } j \in [N]. \end{aligned}$$

Given the above construction, we can show that

$$\begin{aligned} \mathbf{Q}_2^{pow,3} \tilde{\mathbf{H}}^{pow,2} &= -\mathbf{Q}_1^{pow,3} \tilde{\mathbf{H}}^{pow,2} = \begin{bmatrix} \mathbf{0}_{(3d+1) \times N} \\ \mathbf{X} \mathbf{X}^\top \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{K}_2^{pow,3} \tilde{\mathbf{H}}^{pow,2} = \mathbf{K}_1^{pow,3} \tilde{\mathbf{H}}^{pow,2}, \\ \left\| \mathbf{K}_2^{pow,3} \tilde{\mathbf{H}}^{pow,2} - \begin{bmatrix} \mathbf{0}_{(3d+1) \times N} \\ \frac{\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2}, \mathbf{0} \\ \mathbf{0} \end{bmatrix} \right\|_2 &\leq \epsilon \|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2. \end{aligned}$$

Then, using the fact that given $\mathbf{x}_1, \mathbf{x}_2$ with $\|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq \delta_0$, we have $\|\mathbf{X}\mathbf{X}^\top(\mathbf{x}_1 - \mathbf{x}_2)\|_2 \leq \|\mathbf{X}\mathbf{X}^\top\|_2 \delta_0$. Hence, collecting the above pieces, we have

$$\left\| \sum_{m=1}^2 \mathbf{V}_m^{pow,3} \tilde{\mathbf{H}}^{pow,2} \sigma \left((\mathbf{Q}_2^{pow,3} \tilde{\mathbf{H}}^{pow,2})^\top \mathbf{K}_2^{pow,3} \tilde{\mathbf{H}}^{pow,2} \right) - \begin{bmatrix} \mathbf{0}_{(3d+1) \times N} \\ \frac{(\mathbf{X}\mathbf{X}^\top)^2 \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2} - \frac{\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2}, \mathbf{0} \\ \mathbf{0} \end{bmatrix} \right\|_2 \leq \epsilon \|\mathbf{X}\mathbf{X}^\top\|_2 \|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2.$$

Henceforth, one can further show that

$$\left\| \tilde{\mathbf{H}}^{pow,3} - \begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{y}} \\ \mathbf{X}\mathbf{X}^\top, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \frac{(\mathbf{X}\mathbf{X}^\top)^2 \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2}, \mathbf{0} \\ \mathbf{p}_{5,1}, \dots, \mathbf{p}_{5,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_2 \leq$$

$\epsilon \|\mathbf{X}\mathbf{X}^\top\|_2 \|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2$. Consider we are doing in total of τ power iterations, we can set for all $\tau \in \mathbb{N}^*$,

$$\begin{aligned} \mathbf{V}_m^{pow,2\tau+1} &= \mathbf{V}_m^{pow,3}, & \mathbf{Q}_m^{pow,2\tau+1} &= \mathbf{Q}_m^{pow,3}, & \mathbf{K}_m^{pow,2\tau+1} &= \mathbf{K}_m^{pow,3}, \\ \mathbf{V}_m^{pow,2\tau+2} &= \mathbf{V}_m^{pow,4}, & \mathbf{Q}_m^{pow,2\tau+2} &= \mathbf{Q}_m^{pow,4}, & \mathbf{K}_m^{pow,2\tau+2} &= \mathbf{K}_m^{pow,4}. \end{aligned}$$

Therefore, taking another layer of normalization, we can show that

$$\left\| \tilde{\mathbf{H}}^{pow,3} - \begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{y}} \\ \mathbf{X}\mathbf{X}^\top, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \frac{(\mathbf{X}\mathbf{X}^\top)^2 \tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}\|_2}, \mathbf{0} \\ \mathbf{p}_{5,1}, \dots, \mathbf{p}_{5,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_2 \leq 2\epsilon \|\mathbf{X}\mathbf{X}^\top\|_2.$$

Then, using the sublinearity of errors, we can show that for $\tau \in \mathbb{N}$,

$$\left\| \tilde{\mathbf{H}}^{pow,2\tau+2} - \begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{y}} \\ \mathbf{X}\mathbf{X}^\top, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \tilde{\mathbf{p}}_{5,1}, \dots, \tilde{\mathbf{p}}_{5,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_\infty \leq \tau \epsilon \|\mathbf{X}\mathbf{X}^\top\|_2, \quad \tilde{\mathbf{p}}_{3,1}^{(\tau)} = \frac{\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau-1)}}{\|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau-1)}\|_2}, \quad \tilde{\mathbf{p}}_{3,1}^{(0)} = \tilde{\mathbf{p}}_{3,1}.$$

If we denote \mathbf{v}_i as the eigenvector corresponds to the i th largest eigenvalue of $\mathbf{X}\mathbf{X}^\top$. Let the eigenvalues of $\mathbf{X}\mathbf{X}^\top$ be denoted by $\lambda_1 > \lambda_2 > \dots > \lambda_n$. Given $|\tilde{\mathbf{p}}_{3,1}^\top \mathbf{v}_1| > \delta$ and $|\sqrt{\lambda_1} - \sqrt{\lambda_2}| = \Omega(1)$. Theorem 3.11 in (Blum et al., 2020) page 53 shows that given $k = \frac{\log(1/\epsilon_0 \delta)}{2\epsilon_0}$ and $\|\tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 = \|\mathbf{v}_1\|_2 = 1$, one immediately obtains that

$$\tilde{\mathbf{p}}_{3,1}^{(\tau)\top} \mathbf{v}_1 \geq 1 - \epsilon_0, \quad \|\tilde{\mathbf{p}}_{3,1}^{(\tau)} - \mathbf{v}_1\|_2 = \sqrt{2 - 2\mathbf{v}_1^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}} = \sqrt{2\epsilon_0}.$$

And we also consider the approximation of the maximum eigenvalue. Note that using $\|\mathbf{v}_1\|_2 = 1$, we have

$$\begin{aligned}\|\mathbf{X}\mathbf{X}^\top\|_2 &= \|\mathbf{X}\mathbf{X}^\top\mathbf{v}_1\|_2 = \left\| \mathbf{X}\mathbf{X}^\top\tilde{\mathbf{p}}_{3,1}^{(\tau)} + \mathbf{X}\mathbf{X}^\top(\mathbf{v}_1 - \tilde{\mathbf{p}}_{3,1}^{(\tau)}) \right\|_2 \\ &\leq \|\mathbf{X}\mathbf{X}^\top\tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 + \|\mathbf{X}\mathbf{X}^\top(\mathbf{v}_1 - \tilde{\mathbf{p}}_{3,1}^{(\tau)})\|_2 \\ &\leq \|\mathbf{X}\mathbf{X}^\top\tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 + \|\mathbf{X}\mathbf{X}^\top\|_2\|\mathbf{v}_1 - \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2.\end{aligned}$$

Similarly we can also derive that $\|\mathbf{X}\mathbf{X}^\top\|_2 \geq \|\mathbf{X}\mathbf{X}^\top\tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 - \|\mathbf{X}\mathbf{X}^\top\|_2\|\mathbf{v}_1 - \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2$. Then we show that

$$\left| \|\mathbf{X}\mathbf{X}^\top\|_2 - \|\mathbf{X}\mathbf{X}^\top\tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 \right| \leq \|\mathbf{X}\mathbf{X}^\top\|_2\|\mathbf{v}_1 - \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 \leq \sqrt{2\epsilon_0}\|\mathbf{X}\mathbf{X}^\top\|_2.$$

3. The Removal of Principle Eigenvectors.

After τ iterates on the power method, we need to remove the principle term from the matrix $\mathbf{X}\mathbf{X}^\top$, achieved through two important steps: (1) The computation of the estimated eigenvalue $\|\mathbf{X}\mathbf{X}^\top\tilde{\mathbf{p}}_{3,1}\|_2$. (2) The construction of the low rank update $\tilde{\mathbf{p}}_{3,1}\tilde{\mathbf{p}}_{3,1}^\top$. For step (1), we consider the following construction:

$$\begin{aligned}\mathbf{V}_1^{rpe,1} &= -\mathbf{V}_2^{rpe,1} = \begin{bmatrix} \mathbf{0}_{(3d+1)\times(2d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(2d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{Q}_1^{rpe,1} = -\mathbf{Q}_2^{rpe,1} = \begin{bmatrix} \mathbf{0}_{(d+1)\times(d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_1^{rpe,1} &= \mathbf{K}_2^{rpe,1} = \begin{bmatrix} \mathbf{0}_{(4d+1)\times(4d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(4d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}.\end{aligned}$$

Note that the above construction is similar to the first layer of the power method. Under this construction, we can show that

$$\begin{aligned}\tilde{\mathbf{H}}^{rpe,1} &= \tilde{\mathbf{H}}^{pow,2\tau+2} + \sum_{m \in \{1,2\}} \mathbf{V}_m^{rpe,1} \sigma((\mathbf{Q}_m^{rpe,1} \tilde{\mathbf{H}}^{pow,2\tau+2})^\top (\mathbf{K}_m^{rpe,1} \tilde{\mathbf{H}}^{pow,2\tau+2})), \\ \left\| \tilde{\mathbf{H}}^{rpe,1} - \underbrace{\begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{y}} \\ \mathbf{X}\mathbf{X}^\top, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \mathbf{X}\mathbf{X}^\top\tilde{\mathbf{p}}_{3,N}^{(\tau)}, \mathbf{0} \\ \tilde{\mathbf{p}}_{6,1}, \dots, \tilde{\mathbf{p}}_{6,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix}}_{=: \mathbf{H}^{rpe,1}} \right\|_2 &\leq C\tau\epsilon\|\mathbf{X}\mathbf{X}^\top\|_2^2, \quad \tilde{\mathbf{p}}_{5,i} = \mathbf{0}, \quad \forall i \in [N].\end{aligned}\quad (3)$$

Then, we construct the next layer, using the notations in lemma B.2, for $M \geq \|\mathbf{X}\mathbf{X}^\top\|_2^d \frac{C(d)}{\epsilon^2}$ for all $m \in [M]$ we have

$$\begin{aligned}\mathbf{V}_m^{rpe,2} &= \begin{bmatrix} \mathbf{0}_{d\times(4d+1)} & d_m I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{Q}_m^{rpe,2} = \begin{bmatrix} \mathbf{0}_{d\times(2d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_m^{rpe,2} &= \begin{bmatrix} \mathbf{0}_{1\times(5d+1)} & \mathbf{b}_m^\top & \mathbf{0} \\ \vdots & & \\ \mathbf{0}_{1\times(5d+1)} & \mathbf{b}_m^\top & \mathbf{0} \\ \mathbf{0}_{(D-d)\times(5d+1)} & \mathbf{0} & \mathbf{0} \end{bmatrix}.\end{aligned}$$

Given the above construction, we subsequently show that

$$(\mathbf{Q}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1})^\top = \begin{bmatrix} I_{d \times d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{K}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1} = \begin{bmatrix} \mathbf{b}_m^\top \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)} & \mathbf{0} \\ \vdots & \vdots \\ \mathbf{b}_m^\top \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)} & \mathbf{0} \\ \mathbf{0}_{(D-d) \times 1} & \mathbf{0} \end{bmatrix}.$$

Hence, given the construction of $\mathbf{V}_m^{rpe,2}$, we can show that $\tilde{\mathbf{H}}^{rpe,2}$ satisfies

$$\begin{aligned} \tilde{\mathbf{H}}^{rpe,2} &= \tilde{\mathbf{H}}^{rpe,1} + \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1} \times \sigma \left((\mathbf{K}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1})^\top (\mathbf{Q}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1}) \right) \\ &= \underbrace{\mathbf{H}^{rpe,1} + \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \mathbf{H}^{rpe,1} \times \sigma \left((\mathbf{K}_m^{rpe,2} \mathbf{H}^{rpe,1})^\top (\mathbf{Q}_m^{rpe,2} \mathbf{H}^{rpe,1}) \right)}_{=: \tilde{\mathbf{H}}^{rpe,1}} \\ &\quad + \left(\tilde{\mathbf{H}}^{rpe,1} - \mathbf{H}^{rpe,1} \right) + \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1} \times \sigma \left((\mathbf{K}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1})^\top \mathbf{Q}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1} \right) \\ &\quad - \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1} \times \sigma \left((\mathbf{K}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1})^\top \mathbf{Q}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1} \right). \end{aligned}$$

We note that by lemma B.2 we can show that

$$\left\| \widehat{\mathbf{H}}^{rpe,1} - \begin{bmatrix} \mathbf{X} \\ \tilde{\mathbf{y}} \\ \mathbf{X} \mathbf{X}^\top, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,N}^{(\tau)}\|_2^{\frac{1}{2}} \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_2 \leq C\tau\epsilon \|\mathbf{X} \mathbf{X}^\top\|_2^2.$$

Then the rest of the proof focuses on showing that the rest of the terms are small. Note that using equation 3, we show that

$$\left\| \tilde{\mathbf{H}}^{rpe,1} - \mathbf{H}^{rpe,1} \right\|_2 \leq \tau\epsilon \|\mathbf{X} \mathbf{X}^\top\|_2^2.$$

And for the last term, we can show that

$$\begin{aligned} &\left\| \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1} \times \sigma \left((\mathbf{K}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1})^\top (\mathbf{Q}_m^{rpe,2} \tilde{\mathbf{H}}^{rpe,1}) \right) \right. \\ &\quad \left. - \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \mathbf{H}^{rpe,1} \times \sigma \left((\mathbf{K}_m^{rpe,2} \mathbf{H}^{rpe,1})^\top (\mathbf{Q}_m^{rpe,2} \mathbf{H}^{rpe,1}) \right) \right\|_2 \\ &\leq C\tau\epsilon \|\mathbf{X} \mathbf{X}^\top\|_2^2. \end{aligned}$$

Collecting the above pieces, we finally show that

$$\left\| \tilde{\mathbf{H}}^{rpe,2} - \begin{bmatrix} \mathbf{X} \\ \mathbf{X} \mathbf{X}^\top, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \|\mathbf{X} \mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2^{\frac{1}{2}} \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \tilde{\mathbf{p}}_{6,1}, \dots, \tilde{\mathbf{p}}_{6,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_2 \leq C\tau\epsilon \|\mathbf{X} \mathbf{X}^\top\|_2^2.$$

Then we construct another layer to remove the principle components from the matrix $\mathbf{X}\mathbf{X}^\top$, given by

$$\begin{aligned} -\mathbf{V}_1^{rpe,3} = \mathbf{V}_2^{rpe,3} &= \begin{bmatrix} \mathbf{0}_{(d+1)\times(4d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, & \mathbf{Q}_1^{rpe,3} = -\mathbf{Q}_2^{rpe,3} &= \begin{bmatrix} \mathbf{0}_{d\times(4d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_1^{rpe,3} = \mathbf{K}_2^{rpe,3} &= \begin{bmatrix} \mathbf{0}_{d\times(4d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{aligned}$$

Then we can show that

$$(\mathbf{Q}_1^{rpe,3} \tilde{\mathbf{H}}^{rpe,2})^\top = \begin{bmatrix} \|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2^{\frac{1}{2}} \tilde{\mathbf{p}}_{3,1}^{(\tau),\top} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{K}_1^{rpe,3} \tilde{\mathbf{H}}^{rpe,2} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

Then it is further noted that $-(\mathbf{Q}_2^{rpe,3} \tilde{\mathbf{H}}^{rpe,2})^\top \mathbf{K}_2^{rpe,3} \tilde{\mathbf{H}}^{rpe,2} = (\mathbf{Q}_1^{rpe,3} \tilde{\mathbf{H}}^{rpe,2})^\top \mathbf{K}_1^{rpe,3} \tilde{\mathbf{H}}^{rpe,2}$ satisfies

$$\left\| (\mathbf{Q}_1^{rpe,3} \tilde{\mathbf{H}}^{rpe,2})^\top \mathbf{K}_1^{rpe,3} \tilde{\mathbf{H}}^{rpe,2} - \begin{bmatrix} \|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2^{\frac{1}{2}} \tilde{\mathbf{p}}_{3,1}^{(\tau),\top} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right\|_2 \leq C\tau\epsilon \|\mathbf{X}\mathbf{X}^\top\|_2^2.$$

And therefore, combining our construction for \mathbf{V}_m , it is noted that

$$\left\| \sum_{m=1}^2 \mathbf{V}_m \tilde{\mathbf{H}}^{rpe,2} \times \sigma((\mathbf{Q}_1^{rpe,3} \tilde{\mathbf{H}}^{rpe,2})^\top \mathbf{K}_1^{rpe,3} \tilde{\mathbf{H}}^{rpe,2}) - \begin{bmatrix} \mathbf{0}_{(d+1)\times N} \\ -\|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2^{\frac{1}{2}} \tilde{\mathbf{p}}_{3,1}^{(\tau),\top} \\ \mathbf{0} \end{bmatrix} \right\|_2 \leq C\tau\epsilon \|\mathbf{X}\mathbf{X}^\top\|_2^2.$$

Therefore, we can further show that

$$\tilde{\mathbf{H}}^{rpe,3} = \tilde{\mathbf{H}}^{rpe,2} + \sum_{m=1}^2 \mathbf{V}_m^{rpe,3} \tilde{\mathbf{H}}^{rpe,2} \times \sigma((\mathbf{Q}_m^{rpe,3} \tilde{\mathbf{H}}^{rpe,2})^\top \mathbf{K}_m^{rpe,3} \tilde{\mathbf{H}}^{rpe,2})$$

satisfies

$$\left\| \tilde{\mathbf{H}}^{rpe,3} - \begin{bmatrix} \mathbf{X} \\ \mathbf{X}\mathbf{X}^\top - \|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2^{\frac{1}{2}} \tilde{\mathbf{p}}_{3,1}^{(\tau),\top}, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \tilde{\mathbf{p}}_{5,1}, \dots, \tilde{\mathbf{p}}_{5,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_2 \leq C\tau\epsilon \|\mathbf{X}\mathbf{X}^\top\|_2^2.$$

And we can construct another layer to remove the term $\|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2^{\frac{1}{2}} \tilde{\mathbf{p}}_{3,1}^{(\tau),\top}$, which is achieved by

$$\begin{aligned} -\mathbf{V}_1^{rpe,4} = \mathbf{V}_2^{rpe,4} &= \begin{bmatrix} \mathbf{0}_{(4d+1)\times(4d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(4d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{Q}_1^{rpe,4} = -\mathbf{Q}_2^{rpe,4} &= \begin{bmatrix} \mathbf{0}_{(3d+1)\times D} & & \\ \mathbf{0}_{d\times(2d+1)} & I_d & \mathbf{0} \\ & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_1^{rpe,4} = \mathbf{K}_2^{rpe,4} &= \begin{bmatrix} \mathbf{0}_{(3d+1)\times D} & & \\ \mathbf{0}_{d\times(2d+1)} & I_d & \mathbf{0} \\ & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{aligned}$$

Using the above construction, we can further show that

$$\left\| \tilde{\mathbf{H}}^{rpe,4} - \begin{bmatrix} \mathbf{X} \\ \mathbf{X}\mathbf{X}^\top - \|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 \tilde{\mathbf{p}}_{3,1}^{(\tau)} \tilde{\mathbf{p}}_{3,1}^{(\tau)\top}, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \tilde{\mathbf{p}}_{5,1}, \dots, \tilde{\mathbf{p}}_{5,N} \\ \vdots \\ \tilde{\mathbf{p}}_{\ell,1}, \dots, \tilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_2 \leq C\tau\epsilon \|\mathbf{X}\mathbf{X}^\top\|_2^2.$$

And then we proceed to recover the rest of the k principle eigenvectors using similar model architecture given by the ones used by the Power Iterations. For the computation over the τ -th eigenvector, we denote $\tilde{\mathbf{H}}^{pow,\eta,1}$ till $\tilde{\mathbf{H}}^{pow,\eta,\tau}$ to be the intermediate states corresponding to the η -th power iteration. We denote $\tilde{\mathbf{H}}^{rpe,\eta,\tau_0}$ to be the output of η -th removal of principle eigenvector layers for the τ -th eigenvector. Furthermore, we iteratively define

$$\mathbf{A}_1 = \mathbf{X}\mathbf{X}^\top - \|\mathbf{X}\mathbf{X}^\top \tilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 \tilde{\mathbf{p}}_{3,1}^{(\tau)} \tilde{\mathbf{p}}_{3,1}^{(\tau)\top}, \quad \mathbf{A}_{i+1} = \mathbf{A}_i - \|\mathbf{A}_i \tilde{\mathbf{p}}_{3,i}^{(\tau)}\|_2 \tilde{\mathbf{p}}_{3,i}^{(\tau)} \tilde{\mathbf{p}}_{3,i}^{(\tau)\top}, \quad \forall i \in [k].$$

Then, applying the subadditivity of the 2-norm, we can show that

$$\left\| \tilde{\mathbf{H}}^{rpe,4,k} - \begin{bmatrix} \mathbf{X} \\ \mathbf{A}_{k+1}, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \\ \tilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \tilde{\mathbf{p}}_{3,2}^{(\tau)}, \mathbf{0} \\ \vdots \\ \tilde{\mathbf{p}}_{3,k}^{(\tau)}, \mathbf{0} \end{bmatrix} \right\|_2 \leq C\tau k\epsilon \|\mathbf{X}\mathbf{X}^\top\|_2^2.$$

For simplicity, we denote $\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{X} \\ \mathbf{A}_{k+1}, \mathbf{0} \\ \tilde{\mathbf{p}}_{2,1}, \dots, \tilde{\mathbf{p}}_{2,N} \\ \tilde{\mathbf{p}}_{3,1}, \dots, \tilde{\mathbf{p}}_{3,N} \end{bmatrix}$ and $\tilde{\mathbf{P}} = \begin{bmatrix} \tilde{\mathbf{p}}_{3,1}^{(\tau)} \\ \tilde{\mathbf{p}}_{3,2}^{(\tau)} \\ \vdots \\ \tilde{\mathbf{p}}_{3,k}^{(\tau)} \end{bmatrix}$ from here.

4. Finishing Up.

The finishing up phase considers constructing $\tilde{\mathbf{W}}_0$ and $\tilde{\mathbf{W}}_1$ that adjust the final output format. Our construction gives the following

$$\tilde{\mathbf{W}}_0 = [\mathbf{0}, I_{kd}], \quad \tilde{\mathbf{W}}_1 = \begin{bmatrix} 1 \\ \mathbf{0}_{N-1} \end{bmatrix}.$$

And we can show that

$$\left\| \tilde{\mathbf{W}}_0 \tilde{\mathbf{H}}^{rpe,4,k} \tilde{\mathbf{W}}_1 - \begin{bmatrix} \tilde{\mathbf{p}}_{3,1}^{(\tau)} \\ \tilde{\mathbf{p}}_{3,2}^{(\tau)} \\ \vdots \\ \tilde{\mathbf{p}}_{3,k}^{(\tau)} \end{bmatrix} \right\|_2 \leq C\tau k\epsilon \|\mathbf{X}\mathbf{X}^\top\|_2^2.$$

We further use the result given by lemma B.1, denote $a_\eta := \|\mathbf{v}_\eta - \tilde{\mathbf{p}}_{3,\eta}^{(\tau)}\|_2$, $\hat{\lambda}_\eta = \|\mathbf{A}_\eta \tilde{\mathbf{p}}_{3,\eta}^{(\tau)}\|_2$, and $b_\eta := |\lambda_\eta - \hat{\lambda}_\eta|$ for $\eta \in [k]$, we obtain that for all $\eta \geq 1$, given the number of iterations $\tau \geq C \frac{\log(1/\epsilon_0\delta)}{2\epsilon_0}$ where the constant value C depends on d ,

$$a_{\eta+1} \leq \frac{\max_{i \in [\eta]} b_i + \sum_{i=1}^{\eta} 2\lambda_i a_i}{\Delta}, \quad b_{\eta+1} \leq \frac{2\lambda_{\eta+1}}{\Delta} \left(\max_{i \in [\eta]} b_i + \sum_{i=1}^{\eta} 2\lambda_i a_i \right) + \lambda_{\eta+1} \sqrt{2\epsilon_0}.$$

Further note that the starting point is given by $a_1 \leq \sqrt{2\epsilon_0}$, $b_1 \leq \lambda_1 \sqrt{2\epsilon_0}$. Introducing $A_\eta = \sum_{i=1}^\eta 2\lambda_i a_i$, we obtain that $A_{\eta+1} = \sum_{i=1}^{\eta+1} 2\lambda_i a_i = A_\eta + 2\lambda_{\eta+1} a_{\eta+1}$ which alternatively implies that

$$\frac{1}{2\lambda_{\eta+1}}(A_{\eta+1} - A_\eta) \leq \frac{\max_{i \in [\eta]} b_i + A_\eta}{\Delta}, \quad b_{\eta+1} \leq \frac{2\lambda_{\eta+1}}{\Delta} \left(\max_{i \in [\eta]} b_i + A_\eta \right) + \lambda_{\eta+1} \sqrt{2\epsilon_0}.$$

We use the fact $\frac{\lambda_\eta}{\Delta} > 1$ for all $\eta \in [k]$ to show the following

$$A_{\eta+1} + \max_{i \in [\eta+1]} b_i \leq \frac{5\lambda_{\eta+1}}{\Delta} \left(A_\eta + \max_{i \in [\eta]} b_i \right) + \lambda_1 \sqrt{2\epsilon_0}, \quad A_1 + b_1 = 2\lambda_1 \sqrt{2\epsilon_0},$$

which implies that

$$\begin{aligned} A_{\eta+1} + \max_{i \in [\eta+1]} b_i + \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1} &\leq \frac{5\lambda_1}{\Delta} \left(A_\eta + \max_{i \in [\eta]} b_i + \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1} \right), \\ A_{\eta+1} + \max_{i \in [\eta+1]} b_i + \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1} &\leq \left(A_1 + b_1 + \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1} \right) \prod_{i=1}^\eta \left(\frac{5\lambda_{i+1}}{\Delta} \right) \\ &= \lambda_1 \sqrt{2\epsilon_0} \left(2 + \frac{1}{\frac{5\lambda_1}{\Delta} - 1} \right) \prod_{i=1}^\eta \left(\frac{5\lambda_{i+1}}{\Delta} \right). \end{aligned} \quad (4)$$

Therefore, applying the inequality given by equation 4 we can show that, for $\eta \leq k$, we have for all $\eta \in [k-1]$,

$$\begin{aligned} a_{\eta+1} &\leq \frac{1}{\Delta} \left(\lambda_1 \sqrt{2\epsilon_0} \left(2 + \frac{1}{\frac{5\lambda_1}{\Delta} - 1} \right) \prod_{i=1}^\eta \left(\frac{5\lambda_{i+1}}{\Delta} \right) - \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1} \right), \\ b_{\eta+1} &\leq \frac{2\lambda_\eta \lambda_1 \sqrt{2\epsilon_0}}{\Delta} \left(2 + \frac{1}{\frac{5\lambda_1}{\Delta} - 1} \right) \prod_{i=1}^\eta \left(\frac{5\lambda_{i+1}}{\Delta} \right) + \lambda_{\eta+1} \sqrt{2\epsilon_0}. \end{aligned}$$

Therefore collecting pieces, we conclude that there exists a transformer with number of layers $2\tau + 4k + 1$ and number of heads $M \leq \lambda_1^{\frac{d}{C(d) \epsilon^2}}$ such that the final output $\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_k$ given by the Transformer model satisfy $\forall \eta \in [k-1]$,

$$\|\hat{\mathbf{v}}_{\eta+1} - \mathbf{v}_{\eta+1}\|_2 \leq C\tau \epsilon \lambda_1^2 + \frac{1}{\Delta} \left(\lambda_1 \sqrt{2\epsilon_0} \left(2 + \frac{1}{\frac{5\lambda_1}{\Delta} - 1} \right) \prod_{i=1}^\eta \left(\frac{5\lambda_{i+1}}{\Delta} \right) - \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1} \right).$$

And the rest of the result directly follows. \square

B.3 PROOF OF LEMMA 3.1

Proof. To prove the above result, we consider two events $A_1 = \left\{ \|\mathbf{y}\|_2 \geq \sqrt{\frac{1}{\epsilon}} \right\}$, $A_2 = \left\{ |\mathbf{y}^\top \mathbf{v}| \leq \sqrt{\epsilon} \right\}$, then we can show that

$$\left\{ |\mathbf{v}^\top \mathbf{x}| \leq \frac{1}{\sqrt{\epsilon}} \right\} \subset A_1 \cup A_2 \quad \Rightarrow \quad \mathbb{P} \left(|\mathbf{v}^\top \mathbf{x}| \leq \frac{1}{\sqrt{\epsilon}} \right) \leq \mathbb{P}(A_1) + \mathbb{P}(A_2).$$

And we use the tail bound for Chi-square given by (Laurent & Massart, 2000) to obtain that as $\epsilon < d^{-1}$,

$$\mathbb{P}(A_1) = \mathbb{P} \left(\|\mathbf{y}\|_2^2 \geq \epsilon^{-1} \right) \leq \exp(-C\epsilon^{-1}).$$

And similarly, consider the event A_2 , note that $\mathbf{y}^\top \mathbf{v} \sim N(0, 1)$, we use the cdf of the folded normal distribution to obtain that

$$\mathbb{P}(A_2) = \mathbb{P} \left(|\mathbf{v}^\top \mathbf{y}| \leq \sqrt{\epsilon} \right) = \text{erf} \left(\frac{\sqrt{\epsilon}}{\sqrt{2}} \right) = \frac{2}{\sqrt{\pi}} \left(\sqrt{\epsilon} - \frac{(\sqrt{\epsilon})^3}{3} + \frac{(\sqrt{\epsilon})^5}{10} - \frac{(\sqrt{\epsilon})^7}{42} \right) \leq \frac{\sqrt{\epsilon}}{\sqrt{\pi}}.$$

Then we obtain that

$$\mathbb{P}\left(|\mathbf{v}^\top \mathbf{x}| \leq \sqrt{\epsilon}\right) \leq \frac{\sqrt{\epsilon}}{\sqrt{\pi}} + \exp(-C\epsilon^{-1}).$$

Consider in total of k independent random vectors $\mathbf{x}_1, \dots, \mathbf{x}_k$, and arbitrary k vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$, we can show that

$$\mathbb{P}\left(\exists i \text{ such that } \mathbf{x}_i^\top \mathbf{v}_i \leq \epsilon\right) \leq k\mathbb{P}\left(\mathbf{x}_1^\top \mathbf{v}_1 \leq \epsilon\right) \leq \frac{k\sqrt{\epsilon}}{\sqrt{\pi}} + k \exp(-C\epsilon^{-1}).$$

□

B.4 PROOF OF LEMMA B.1

Lemma B.1. *Assume that the correlation matrix $\mathbf{X}\mathbf{X}^\top$ has eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_k$. Assume that the eigenvectors are given by $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ and the eigenvalues satisfy $\inf_{i \neq j} |\lambda_i - \lambda_j| = \Delta$. Then, given that the estimate for the first τ eigenvectors satisfy $\mathbf{v}_i^\top \hat{\mathbf{v}}_i \geq 1 - \epsilon_i$ and the eigenvalues satisfy $|\lambda_i - \hat{\lambda}_i| \leq \delta_i$, the principle eigenvector of $\mathbf{X}\mathbf{X}^\top - \sum_{i=1}^{\tau} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top$ denoted by $\tilde{\mathbf{v}}_{\tau+1}$ satisfies*

$$\|\tilde{\mathbf{v}}_{\tau+1} - \mathbf{v}_{\tau+1}\|_2 \leq \frac{\max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8\lambda_i \sqrt{\epsilon_i}}}{\Delta}.$$

Alternatively, we can also show that the eigenvector $\hat{\mathbf{v}}_{\tau+1}$ returned by power method with $k = \frac{\log(1/\epsilon_0 \delta)}{2\epsilon_0}$ that is initialized by satisfies

$$\hat{\mathbf{v}}_{\tau+1}^\top \mathbf{v}_{\tau+1} \geq 1 - \epsilon_{\tau+1} := 1 - \frac{1}{2} \left(\frac{\max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8\lambda_i \sqrt{\epsilon_i}}}{\Delta} + \sqrt{2\epsilon_0} \right)^2,$$

Proof. Our proof is given by inductive arguments. Consider our obtained estimates $\{\hat{\mathbf{v}}_i\}_{i \in [k]}$ for the eigenvectors $\{\mathbf{v}_i\}_{i \in [k]}$ satisfy

$$\mathbf{v}_i^\top \hat{\mathbf{v}}_i \geq 1 - \epsilon_i \quad \forall i \in [\tau], \quad |\lambda_i - \hat{\lambda}_i| \leq \delta_i.$$

We note that for the eigenvectors, we have for a vector \mathbf{v}_0 ,

$$\begin{aligned} \|\mathbf{v}_i \mathbf{v}_i^\top - \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top\|_2 &= \sup_{\mathbf{v}_0 \in \mathbb{S}^{d-1}} \mathbf{v}_0^\top (\mathbf{v}_i \mathbf{v}_i^\top - \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top) \mathbf{v}_0 = \sup_{\mathbf{v}_0 \in \mathbb{S}^{d-1}} (\mathbf{v}_0^\top \mathbf{v}_i)^2 - (\mathbf{v}_0^\top \hat{\mathbf{v}}_i)^2 \\ &= \sup_{\mathbf{v}_0 \in \mathbb{S}^{d-1}} (\mathbf{v}_0^\top (\mathbf{v}_i - \hat{\mathbf{v}}_i)) (\mathbf{v}_0^\top (\mathbf{v}_i + \hat{\mathbf{v}}_i)) \\ &\leq 2\|\mathbf{v}_i - \hat{\mathbf{v}}_i\|_2 = 2\sqrt{\|\mathbf{v}_i - \hat{\mathbf{v}}_i\|_2^2} = 2\sqrt{\|\mathbf{v}_i\|_2^2 + \|\hat{\mathbf{v}}_i\|_2^2 - 2\mathbf{v}_i^\top \hat{\mathbf{v}}_i} = 2\sqrt{2\epsilon_i}. \end{aligned}$$

Then, we can show by the subadditivity of the spectral norm,

$$\begin{aligned} \left\| \mathbf{X}\mathbf{X}^\top - \sum_{i=1}^{\tau} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 &= \left\| \sum_{i=1}^k \lambda_i \mathbf{v}_i \mathbf{v}_i^\top - \sum_{i=1}^{\tau} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 \\ &\leq \left\| \sum_{i=1}^k \lambda_i \mathbf{v}_i \mathbf{v}_i^\top - \sum_{i=1}^{\tau} \lambda_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 + \left\| \sum_{i=1}^{\tau} \delta_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 \\ &\leq \left\| \sum_{i=\tau+1}^k \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \right\|_2 + \left\| \sum_{i=1}^{\tau} \delta_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 + \left\| \sum_{i=1}^{\tau} \lambda_i \mathbf{v}_i \mathbf{v}_i^\top - \sum_{i=1}^{\tau} \lambda_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 \\ &\leq \lambda_{\tau+1} + \max_{i \in [\tau]} \delta_i + \left\| \sum_{i=1}^{\tau} \lambda_i (\mathbf{v}_i \mathbf{v}_i^\top - \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top) \right\|_2 \\ &\leq \lambda_{\tau+1} + \max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \lambda_i \left\| \mathbf{v}_i \mathbf{v}_i^\top - \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 \\ &\leq \lambda_{\tau+1} + \max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8\lambda_i \sqrt{\epsilon_i}}. \end{aligned}$$

By similar argument, we can also show that

$$\left\| \mathbf{X} \mathbf{X}^\top - \sum_{i=1}^{\tau} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 \geq \lambda_{\tau+1} - \max_{i \in [\tau]} \delta_i - \sum_{i=1}^{\tau} \sqrt{8} \lambda_i \sqrt{\epsilon_i}.$$

To study the convergence of the eigenvectors, we notice that by Davis-Kahan Theorem by (Yu et al., 2015) we can show that the principle eigenvector $\tilde{\mathbf{v}}_{\tau+1} = \arg \max_{\mathbf{v} \in \mathbb{S}^{d-1}}$ satisfies

$$\|\tilde{\mathbf{v}}_{\tau+1} - \mathbf{v}_{\tau+1}\|_2 \leq \frac{\left| \left\| \mathbf{X} \mathbf{X}^\top - \sum_{i=1}^{\tau} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\| - \lambda_{\tau+1} \right|}{\max\{|\lambda_{\tau+1} - \lambda_{\tau}|, |\lambda_{\tau-1} - \lambda_{\tau}|\}} \leq \frac{\max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8} \lambda_i \sqrt{\epsilon_i}}{\Delta}.$$

Consider the eigenvector returned by the power method, we can show by the subadditivity of L_2 norm, we obtain that $\|\tilde{\mathbf{v}}_{\tau+1} - \mathbf{v}_{\tau+1}\|_2 \leq \|\tilde{\mathbf{v}}_{\tau+1} - \hat{\mathbf{v}}_{\tau+1}\|_2 + \|\hat{\mathbf{v}}_{\tau+1} - \mathbf{v}_{\tau+1}\|_2 \leq \frac{\max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8} \lambda_i \sqrt{\epsilon_i}}{\Delta} + \sqrt{2\epsilon_0}$

$$\begin{aligned} \hat{\mathbf{v}}_{\tau+1}^\top \mathbf{v}_{\tau+1} &= \frac{1}{2} (2 - \|\mathbf{v}_{\tau+1} - \hat{\mathbf{v}}_{\tau+1}\|_2^2) \geq \frac{1}{2} (2 - (\|\tilde{\mathbf{v}}_{\tau+1} - \hat{\mathbf{v}}_{\tau+1}\|_2 + \|\mathbf{v}_{\tau+1} - \tilde{\mathbf{v}}_{\tau+1}\|_2)^2) \\ &= 1 - \frac{1}{2} \left(\frac{\max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8} \lambda_i \sqrt{\epsilon_i}}{\Delta} + \sqrt{2\epsilon_0} \right)^2. \end{aligned}$$

Moreover, consider the estimate of the eigenvalue, we have

$$\begin{aligned} &\left\| \left(\mathbf{X} \mathbf{X}^\top - \sum_{i=1}^{\tau} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right) \hat{\mathbf{v}}_{\tau+1} \right\|_2 \\ &\leq \left\| \left(\mathbf{X} \mathbf{X}^\top - \sum_{i=1}^{\tau} \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \right) \hat{\mathbf{v}}_{\tau+1} \right\|_2 + \left\| \sum_{i=1}^{\tau} \lambda_i \mathbf{v}_i \mathbf{v}_i^\top - \sum_{i=1}^{\tau} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 \\ &\leq \left\| \left(\mathbf{X} \mathbf{X}^\top - \sum_{i=1}^{\tau} \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \right) \hat{\mathbf{v}}_{\tau+1} \right\|_2 + \left\| \sum_{i=1}^{\tau} \lambda_i \mathbf{v}_i \mathbf{v}_i^\top - \sum_{i=1}^{\tau} \lambda_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 + \left\| \sum_{i=1}^{\tau} (\lambda_i - \hat{\lambda}_i) \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right\|_2 \\ &\leq \left\| \left(\mathbf{X} \mathbf{X}^\top - \sum_{i=1}^{\tau} \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \right) \mathbf{v}_{\tau+1} \right\|_2 + \left\| \left(\mathbf{X} \mathbf{X}^\top - \sum_{i=1}^{\tau} \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \right) \right\|_2 \|\hat{\mathbf{v}}_{\tau+1} - \mathbf{v}_{\tau+1}\|_2 \\ &\quad + \max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8} \lambda_i \sqrt{\epsilon_i} \\ &= \lambda_{\tau+1} + \lambda_{\tau+1} \left(\frac{\max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8} \lambda_i \sqrt{\epsilon_i}}{\Delta} + \sqrt{2\epsilon_0} \right) + \max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8} \lambda_i \sqrt{\epsilon_i}. \end{aligned}$$

Therefore, by similar arguments, we can show that

$$\left| \left\| \left(\mathbf{X} \mathbf{X}^\top - \sum_{i=1}^{\tau} \hat{\lambda}_i \hat{\mathbf{v}}_i \hat{\mathbf{v}}_i^\top \right) \hat{\mathbf{v}}_{\tau+1} \right\|_2 - \lambda_{\tau+1} \right| \leq \frac{2\lambda_{\tau+1}}{\Delta} \left(\max_{i \in [\tau]} \delta_i + \sum_{i=1}^{\tau} \sqrt{8} \lambda_i \sqrt{\epsilon_i} \right) + \lambda_{\tau+1} \sqrt{2\epsilon_0}.$$

□

B.5 PROOF OF LEMMA B.2

Lemma B.2 (Approximation of norm by sum of Relu activations by Transformer networks). *Assume that there exists a constant C with $\|\mathbf{v}\|_2 \leq C$. There exists a multihead Relu attention layer with number of heads $M < \left(\frac{\bar{R}}{\underline{R}}\right)^d \frac{C(d)}{\epsilon^2} \log(1 + C/\epsilon)$ such that there exists $\{\mathbf{a}_m\}_{m \in [M]} \subset \mathbb{S}^{N-1}$ and $\{c_m\}_{m \in [M]} \subset \mathbb{R}$ where for all \mathbf{v} with $\bar{R} \geq \|\mathbf{v}\|_2 \geq \underline{R}$, we have*

$$\left| \sum_{m=1}^M c_m \sigma(\mathbf{a}_m^\top \mathbf{v}) - \frac{1}{\|\mathbf{v}\|_2} + 1 \right| \leq \epsilon.$$

1296 Similarly, there exists a multihead Relu attention layer with number of heads $M \leq$
 1297 $\bar{R}^{\frac{d}{2}} \frac{C(d)}{\epsilon^2} \log(1 + C/\epsilon)$, a set of vectors $\{\mathbf{b}_m\}_{m \in [M]} \subset \mathbb{S}^{N-1}$ and $\{d_m\}_{m \in [M]} \subset \mathbb{R}$ such that
 1298

$$1299 \left| \sum_{m=1}^M d_m \sigma(\mathbf{b}_m^\top \mathbf{v}) - \|\mathbf{v}\|_2^{1/2} + 1 \right| \leq \epsilon.$$

1300
 1301
 1302 *Proof.* Consider a set $C^d(\bar{R}) := B_\infty^d(\bar{R}) \setminus B_2^d(\underline{R})$, then it is not hard to check that given $\|\mathbf{v}\|_2 > C$
 1303 with some $C(d) > 0$ depending on d such that we have
 1304

$$1305 \sup_{\mathbf{v} \in C^d(\bar{R})} \partial_{v_{j_1}, \dots, v_{j_i} \in [d]} \left(\frac{1}{\|\mathbf{v}\|_2} \right) \leq \frac{C(d)}{\|\mathbf{v}\|_2^d} \leq \frac{C(d)}{\underline{R}^d}.$$

1306
 1307 Therefore, consider the definition 5, we have $C_\ell = \left(\frac{\bar{R}}{\underline{R}}\right)^d C(d)$. Note that by proposition A.1 in
 1308 (Bai et al., 2024) shows that for a function that is (R, C_ℓ) smooth with $R \geq 1$ is $(\epsilon_{approx}, R, M, C)$
 1309 approximable with $M \leq C(d)C_\ell \log(1 + C_\ell/\epsilon_{approx})/\epsilon_{approx}^2$, we complete the proof.
 1310
 1311
 1312

1313 Then we consider the function $\|\mathbf{v}\|_2^{\frac{1}{2}}$, note that
 1314

$$1315 \sup_{\mathbf{v} \in C^d(\bar{R})} \partial_{v_{j_1}, \dots, v_{j_i} \in [d]} \|\mathbf{v}\|_2^{\frac{1}{2}} \leq C \|\mathbf{v}\|_2^{-\frac{1}{2}} \leq C \bar{R}^{-\frac{1}{2}}.$$

1316
 1317 And the rest of the proof follows similarly to the previous step. □
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

C EXPERIMENTAL DETAILS

C.1 SETUP.

We run all our experiments on RTX 2080 Ti GPUs. We use PyTorch to construct our models and training process. We use sklearn for data generation. A training process with 2k steps roughly takes 0.5 hours.

C.2 DATA.

Synthetic Dataset. For each $\mathbf{X}_i \in \mathbb{R}^D$, we sample $Z_i \sim N(0, I) \in \mathbb{R}^D$. We then form $Z = [Z_1, \dots, Z_N]$ and transform it using an invertible matrix $L \sim N(0, I) \in \mathbb{R}^{D \times D}$, yielding the desired training sample \mathbf{X} . To speed up the training process, we set $N < D$ in all our experiment setting. With this design, the rank of the covariance matrix $\mathbf{X}^T \mathbf{X}$ is at most N , meaning there are at least $D - N$ zero eigenvalues. The eigenvectors corresponding to these zero eigenvalues are less meaningful. Thus, to ensure predictions focus on meaningful eigenvectors, we increase N to 10 when predicting multiple eigenvectors. We also adjust the data generation process to ensure the magnitude of eigenvalue across different D to be at a similar level.

Real-world Dataset. For both the MNIST and FMNIST, we first normalize the images to zero mean. Next, we perform SVD to extract the top- D principal components and project the data onto these components, reducing feature dimension to $D = 10, 20$, and use $N = 10, 50$ for eigenvalue and eigenvector prediction respectively. Last, we rescale the resulting matrix to ensure its magnitude is roughly the same level as training data (transformers are trained on synthetic data). The rescaling process is critical to transformers as some images after SVD contains entries large as $7e3$. This will largely degrade transformer’s performance as it changes the input domain by a large margin.

C.3 HYPERPARAMETERS.

We list the hyperparameters in our experiments as below (table 2). We separate the hyperparameters used in predicting (1) eigenvalues and single eigenvectors, and (2) multiple eigenvectors.

Table 2: Hyperparameters for Eigenvalue and Eigenvector Prediction.

parameter	$N = 5$	$N = 10$	$N = 20$
steps (eigenvalue)	20k	20k	20k
steps (eigenvector)	20k	20k	60k
learning rate	1e-3	5e-3	5e-3
Optimizer	Adam	Adam	Adam
batch size	64	64	64
number of layers	3	3	3
hidden dimension	64	64	64
number of heads	2	2	2

C.4 ADDITIONAL EXPERIMENTAL RESULTS

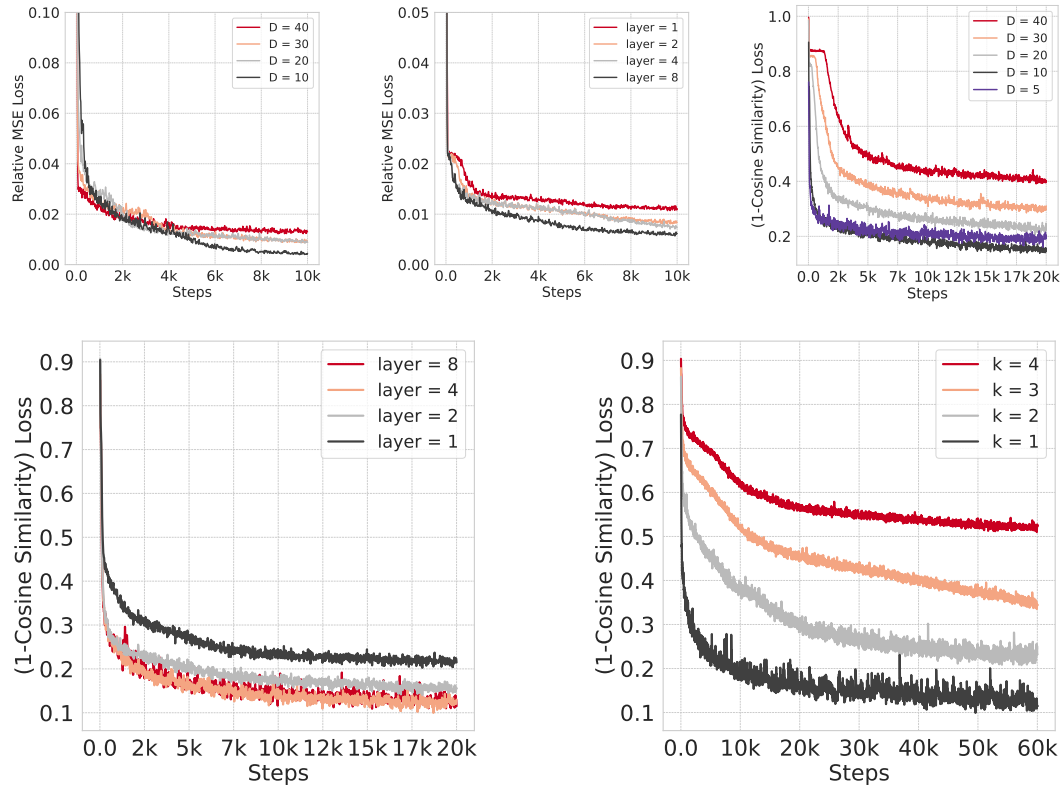
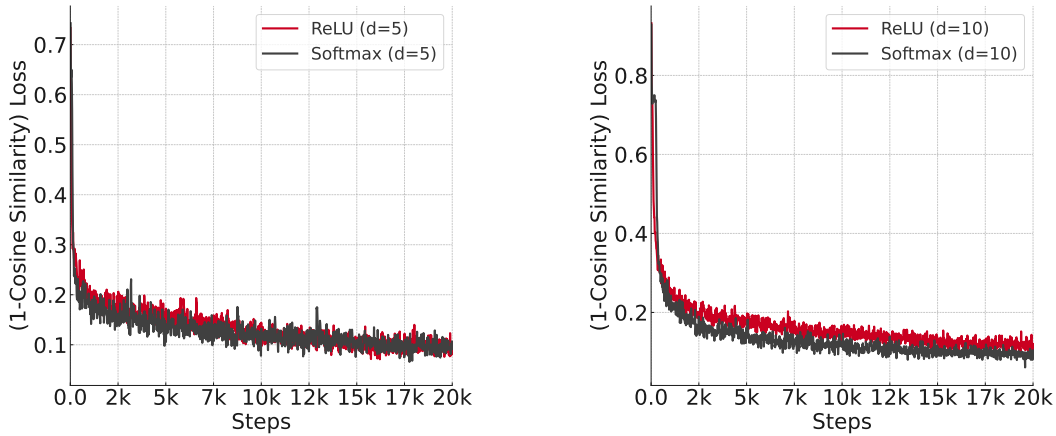


Figure 5: **Convergence Results on Eigenvalue, Eigenvector Prediction with Different Parameters.** (1) Top left: Loss curve on eigenvalue prediction with different size of D (2) Top middle: Loss curve on eigenvalue prediction with different number of layers (3) Top right: Loss curve on eigenvector prediction with different size of D (4) Bottom left: Top right: Loss curve on eigenvector prediction with different number of layers (5) Bottom left: Loss curve on eigenvector prediction with different number of k_{train} . For (1), we observe that smaller D is easier for transformers as they present lower loss. For (2), we see that with more layers, transformers are also capable of predicting eigenvalues more accurately. For (3), transformers also predict eigenvectors better when D is small. For (4), similar to (2), transformers with more layers shows improved performance. For (5), we want to highlight that the loss value is mainly affected by the fact that predicting 3rd or 4th eigenvectors are significantly harder, which contributes to higher loss value.

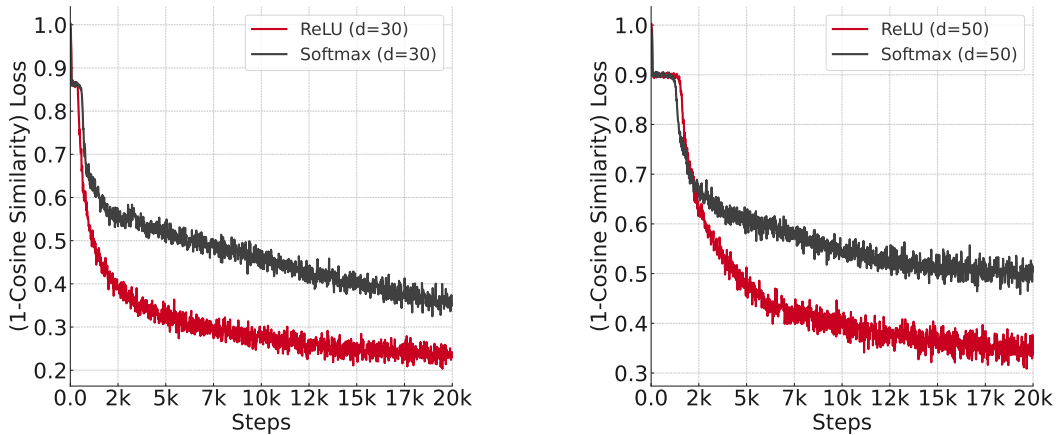
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475



1476
1477
1478
1479
1480
1481

Figure 6: **Loss Curve Comparison between Softmax and ReLU Transformers (Top-1 Eigenvector Prediction)**. *Left: $D = 5$ Right: $D = 10$* We use a 3-layer, 2 head, 64 hidden dimension transformer to predict top-1 eigenvector across all experiments in this figure. An explanation for the superior performance of ReLU transformers is that the normalizing behavior of Softmax can potentially hinder the PCA process.

1482
1483
1484
1485
1486
1487
1488
1489



1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Figure 7: **Loss Curve Comparison between Softmax and ReLU Transformers (Top-1 Eigenvector Prediction)**. *Left: $D = 30$ Right: $D = 50$* We use a 3-layer, 2 head, 64 hidden dimension transformer to predict top-1 eigenvector across all experiments in this figure. We also observe that the performance gap enlarges as D increases, likely because the difference between eigenvectors becomes larger with increasing D , making the normalizing nature of Softmax unsuitable for PCA.