
Supplementary Materials of "CosNet: A Generalized Spectral Kernel Network"

Yanfang Xue^{1,2}, Pengfei Fang^{1,2}, Jinyue Tian^{1,2}, Shipeng Zhu^{1,2}, Hui Xue^{1,2*}

¹School of Computer Science and Engineering, Southeast University, Nanjing, 210096, China

²Key Laboratory of New Generation Artificial Intelligence Technology and Its Interdisciplinary Applications (Southeast University), Ministry of Education, China

{230218795, fangpengfei, 220222083, shipengzhu, hxue}@seu.edu.cn

In the supplementary material, we provide:

- Detailed proof of theoretical results (mentioned in the CosNet analysis part of the main paper).
- More explanation of the initialization scheme (mentioned in the complex-valued spectral kernel network part of the main paper)
- More details of the experiment and ablation studies (mentioned in the ablation study part of the main paper).

1 Proof of Theoretical Results

In this section, we collect the proof of the error bound omitted from the main paper. Before the proof, we introduce the necessary preliminary knowledge and notations.

Foundation setting The proposed network with L layers consists of two modules, i.e. SKMG module and CSKE module. The SKMG module includes 1 layer and CSKE model includes $L - 1$ layers. The input $\mathbf{X} \in \mathbb{R}^{d^x \times n}$ has n samples and each sample x_i is d^x -dimension.

For the SKMG module, the weight matrix of the first linear transformation is denoted as

$$\Omega^1 = \begin{bmatrix} \omega \\ \omega' \\ \omega \\ \omega' \end{bmatrix} \in \mathbb{R}^{4d^0 \times d^x}, \text{ where } \omega, \omega' \in \mathbb{R}^{d^0 \times d^x}. \text{ The second weight matrix is } \Omega^2 = \begin{bmatrix} \frac{1}{\sqrt{4d^0}} \mathbf{I}_{d^0} & \frac{1}{\sqrt{4d^0}} \mathbf{I}_{d^0} & 0 * \mathbf{J} & 0 * \mathbf{J} \\ 0 * \mathbf{J} & 0 * \mathbf{J} & \frac{1}{\sqrt{4d^0}} \mathbf{I}_{d^0} & \frac{1}{\sqrt{4d^0}} \mathbf{I}_{d^0} \end{bmatrix} \in \mathbb{R}^{2d^0 \times 4d^0}, \text{ where } \mathbf{I}_{d^0} \in \mathbb{R}^{d^0 \times d^0} \text{ is the identity matrix and } \mathbf{J} \in \mathbb{R}^{d^0 \times d^0} \text{ is a matrix of ones. The output of this module is denoted as } \mathbf{X}^0 \in \mathbb{R}^{d^0 \times n}.$$

For the CSKE module, the input is denoted as \mathbf{X}^{l-1} with d^{l-1} -dimension and the output is denoted as \mathbf{X}^l in the l^{th} layer. Note that: $\mathbf{X}^1 = \mathbf{X}^0$, $d^1 = d^0$, which means the output of the SKMG module is used as the input of CSKE module. The weight matrix of l^{th} ($2 \leq l \leq L$) layer is denoted as

$$\mathbf{W}^l = \begin{bmatrix} \cos(\mathbf{A}^l) & -\sin(\mathbf{A}^l) \\ \sin(\mathbf{A}^l) & \cos(\mathbf{A}^l) \end{bmatrix} \in \mathbb{R}^{2d^l \times 2d^{l-1}},$$

where $\mathbf{A}^l \in \mathbb{R}^{d^l \times d^{l-1}}$.

Definition 1. (ϵ -net) A 's subset \tilde{A} is an ϵ -net of A under the metric d if for any $a \in A$ there exists $\tilde{a} \in \tilde{A}$ that $d(a, \tilde{a}) \leq \epsilon$.

*Corresponding author

Definition 2. (covering number) The covering number $N_d(A, \epsilon)$ is the size of the smallest ϵ -net of A .

Lemma 1. (Maury's sparsification lemma Pisier). Fix Hilbert space \mathcal{H} with norm $\|\cdot\|$. Let $\mathbf{U} \in \mathcal{H}$ be given with representation $\mathbf{U} = \sum_{i=1}^d \alpha_i \mathbf{V}_i$, where $\mathbf{V}_i \in \mathcal{H}$ and $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_d] \in \mathbb{R}_{\geq 0}^d \setminus \{0\}$. Then for any positive integer k , there exists a choice of non-negative integers (k_1, \dots, k_d) , $\sum_{i=1}^d k_i = k$, such that

$$\|\mathbf{U} - \frac{\|\boldsymbol{\alpha}\|_1}{k} \sum_{i=1}^d k_i \mathbf{V}_i\|^2 \leq \frac{\|\boldsymbol{\alpha}\|_1}{k} \sum_{i=1}^d \alpha_i \|\mathbf{V}_i\|^2 \leq \frac{\|\boldsymbol{\alpha}\|_1^2}{k} \max_i \|\mathbf{V}_i\|^2.$$

Proof. Set $\beta = \|\boldsymbol{\alpha}\|_1$, and let $(\mathbf{B}_1, \mathbf{B}_2, \dots, \mathbf{B}_k)$ denote k iid random variables where $P(\mathbf{B}_1 = \beta \mathbf{V}_i) = \frac{\alpha_i}{\beta}$. Define $\mathbf{B} = \frac{1}{k} \sum_{i=1}^k \mathbf{B}_i$, whereby

$$\mathbb{E}[\mathbf{B}] = \mathbb{E}\left[\frac{1}{k} \sum_{i=1}^k \mathbf{B}_i\right] = \frac{1}{k} \mathbb{E}[k\mathbf{B}_1] = \mathbb{E}[\mathbf{B}_1] = \sum_{i=1}^d (\beta \mathbf{V}_i) \frac{\alpha_i}{\beta} = \sum_{i=1}^d \alpha_i \mathbf{V}_i = \mathbf{U}.$$

Consequently,

$$\begin{aligned} \mathbb{E}[\|\mathbf{U} - \mathbf{B}\|^2] &= \frac{1}{k^2} \mathbb{E}\left[\left\|\sum_{i=1}^d (\mathbf{U} - \mathbf{B}_i)\right\|^2\right] = \frac{1}{k^2} \mathbb{E}\left[\sum_i \|\mathbf{U} - \mathbf{B}_i\|^2 + \sum_{i \neq j} \langle \mathbf{U} - \mathbf{B}_i, \mathbf{U} - \mathbf{B}_j \rangle\right] \\ &= \frac{1}{k} \mathbb{E}[\|\mathbf{U} - \mathbf{B}_1\|^2] = \frac{1}{k} (\mathbb{E}[\|\mathbf{B}_1\|^2] - \|\mathbf{U}\|^2) \leq \frac{1}{k} \mathbb{E}[\|\mathbf{B}_1\|^2] \\ &= \frac{1}{k} \sum_{i=1}^d \frac{\alpha_i}{\beta} \|\beta \mathbf{V}_i\|^2 = \frac{\beta}{k} \sum_{i=1}^d \alpha_i \|\mathbf{V}_i\|^2 \\ &\leq \frac{\beta^2}{k} \max_i \|\mathbf{V}_i\|^2. \end{aligned}$$

By the probabilistic method, there exists integers $(j_1, \dots, j_k) \in \{1, \dots, d\}^k$ and an assignment $\widetilde{\mathbf{B}}_i = \beta \mathbf{V}_{j_i}$ and $\widetilde{\mathbf{B}} = \frac{1}{k} \sum_{i=1}^k \widetilde{\mathbf{B}}_i$ such that

$$\|\mathbf{U} - \widetilde{\mathbf{B}}\|^2 \leq \mathbb{E}[\|\mathbf{U} - \mathbf{B}\|^2].$$

The result now follows by defining integers (k_1, \dots, k_d) according to $k_i = \sum_{l=1}^k \mathbb{I}_{[j_l=i]}$, where $\mathbb{I}(\cdot)$ denotes the indicator function. \square

Corollary 1. Fix Hilbert space \mathcal{H} with norm $\|\cdot\|$. Let $\mathbf{U} \in \mathcal{H}$ be given with representation $\mathbf{U} = \sum_{i=1}^d \alpha_i \mathbf{V}_i$, where $\mathbf{V}_i \in \mathcal{H}$ and $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_d] \in \mathbb{R}_{\geq 0}^d \setminus \{0\}$. Then for any positive integer k and for any $m \geq \alpha_i (\forall i)$, there exists a choice of non-negative integers (k_1, \dots, k_d) , $\sum_{i=1}^d k_i = k$ such that

$$\|\mathbf{U} - \frac{m}{k} \sum_{i=1}^d k_i \mathbf{V}_i\|^2 \leq \frac{m}{k} \sum_{i=1}^d \alpha_i \|\mathbf{V}_i\|^2 \leq \frac{m^2}{k} \max_i \|\mathbf{V}_i\|^2.$$

Proof. Following the proof of Lemma 1 with $\beta = m$, the result is trivial. \square

Corollary 2. For set $A = \{\sum_{i=1}^d \alpha_i \mathbf{V}_i \mid \boldsymbol{\alpha} \in \mathbb{R}_{\geq 0}^d \setminus \{0\}, \|\boldsymbol{\alpha}\|_1 \leq \tilde{\alpha}\} \in \text{conv}(\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_d)$, its covering number satisfies $N_d(A, \epsilon) \leq d^k$ and $\ln(N_d(A, \epsilon)) \leq k \ln(d)$, where k is an integer and $k \geq \frac{\tilde{\alpha}^2}{\epsilon^2} \max_i \|\mathbf{V}_i\|^2$.

Proof. For any $\mathbf{U} \in A$, there exist α satisfying $\mathbf{U} = \sum_{i=1}^d \alpha_i \mathbf{V}_i$. Since $\|\alpha\|_1 \leq \tilde{\alpha}$, by Corollary 1, for a fixed positive integer k satisfying $k \geq \frac{\tilde{\alpha}^2}{\epsilon^2} \max_i \|\mathbf{V}_i\|^2$, there exists a choice of non-negative integers (k_1, \dots, k_d) , $\sum_{i=1}^d k_i = k$, such that

$$\|\mathbf{U} - \frac{\tilde{\alpha}}{k} \sum_{i=1}^d k_i \mathbf{V}_i\|^2 \leq \frac{\tilde{\alpha}}{k} \sum_{i=1}^d \alpha_i \|\mathbf{V}_i\|^2 \leq \frac{\tilde{\alpha}_1^2}{k} \max_i \|\mathbf{V}_i\|^2 \leq \epsilon^2.$$

By Definition 1, $\{\frac{\tilde{\alpha}}{k} \sum_{i=1}^d k_i \mathbf{V}_i | k_i (i = 1, \dots, d), \sum_{i=1}^d k_i = k\}$ is a ϵ -net of A . With k_i being non-negative integers and $\sum_{i=1}^d k_i = k$, the cardinality of this set satisfies

$$|\{\frac{\tilde{\alpha}}{k} \sum_{i=1}^d k_i \mathbf{V}_i\}| \leq d^k.$$

By Definition 2, $N_d(A, \epsilon) \leq d^k$ and $\ln(N_d(A, \epsilon)) \leq k \ln(d)$, where k is an integer and $k \geq \frac{\tilde{\alpha}^2}{\epsilon^2} \max_i \|\mathbf{V}_i\|^2$. □

Lemma 2. For the non-empty set $A \subset \mathbb{R}^n$, define a projection $f(\mathbf{a}) = \begin{bmatrix} \sigma(a_1) \\ \sigma(a_2) \\ \vdots \\ \sigma(a_n) \end{bmatrix}$, where $\mathbf{a} \in A$,

$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$, σ is a function. And the distance metric of this space is defined by p -norm. If σ is

l -Lipschitz and \tilde{A} is a ϵ -net of A , then $f(\tilde{A})$ is a ϵ -net of $f(A)$.

Proof. Denote the distance metric in the space as $d(\mathbf{a}, \mathbf{a}') = \|\mathbf{a} - \mathbf{a}'\|_p$.

Since \tilde{A} is a ϵ -net of A , then for any $\mathbf{a} \in A$, there exists $\tilde{\mathbf{a}} \in \tilde{A}$ satisfying that $d(\mathbf{a}, \tilde{\mathbf{a}}) \leq \epsilon$.

For any $\mathbf{f} \in f(A)$, there exists $\mathbf{a} \in A$ so that $f(\mathbf{a}) = \mathbf{f}$ and corresponding $\tilde{\mathbf{a}}$ that $\tilde{\mathbf{f}} = f(\tilde{\mathbf{a}}) \in f(\tilde{A})$ and $d(\mathbf{a}, \tilde{\mathbf{a}}) \leq \epsilon$. Then,

$$d(\mathbf{f}, \tilde{\mathbf{f}}) = d(f(\mathbf{a}), f(\tilde{\mathbf{a}})) = (\sum_{i=1}^n (\sigma(a_i) - \sigma(\tilde{a}_i))^p)^{1/p} = l d(\mathbf{a}, \tilde{\mathbf{a}}) \leq l\epsilon \triangleq \epsilon'.$$

□

Corollary 3. Following the definitions in Lemma 2, if there exists a covering number for set A , then $N_d(f(A), \epsilon') \leq N_d(A, \epsilon)$.

Proof. Assume \tilde{A} is an ϵ -net of A .

By Lemma 2, $f(\tilde{A})$ is an ϵ -net of $f(A)$.

By Definition 2, since the cardinality: $|f(\tilde{A})| \leq |\tilde{A}|$, so $N_d(f(A), \epsilon') \leq N_d(A, \epsilon)$. □

Theorem 1. Denote the covering number of set S as $N_d(S, \epsilon)$. $\mathbf{X} \in \mathbb{R}^{d^x \times n}$ is the input of n samples and each sample is d^x -dimensioned. $\mathbf{X}^l \in \mathbb{R}^{d^l \times n}$ is the input of layer l ($l > 1$) and \mathbf{W}^l is the weight matrix of layer l ($l \geq 2$). The other notations remain the same as mentioned above. For different layers, their covering numbers satisfy that

1. In the first layer (i.e., the SKMG module), $N_d(\Omega^1 \mathbf{X}, \epsilon) \leq (4d^0 d^x)^k$, where $k \geq \frac{\|\Omega_{i,j}^1\|_1^2}{\epsilon^2} \max_{i,j} \|\mathbf{x}_{i,j}\|^2$.

2. In layer l ($l > 1$) (i.e., the CSKE module), $N_d(\mathbf{W}^l \mathbf{X}^{l-1}, \epsilon) \leq (2d^l d^{l-1} + 1)^k$, where $k \geq \frac{\|\mathbf{A}_{ij}\|_1^2}{\epsilon^2} \frac{\pi}{2} d^l \|\mathbf{X}^{l-1}\|_1^2$.

Proof. Part 1

We denote $\mathbf{P}_{kq} \in \mathbb{R}^{d^0 \times d^x}$ as a matrix where only the element of row k column q is 1, other element is 0. Trivially, any $\mathbf{\Omega}^1 \in \mathbb{R}^{4d^0 \times d^x}$ can be written as:

$$\mathbf{\Omega}^1 = \alpha \sum_{1 \leq k \leq d^0, 1 \leq q \leq d^x} \pm \begin{bmatrix} \mathbf{P}_{kq} \\ 0 * \mathbf{J} \\ \mathbf{P}_{kq} \\ 0 * \mathbf{J} \end{bmatrix} \pm \begin{bmatrix} 0 * \mathbf{J} \\ \mathbf{P}_{kq} \\ 0 * \mathbf{J} \\ \mathbf{P}_{kq} \end{bmatrix}.$$

For simplicity, denote the $4d^0 d^x$ components shown above as $\{\mathbf{V}_i\}_{i=1}^{4d^0 d^x}$, then $\mathbf{\Omega}^1 = \sum_{i=1}^{4d^0 d^x} \alpha_i \mathbf{V}_i$, where $\alpha \in \mathbb{R}_{\geq 0}^d \setminus \{0\}$. And $\mathbf{\Omega}^1 \mathbf{X} = \sum_{i=1}^{4d^0 d^x} \alpha_i \mathbf{V}_i \mathbf{X} = \sum_{i=1}^{4d^0 d^x} \alpha_i \mathbf{V}'_i$, where $\mathbf{V}'_i = \mathbf{V}_i \mathbf{X}$.

By Corollary 2, $N_d(\mathbf{\Omega} \mathbf{X}, \epsilon) \leq (4d^0 d^x)^k$, where $k \geq \frac{\|\mathbf{\Omega}_{ij}\|_1^2}{\epsilon^2} \max_{i,j} \|x_{ij}\|^2$.

Part 2

To analyze the covering number of $\mathbf{W}^l \mathbf{X}^{l-1}$, first begin with the covering number of $\mathbf{B}^l \mathbf{X}^{l-1}$, and \mathbf{B}^l is defined as

$$\mathbf{B}^l = \begin{bmatrix} \mathbf{A}^l & \mathbf{A}^l + \frac{\pi}{2} * \mathbf{J} \\ \mathbf{A}^l - \frac{\pi}{2} * \mathbf{J} & \mathbf{A}^l \end{bmatrix}.$$

We rewrite \mathbf{B}^l as

$$\mathbf{B}^l = \alpha \left(\sum_{1 \leq k \leq d^l, 1 \leq q \leq d^{l-1}} \pm \begin{bmatrix} \mathbf{P}_{kq} & \mathbf{P}_{kq} \\ \mathbf{P}_{kq} & \mathbf{P}_{kq} \end{bmatrix} \right) + 1 \cdot \begin{bmatrix} 0 * \mathbf{J} & \frac{\pi}{2} * \mathbf{J} \\ -\frac{\pi}{2} * \mathbf{J} & 0 * \mathbf{J} \end{bmatrix},$$

where \mathbf{P}_{kq} is in $\mathbb{R}^{d^l \times d^{l-1}}$.

For simplicity, denote the $2d^l d^{l-1} + 1$ components shown above as $\{\mathbf{V}_i\}_{i=1}^{2d^l d^{l-1} + 1}$, then $\mathbf{B}^l = \sum_{i=1}^{2d^l d^{l-1} + 1} \alpha_i \mathbf{V}_i$, where $\alpha \in \mathbb{R}_{\geq 0}^d \setminus \{0\}$. And $\mathbf{B}^l \mathbf{X}^l = \sum_{i=1}^{2d^l d^{l-1} + 1} \alpha_i \mathbf{V}_i \mathbf{X}^{l-1} = \sum_{i=1}^{2d^l d^{l-1} + 1} \alpha_i \mathbf{V}'_i$, where $\mathbf{V}'_i = \mathbf{V}_i \mathbf{X}^l$.

If \mathbf{V}_i is in the form of $\begin{bmatrix} \mathbf{P}_{kq} & \mathbf{P}_{kq} \\ \mathbf{P}_{kq} & \mathbf{P}_{kq} \end{bmatrix}$, then

$$\|\mathbf{V}'_i\|_1 \leq \max_{1 \leq k \leq d^{l-1}} 2 \sum_{j=1}^n |\mathbf{X}_{kj}^{l-1} + \mathbf{X}_{k+d^{l-1},j}^{l-1}|.$$

If \mathbf{V}_i is in the form of $\begin{bmatrix} 0 * \mathbf{J} & \frac{\pi}{2} * \mathbf{J} \\ -\frac{\pi}{2} * \mathbf{J} & 0 * \mathbf{J} \end{bmatrix}$, then

$$\|\mathbf{V}'_i\|_1 \leq \frac{\pi}{2} d^l \sum_{j=1}^n \left[\left| \sum_{k=1}^{d^{l-1}} \mathbf{X}_{kj}^{l-1} \right| + \left| \sum_{k=d^{l-1}+1}^{2d^{l-1}} \mathbf{X}_{kj}^{l-1} \right| \right] \leq \frac{\pi}{2} d^l \|\mathbf{X}^{l-1}\|_1.$$

In summary, since $\frac{\pi}{2} d^l \|\mathbf{X}^{l-1}\|_1 > \max_{1 \leq k \leq d^{l-1}} 2 \sum_{j=1}^n |\mathbf{X}_{kj}^{l-1} + \mathbf{X}_{k+d^{l-1},j}^{l-1}|$, by Corollary 2, $N_d(\mathbf{B}^l \mathbf{X}^{l-1}, \epsilon) \leq (2d^l d^{l-1} + 1)^k$, where $k \geq \frac{\|\mathbf{A}_{ij}\|_1^2}{\epsilon^2} \frac{\pi}{2} d^l \|\mathbf{X}^{l-1}\|_1^2$.

Since $\mathbf{W}^l = \cos(\mathbf{B}^l)$, by Corollary 3, $N_d(\mathbf{B}^l \mathbf{X}^{l-1}, \epsilon) \leq (2d^l d^{l-1} + 1)^k$, where $k \geq \frac{\|\mathbf{A}_{ij}\|_1^2}{\epsilon^2} \frac{\pi}{2} d^l \|\mathbf{X}^{l-1}\|_1^2$. \square

Then we introduce two important lemmas about Rademacher complexity and empirical Rademacher complexity respectively.

Lemma 3. Mohri et al. [2018] Let $F|_S$ be a real-valued function class taking values in $[0, 1]$ given the dataset S , and assume that $0 \in F|_S$. Then the Rademacher complexity given the dataset S satisfies that

$$R(F|_S) \leq \inf_{\alpha > 0} \left(\frac{4\alpha}{\sqrt{n}} + \frac{12}{n} \int_{\alpha}^{\sqrt{n}} \sqrt{\ln N_d(F|_S, \epsilon, \|\cdot\|_2)} d\epsilon \right).$$

Lemma 4. Mohri et al. [2018] Let $\mathcal{L} : X \times Y \rightarrow \mathbb{R}$ be an L_p loss function bounded by $M > 0$, \mathcal{F} be the hypothesis set, family $\mathcal{G} = \{(\mathbf{x}, \mathbf{y}) \rightarrow \mathcal{L}(F(\mathbf{x}), \mathbf{y}) : F \in \mathcal{F}\}$, then for any δ , with probability at least $1 - \delta$, the following inequality holds:

$$E_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[\mathcal{L}(F(\mathbf{x}), \mathbf{y})] \leq \frac{1}{n} \sum_{i=1}^n l(F(\mathbf{x}_i), \mathbf{y}_i) + 2\hat{R}_S(\mathcal{G}) + 3M \sqrt{\frac{\ln(\frac{2}{\delta})}{2n}}.$$

where $\hat{R}_S(\mathcal{G})$ is the empirical Rademacher complexity of the family \mathcal{G} given the dataset S .

Lemma 5. Let $(\epsilon_1, \dots, \epsilon_L)$ be given, along with operator norm bounds (c_1, \dots, c_L) . Suppose the matrix $\Theta = (\Omega^1, \Omega^2, \mathbf{W}^2, \dots, \mathbf{W}^L)$ lie within $B^1 \times B^2 \times \dots \times B^{L+1}$ where B^l are arbitrary classes with the property that each $\mathbf{W}^l \in B^l$ has $\sup_{\|\mathbf{x}\| \leq 1} \|\mathbf{W}^l \mathbf{x}\| = c_l$. Lastly, let data \mathbf{X} be given with $\|\mathbf{X}\|_1 \leq B$. Then, letting $\tau = \sum_{j \leq L} \epsilon_j \prod_{l=j+1}^L c_l$, the neural network hypothesis space $H_{\mathbf{X}} = \{F_{\Theta}(\mathbf{X}) : \Theta \in B^1 \times B^2 \times \dots \times B^{L+1}\}$ has covering number bound

$$N_d(H_{\mathbf{X}}, \tau, \|\cdot\|_L) \leq \prod_{l=1}^L \sup_{(\Omega^1, \Omega^2, \mathbf{W}^2, \dots, \mathbf{W}^L), \forall j < i} N_d(\{W^l F_{(\Omega^1, \Omega^2, \mathbf{W}^2, \dots, \mathbf{W}^L)}(\mathbf{X})\}, \epsilon_l, \|\cdot\|_{l+1}).$$

Proof. It is proved by mathematical induction. Inductively construct covers \mathcal{C}_l of $\Omega^1 \mathbf{X}, \dots, \mathbf{W}^L \mathbf{X}^{L-1}$.

- When $l = 1$, since Ω^2 is fixed once the output dimension is chosen, it is trivial that the lemma holds.
- When $l = 2$

Denote \mathcal{C}_1 as an ϵ -net of $\Omega^1 \mathbf{X}$, then

$$|\mathcal{C}_1| \leq N_d(\{\Omega^1 \mathbf{X} : \Omega^1 \in B^1\}, \epsilon_1, \|\cdot\|_2) \triangleq N_1.$$

For a fixed $C \in \mathcal{C}_1$, there exists an ϵ -net $G(C)$ that

$$|G(C)| \leq N_d(\{\Omega^2(\sigma(\Omega^1 \mathbf{X})) : \Omega^2 \in B^2\}, \epsilon_2, \|\cdot\|_2) \triangleq N_2.$$

Set $\mathcal{C}_2 = \cup_{C \in \mathcal{C}_1} G(C)$, then \mathcal{C}_2 is an ϵ -net of $\{\Omega^2(\sigma(\Omega^1 \mathbf{X})) : \Omega^2 \in B^2\}$.

Then, $|\mathcal{C}_2| \leq N_1 N_2$.

And for any $\mathbf{X}^0 \in \Omega^2(\sigma(\Omega^1 \mathbf{X}))$, there exists $\hat{\mathbf{X}}^0 \in \mathcal{C}_2$ that:

$$\begin{aligned} |\mathbf{X}^0 - \hat{\mathbf{X}}^0| &= |\Omega^2(\sigma(\Omega^1(\mathbf{X}))) - \Omega^2(\sigma(\widetilde{\Omega^1(\mathbf{X})}))| \\ &\leq |\Omega^2(\sigma(\Omega^1(\mathbf{X}))) - \Omega^2(\sigma(\widetilde{\Omega^1(\mathbf{X})}))| + |\Omega^2(\sigma(\widetilde{\Omega^1(\mathbf{X})})) - \Omega^2(\sigma(\widetilde{\widetilde{\Omega^1(\mathbf{X})}}))| \\ &\leq |\Omega^2| |\sigma(\Omega^1(\mathbf{X})) - \sigma(\widetilde{\Omega^1(\mathbf{X})})| + \epsilon_2 \\ &\leq c_2 \epsilon_1 + \epsilon_2. \end{aligned}$$

Note that $\hat{\mathbf{X}}^0 = \Omega^2(\sigma(\widetilde{\widetilde{\Omega^1(\mathbf{X})}}))$ is in an ϵ -net of $\Omega^2(\sigma(\widetilde{\Omega^1(\mathbf{X})}))$, and $\sigma(\widetilde{\widetilde{\Omega^1(\mathbf{X})}})$ is in an ϵ -net of $\sigma(\widetilde{\Omega^1(\mathbf{X})})$.

The lemma holds under this condition.

- Assume the lemma holds when $1 \leq l < L$.
- When $l = L$, use the same notation as above, set $\mathcal{C}_{L+1} = \cup_{C \in \mathcal{C}_L} G(C)$, then $|\mathcal{C}_{L+1}| \leq \prod_{l=1}^L N_l$.

And for any \mathbf{X}^L , there exists $\hat{\mathbf{X}}^L \in \mathcal{C}_{L+1}$ that:

$$\begin{aligned}
|\mathbf{X}^L - \hat{\mathbf{X}}^L| &= |\mathbf{W}^L \mathbf{X}^{L-1} - \widetilde{\mathbf{W}^L \mathbf{X}^{L-1}}| \\
&\leq |\mathbf{W}^L \mathbf{X}^{L-1} - \widetilde{\mathbf{W}^L \mathbf{X}^{L-1}}| + |\widetilde{\mathbf{W}^L \mathbf{X}^{L-1}} - \widetilde{\mathbf{W}^L \mathbf{X}^{L-1}}| \\
&\leq |\mathbf{W}^L| |\mathbf{X}^{L-1} - \widetilde{\mathbf{X}^{L-1}}| + \epsilon_{L+1} \\
&\leq c_L \left(\sum_{j \leq i} \epsilon_j \prod_{l=j+1}^{L-1} c_l \right) + \epsilon_L \\
&= \sum_{j \leq L} \epsilon_j \prod_{l=j+1}^L c_l.
\end{aligned}$$

Note that $\hat{\mathbf{X}}^L = \widetilde{\mathbf{W}^L \mathbf{X}^{L-1}}$ is in an ϵ -net of $\mathbf{W}^L \mathbf{X}^{L-1}$, and $\widetilde{\mathbf{X}^{L-1}}$ is in an ϵ -net of \mathbf{X}^{L-1} .

By induction, the lemma holds. \square

Theorem 2. Let $S = \{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ be a sample data of size n from distribution \mathcal{D} . Given the weight matrices defined before $(\Omega^1, \Omega^2, \mathbf{W}^2, \dots, \mathbf{W}^L)$, and they satisfy that $\|\mathbf{W}^l\| \leq c_l (l > 2)$, $\|\Omega^1\| \leq a_1$, $\|\Omega^2\| \leq c_1$, $\|\mathbf{A}^l\| \leq b_l$, $\|\mathbf{X}\|_1 \leq B$, $d^l \leq W$ and $T = (\sum_{l=2}^L (\frac{b_l}{c_l})^{2/3})^{3/2} \prod_{l=1}^L c_l$. And the loss function $\mathcal{L}(F(\mathbf{x}), \mathbf{y}) \leq M$. Then, with the probability at least $1 - \delta$, the proposed network F satisfy:

$$\begin{aligned}
\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [(\mathcal{L}(F_{\Theta}(\mathbf{x}), \mathbf{y}))] &\leq \frac{1}{n} \sum_{i=1}^n \mathcal{L}(F_{\Theta}(\mathbf{x}_i), \mathbf{y}_i) \\
&+ \mathcal{O}\left(\frac{8M}{n^{3/2}} + M \sqrt{\frac{\ln(1/\delta)}{n}} + \ln(n) \frac{\sqrt{\ln(W)W \|\mathbf{X}^0\|^2 T^2 + \ln(W) a_1^2 \|\mathbf{X}\|^2}}{n}\right).
\end{aligned}$$

Proof. Follow the notaion above, by Lemma 5 and Theorem 1, the covering number of the whole network has:

$$\begin{aligned}
\ln(N_d(H_X, \tau, |\cdot|_{L+1})) &\leq \sum_{l=1}^L \sup_{(\Omega^1, \Omega^2, \mathbf{W}^2, \dots, \mathbf{W}^j), \forall j < l} \ln(N_d(\{\mathbf{W}^l F_{\Omega^1, \Omega^2, \mathbf{W}^2, \dots, \mathbf{W}^{l-1}}(\mathbf{X})\}, \epsilon_l, \|\cdot\|_{l+1})). \\
&\leq \sum_{l=1}^L \sup_{(\Omega^1, \Omega^2, \mathbf{W}^2, \dots, \mathbf{W}^L), \forall j < l} \ln(2d^l d^{l-1} + 1)k^l + \ln(4d^0 d^x)k^0,
\end{aligned}$$

where, k^l is an integer and satisfies that $k^l \geq \frac{\|\mathbf{A}_{ij}\|_1^2}{\epsilon^2} \frac{\pi}{2} d^l \|\mathbf{X}^{l-1}\|_1^2$, $k^0 \geq \frac{\|\Omega_{ij}\|_1^2}{\epsilon^2} \max_{i,j} \|x_{ij}\|^2$. And

$$\|\mathbf{X}^l\| = \|\mathbf{W}^l \mathbf{X}^{l-1}\| \leq \|\mathbf{W}^l\| \|\mathbf{X}^{l-1}\| \leq \dots \leq \prod_{j=1}^l \|\mathbf{W}^j\| \|\mathbf{X}^0\|.$$

In order to satisfy that $\tau < \epsilon$, set

$$\epsilon_i = \frac{\alpha_i \epsilon}{\prod_{j>i} c_j}, \quad \alpha_i = \frac{1}{\bar{\alpha}} \left(\frac{b_i}{c_i}\right)^{2/3}, \quad \bar{\alpha} = \sum_{i=1}^L \left(\frac{b_i}{c_i}\right)^{2/3}.$$

Then, by Lemma 5

$$\tau = \sum_{j \leq L} \epsilon_j \prod_{l=j+1}^L c_l = \sum_{j=1}^L \alpha_j \epsilon_j = \epsilon.$$

And

$$\begin{aligned}
\ln(N_d(H_X, \tau, |\cdot|_{L+1})) &\leq \sum_{l=1}^L \ln(2d^l d^{l-1} + 1) \cdot \frac{b_l^2 \pi}{\epsilon_l^2} d^l \cdot \|\mathbf{X}^0\|^2 \prod_{j<l} c_j^2 + \ln(4d^0 d^x) \frac{a_1^2}{\epsilon^2} B^2 \\
&= \sum_{l=1}^L \ln(2d^l d^{l-1} + 1) \cdot \frac{b_l^2 \pi}{\epsilon^2 c_l^2 \alpha_l^2} d^l \cdot \|\mathbf{X}^0\|^2 \prod_{j=1}^L c_j^2 + \ln(4d^0 d^x) \frac{a_1^2}{\epsilon^2} B^2 \\
&= \ln(2W^2 + 1) \frac{\pi}{2\epsilon^2} W \|\mathbf{X}^0\|^2 \sum_{l=1}^L \frac{b_l^2}{c_l^2 \alpha_l^2} \prod_{j=1}^L c_j^2 + \ln(4W^2) \frac{a_1^2}{\epsilon^2} B^2 \quad (d^l \leq W) \\
&= \ln(2W^2 + 1) \frac{\pi}{2\epsilon^2} W \|\mathbf{X}^0\|^2 (\bar{\alpha})^3 \prod_{j=1}^L c_j^2 + \ln(4W^2) \frac{a_1^2}{\epsilon^2} B^2 \triangleq \frac{R}{\epsilon^2}.
\end{aligned}$$

Consider the class of networks \mathcal{F} obtained by affixing the loss $\mathcal{L}(F(x), y)$ and $\mathcal{L}(F(x), y) \leq M$. When \mathcal{L} is fixed, the covering number of the obtained network \mathcal{F} is not larger than the original network.

Then, by Lemma 3

$$\begin{aligned}
R(\frac{\mathcal{F}_S}{M}) &\leq \inf_{\alpha>0} \left(\frac{4\alpha}{\sqrt{n}} + \frac{12}{n} \int_{\alpha}^{\sqrt{n}} \sqrt{\ln N_d(\frac{F_{\Theta}}{M}, \epsilon, \|\cdot\|_2) d\epsilon} \right) = \inf_{\alpha>0} \left(\frac{4\alpha}{\sqrt{n}} + \frac{12}{n} \int_{\alpha}^{\sqrt{n}} \sqrt{\ln N(\frac{H_X}{M}, \frac{\tau}{M}, \|\cdot\|_2) d\epsilon} \right) \\
&\leq \inf_{\alpha>0} \left(\frac{4\alpha}{\sqrt{n}} + \frac{12}{n} \int_{\alpha}^{\sqrt{n}} \sqrt{\frac{R}{M^2 \epsilon^2}} d\epsilon \right) = \inf_{\alpha>0} \left(\frac{4\alpha}{\sqrt{n}} + \ln(\sqrt{n}/\alpha) \frac{12\sqrt{R}}{Mn} \right) \\
&\leq \frac{4}{n^{3/2}} + \ln(n^{3/2}) \frac{12\sqrt{R}}{Mn} \quad (\alpha = 1/n) \\
&= \frac{4}{n^{3/2}} + \ln(n^{3/2}) \frac{12\sqrt{\ln(2W^2 + 1) \frac{\pi}{2} W \|\mathbf{X}^0\|^2 (\bar{\alpha})^3 \prod_{l=1}^L c_l^2 + \ln(4W^2) a_1^2 B^2}}{Mn}.
\end{aligned}$$

And the empirical Rademacher complexities of $\frac{\mathcal{F}}{M}$ follows:

$$\hat{R}_S(\frac{\mathcal{F}}{M}) = E_{\sigma} \left[\sup_{\frac{f}{M} \in \frac{\mathcal{F}}{M}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{f(\mathbf{x}_i)}{M} \right] = \frac{1}{M} E_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i) \right] = \frac{1}{M} \hat{R}_S(\mathcal{F}).$$

By Lemma 4

$$\begin{aligned}
\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\mathcal{L}(F_{\Theta}(\mathbf{x}), \mathbf{y})] &\leq \frac{1}{n} \sum_{i=1}^n \mathcal{L}(F_{\Theta}(\mathbf{x}_i), \mathbf{y}_i) + 2\hat{R}_S(\mathcal{F}) + 3M \sqrt{\frac{\ln \frac{2}{\delta}}{2n}} \\
&\leq \frac{1}{n} \sum_{i=1}^n \mathcal{L}(F_{\Theta}(\mathbf{x}_i), \mathbf{y}_i) + 3M \sqrt{\frac{\ln \frac{2}{\delta}}{2n}} + \frac{8M}{n^{3/2}} \\
&\quad + \ln(n^{3/2}) \frac{24\sqrt{\ln(2W^2 + 1) \frac{\pi}{2} W \|\mathbf{X}^0\|^2 (\bar{\alpha})^3 \prod_{l=1}^L c_l^2 + \ln(4W^2) a_1^2 B^2}}{n} \\
&\leq \frac{1}{n} \sum_{i=1}^n \mathcal{L}(F_{\Theta}(\mathbf{x}_i), \mathbf{y}_i) \\
&\quad + \mathcal{O}\left(\frac{8M}{n^{3/2}} + M \sqrt{\frac{\ln 1/\delta}{n}} + \ln(n) \frac{\sqrt{\ln(W) W \|\mathbf{X}^0\|^2 T^2 + \ln(W) a_1^2 \|\mathbf{X}\|^2}}{n}\right),
\end{aligned}$$

where $T^2 = (\bar{\alpha})^3 \prod_{l=1}^L c_l^2$. \square

2 Initialization

As mentioned in the main paper, we initialize the complex-valued weight matrix as $\mathbf{W} = \cos(\mathbf{A}) + i\sin(\mathbf{A})$. This design ensures CosNet retains the property of non-stationary spectral kernels and

takes the relative distance of data in the complex number domain without increasing the number of parameters. In this section, we further discuss the initialization, including non-stationary and multi-kernel learning.

Non-stationary ensuring Compared with the CSKE module of CosNet, sampling stack the complex-valued spectral kernel mapping in the neural networks cannot ensure that the model retains the non-stationary of the spectral kernel. In this section, we explain that in a two-dimensional case.

For the complex-valued spectral mapping $z = \begin{bmatrix} \cos(u_{11}) + \cos(u'_{11}) \\ \cos(u_{21}) + \cos(u'_{21}) \end{bmatrix} + i \begin{bmatrix} \sin(v_{11}) + \sin(v'_{11}) \\ \sin(v_{21}) + \sin(v'_{21}) \end{bmatrix} \in \mathbb{C}^2$, following the commonly used setting, the weight matrix is defined as $\mathbf{W} = \mathbf{A} + i\mathbf{B}$. The complex-valued transformation with the matrix formula can be defined as:

$$\begin{aligned} \Psi(z) &= \begin{bmatrix} \mathbf{A} & -\mathbf{B} \\ \mathbf{B} & \mathbf{A} \end{bmatrix} * \begin{bmatrix} \Re(z) \\ \Im(z) \end{bmatrix} = \begin{bmatrix} a_{11} & b_{11} \\ a_{12} & b_{12} \\ -b_{11} & a_{11} \\ -b_{12} & a_{12} \end{bmatrix}^\top * \begin{bmatrix} \cos(u_{11}) + \cos(u'_{11}) \\ \cos(u_{21}) + \cos(u'_{21}) \\ \sin(v_{11}) + \sin(v'_{11}) \\ \sin(v_{21}) + \sin(v'_{21}) \end{bmatrix} \\ &= \begin{bmatrix} a_{11}(\cos(u_{11}) + \cos(u'_{11})) + a_{12}(\cos(u_{21}) + \cos(u'_{21})) - b_{11}(\sin(v_{11}) + \sin(v'_{11})) - b_{12}(\sin(v_{21}) + \sin(v'_{21})) \\ b_{11}(\cos(u_{11}) + \cos(u'_{11})) + b_{12}(\cos(u_{21}) + \cos(u'_{21})) + a_{11}(\sin(v_{11}) + \sin(v'_{11})) + a_{12}(\sin(v_{21}) + \sin(v'_{21})) \end{bmatrix}, \end{aligned}$$

Obviously, for a general non-stationary spectral kernel $k(z, z')$, it cannot be defined as the inner product of two complex-valued mappings, $\Psi(z)$ and $\overline{\Psi(z')}$.

Multi-kernels learning In addition to the mentioned property in the main paper, our initialization enables each unit can be considered a combination of multiple spectral kernels. In this section, we explain that in two-dimensional complex input space.

Example 1. *Let*

$$z = \begin{bmatrix} \cos(u_{11}) + \cos(u'_{11}) \\ \cos(u_{21}) + \cos(u'_{21}) \end{bmatrix} + i \begin{bmatrix} \sin(v_{11}) + \sin(v'_{11}) \\ \sin(v_{21}) + \sin(v'_{21}) \end{bmatrix} \in \mathbb{C}^2$$

be a complex-valued vector. The complex-valued weight matrix is defined as $\mathbf{W} = \cos(\mathbf{A}) + i\sin(\mathbf{A})$, where $\mathbf{A} = [a_{11}, a_{12}] \in \mathbb{R}^{1 \times 2}$ is a real-valued matrix. The complex-valued mapping can be defined as:

$$\begin{aligned} \Psi(z) &= \mathbf{W} * z \\ &= (\cos(\mathbf{A}) + i\sin(\mathbf{A})) \\ &\quad * \left(\begin{bmatrix} \cos(u_{11}) + \cos(u'_{11}) \\ \cos(u_{21}) + \cos(u'_{21}) \end{bmatrix} + i \begin{bmatrix} \sin(v_{11}) + \sin(v'_{11}) \\ \sin(v_{21}) + \sin(v'_{21}) \end{bmatrix} \right) \end{aligned}$$

Without loss of generality, we formalize the complex-valued mapping as the following matrix notation:

$$\begin{aligned} \Psi(z) &= \begin{bmatrix} \cos(\mathbf{A}) & -\sin(\mathbf{A}) \\ \sin(\mathbf{A}) & \cos(\mathbf{A}) \end{bmatrix} * \begin{bmatrix} \Re(z) \\ \Im(z) \end{bmatrix} \\ &= \begin{bmatrix} \cos(a_{11}) & \sin(a_{11}) \\ \cos(a_{12}) & \sin(a_{12}) \\ -\sin(a_{11}) & \cos(a_{11}) \\ -\sin(a_{12}) & \cos(a_{12}) \end{bmatrix}^\top * \begin{bmatrix} \cos(u_{11}) + \cos(u'_{11}) \\ \cos(u_{21}) + \cos(u'_{21}) \\ \sin(u_{11}) + \sin(u'_{11}) \\ \sin(u_{21}) + \sin(u'_{21}) \end{bmatrix} \\ &= \begin{bmatrix} \Psi_{a_{11}, u_{11}, u'_{11}} + \Psi_{a_{12}, u_{21}, u'_{21}} \\ \Psi'_{a_{11}, v_{11}, v'_{11}} + \Psi'_{a_{12}, v_{21}, v'_{21}} \end{bmatrix} \end{aligned}$$

where $\Psi_{a_{11}, u_{11}, u'_{11}} = \cos(a_{11} + u_{11}) + \cos(a_{11} + u'_{11})$, $\Psi_{a_{12}, u_{21}, u'_{21}} = \cos(a_{12} + u_{21}) + \cos(a_{12} + u'_{21})$, $\Psi'_{a_{11}, v_{11}, v'_{11}} = \sin(a_{11} + v_{11}) + \sin(a_{11} + v'_{11})$, and $\Psi'_{a_{12}, v_{21}, v'_{21}} = \sin(a_{12} + v_{21}) + \sin(a_{12} + v'_{21})$.

We can observe that $\Psi_{a_{11}, u_{11}, u'_{11}}$, $\Psi_{a_{12}, u_{21}, u'_{21}}$, $\Psi'_{a_{11}, v_{11}, v'_{11}}$, and $\Psi'_{a_{12}, v_{21}, v'_{21}}$ can be seen as two separate parts of two different spectral kernel mappings. Hence, the proposed CosNet can be regarded as a linear combination of different kernels (the number of kernels is restricted by the feature numbers), which indicates that our method naturally has a close relation with multi-kernel learning.

3 Experiment

In this section, we include more details of the experiment section in the main paper, including the information on the involved datasets (shown in Table 1), the detailed setting for each dataset (shown in Table 2), and extra experiments.

Table 1: The detailed information of the involved dataset. Specifically, the input size denotes the number of time points and features for the time-series classification task and regression task, respectively.

Dataset	Type	Input size	Train.Data	Test.Data	Class
FordA	Sensor	500	3601	1320	2
FordB	Sensor	500	3636	810	2
PhalangesOutlinesCorrect	Image	80	1800	858	2
Wine	Spectro	234	57	54	2
ECG200	ECG	96	100	100	2
ECG5000	ECG	140	500	4500	5
Herring	Image	512	64	64	2
Ham	Spectro	431	109	105	2
ProximalPhalanxOutlineAgeGroup	Image	80	400	139	6
Earthquakes	Sensor	512	322	139	2
DistalPhalanxTW	Image	80	400	139	6
Strawberry	Spectro	235	613	370	2
power	-	4	7654	1914	-
concreat	-	8	824	206	-
yacht	-	6	246	62	-

Table 2: The detailed settings on different datasets. Specifically, Init denotes that the weight matrix is sampled from $\mathcal{N}(0, p)$. Networks denote the unified architecture, where, the first number is the input size, the last number is the class, and the others denote the neuron numbers of the hidden layers.

	Networks
FordA	$500 \times 500 \times 256 \times 64 \times 2$
FordB	$500 \times 500 \times 256 \times 64 \times 2$
PhalangesOutlinesCorrect	$80 \times 80 \times 80 \times 64 \times 2$
Wine	$234 \times 234 \times 128 \times 64 \times 2$
ECG200	$96 \times 96 \times 96 \times 32 \times 2$
ECG5000	$140 \times 140 \times 64 \times 32 \times 5$
Herring	$512 \times 512 \times 128 \times 64 \times 2$
Ham	$431 \times 431 \times 128 \times 64 \times 2$
ProximalPhalanxOutlineAgeGroup	$80 \times 80 \times 80 \times 32 \times 3$
Earthquakes	$512 \times 512 \times 128 \times 32 \times 2$
DistalPhalanxTW	$80 \times 80 \times 80 \times 32 \times 6$
Strawberry	$235 \times 235 \times 128 \times 32 \times 2$
power	$4 \times 4 \times 4 \times 4 \times 1$
concreat	$8 \times 8 \times 8 \times 4 \times 1$
yacht	$6 \times 6 \times 6 \times 3 \times 1$

Image classification Addition to the time-sequential data, complex-valued representation is commonly used in the image processing. The phase describes objects in an image in terms of edges, shapes and their orientation. To explore the capability of CosNet on image-related tasks, we extend CosNet to the convolutional neural networks (CNNs), namely complex-valued spectral convolutional networks (CosCNet). Similar with CosNet, CosCNet also include two modules, including complex-valued representation learning (CRL) module and complex-valued convolutional (CC) module. The CRL module is used to transform the image in the real number domain to the complex number domain, and the CC module is used to explore the inherently complex-valued representation and further explore the detailed information of edges and shape.

Table 3: Classification accuracy (%) under different hyper-parameters. The best results are highlighted in **bold**.

lr	init (p)	SRFF	DSKN	DCN^1	DCN^2	ASKL	CosNet
0.1	1	64	60.85	90.30	83.15	61.00	86.05
0.1	0.1	85.55	61.50	90.30	83.15	80.20	85.85
0.1	0.01	78.25	62.80	90.30	83.15	74.10	85.05
0.01	1	52.60	64.00	88.10	84.05	75.00	90.05
0.01	0.1	85.40	67.95	88.10	84.05	89.75	91.30
0.01	0.01	73.40	77.80	88.10	84.05	87.53	90.10
0.001	1	50.85	62.75	80.35	79.25	72.55	89.25
0.001	0.1	83.50	73.00	80.35	79.25	90.90	90.25
0.001	0.01	64.00	84.40	80.35	79.25	88.90	90.45

Specifically, the CRL module is defined as :

$$\Phi(\mathbf{x}) = \sqrt{\frac{1}{4M}} \left[(\cos(\mathbf{\Omega} * \mathbf{x}) + \cos(\mathbf{\Omega}' * \mathbf{x})) + i(\sin(\mathbf{\Omega} * \mathbf{x}) + \sin(\mathbf{\Omega}' * \mathbf{x})) \right],$$

and the convolution operation of CC module with the matrix notation is defined as:

$$\begin{bmatrix} \Re(\Psi(\mathbf{h})) \\ \Im(\Psi(\mathbf{h})) \end{bmatrix} = \begin{bmatrix} \cos(\mathbf{A}) & -\sin(\mathbf{A}) \\ \sin(\mathbf{A}) & \cos(\mathbf{A}) \end{bmatrix} * \sqrt{\frac{1}{4M}} \begin{bmatrix} \cos(\mathbf{\Omega} * \mathbf{x}) + \cos(\mathbf{\Omega}' * \mathbf{x}) \\ \sin(\mathbf{\Omega} * \mathbf{x}) + \sin(\mathbf{\Omega}' * \mathbf{x}) \end{bmatrix},$$

where, $\mathbf{\Omega}$, $\mathbf{\Omega}'$, and \mathbf{A} are filters. Moreover, the CosCNet with l layers is defined as:

$$CosCNet(\mathbf{x}) = \Psi^{l-1}(\dots \Psi^1(\Phi^1(\mathbf{x}))).$$

Generalizaion of CosNet Furthermore, to evaluate the generalization of our CosNet, we explore the influence of varying learning rates and distribution of weight matrices on the result based on ECG200 dataset. The results are shown in Table 3 The results show the superior performance and stability of our CosNet.

References

- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- Gilles Pisier. Remarques sur un résultat non publié de b. maurey. *Séminaire d'Analyse fonctionnelle (dit "Maurey-Schwartz")*, pages 1–12.