

The Appendix is divided into the following six sections,

- A. *Preliminaries and Structural Results*
- B. *Equivalence of Optimization problem in \mathcal{H} and $L^2(\mathcal{X}, \rho)$*
- C. *Proof of the main Theorem*
- D. *Examples of ESL*
- E. *Experiments*
- F. *Auxiliary Results*

A Preliminaries and Structural Results

We begin this section by giving some definitions and structural theorems that are required for the proofs. We first introduce a few additional notation. Consider two separable Hilbert spaces \mathcal{H}_1 and \mathcal{H}_2 . For an operator $\mathcal{D} : \mathcal{H}_1 \rightarrow \mathcal{H}_2$, we use $\|\mathcal{D}\|_{\mathcal{L}^p(\mathcal{H}_1, \mathcal{H}_2)}$ to denote its p^{th} Schatten norm, assuming that it is finite. We omit the p when we talk about about Hilbert-Schmidt norm i.e. $p = 2$.

A.1 Covariance operators C and C_m

The covariance operators in the RKHS and the random feature space are defined as follows:

Definition A.1. $C : \mathcal{H} \rightarrow \mathcal{H}$ is the co-variance operator of the random variables $k(\mathbf{x}, \cdot)$ with measure ρ , defined as:

$$Cf := \int_{\mathcal{X}} k(\mathbf{x}, \cdot) f(\mathbf{x}, t) d\rho(\mathbf{x}, t)$$

C is compact and self-adjoint, which implies C has a spectral decomposition as follows:

$$C = \sum_{i=1}^{\infty} \bar{\lambda}_i \bar{\phi}_i \otimes_{\mathcal{H}} \bar{\phi}_i$$

where $\bar{\lambda}_i, \bar{\phi}_i$'s are the eigenvalues and eigenfunctions of C . Also, $\bar{\phi}_i$ are a unitary basis for \mathcal{H} .

Definition A.2. $C_m : \mathcal{F} \rightarrow \mathcal{F}$ is the covariance operator in the random feature space, defined as

$$C_m := \mathbb{E}_{\rho} [z(\mathbf{x}, t) \otimes_{\mathcal{F}} z(\mathbf{x}, t)]$$

Equivalently, for any $\beta \in \mathcal{F}$, $C_m \beta = \int_{\mathcal{X}} \langle z(\mathbf{x}, t), \beta \rangle z(\mathbf{x}, t) d\rho(\mathbf{x}, t)$.

C_m is compact and self-adjoint which implies that C_m has a spectral decomposition as follows:

$$C_m = \sum_{i=1}^m \lambda_i \phi_i \otimes_{\mathcal{F}} \phi_i$$

The kernel integral operators and its approximation based on random features are defined as follows:

Definition A.3. The kernel integral operator $L : L^2(\mathcal{X}, \rho) \rightarrow L^2(\mathcal{X}, \rho)$ is defined as follows:

$$Lg = \int_{\mathcal{X}} k(\mathbf{x}, \cdot) g(\mathbf{x}) d\rho(\mathbf{x}) \quad \forall g \in L^2(\mathcal{X}, \rho)$$

Definition A.4. $L_m : L^2(\mathcal{X}, \rho) \rightarrow L^2(\mathcal{X}, \rho)$ is the (approximated) kernel integral operator, defined as:

$$(L_m g)(\cdot) = \int_{\mathcal{X}} k_m(\mathbf{x}, \cdot) g(\mathbf{x}) d\rho(\mathbf{x})$$

We state the classical Mercer's and Bochner's theorems for completeness.

Theorem A.5 (Mercer's Theorem). *For every positive definition kernel $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, there exists a set Ω with measure π , and functions $\varpi(\cdot) : \mathcal{X} \times \Omega \rightarrow \mathbb{R}$ such that the kernel has an integral representation of the following form,*

$$k(\mathbf{x}, \mathbf{y}) = \int_{\Omega} z(\mathbf{x}, \omega) z(\mathbf{y}, \omega) d\pi(\omega) \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$$

In particular, for shift-invariant kernels, we have

Theorem A.6 (Bochner's Theorem [Rudin \[2017\]](#)). *A continuous, real-valued, symmetric shift-invariant kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a positive-definite kernel if and only if there exists a non-negative measure $\pi(\omega)$ such that $k(x-y) = \int_{\mathcal{X}} e^{i\omega^\top(x-y)} d\pi(\omega)$ i.e the inverse Fourier transform of $k(x-y)$*

For a comprehensive list of kernels with their Fourier transform see Table 1 of [Xie et al. \[2015\]](#).

We now define operators to lift functions and operators from and into different spaces. These are crucially used in the analysis of the algorithm. See Figure 3 in for a schematic of the lifting operators.

A.2 Inclusions operators I, \mathfrak{I}

We first recall the definitions of Inclusion operators I, \mathfrak{I} .

Definition A.7. [Inclusion Operators I and \mathfrak{I}] The inclusion operator is defined

$$I : \mathcal{H} \rightarrow L^2(\mathcal{X}, \rho), (If) = f, \text{ where } f \in \mathcal{H}$$

Also, for an operator $D \in HS(\mathcal{H})$ with spectral decomposition $D = \sum_{i \in I \subset \mathbb{R}} \mu_i \psi_i \otimes \psi_i$,

$$\mathfrak{I} : HS(\mathcal{H}) \rightarrow HS(\rho), \mathfrak{I}D := \sum_{i \in I \subset \mathbb{R}} \mu_i \frac{I\psi_i}{\sqrt{\langle C\psi_i, \psi_i \rangle_{\mathcal{H}}}} \otimes \frac{I\psi_i}{\sqrt{\langle C\psi_i, \psi_i \rangle_{\mathcal{H}}}}$$

In Proposition A.8, we show that the adjoint of the Inclusion operator I is

$$I^* : L^2(\mathcal{X}, \rho) \rightarrow \mathcal{H}, (I^*g)(\cdot) = \int k(x, \cdot)g(x)d\rho(x).$$

Moreover, In Proposition A.9 we show that the covariance operator and the kernel integral operator can be expressed in terms of I and I^* as $C = I^*I$ and $L = II^*$.

Proposition A.8. *The following holds with regard to the inclusion operator,*

- (a). *The adjoint of the Inclusion operator I is given by $(I^*g)(\cdot) = \int_{\mathcal{X}} k(x, \cdot)g(x)d\rho(x)$.*
- (b). *I and I^* are Hilbert-Schmidt.*

Proof of Proposition A.8. (a). We first show that the adjoint of the Inclusion operator I is given by $(I^*g)(\cdot) = \int_{\mathcal{X}} k(x, \cdot)g(x)d\rho(x)$. For $f \in \mathcal{H}$ and $g \in L^2(\mathcal{X}, \rho)$, we have that

$$\begin{aligned} \langle If, g \rangle_{\rho} &= \langle f, g \rangle_{\rho} && \text{(Definition of } I) \\ &= \int_{\mathcal{X}} f(x)g(x)d\rho(x) \\ &= \int_{\mathcal{X}} \langle k(x, \cdot), f \rangle_{\mathcal{H}} g(x)d\rho(x) && \text{(Reproducing property)} \\ &= \int_{\mathcal{X}} \langle k(x, \cdot)g(x), f \rangle_{\mathcal{H}} d\rho(x) && \text{(Linearity of inner product)} \\ &= \left\langle \int_{\mathcal{X}} k(x, \cdot)g(x)d\rho(x), f \right\rangle_{\mathcal{H}} && \text{(Fubini's Theorem)} \\ &= \langle I^*g, f \rangle_{\mathcal{H}} \end{aligned}$$

(b). Let $\{\bar{e}_i\}_{i=1}^{\infty}$ be an orthonormal basis for \mathcal{H} . We have,

$$\begin{aligned}
\|I\|_{\mathcal{L}(\mathcal{H}, \rho)}^2 &= \sum_{i=1}^{\infty} \|I\bar{e}_i\|_{\rho}^2 && \text{(Pythagoras Theorem)} \\
&= \sum_{i=1}^{\infty} \|\bar{e}_i\|_{\rho}^2 && \text{(Definition 3.5)} \\
&= \sum_{i=1}^{\infty} \int_{\mathcal{X}} \langle k(x, \cdot), \bar{e}_i \rangle_{\mathcal{H}}^2 d\rho(x) && \text{(Reproducing Property)} \\
&= \int_{\mathcal{X}} \sum_{i=1}^{\infty} \langle k(x, \cdot), \bar{e}_i \rangle_{\mathcal{H}}^2 d\rho(x) && \text{(Fubini's Theorem)} \\
&= \int_{\mathcal{X}} k(x, x) d\rho(x) \leq \tau^2 < \infty && \text{(Assumption 3.1)}
\end{aligned}$$

For the adjoint I^* , we have $\|I^*\|_{\mathcal{L}(\rho, \mathcal{H})} = \|I\|_{\mathcal{L}(\mathcal{H}, \rho)} < \infty$ \square

Proposition A.9. *The following properties hold,*

- (a). *The covariance operator and the kernel integral operator satisfy $C = I^*I$ and $L = \Pi^*$ respectively.*
- (b). *C and L are trace-class*

Proof of Proposition A.9. (a). We first show that $C = I^*I$. For any $f \in L^2(\mathcal{X}, \rho)$, we have

$$\begin{aligned}
I^*I f &= I^* f && \text{(Definition A.7)} \\
&= \int_{\mathcal{X}} k(x, \cdot) f(x) d\rho(x) && \text{(Proposition A.8)} \\
&= C f && \text{(Definition A.1)}
\end{aligned}$$

We now show that $L = \Pi^*$. For any $g \in L^2(\mathcal{X}, \rho)$, we have

$$\begin{aligned}
\Pi^* g &= I \left(\int_{\mathcal{X}} k(x, \cdot) g(x) d\rho(x) \right) && \text{(Proposition A.8)} \\
&= \int_{\mathcal{X}} k(x, \cdot) g(x) d\rho(x) && \text{(Definition A.7)} \\
&= L g
\end{aligned}$$

(b). Now we show that C and L are trace-class.

$$\begin{aligned}
\|C\|_{\mathcal{L}^1(\mathcal{H})} &= \|I^*I\|_{\mathcal{L}^1(\mathcal{H})} && \text{(Proposition A.9)} \\
&= \|I\|_{\mathcal{L}^2(\mathcal{H})}^2 < \infty && \text{(Proposition A.8)}
\end{aligned}$$

Similarly,

$$\begin{aligned}
\|L\|_{\mathcal{L}^1(\rho)} &= \|\Pi^*\|_{\mathcal{L}^1(\rho)} && \text{(Proposition A.9)} \\
&= \|I\|_{\mathcal{L}^2(\mathcal{H})}^2 < \infty && \text{(Proposition A.8)}
\end{aligned}$$

\square

A.3 Approximation operators A and \mathfrak{A}

We first recall the definitions of approximation operators A and \mathfrak{A} .

Definition A.10. [Approximation Operators A and \mathfrak{A}] The Approximation operator A is defined as

$$A : \mathcal{F} \rightarrow L^2(X, \rho), (Av)(\cdot) = \langle z(\cdot), v \rangle, \text{ where } v \in \mathcal{F}$$

For an operator $D \in HS(\mathcal{F})$ with rank k with spectral decomposition $D = \sum_{i=1}^{\infty} \mu_i \psi_i \otimes \psi_i$, let Ψ be the matrix with eigenvectors ψ_i as columns and Φ be the matrix with eigenvectors of C_m as columns (see Definition 3.4). Define

$$R^* = \arg \min_{R^T R = R R^T = I} \|\Psi R - \Phi\|_{\mathcal{F}}^2, \quad \tilde{\Psi} := \Psi R^*$$

Let $\tilde{\psi}_i$ be the columns of $\tilde{\Psi}$, define

$$\mathfrak{A} : HS(\mathcal{F}) \rightarrow HS(\rho), \quad \mathfrak{A}D := \sum_{i=1}^k \mu_i \frac{A\tilde{\psi}_i}{\sqrt{\langle C_m \tilde{\psi}_i, \tilde{\psi}_i \rangle_{\mathcal{F}}}} \otimes_{\rho} \frac{A\tilde{\psi}_i}{\sqrt{\langle C_m \tilde{\psi}_i, \tilde{\psi}_i \rangle_{\mathcal{F}}}}$$

Note that the definition of the approximation operator \mathfrak{A} requires knowledge of the co-variance matrix C_m to find the optimal rotation matrix R^* , but this is solely for the purpose of analysis and is not used in the algorithm in any form.

In Proposition A.11, we show that the adjoint of the Approximation Operator is

$$A^* : L^2(X, \rho) \rightarrow \mathcal{F}, \quad (A^* f)_i = \int_{\mathcal{X}} f(x) z_{\omega_i}(x) d\rho(x).$$

Moreover, in Proposition A.12 we show that the approximate covariance operator and the approximate kernel integral operator can be expressed in terms of the Approximation operator A as $C_m = A^* A$ and $L_m = A A^*$.

Proposition A.11. *The approximation operator satisfies the following properties,*

- (a). *The adjoint of A is $(A^* f)_i = \frac{1}{\sqrt{m}} \int_{\mathcal{X}} f(x) z_{\omega_i}(x) d\rho(x)$*
- (b). *A and A^* are Hilbert-Schmidt.*

Proof of Proposition A.11. (a). First we show that the adjoint of A is $(A^* f)_i = \frac{1}{\sqrt{m}} \int_{\mathcal{X}} f(x) z_{\omega_i}(x) d\rho(x)$. For $v \in \mathcal{F}, f \in L^2(\mathcal{X}, \rho)$, we have,

$$\begin{aligned} \langle Av, f \rangle_{\rho} &= \int_{\mathcal{X}} (Av)(x) f(x) d\rho(x) \\ &= \int_{\mathcal{X}} \langle z(x), v \rangle_{\mathcal{F}} f(x) d\rho(x) && \text{(Definition A.10)} \\ &= \left\langle \int_{\mathcal{X}} z(x) f(x) d\rho(x), v \right\rangle_{\mathcal{F}} && \text{(Fubini's Theorem)} \\ &= \langle A^* f, v \rangle_{\mathcal{F}} \end{aligned}$$

(b). Let $\{e_i\}$ be an orthonormal basis for \mathcal{F} .

$$\begin{aligned} \|A\|_{\mathcal{L}(\mathcal{F}, \rho)}^2 &= \sum_{i=1}^m \|Ae_i\|_{\rho}^2 && \text{(Pythagoras Theorem)} \\ &= \sum_{i=1}^m \|\langle z(\cdot), e_i \rangle_{\mathcal{F}}\|_{\rho}^2 && \text{(Definition 3.6)} \\ &= \sum_{i=1}^m \int_{\mathcal{X}} (\langle z(x), e_i \rangle_{\mathcal{F}})^2 d\rho(x) \\ &\leq \sum_{i=1}^m \int_{\mathcal{X}} \|z(x)\|_{\mathcal{F}}^2 d\rho(x) < m\tau^2 < \infty \end{aligned}$$

where third last and second inequality follows from Cauchy Schwartz inequality and Assumption 3.1 respectively.

Similarly, to show A^* is Hilbert-Schmidt, we note that $\|A^*\|_{\mathcal{L}^2(\rho, \mathcal{F})} = \|A\|_{\mathcal{L}^2(\mathcal{F}, \rho)} < \infty$ \square

In the following proposition, we show how the Covariance operator C_m and kernel Integral operator L_m are related.

Proposition A.12. *The following properties hold,*

- (a). C_m and L_m satisfy that $C_m = A^*A$, $L_m = AA^*$
- (b). C_m and L_m are trace-class.

Proof of Proposition A.12. (a). We first show the first part of the Proposition. For any $v \in \mathcal{F}$, we have,

$$\begin{aligned} A^*Av &= \int_{\mathcal{X}} \langle z(x), v \rangle_{\mathcal{F}} z(x) d\rho(x) && \text{(Definition A.10 and Proposition A.11)} \\ &= \mathbb{E}_{\rho} [z(x) \otimes_{\mathcal{F}} z(x)] v \\ &= C_m v && \text{(Definition A.2)} \end{aligned}$$

For any $g \in L^2(\mathcal{X}, \rho)$,

$$\begin{aligned} AA^*g &= \frac{1}{m} \sum_{i=1}^m z_{\omega_i}(\cdot) \int_{\mathcal{X}} z_{\omega_i}(x) g(x) d\rho(x) && \text{(Definition A.10 and Proposition A.11)} \\ &= \int_{\mathcal{X}} \sum_{i=1}^m \frac{1}{\sqrt{m}} z_{\omega_i}(\cdot) \frac{1}{\sqrt{m}} z_{\omega_i}(x) g(x) d\rho(x) && \text{(Fubini's Theorem)} \\ &= \int_{\mathcal{X}} \langle z(x), z(\cdot) \rangle_{\mathcal{F}} g(x) d\rho(x) \\ &= \int_{\mathcal{X}} k_m(x, \cdot) g(x) d\rho(x, t) && \text{(Definition of the approximate kernel mapping)} \\ &= L_m g && \text{(Definition A.4)} \end{aligned}$$

(b). Now we show that C_m and L_m are trace-class.

$$\begin{aligned} \|C_m\|_{\mathcal{L}^1(\mathcal{F})} &= \|A^*A\|_{\mathcal{L}^1(\mathcal{F})} && \text{(Proposition A.12)} \\ &= \|A\|_{\mathcal{L}^2(\mathcal{F})}^2 < \infty && \text{(Proposition A.11)} \end{aligned}$$

Similarly,

$$\begin{aligned} \|L_m\|_{\mathcal{L}^1(\rho)} &= \|AA^*\|_{\mathcal{L}^1(\rho)} && \text{(Proposition A.12)} \\ &= \|A\|_{\mathcal{L}^2(\rho)}^2 < \infty && \text{(Proposition A.11)} \end{aligned}$$

□

A.4 Kernel integral operator L and its approximation L_m

We first recall the definition of Kernel integral operator L and its approximation L_m .

Definition A.13. The kernel integral operator $L : L^2(\mathcal{X}, \rho) \rightarrow L^2(\mathcal{X}, \rho)$ is defined as follows:

$$Lg = \int_{\mathcal{X}} k(x, \cdot) g(x) d\rho(x) \quad \forall g \in L^2(\mathcal{X}, \rho)$$

Definition A.14. $L_m : L^2(\mathcal{X}, \rho) \rightarrow L^2(\mathcal{X}, \rho)$ is the (approximated) kernel integral operator, defined as:

$$(L_m g)(\cdot) = \int_{\mathcal{X}} k_m(x, \cdot) g(x) d\rho(x)$$

We now show in Proposition A.15 that spectral decomposition of the kernel integral operator L can be given in terms of the eigenfunctions and the eigenvalues of the covariance operator C .

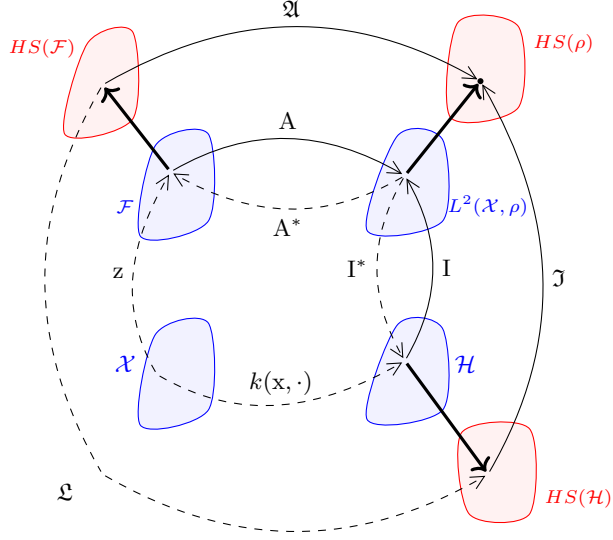


Figure 3: Maps between the data domain (\mathcal{X}), space of square integrable functions on \mathcal{X} ($L^2(\mathcal{X}, \rho)$), the RKHS of kernel $k(\cdot, \cdot)$, and RKHS of the approximate feature map, as well as maps between Hilbert-Schmidt operators on these spaces.

Proposition A.15. *The spectral decomposition of L is:*

$$L = \sum_{i=1}^{\infty} \bar{\lambda}_i \frac{I\bar{\phi}_i}{\sqrt{\bar{\lambda}_i}} \otimes_{\rho} \frac{I\bar{\phi}_i}{\sqrt{\bar{\lambda}_i}}$$

where $\frac{I\bar{\phi}_i}{\sqrt{\bar{\lambda}_i}}$ are the (unit norm) eigenfunctions of L with eigenvalues $\bar{\lambda}_i$

Proof of Proposition A.15. First we show that the operators L and $\sum_{i=1}^{\infty} \bar{\lambda}_i \frac{I\bar{\phi}_i}{\sqrt{\bar{\lambda}_i}} \otimes_{\rho} \frac{I\bar{\phi}_i}{\sqrt{\bar{\lambda}_i}}$ agree on any function $g \in L^2(\mathcal{X}, \rho)$.

$$\begin{aligned} \sum_{i=1}^{\infty} \bar{\lambda}_i \frac{I\bar{\phi}_i}{\sqrt{\bar{\lambda}_i}} \otimes_{\rho} \frac{I\bar{\phi}_i}{\sqrt{\bar{\lambda}_i}} g &= \sum_{i=1}^{\infty} I\bar{\phi}_i \langle I\bar{\phi}_i, g \rangle_{\rho} && \text{(Definition of outer product)} \\ &= \sum_{i=1}^{\infty} \bar{\phi}_i \langle \bar{\phi}_i, g \rangle_{\rho} && \text{(Definition A.7)} \\ &= \sum_{i=1}^{\infty} \bar{\phi}_i \int_{\mathcal{X}} \langle \bar{\phi}_i, k(x, \cdot) \rangle_{\mathcal{H}} g(x, t) d\rho(x, t) && \text{(Reproducing property)} \\ &= \sum_{i=1}^{\infty} \bar{\phi}_i \left\langle \bar{\phi}_i, \int_{\mathcal{X}} k(x, \cdot) g(x, t) d\rho(x, t) \right\rangle_{\mathcal{H}} && \text{(Fubini's Theorem)} \\ &= \int_{\mathcal{X}} k(x, \cdot) g(x, t) d\rho(x, t) && \text{(Pythagoras theorem)} \\ &= Lg \end{aligned}$$

where the second to last equality holds by Pythagoras theorem, since $\bar{\phi}_i$'s are basis for \mathcal{H} . Now we show that $\frac{I\bar{\phi}_i}{\sqrt{\bar{\lambda}_i}}$ are eigenfunctions of L with eigenvalues $\bar{\lambda}_i$.

$$\begin{aligned}
L \frac{\bar{\phi}_i}{\sqrt{\lambda_i}} &= \Pi^* \frac{\bar{\phi}_i}{\sqrt{\lambda_i}} && \text{(Proposition A.9)} \\
&= IC \frac{\bar{\phi}_i}{\sqrt{\lambda_i}} && \text{(Proposition A.9)} \\
&= \bar{\lambda}_i \frac{\bar{\phi}_i}{\sqrt{\lambda_i}} && (\bar{\phi}_i \text{ is an eigenfunction of } C)
\end{aligned}$$

Moreover, they have unit norms and are mutually orthogonal.

$$\begin{aligned}
\left\langle \frac{\bar{\phi}_i}{\sqrt{\lambda_i}}, \frac{\bar{\phi}_j}{\sqrt{\lambda_j}} \right\rangle_\rho &= \frac{1}{\sqrt{\lambda_i \lambda_j}} \langle \Pi^* \bar{\phi}_i, \bar{\phi}_j \rangle_{\mathcal{H}} && \text{(Proposition A.8)} \\
&= \frac{1}{\sqrt{\lambda_i \lambda_j}} \langle C \bar{\phi}_i, \bar{\phi}_j \rangle_{\mathcal{H}} && \text{(Proposition A.9)} \\
&= \frac{1}{\sqrt{\lambda_i \lambda_j}} \lambda_i \delta_{ij} = \delta_{ij}
\end{aligned}$$

□

Similarly, we characterize the relationship between the spectral decomposition of the approximate kernel integral operator L_m and the approximate covariance operator C_m in Proposition A.16.

Proposition A.16. *The spectral decomposition of L_m is:*

$$L_m = \sum_{i=1}^m \lambda_i \frac{A\phi_i}{\sqrt{\lambda_i}} \otimes_\rho \frac{A\phi_i}{\sqrt{\lambda_i}}$$

where $\frac{A\phi_i}{\sqrt{\lambda_i}}$ are the (unit norm) eigenfunctions of L_m with eigenvalues λ_i

Proof of Proposition A.16. First we prove that L_m and $\sum_{i=1}^m \lambda_i \frac{A\phi_i}{\sqrt{\lambda_i}} \otimes_\rho \frac{A\phi_i}{\sqrt{\lambda_i}}$ agree on any function g in $L^2(\mathcal{X}, \rho)$:

$$\begin{aligned}
\sum_{i=1}^m \lambda_i \frac{A\phi_i}{\sqrt{\lambda_i}} \otimes_\rho \frac{A\phi_i}{\sqrt{\lambda_i}} g &= \sum_{i=1}^m \langle A\phi_i, g \rangle_\rho A\phi_i && \text{(Definition of outer product)} \\
&= \sum_{i=1}^m \int_{\mathcal{X}} \langle z(x, t), \phi_i \rangle_{\mathcal{F}} g(x, t) d\rho(x, t) \langle z(\cdot), \phi_i \rangle_{\mathcal{F}} && \text{(Definition A.10)} \\
&= \int_{\mathcal{X}} \sum_{i=1}^m \langle z(x, t), \phi_i \rangle_{\mathcal{F}} \langle z(\cdot), \phi_i \rangle_{\mathcal{F}} g(x, t) d\rho(x, t) && \text{(Fubini's Theorem)} \\
&= \int_{\mathcal{X}} \sum_{i=1}^m \frac{1}{\sqrt{m}} z_{\omega_i}(x, t) \frac{1}{\sqrt{m}} z_{\omega_i}(\cdot) g(x, t) d\rho(x, t) && \text{(Change of basis)} \\
&= \int_{\mathcal{X}} \langle z(x, t), z(\cdot) \rangle_{\mathcal{F}} g(x, t) d\rho(x, t) \\
&= \int_{\mathcal{X}} k_m(x, \cdot) g(x, t) d\rho(x, t) \\
&= L_m g
\end{aligned}$$

The fourth equality follows from the fact that inner products are invariant under orthogonal change of basis. Now we show that $\frac{A\phi_i}{\sqrt{\lambda_i}}$ are eigenfunctions of L_m with the corresponding eigenvalue λ_i

$$L_m \frac{A\phi_i}{\sqrt{\lambda_i}} = AA^* \frac{A\phi_i}{\sqrt{\lambda_i}} \quad (\text{Proposition A.12})$$

$$= \frac{AC_m\phi_i}{\sqrt{\lambda_i}} \quad (\text{Proposition A.12})$$

$$= \lambda_i \frac{A\phi_i}{\sqrt{\lambda_i}} \quad (\phi_i \text{ is an eigenvector of } C_m)$$

Moreover, they have unit norms and are mutually orthogonal.

$$\left\langle \frac{A\phi_i}{\sqrt{\lambda_i}}, \frac{A\phi_j}{\sqrt{\lambda_j}} \right\rangle_\rho = \frac{1}{\sqrt{\lambda_i\lambda_j}} \langle A^*A\phi_i, \phi_j \rangle_{\mathcal{F}} \quad (\text{Proposition A.11})$$

$$= \frac{1}{\sqrt{\lambda_i\lambda_j}} \langle C_m\phi_i, \phi_j \rangle_{\mathcal{F}} \quad (\text{Proposition A.12})$$

$$= \frac{1}{\sqrt{\lambda_i\lambda_j}} \lambda_i \delta_{ij} = \delta_{ij}$$

which completes the proof. \square

The following is an important lemma which shows that the kernel integral operator and its approximation can be seen as true and empirical covariance operators in $HS(\rho)$ associated with the random variable z_ω . This allows us to use concentration of measure tools to bound the approximation error in $H(\rho)$.

Lemma A.17. $L = \mathbb{E}_\omega [z_\omega \otimes_\rho z_\omega], L_m = \frac{1}{m} \sum_{i=1}^m z_{\omega_i} \otimes_\rho z_{\omega_i}$

Proof of Lemma A.17. For any $f, g \in L^2(\mathcal{X}, \rho)$ it holds that

$$\langle Lf, g \rangle_\rho = \left\langle \int_{\mathcal{X}} k(x, \cdot) f(x, t) d\rho(x, t), g \right\rangle_\rho \quad (\text{Definition A.3})$$

$$= \left\langle \int_{\mathcal{X}} \int_{\Omega} z_\omega(x, t) z_\omega(\cdot) f(x, t) d\pi(\omega) d\rho(x, t), g \right\rangle_\rho \quad (\text{Theorem A.6})$$

$$= \left\langle \int_{\Omega} \int_{\mathcal{X}} z_\omega(x, t) f(x, t) d\rho(x, t) z_\omega(\cdot) d\pi(\omega), g \right\rangle_\rho \quad (\text{Fubini's theorem})$$

$$= \left\langle \int_{\Omega} \langle z_\omega(\cdot), f \rangle_\rho z_\omega(\cdot) d\pi(\omega), g \right\rangle_\rho$$

$$= \left\langle \int_{\Omega} z_\omega(\cdot) \otimes_\rho z_\omega(\cdot) d\pi(\omega) f, g \right\rangle_\rho \quad (\text{Definition of outer product})$$

$$= \langle \mathbb{E}_\omega [z_\omega \otimes_\rho z_\omega] f, g \rangle_\rho$$

Similarly, for any $f, g \in L^2(\mathcal{X}, \rho)$ we have that

$$\begin{aligned}
\langle L_m f, g \rangle_\rho &= \left\langle \int_{\mathcal{X}} k_m(x, \cdot) f(x, t) d\rho(x, t), g \right\rangle_\rho && \text{(Definition A.4)} \\
&= \left\langle \int_{\mathcal{X}} \langle z(x, t), z(\cdot) \rangle_{\mathcal{F}} f(x, t) d\rho(x, t), g \right\rangle_\rho && \text{(Definition of } k_m) \\
&= \left\langle \int_{\mathcal{X}} \frac{1}{m} \sum_{i=1}^m z_{\omega_i}(x, t) z_{\omega_i}(\cdot) f(x, t) d\rho(x, t), g \right\rangle_\rho \\
&= \left\langle \frac{1}{m} \sum_{i=1}^m \int_{\mathcal{X}} z_{\omega_i}(x, t) f(x, t) d\rho(x, t) z_{\omega_i}(\cdot), g \right\rangle_\rho \\
&= \left\langle \frac{1}{m} \sum_{i=1}^m \langle z_{\omega_i}(\cdot), f \rangle_\rho z_{\omega_i}(\cdot), g \right\rangle_\rho \\
&= \left\langle \frac{1}{m} \sum_{i=1}^m z_{\omega_i}(\cdot) \otimes_\rho z_{\omega_i}(\cdot) f, g \right\rangle_\rho && \text{(Definition of outer product)} \\
&= \langle L_m f, g \rangle_\rho && \text{(Definition A.4)}
\end{aligned}$$

which completes the proof. \square

The following Proposition shows the relation between the outer products in two separable Hilbert spaces; this is useful in the proof of the main theorem.

Proposition A.18. *For any Hilbert-Schmidt Operator $B : \mathcal{H}_1 \rightarrow \mathcal{H}_2$, it holds that $Bu \otimes_{\mathcal{H}_2} Bv = B(u \otimes_{\mathcal{H}_1} v)B^*$, where $u, v \in \mathcal{H}_1$.*

Proof of Proposition A.18. For any $f \in \mathcal{H}_2$, the following equalities hold:

$$\begin{aligned}
(Bu \otimes_{\mathcal{H}_2} Bv)f &= Bu \langle Bv, f \rangle_{\mathcal{H}_2} && \text{(Definition of outer product)} \\
&= Bu \langle v, B^* f \rangle_{\mathcal{H}_1} && \text{(Proposition A.11)} \\
&= B(u \langle v, B^* f \rangle_{\mathcal{H}_1}) \\
&= B(u \otimes_{\mathcal{H}_1} v)B^* f && \text{(Definition of outer product)}
\end{aligned}$$

\square

Finally, we state the assumptions that we make throughout the paper.

Assumption A.19. The kernel function k is a Mercer's kernel (see Theorem A.5) and has the following integral representation, $k(x, y) = \int_{\Omega} z(x, \omega) z(y, \omega) d\pi(\omega) \forall x, y \in \mathcal{X}$ where (\mathcal{H}, k) is a separable RKHS of real-valued functions on \mathcal{X} with a bounded positive definite kernel k . We also assume that there exists $\tau > 1$ such that $|z(x, \omega)| \leq \tau$ for all $x \in \mathcal{X}, \omega \in \Omega$. Furthermore, we assume that the operator $L^{\frac{1}{2}}$ exists.

B Equivalence of optimization problems in \mathcal{H} and $L^2(\mathcal{X}, \rho)$

We now show that Kernel PCA in the RKHS \mathcal{H} and $L^2(\mathcal{X}, \rho)$ are equivalent under some assumptions which we show are naturally satisfied in the case of Kernel PCA with random features. This allows us to transfer our generalization bounds established in $L^2(\mathcal{X}, \rho)$ to \mathcal{H} .

The Kernel PCA problem essentially reduces to solving the following optimization problem:

$$\max_{P \in \mathcal{P}_{\mathcal{H}}^k} \langle P, C \rangle_{HS(\mathcal{H})} \quad (\text{OPT-1})$$

For any $P \in \mathcal{P}_{HS(\mathcal{H})}^k$, by spectral decomposition, P has an eigendecomposition given by $P = \sum_{i=1}^k u_i \otimes_{\rho} u_i$ where $u_i \in \mathcal{H}, i \in [k]$ are a set of orthonormal functions. We define an operator $U : \mathbb{R}^k \rightarrow \mathcal{H}$ such that $Ub = \sum_{i=1}^k b_i u_i$, where $b \in \mathbb{R}^k$.

Proposition B.1. *U satisfies the following properties.*

- (a) *U is Hilbert-Schmidt.*
- (b) *The adjoint of U is $U^* : \mathcal{H} \rightarrow \mathbb{R}^k$ such that $(U^* f)_i = \langle u_i, f \rangle_{\mathcal{H}}$ where $f \in \mathcal{H}$.*
- (c) *$P = UU^*$ and $U^*U = I_k$*

Proof. (a) First we show that the operator U is Hilbert-Schmidt. Let $\{e_i\}_{i=1}^k$ be the canonical basis of \mathbb{R}^k .

$$\begin{aligned} \|U\|_{\mathcal{L}^2(\mathbb{R}^k, \mathcal{H})}^2 &= \sum_{i=1}^k \|Ue_i\|_{\mathcal{H}}^2 && (\text{Pythagoras Theorem}) \\ &= \sum_{i=1}^k \|u_i\|_{\mathcal{H}}^2 = k \end{aligned}$$

(b) Let U^* be the adjoint of U . We now show that $(U^* f)_i = \langle u_i, f \rangle_{\mathcal{H}}$. For any $b \in \mathbb{R}^k, f \in \mathcal{H}$,

$$\begin{aligned} \langle U^* f, b \rangle &= \langle f, Ub \rangle_{\mathcal{H}} \\ &= \left\langle f, \sum_{i=1}^k b_i u_i \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^k \langle f, u_i \rangle_{\mathcal{H}} b_i \\ &= \langle d, b \rangle \end{aligned}$$

where $d \in \mathbb{R}^k, d_i = \langle f, u_i \rangle_{\mathcal{H}}$

(c) For the first part, for any $f \in \mathcal{H}$, we have,

$$\begin{aligned} Pf &= \sum_{i=1}^k (u_i \otimes_{\mathcal{H}} u_i) f \\ &= \sum_{i=1}^k \langle u_i, f \rangle_{\mathcal{H}} u_i \\ &= \sum_{i=1}^k (U^* f)_i u_i \\ &= UU^* f \end{aligned}$$

Now we show that the constraint $P \in \mathcal{P}_{\mathcal{H}}^k$ reduces to $U^*U = I_k$.

For any $\mathbf{b} \in \mathbb{R}^k$,

$$\begin{aligned} \mathbf{U}^* \mathbf{U} \mathbf{b} &= \sum_{i=1}^k \mathbf{b}_i \mathbf{U}^* \mathbf{u}_i \\ &= \sum_{i=1}^k \mathbf{b}_i \mathbf{d}_i \end{aligned}$$

where $\mathbf{d}_i \in \mathbb{R}^k$, $(\mathbf{d}_i)_j = \langle \mathbf{u}_j, \mathbf{u}_i \rangle_{\mathcal{H}}$. Note that since \mathbf{u}_i 's are orthonormal functions, therefore, $\mathbf{d}_i = \mathbf{e}_i$, where \mathbf{e}_i 's is the canonical basis of \mathbb{R}^k . Therefore,

$$\mathbf{U}^* \mathbf{U} \mathbf{b} = \sum_{i=1}^k \mathbf{b}_i \mathbf{e}_i = \mathbf{b}$$

□

We now write the optimization problem **OPT-1** in terms of \mathbf{U} as,

$$\max_{\mathbf{U}^* \mathbf{U} = \mathbf{I}_k} \langle \mathbf{U} \mathbf{U}^*, \mathbf{C} \rangle_{HS(\mathcal{H})} \quad (\text{OPT-2})$$

Now consider $\mathbf{v}_i \in L^2(\mathcal{X}, \rho)$, $i \in [k]$ such that $\mathbf{u}_i = \mathbf{I}^* \mathbf{v}_i$. Note that the existence of \mathbf{v}_i is guaranteed from the construction of RKHS from eigenfunctions of \mathbf{L} (for details, see [Sejdinovic and Gretton, 2012, Theorem 51]). We now define an operator $\mathbf{V} : \mathbb{R}^k \rightarrow L^2(\mathcal{X}, \rho)$ such that $\mathbf{V} \mathbf{b} = \sum_{i=1}^k \mathbf{b}_i \mathbf{v}_i$, where $\mathbf{b} \in \mathbb{R}^k$. We have the following proposition about \mathbf{V} .

Proposition B.2. *\mathbf{V} satisfies the following properties,*

- (a) *\mathbf{V} is Hilbert-Schmidt.*
- (b) *The adjoint of \mathbf{V} is $\mathbf{V}^* : L^2(\mathcal{X}, \rho) \rightarrow \mathcal{H}$, defined as $(\mathbf{V}^* f)_i = \langle \mathbf{v}_i, f \rangle_{\rho}$*
- (c) *$\langle \mathbf{V} \mathbf{V}^*, \mathbf{L}^2 \rangle_{HS(\rho)} = \langle \mathbf{U} \mathbf{U}^*, \mathbf{C} \rangle_{HS(\mathcal{H})}$ and $\mathbf{V}^* \mathbf{L} \mathbf{V} = \mathbf{U}^* \mathbf{U} = \mathbf{I}_k$*

Proof. The proofs of (a) and (b) are similar to that of Proposition B.1.

(c). We start with the first part. The objective in terms of \mathbf{V} is,

$$\begin{aligned} \langle \mathbf{P}, \mathbf{C} \rangle_{HS(\mathcal{H})} &= \left\langle \sum_{i=1}^k \mathbf{u}_i \otimes_{\mathcal{H}} \mathbf{u}_i, \mathbf{C} \right\rangle_{HS(\mathcal{H})} \\ &= \left\langle \sum_{i=1}^k \mathbf{I}^* \mathbf{v}_i \otimes_{\mathcal{H}} \mathbf{I}^* \mathbf{v}_i, \mathbf{C} \right\rangle_{HS(\mathcal{H})} && (\mathbf{u}_i = \mathbf{I}^* \mathbf{v}_i) \\ &= \left\langle \sum_{i=1}^k \mathbf{I}(\mathbf{v}_i \otimes_{\rho} \mathbf{v}_i) \mathbf{I}^*, \mathbf{C} \right\rangle_{HS(\mathcal{H})} && (\text{Proposition A.18}) \\ &= \left\langle \sum_{i=1}^k \mathbf{v}_i \otimes_{\rho} \mathbf{v}_i, \mathbf{I} \mathbf{C} \mathbf{I}^* \right\rangle_{HS(\rho)} && (\text{Definition of adjoint}) \\ &= \left\langle \sum_{i=1}^k \mathbf{v}_i \otimes_{\rho} \mathbf{v}_i, \mathbf{\Pi}^* \mathbf{\Pi} \right\rangle_{HS(\rho)} && (\text{Proposition A.9}) \\ &= \left\langle \sum_{i=1}^k \mathbf{v}_i \otimes_{\rho} \mathbf{v}_i, \mathbf{L}^2 \right\rangle_{HS(\rho)} && (\text{Proposition A.9}) \\ &= \langle \mathbf{V} \mathbf{V}^*, \mathbf{L}^2 \rangle_{HS(\rho)} \end{aligned}$$

For the second part, for any $\mathbf{b} \in \mathbb{R}^k$, we have,

$$\begin{aligned} \mathbf{U}^* \mathbf{U} \mathbf{b} &= \sum_{i=1}^k \mathbf{b}_i \mathbf{U}^* \mathbf{u}_i \\ &= \sum_{i=1}^k \mathbf{b}_i \mathbf{U}^* \mathbf{I}^* \mathbf{v}_i \quad (\mathbf{u}_i = \mathbf{I}^* \mathbf{v}_i) \\ &= \sum_{i=1}^k \mathbf{b}_i \mathbf{d}_i \end{aligned}$$

where $\mathbf{d}_i \in \mathbb{R}^k$, $(\mathbf{d}_i)_j = \langle \mathbf{u}_j, \mathbf{I}^* \mathbf{v}_i \rangle_{\mathcal{H}} = \langle \mathbf{I}^* \mathbf{v}_j, \mathbf{I}^* \mathbf{v}_i \rangle_{\mathcal{H}} = \langle \mathbf{v}_j, \mathbf{I} \mathbf{I}^* \mathbf{v}_i \rangle_{\rho} = \langle \mathbf{v}_j, \mathbf{L} \mathbf{v}_i \rangle_{\rho}$, where the third equality follows from the property of adjoints, and the last equality because $\mathbf{L} = \mathbf{I} \mathbf{I}^*$. Since $\mathbf{U}^* \mathbf{U} \mathbf{b} = \mathbf{b}$, so $\mathbf{d}_i = \mathbf{e}_i$. Therefore, we get $\langle \mathbf{v}_j, \mathbf{L} \mathbf{v}_i \rangle = \delta_{ij}$.

Let us now look at the j^{th} element of $\mathbf{V}^* \mathbf{L} \mathbf{V} \mathbf{b}$,

$$\begin{aligned} (\mathbf{V}^* \mathbf{L} \mathbf{V} \mathbf{b})_j &= (\mathbf{V}^* \mathbf{L} \sum_{i=1}^k \mathbf{b}_i \mathbf{v}_i)_j \quad (\text{Definition of } \mathbf{V}) \\ &= \sum_{i=1}^k (\mathbf{b}_i \mathbf{V}^* \mathbf{L} \mathbf{v}_i)_j \\ &= \sum_{i=1}^k \mathbf{b}_i \langle \mathbf{v}_j, \mathbf{L} \mathbf{v}_i \rangle_{\rho} \\ &= \sum_{i=1}^k \mathbf{b}_i \delta_{ij} = \mathbf{b}_j \end{aligned}$$

□

We can now restate the optimization problem in terms of \mathbf{V} .

$$\max_{\mathbf{V}^* \mathbf{L} \mathbf{V} = \mathbf{I}_k} \langle \mathbf{V} \mathbf{V}^*, \mathbf{L}^2 \rangle_{HS(\rho)} \quad (\text{OPT-3})$$

Now, let $\mathbf{w}_i = \mathbf{L}^{1/2} \mathbf{v}_i$. Note that \mathbf{w}_i is well-defined since we assume that $\mathbf{L}^{1/2}$ exists (See Assumption 3.1). Define $\mathbf{W} : \mathbb{R}^k \rightarrow L^2(\mathcal{X}, \rho)$, such that $\mathbf{W} \mathbf{b} = \sum_{i=1}^k \mathbf{b}_i \mathbf{w}_i$.

Proposition B.3. *\mathbf{W} satisfies the following properties,*

- (a) *\mathbf{W} is Hilbert-Schmidt.*
- (b) *The adjoint of \mathbf{W} is $\mathbf{W}^* : L^2(\mathcal{X}, \rho) \rightarrow \mathbb{R}^k$ $(\mathbf{W}^* f)_i = \langle \mathbf{w}_i, f \rangle_{\rho}$.*
- (c) *$\mathbf{W} = \mathbf{L}^{1/2} \mathbf{V}$, $\langle \mathbf{V} \mathbf{V}^*, \mathbf{L}^2 \rangle_{HS(\rho)} = \langle \mathbf{W} \mathbf{W}^*, \mathbf{L} \rangle_{HS(\rho)}$ and $\mathbf{W}^* \mathbf{W} = \mathbf{V}^* \mathbf{L} \mathbf{V} = \mathbf{I}_k$*

Proof. The proofs of (a) and (b) are similar to that of Proposition B.1.

(c) For the first part, for any $\mathbf{b} \in \mathbb{R}^k$, we have

$$\begin{aligned} \mathbf{W} \mathbf{b} &= \sum_{i=1}^k \mathbf{b}_i \mathbf{w}_i \\ &= \sum_{i=1}^k \mathbf{b}_i \mathbf{L}^{1/2} \mathbf{v}_i \\ &= \mathbf{L}^{1/2} \sum_{i=1}^k \mathbf{b}_i \mathbf{v}_i \\ &= \mathbf{L}^{1/2} \mathbf{V} \mathbf{b} \end{aligned}$$

Note that since L is self-ajoint, $L^{1/2}$ is self-adjoint too. The objective in terms of W is

$$\begin{aligned}\langle VV^*, L^2 \rangle_{HS(\rho)} &= \langle L^{1/2}VV^*L^{1/2}, L \rangle_{HS(\rho)} && \text{(Definition of adjoint)} \\ &= \langle WW^*, L \rangle_{HS(\rho)}\end{aligned}$$

Equivalently, we can restate the constraint in terms of W as,

$$\begin{aligned}V^*LV &= V^*L^{1/2}L^{1/2}V \\ &= (L^{1/2}V)^*(L^{1/2}V) \\ &= W^*W = I_k\end{aligned}$$

□

We now restate the optimization problem in terms of W .

$$\max_{W^*W=I_k} \langle WW^*, L \rangle_{HS(\rho)} \quad (\text{OPT-4})$$

We now state this equivalence of objective in the following Lemma.

Lemma B.4 (Equivalence of Objective).

$$\langle P, C \rangle_{HS(\mathcal{H})} = \langle WW^*, L \rangle_{HS(\rho)}$$

where the relation between P and W is presented via Propositions [B.1](#), [B.2](#) and [B.3](#).

Proof. One direction of implication simply follows from the construction in Propositions [B.1](#), [B.2](#) and [B.3](#). In particular, from Propositions [B.1](#), [B.2](#) and [B.3](#), we conclude that [OPT-1](#) \implies [OPT-2](#) \implies [OPT-3](#) \implies [OPT-4](#). It is easy to see that [OPT-3](#) \implies [OPT-2](#) \implies [OPT-1](#) where the first implication simply follows from the construction of u_i 's and the second from Proposition [B.1](#). However, showing that [OPT-4](#) \implies [OPT-3](#) is conditioned on w_i 's lying in the range of $L^{1/2}$ because otherwise there might not exist v_i 's such that $w_i = L^{1/2}v_i$. In Lemma [B.5](#), we show that when using random features, with the approximation operator defined in Definition [3.6](#), the functions obtained via random feature approximation lies in the range of $L^{1/2}$ with probability 1 on the support of π . This establishes the equivalence claimed. □

We now formally show that vectors from \mathcal{F} lifted to $L^2(\mathcal{X}, \rho)$ via the approximation operator A lie in the range of $L^{1/2}$ almost surely with respect to measure π . The proof of the following lemma closely follows [[Rudi and Rosasco, 2017](#), Lemma 2].

Lemma B.5. *For every $v \in \mathcal{F}$, $Av \in L^2(\mathcal{X}, \rho)$ lies in the range of $L^{1/2}$ almost surely on the support of π .*

Proof. Let $\Pi \in HS(\rho)$ denote the projection operator projecting to the range of $L^{1/2}$. Then $(I_\rho - \Pi)L^{1/2}f = 0 \forall f \in L^2(\mathcal{X}, \rho)$ as $(I_\rho - \Pi)$ is the projection to the orthogonal complement to the range of $L^{1/2}$. From this, we have, $\text{Tr}((I_\rho - \Pi)L^{1/2}L^{1/2}(I_\rho - \Pi)) = \text{Tr}((I_\rho - \Pi)L(I_\rho - \Pi)) = 0$.

$$\begin{aligned}\text{Tr}((I_\rho - \Pi)L(I_\rho - \Pi)) &= \text{Tr}\left((I_\rho - \Pi) \int_{\Omega} z_\omega \otimes_{\rho} z_\omega d\pi(\omega)(I_\rho - \Pi)\right) \\ &= \int_{\Omega} \text{Tr}((I_\rho - \Pi)(z_\omega \otimes_{\rho} z_\omega)(I_\rho - \Pi)) d\pi(\omega) \\ &= \int_{\Omega} \text{Tr}((I_\rho - \Pi)z_\omega \otimes_{\rho} (I_\rho - \Pi)z_\omega) d\pi(\omega) \\ &= \int_{\Omega} \|(I_\rho - \Pi)z_\omega\|_{\rho}^2 d\pi(\omega) = 0\end{aligned}$$

From the above equation, we see that $\|(I_\rho - \Pi)z_\omega\|_{\rho} = 0$ almost surely on the support of π . This implies that $(I_\rho - \Pi)z_\omega = 0$ a.s. on the support of π .

Now we show that all functions of interest i.e. anything lifted from \mathcal{F} to $L^2(\mathcal{X}, \rho)$ lie in the range of $L^{1/2}$. Let $v \in \mathcal{F}$, $f \in L^2(\mathcal{X}, \rho)$.

$$\begin{aligned} \langle (I_\rho - \Pi)f, Av \rangle_\rho &= \langle A^*(I_\rho - \Pi)f, v \rangle_{\mathcal{F}} \\ &= \sum_{i=1}^m ((I_\rho - \Pi)z_{\omega_i} f) v_i = 0 \end{aligned}$$

This is because ω_i 's are drawn from π , and we already argued that $(I_\rho - \Pi)z_\omega = 0$ a.s. on the support of π . Since this holds for any $v \in \mathcal{F}$ and $f \in L^2(\mathcal{X}, \rho)$, this implies that Av lies in the range of $L^{1/2}$ for all $v \in \mathcal{F}$. \square

Moreover, note that since $L_m = \mathfrak{A}C_m$ (See Proposition A.16), the eigenfunctions of L_m are the lifted eigenvectors of C_m from \mathcal{F} to $L^2(\mathcal{X}, \rho)$. This equivalence entails that any candidate solution of **OPT-2** has an equivalent candidate solution for **OPT-4**, and they would both have the same objective.

As already hinted, the solution of Kernel PCA with random features might not lie in the constraint set of rank k projection operators over $L^2(\mathcal{X}, \rho)$. By the equivalence of the optimization problems **OPT-1** and **OPT-4**, we violate the constraint in \mathcal{H} as well. We, however, remarked that we counter this problem by showing a fast $O(1/\sqrt{n})$ speed of convergence to the constraint set in $L^2(\mathcal{X}, \rho)$. A natural question to ask is does the speed of convergence to the constraint set preserved too? We answer affirmatively as shown below.

We now use this equivalence to give a reduction from a candidate solution **OPT-4** to **OPT-2**. Let $\tilde{P} = \sum_{i=1}^k \tilde{p}_i \otimes_\rho \tilde{p}_i$ be the output of some algorithm for Kernel PCA with random features, lifted through the approximation operator \mathfrak{A} . We have show in Theorem C.1 that $\mathfrak{A}P_{C_m}^k = P_{L_m}^k$ is a rank k projection over $L^2(\mathcal{X}, \rho)$. Let $P_{L_m} = \sum_{i=1}^k \tilde{q}_i \otimes \tilde{q}_i$. Since \tilde{p}_i and \tilde{q}_i lie in the range of $L^{1/2}$, $\forall i \in [k]$ (See Lemma B.5), there exists p_i 's and q_i 's such that $\tilde{p}_i = L^{1/2}p_i$ and $\tilde{q}_i = L^{1/2}q_i$, $i \in [k]$. Define $P := \sum_{i=1}^k I^*p_i \otimes_{\mathcal{H}} I^*p_i$, and $Q := \sum_{i=1}^k I^*q_i \otimes_{\mathcal{H}} I^*q_i$.

First we show that Q is a projection operator in $HS(\mathcal{H})$.

$$\begin{aligned} \langle I^*q_i, I^*q_j \rangle_{\mathcal{H}} &= \langle q_i, \Pi^*q_j \rangle_\rho && \text{(Definition of adjoints)} \\ &= \langle q_i, Lq_j \rangle_\rho && \text{(Proposition A.9)} \\ &= \left\langle L^{1/2}q_i, L^{1/2}q_j \right\rangle_\rho && \text{(Definition of adjoints)} \\ &= \langle \tilde{q}_i, \tilde{q}_j \rangle_\rho = \delta_{ij} \end{aligned}$$

Now, let us look at the rate of convergence P to $\mathcal{P}_{HS(\mathcal{H})}^k$.

Lemma B.6 (Equivalence of convergence to the constraint set).

$$d(\bar{P}, \mathcal{P}_{HS(\mathcal{H})}^k) \leq \left\| \tilde{P} - \mathfrak{A}C_m \right\|_{HS(\rho)}$$

Proof.

$$\begin{aligned}
d(\tilde{P}, \mathcal{P}_{HS(\mathcal{H})}^k) &\leq \|P - Q\|_{HS(\mathcal{H})} \\
&= \left\| \sum_{i=1}^k I^* p_i \otimes_{\mathcal{H}} I^* p_i - I^* q_i \otimes_{\mathcal{H}} I^* q_i \right\|_{HS(\mathcal{H})} \\
&= \left\| I \left(\sum_{i=1}^k p_i \otimes_{\rho} p_i - q_i \otimes_{\rho} q_i \right) I^* \right\|_{HS(\mathcal{H})} && \text{(Proposition A.18)} \\
&= \left\| L \left(\sum_{i=1}^k p_i \otimes_{\rho} p_i - q_i \otimes_{\rho} q_i \right) \right\|_{HS(\rho)} \\
&= \left\| L^{1/2} \left(\sum_{i=1}^k p_i \otimes_{\rho} p_i - q_i \otimes_{\rho} q_i \right) L^{1/2} \right\|_{HS(\rho)} && \text{(Cyclic property)} \\
&= \left\| \sum_{i=1}^k L^{1/2} p_i \otimes_{\rho} L^{1/2} p_i - L^{1/2} q_i \otimes_{\rho} L^{1/2} q_i \right\|_{HS(\rho)} && \text{(Proposition A.18)} \\
&= \left\| \sum_{i=1}^k \tilde{p}_i \otimes_{\rho} \tilde{p}_i - \tilde{q}_i \otimes_{\rho} \tilde{q}_i \right\|_{HS(\rho)} \\
&= \left\| \tilde{P} - \mathfrak{A}C_m \right\|_{HS(\rho)}
\end{aligned}$$

□

In Lemma C.5, we will bound $\left\| \tilde{P} - \mathfrak{A}C_m \right\|_{HS(\rho)}$ which implies the bound given in the main theorem 4.2.

We now combine the above relations into a definition to lift operators from $HS(\mathcal{F})$ to $HS(\mathcal{H})$ and then discuss that the operator is well-defined.

Definition B.7 (Operator \mathfrak{L}). Let $\tilde{P} \in HS(\mathcal{F})$ and $\mathfrak{A}\tilde{P} = \sum_{i=1}^k \tilde{p}_i \otimes_{\rho} \tilde{p}_i$ be \tilde{P} lifted to $L^2(\mathcal{X}, \rho)$. Consider the equivalence relation $p_i \sim p_j$ if $L^{1/2}p_i = L^{1/2}p_j$. Let $[p_i]$ be the equivalence class such that $L^{1/2}p_i = \tilde{p}_i$. The operator $\mathfrak{L} : HS(\mathcal{F}) \rightarrow HS(\mathcal{H})$ is defined as,

$$\mathfrak{L}\hat{P} = \sum_{i=1}^k I^* p_i \otimes_{\mathcal{H}} I^* p_i$$

Here I^* is the restriction of the operator I^* to the quotient space $L^2(\mathcal{X}, \rho) / \sim$.

We now discuss that the operator \mathfrak{L} is indeed a well defined operator. We guarantee by Lemma B.5 that there is at least one element in $[p_i]$ such that $\tilde{p}_i = L^{1/2}p_i$. It remains to argue that all the elements in the equivalence class $[p_i]$ are being mapped to the same element in \mathcal{H} through I^* . Let p_i and p_j be two elements of $[p_i]$. Since $L^{1/2}p_i = L^{1/2}p_j$, therefore $L^{1/2}(p_i - p_j) = 0$. This implies that $p_j = p_i + \text{Ker}(L^{1/2})$. Note that any $r_i \in \text{Ker}(L^{1/2})$ will be mapped by I^* to 0, i.e. $I^*r_i = 0$. It follows from linearity of I^* that $I^*\tilde{p}_j = I^*p_i$. Thus this maps an equivalence class to a single element in \mathcal{H} .

C Proof of the main Theorem

From we have already established the problems in $HS(\rho)$ and $HS(\mathcal{H})$, we focus on error decomposition and bounding the corresponding error terms in $L^2(\mathcal{X}, \rho)$. Solving KPCA by using a kernel approximation, one needs to consider two different sources of error. First, the error coming from approximating the true kernel operator by random features. Second, the statistical error due to estimating the covariance operator using iid samples from the unknown distribution. Thus, we distinguish between and base our proof around these two sources of error, namely *approximation error* and *estimation error*. In particular we decompose our objective as:

$$\begin{aligned} \langle \mathfrak{P}_C^k, \mathfrak{P}_C \rangle_{HS(\rho)} - \langle \mathfrak{A}\hat{\mathfrak{P}}, \mathfrak{P}_C \rangle_{HS(\rho)} &= \underbrace{\langle \mathfrak{P}_C^k, \mathfrak{P}_C \rangle_{HS(\rho)} - \langle \mathfrak{A}\mathfrak{P}_{C_m}^k, \mathfrak{P}_C \rangle_{HS(\rho)}}_{\epsilon_a: \text{Approximation Error}} \\ &\quad + \underbrace{\langle \mathfrak{A}\mathfrak{P}_{C_m}^k, \mathfrak{P}_C \rangle_{HS(\rho)} - \langle \mathfrak{A}\hat{\mathfrak{P}}, \mathfrak{P}_C \rangle_{HS(\rho)}}_{\epsilon_e: \text{Estimation Error}}. \end{aligned}$$

The first term in the decomposition is interpreted as approximation error because it essentially captures the error incurred by approximating the kernel function with random features. The second term in the decomposition is interpreted as estimation error as it is the error incurred in the original statistical estimation problem. In what follows, we give a bound on each of the error terms and provide a detailed analysis. Throughout this section, we use the following Lemma that shows the relation between different projection operators.

Lemma C.1. \mathfrak{P}_C^k and $\mathfrak{A}\mathfrak{P}_L^k$ are rank k projection operators in $L^2(\mathcal{X}, \rho)$. Furthermore, it holds that $\mathfrak{P}_C^k = \mathfrak{P}_L^k$ and $\mathfrak{A}\mathfrak{P}_{C_m}^k = \mathfrak{P}_{L_m}^k$.

Proof of Lemma C.1. We have

$$\mathfrak{P}_C^k = \sum_{i=1}^k \frac{\mathbf{I}\bar{\phi}_i}{\sqrt{\lambda_i}} \otimes_{\rho} \frac{\mathbf{I}\bar{\phi}_i}{\sqrt{\lambda_i}} = \mathfrak{P}_L^k$$

where the second inequality follows from Lemma A.9. Similarly,

$$\mathfrak{A}\mathfrak{P}_{C_m}^k = \sum_{i=1}^k \frac{\mathbf{A}\phi_i}{\sqrt{\lambda_i}} \otimes_{\rho} \frac{\mathbf{A}\phi_i}{\sqrt{\lambda_i}} = \mathfrak{P}_{L_m}^k$$

where the second inequality follows from Lemma A.12. and $\mathfrak{A}\mathfrak{P}_C^k$ □

C.1 Approximation Error

The main idea behind controlling the approximation error is to use the local Rademacher complexity of the kernel class Massart [2000], Bartlett et al. [2002]. More precisely, we use the following result in [Blanchard et al., 2007], which allows us to get rates depending both on the number of features used and the decay of the spectrum of the operator C_2 .

Theorem C.2 (Blanchard et al. [2007]). Assume $\|\zeta\|^2 \leq M$ almost surely, and let (λ_i) denote the ordered eigenvalues of $C := E[\zeta\zeta^\top]$, and further assume that (λ_i) are distinct. Let $B_k := \frac{\sqrt{E[(\langle \zeta, \zeta' \rangle)^4]}}{\lambda_k - \lambda_{k+1}}$, where ζ' is an iid copy of ζ . Then for all δ , with overwhelming probability of at least $1 - e^{-\delta}$ it holds that

$$\langle P_{C^\perp}^k, C \rangle - \langle P_{C^\perp}^k, C \rangle \leq 24\kappa(B_k, k, n) + \frac{11\delta(M + B_k)}{n}$$

where κ is defined as follows:

$$\kappa(B_k, k, n) = \inf_{h \geq 0} \left\{ \frac{B_k h}{n} + \sqrt{\frac{k}{n} \sum_{j > h} \lambda_j(C')} \right\}$$

Lemma C.3 (Approximation Error). With probability at least $1 - \frac{\delta}{2}$, we have

$$\langle P_{L_m^\perp}^k, L \rangle - \langle P_{L^\perp}^k, L \rangle \leq 24\kappa(B_k, k, m) + \frac{11 \log(\delta/2) \tau^2 + 7B_k}{m}$$

Proof. We first note that $\mathfrak{A}P_{C_m}^k = \sum_{i=1}^k \frac{A\phi_i}{\sqrt{\lambda_i}} \otimes \frac{A\phi_i}{\sqrt{\lambda_i}}$ (see definition of the approximation operator in A.10). The following holds for the approximation error:

$$\begin{aligned}
\epsilon_a &= \langle \mathfrak{I}P_C^k, \mathfrak{I}C \rangle_{HS(\rho)} - \langle \mathfrak{A}P_{C_m}^k, \mathfrak{I}C \rangle_{HS(\rho)} \quad (5) \\
&= \left\langle \sum_{i=1}^k \frac{I\bar{\phi}_i}{\sqrt{\lambda_i}} \otimes_{\rho} \frac{I\bar{\phi}_i}{\sqrt{\lambda_i}}, \sum_{i \in I \subset \mathbb{R}} \bar{\lambda}_i \left(\frac{I\bar{\phi}_i}{\sqrt{\lambda_i}} \otimes_{\rho} \frac{I\bar{\phi}_i}{\sqrt{\lambda_i}} \right) \right\rangle_{HS(\rho)} \quad (\text{definition A.7}) \\
&\quad - \left\langle \sum_{i=1}^k \frac{A\phi_i}{\sqrt{\lambda_i}} \otimes_{\rho} \frac{A\phi_i}{\sqrt{\lambda_i}}, \sum_{i \in I \subset \mathbb{R}} \bar{\lambda}_i \left(\frac{I\bar{\phi}_i}{\sqrt{\lambda_i}} \otimes_{\rho} \frac{I\bar{\phi}_i}{\sqrt{\lambda_i}} \right) \right\rangle_{HS(\rho)} \quad (\text{definition A.7, A.10}) \\
&= \langle P_L^k, L \rangle_{HS(\rho)} - \langle P_{L_m}^k, L \rangle_{HS(\rho)} \quad (\text{Lemma A.15 and Lemma C.1}) \\
&= \langle P_{L_m^\perp}^k, L \rangle_{HS(\rho)} - \langle P_{L^\perp}^k, L \rangle_{HS(\rho)} \quad (\text{properties of the orthogonal subspace})
\end{aligned}$$

We have already showed that L and L_m in the right hand side of the equation (5) are true and empirical covariance operators respectively (see Lemma A.17). As required by Theorem C.2, we need to show that norm of the random variables z_ω are bounded. We have

$$\begin{aligned}
\|z_\omega\|^2 &= \langle z_\omega, z_\omega \rangle_{\rho} \\
&= \int_{\mathcal{X}} z_\omega(x, t)^2 d\rho(x, t) \\
&\leq \tau^2
\end{aligned}$$

where the last inequality follows Assumption 3.1.

Invoking Theorem C.2, we have with probability at least $1 - \delta$,

$$\langle P_{L_m^\perp}^k, L \rangle - \langle P_{L^\perp}^k, L \rangle \leq 24\kappa(B_k, k, m) + \frac{11 \log(\delta) \tau^2 + 7B_k}{m}$$

$$\text{where } \kappa(B_k, k, m) = \inf_{h \geq 0} \left\{ \frac{B_k h}{m} + \sqrt{\frac{k \sum_{j>h} \lambda_j(C'_2)}{m}} \right\}. \quad \square$$

Lemma C.4 (Approximation Error - Good decay). *When the spectrum of operator C'_2 has an exponential decay, i.e. $\lambda_j(C'_2) = \alpha^j$ for some $\alpha < 1$, then with probability at least $1 - \delta$, we have*

$$\langle P_{L_m^\perp}^k, L \rangle - \langle P_{L^\perp}^k, L \rangle \leq \frac{24B_k \log(m)}{\log(1/\alpha)m} + \frac{k + (1 - \alpha)(11 \log(\delta) \tau^2 + 7B_k)}{(1 - \alpha)m}$$

Proof. When $\lambda_i(C'_2)$ have an exponential decay, i.e $\lambda_j(C'_2) = \alpha^j$ for some $\alpha < 1$, we have

$$\sum_{j>h} \lambda_j(C'_2) = \frac{\alpha^{h+1}}{1 - \alpha}$$

Set $h = \lceil -\log_\alpha(m) \rceil - 1$, we get

$$\sum_{j>h} \lambda_j(C'_2) \leq \frac{1}{(1 - \alpha)m}$$

Now,

$$\begin{aligned}
\kappa(B_k, k, m) &= \inf_{h \geq 0} \left\{ \frac{B_k h}{m} + \sqrt{\frac{k \sum_{j>h} \lambda_j(C'_2)}{m}} \right\} \\
&\leq \frac{-B_k \log_\alpha m}{m} + \frac{k}{(1 - \alpha)m} \\
&= \frac{B_k \log(m)}{\log(1/\alpha)m} + \frac{k}{(1 - \alpha)m}
\end{aligned}$$

where the last equality follows from the identity $\log_b a = \frac{\log_d(a)}{\log_d(b)}$

So essentially, $\kappa(B_k, k, m) = O\left(\frac{\log(m)}{m}\right)$. Therefore, we get

$$\begin{aligned}\epsilon_a &= \langle \mathfrak{I}P_C^k, \mathfrak{I}C \rangle_{HS(\rho)} - \langle \mathfrak{A}P_{C_m}^k, \mathfrak{I}C \rangle_{HS(\rho)} \leq \frac{24B_k \log(m)}{\log(1/\alpha)m} + \frac{k}{(1-\alpha)m} + \frac{11 \log(\delta) \tau^2 + 7B_k}{m} \\ &= \frac{24B_k \log(m)}{\log(1/\alpha)m} + \frac{k + (1-\alpha)(11 \log(\delta) \tau^2 + 7B_k)}{(1-\alpha)m}\end{aligned}$$

which completes the proof. \square

C.2 Estimation Error

We first remind the reader that $\mathfrak{A}P_{C_m}^k$ is a projection operator in $L^2(\mathcal{X}, \rho)$ (See Lemma C.1). However, the problem we face is that $\mathfrak{A}\hat{P}$ might not be a projection operator in $L^2(\mathcal{X}, \rho)$. This is because the lifting is accomplished by lifting a particular set of eigenvectors of $\hat{P}_{\mathcal{A}}$ through A , and we remark that A doesn't necessarily preserve norms and angles between elements. To get around this predicament, we show that lifted operator converges to a projection operator, i.e the lifted set of eigenvectors go to an orthogonal set of functions in $L^2(\mathcal{X}, \rho)$. Moreover, from Lemma B.6, we have that this convergence in $HS(\rho)$ is equivalent to convergence in $HS(\mathcal{H})$.

Lemma C.5. *When the number of samples $n \geq \frac{2\lambda_1^2 q_{\mathcal{A}}(1/\delta, \log(m), \log(n))^2}{\lambda_k^2(\sqrt{2}-1)}$, with probability at least $1 - \frac{\delta}{2}$, we have*

$$\begin{aligned}d(\mathfrak{A}C_m, \mathcal{P}_{HS(\rho)}^k) &\leq \left\| \mathfrak{A}P_{C_m}^k - \mathfrak{A}\hat{P}_{\mathcal{A}} \right\|_{HS(\rho)} \leq \left\| \mathfrak{A}P_{C_m}^k - \mathfrak{A}\hat{P}_{\mathcal{A}} \right\|_{\mathcal{L}^1(\rho)} \\ &\leq \frac{\lambda_1}{(\sqrt{2}-1)} \sqrt{\sum_{i=1}^k \left(\frac{2\lambda_i + 4\lambda_1}{\lambda_i^2} \right)^2} \frac{q_{\mathcal{A}}(1/\delta, \log(m), \log(n))}{n}\end{aligned}$$

Proof of Lemma C.5. Since $\mathfrak{A}C_m$ is a rank k projection operator in $HS(\rho)$ (from Lemma C.1), the first inequality follows trivially. The second inequality is just from the property of norms that Schatten norms $\|D\|_{\mathcal{L}^p(\rho)}$ decreases with increasing p . We focus on proving the third inequality below.

Let $\hat{P}_{\mathcal{A}} = \tilde{\Phi}\tilde{\Phi}^\top$ be an eigendecomposition of the output $\hat{P}_{\mathcal{A}}$. Let

$$R^* = \arg \min_{R^\top R = R R^\top = I} \left\| \tilde{\Phi}R - \Phi_k \right\|_F^2$$

where Φ_k is the matrix corresponding top k eigenvectors of C_m . Define $\hat{\Phi} := \tilde{\Phi}R^*$. This means that we rotate the eigenvectors of our output to a basis such that it is closest to the truth (in element-wise metric sense). An important point on why we can do this is that this rotation (or any other rotation for that matter) doesn't change the output, i.e. $\hat{\Phi}\hat{\Phi}^\top = \tilde{\Phi}R^*R^{*\top}\tilde{\Phi}^\top = \tilde{\Phi}\tilde{\Phi}^\top = \hat{P}_{\mathcal{A}}$. We now lifting the output by lifting this rotated set of eigenvectors. We have, $\mathfrak{A}\hat{P} = \sum_{i=1}^k \frac{A\hat{\phi}_i}{\sqrt{\hat{\lambda}_i}} \otimes_\rho \frac{A\hat{\phi}_i}{\sqrt{\hat{\lambda}_i}}$, where $\hat{\lambda}_i := \langle \hat{\phi}_i, C_m \hat{\phi}_i \rangle_{\mathcal{F}}$.

$$\begin{aligned}
\|\mathfrak{A}P_{C_m}^k - \mathfrak{A}\hat{P}_A\|_{\mathcal{L}^1(\rho)} &= \left\| \sum_{i=1}^k \frac{A\phi_i}{\sqrt{\lambda_i}} \otimes \frac{A\phi_i}{\sqrt{\lambda_i}} - \sum_{i=1}^k \frac{A\hat{\phi}_i}{\sqrt{\hat{\lambda}_i}} \otimes \frac{A\hat{\phi}_i}{\sqrt{\hat{\lambda}_i}} \right\|_{\mathcal{L}^1(\rho)} \\
&= \left\| A \sum_{i=1}^k \left(\frac{\phi_i}{\sqrt{\lambda_i}} \otimes \frac{\phi_i}{\sqrt{\lambda_i}} - \frac{\hat{\phi}_i}{\sqrt{\hat{\lambda}_i}} \otimes \frac{\hat{\phi}_i}{\sqrt{\hat{\lambda}_i}} \right) A^* \right\|_{\mathcal{L}^1(\rho)} \\
&\leq \|A\| \left\| \sum_{i=1}^k \left(\frac{1}{\lambda_i} \phi_i \otimes \phi_i - \frac{1}{\hat{\lambda}_i} \hat{\phi}_i \otimes \hat{\phi}_i \right) A^* \right\|_{\mathcal{L}^1(\mathcal{F}, \rho)} \\
&\leq \|A\| \|A^*\| \left\| \sum_{i=1}^k \left(\frac{1}{\lambda_i} \phi_i \otimes \phi_i - \frac{1}{\hat{\lambda}_i} \hat{\phi}_i \otimes \hat{\phi}_i \right) \right\|_{\mathcal{L}^1(\mathcal{F})} \\
&\leq \lambda_1 \left\| \sum_{i=1}^k \left(\frac{1}{\lambda_i} \phi_i \otimes \phi_i - \frac{1}{\hat{\lambda}_i} \hat{\phi}_i \otimes \hat{\phi}_i \right) \right\|_{\mathcal{L}^1(\mathcal{F})}
\end{aligned}$$

Where third and fourth inequalities follows from the fact that for trace-class operators $\|AB\|_{\mathcal{L}^1} \leq \|A\|_2 \|B\|_{\mathcal{L}^1}$. See [Reed and Simon, 1972, Exercise 28, Page 218].

Adding and subtracting $\frac{1}{\lambda_i} \hat{\phi}_i \otimes \hat{\phi}_i$ inside the summation to get

$$\begin{aligned}
&\leq \lambda_1 \left\| \sum_{i=1}^k \frac{1}{\lambda_i} \phi_i \otimes \phi_i - \frac{1}{\lambda_i} \hat{\phi}_i \otimes \hat{\phi}_i + \frac{1}{\lambda_i} \hat{\phi}_i \otimes \hat{\phi}_i - \frac{1}{\hat{\lambda}_i} \hat{\phi}_i \otimes \hat{\phi}_i \right\|_{\mathcal{L}^1(\mathcal{F})} \\
&\leq \lambda_1 \left\| \sum_{i=1}^k \left(\frac{1}{\lambda_i} (\phi_i \otimes \phi_i - \hat{\phi}_i \otimes \hat{\phi}_i) + \left(\frac{1}{\lambda_i} - \frac{1}{\hat{\lambda}_i} \right) \hat{\phi}_i \otimes \hat{\phi}_i \right) \right\|_{\mathcal{L}^1(\mathcal{F})} \\
&\leq \lambda_1 \sum_{i=1}^k \frac{1}{\lambda_i} \left\| \phi_i \otimes \phi_i - \hat{\phi}_i \otimes \hat{\phi}_i \right\|_{\mathcal{L}^1(\mathcal{F})} + \left| \frac{1}{\lambda_i} - \frac{1}{\hat{\lambda}_i} \right| \\
&\leq \lambda_1 \sum_{i=1}^k \frac{1}{\lambda_i} \left\| \phi_i \otimes \phi_i - \hat{\phi}_i \otimes \hat{\phi}_i \right\|_{\mathcal{L}^1(\mathcal{F})} + \left| \frac{\lambda_i - \hat{\lambda}_i}{\lambda_i \hat{\lambda}_i} \right| \\
&\leq \lambda_1 \sum_{i=1}^k \frac{2}{\lambda_i} \left\| \phi_i - \hat{\phi}_i \right\|_2 + \frac{4\lambda_1}{\lambda_i^2} \left\| \phi_i - \hat{\phi}_i \right\|_2 \\
&\leq \lambda_1 \sum_{i=1}^k \left(\frac{2\lambda_i + 4\lambda_1}{\lambda_i^2} \right) \left\| \phi_i - \hat{\phi}_i \right\|_2 \\
&\leq \lambda_1 \sqrt{\sum_{i=1}^k \left(\frac{2\lambda_i + 4\lambda_1}{\lambda_i^2} \right)^2} \left\| \Phi_k - \hat{\Phi} \right\|_{\mathcal{F}} \\
&\leq \frac{\lambda_1}{2(\sqrt{2}-1)} \sqrt{\sum_{i=1}^k \left(\frac{2\lambda_i + 4\lambda_1}{\lambda_i^2} \right)^2} \left\| P_{C_m}^k - \hat{P} \right\|_F^2 \\
&\leq \frac{\lambda_1}{(\sqrt{2}-1)} \sqrt{\sum_{i=1}^k \left(\frac{2\lambda_i + 4\lambda_1}{\lambda_i^2} \right)^2} \frac{q_A(1/\delta, \log(m), \log(n))}{n}
\end{aligned}$$

The second to last inequality follows from Lemma C.6 and Lemma C.7. \square

Lemma C.6. $\|\phi_i \otimes \phi_i - \hat{\phi}_i \otimes \hat{\phi}_i\|_{\mathcal{L}^1(\mathcal{F})} \leq 2\|\phi_i - \hat{\phi}_i\|_2 \forall i \in [k]$

Proof.

$$\begin{aligned}
\|\phi_i \otimes \phi_i - \widehat{\phi}_i \otimes \widehat{\phi}_i\|_{\mathcal{L}^1(\mathcal{F})} &= \|\phi_i \otimes \phi_i - \widehat{\phi}_i \otimes \phi_i + \widehat{\phi}_i \otimes \phi_i - \widehat{\phi}_i \otimes \widehat{\phi}_i\|_{\mathcal{L}^1(\mathcal{F})} \\
&\leq \|(\phi_i - \widehat{\phi}_i) \otimes \phi_i\|_{\mathcal{L}^1(\mathcal{F})} + \|\widehat{\phi}_i \otimes (\phi_i - \widehat{\phi}_i)\|_{\mathcal{L}^1(\mathcal{F})} \\
&= 2\|\phi_i - \widehat{\phi}_i\|_2
\end{aligned}$$

□

Lemma C.7. When the number of samples $n \geq \frac{2q_{\mathcal{A}}(1/\delta, \log(m), \log(n))^2 \lambda_1^2}{\lambda_i^2(\sqrt{2}-1)}$, with probability at least $1 - \delta$, $\forall i \in [k]$ we have,

$$\left| \frac{\lambda_i - \widehat{\lambda}_i}{\lambda_i \widehat{\lambda}_i} \right| \leq \frac{4\lambda_1}{\lambda_i^2} \|\phi_i - \widehat{\phi}_i\|_2$$

where $C_{\mathcal{A}}$ is a constant specific to the algorithm \mathcal{A} .

The numerator is bounded as follows

Proof.

$$\begin{aligned}
|\lambda_i - \widehat{\lambda}_i| &= |\phi_i^\top C_m \phi_i - \widehat{\phi}_i^\top C_m \widehat{\phi}_i| \\
&= |\phi_i^\top C_m \phi_i - \phi_i^\top C_m \widehat{\phi}_i + \phi_i^\top C_m \widehat{\phi}_i - \widehat{\phi}_i^\top C_m \widehat{\phi}_i| \\
&= |\phi_i^\top C_m (\phi_i - \widehat{\phi}_i) + (\phi_i - \widehat{\phi}_i)^\top C_m \widehat{\phi}_i| \\
&\leq \|C_m \phi_i\|_2 \|\phi_i - \widehat{\phi}_i\|_2 + \|\phi_i - \widehat{\phi}_i\|_2 \|C_m \widehat{\phi}_i\|_2 \\
&= (\lambda_i + \widehat{\lambda}_i) \|\phi_i - \widehat{\phi}_i\|_2 \\
&\leq (\lambda_i + \lambda_1) \|\phi_i - \widehat{\phi}_i\|_2 \\
&\leq 2\lambda_1 \|\phi_i - \widehat{\phi}_i\|_2
\end{aligned}$$

where the second inequality holds since $\widehat{\lambda}_i < \lambda_1$ by definition of $\widehat{\lambda}_i$, and the last inequality follows because $\widehat{\lambda}_i \leq \lambda_1$.

The denominator is lower bounded similarly as

$$\begin{aligned}
\lambda_i \widehat{\lambda}_i &\geq \lambda_i (\lambda_i - 2\lambda_1 \|\phi_i - \widehat{\phi}_i\|_2) \\
&\geq \frac{\lambda_i^2}{2}
\end{aligned}$$

where the first inequality follows from the bound on the numerator and the last inequality follows when $2\lambda_1 \|\phi_i - \widehat{\phi}_i\|_2 \leq \frac{\lambda_i}{2}$. From Lemma C.10, we know that $\|\phi_i - \widehat{\phi}_i\|_2^2 \leq \frac{1}{2(\sqrt{2}-1)} \left(\frac{q_{\mathcal{A}}(1/\delta, \log(m), \log(n))}{n} \right)$ with probability at least $1 - \delta$. Combining, we get, with probability at least $1 - \delta$,

$$\|\phi_i - \widehat{\phi}_i\|_2 \leq \sqrt{\frac{1}{2(\sqrt{2}-1)} \left(\frac{q_{\mathcal{A}}(1/\delta, \log(m), \log(n))}{n} \right)} \leq \frac{\lambda_i}{4\lambda_1}$$

The above holds when the number of samples $n \geq \frac{2\lambda_1^2 q_{\mathcal{A}}(1/\delta, \log(m), \log(n))^2}{\lambda_i^2(\sqrt{2}-1)}$. Combining, we get

$$\left| \frac{\lambda_i - \widehat{\lambda}_i}{\lambda_i \widehat{\lambda}_i} \right| \leq \frac{4\lambda_1}{\lambda_i^2} \|\phi_i - \widehat{\phi}_i\|_2$$

with probability at least $1 - \delta$ and when $n \geq \frac{2\lambda_1^2 q_{\mathcal{A}}(1/\delta, \log(m), \log(n))^2}{\lambda_i^2(\sqrt{2}-1)}$. □

Note that in particular since Oja's algorithm has a warm-up phase, the lower bound on the denominator

Lemma C.8. For rank k orthogonal matrices $U \in \mathbb{R}^{m \times k}$ and $V \in \mathbb{R}^{m \times k}$, i.e. $U^\top U = V^\top V = I_k$, the following holds,

$$\|U - V\hat{R}\|_F^2 \leq \frac{1}{2(\sqrt{2} - 1)} \|UU^\top - VV^\top\|_F^2,$$

where

$$\hat{R} = \arg \min_{R^\top R = RR^\top = I_k} \|U - VR\|_F^2$$

Proof. Proof in [Ge et al., 2017, Lemma 6]. \square

Since Φ_k and $\hat{\Phi}$ are rank k orthogonal matrices, from Lemma C.8, we have

$$\|\Phi_k - \hat{\Phi}\|_F^2 \leq \frac{1}{2(\sqrt{2} - 1)} \|P_{C_m}^k - \hat{P}\|_F^2$$

Lemma C.9. For any efficient subspace learner \mathcal{A} , we have $\|P_{C_m}^k - \hat{P}\|_F^2 \leq \frac{2q_{\mathcal{A}}(1/\delta, \log(m), \log(n))}{n}$ with probability at least $1 - \frac{\delta}{2}$.

Proof.

$$\begin{aligned} \|P_{C_m}^k - \hat{P}\|_F^2 &= \|P_{C_m}^k\|_F^2 + \|\hat{P}\|_F^2 - 2\langle \hat{P}, P_{C_m}^k \rangle \\ &= 2\left(k - \langle \hat{P}, P_{C_m}^k \rangle\right) \\ &= 2\left(\langle I - P_{C_m}^k, \hat{P} \rangle\right) \\ &= 2\left\|(\Phi_k^\perp)^\top \hat{\Phi}\right\|_F^2 \\ &\leq \frac{2q_{\mathcal{A}}(1/\delta, \log(m), \log(n))}{n} \end{aligned}$$

where the last inequality follows from the definition of efficient subspace learner. \square

Lemma C.10. With probability at least $1 - \delta$,

$$\|\phi_i - \hat{\phi}_i\|_2 \leq \frac{1}{2(\sqrt{2} - 1)} \left(\frac{q_{\mathcal{A}}(1/\delta, \log(m), \log(n))}{n} \right)$$

where $q_{\mathcal{A}}(1/\delta, \log(m), \log(n))$ is specific to the algorithm \mathcal{A} .

Proof.

$$\begin{aligned} \|\phi_i - \hat{\phi}_i\|_2^2 &\leq \sum_{i=1}^k \|\phi_i - \hat{\phi}_i\|_2^2 \\ &= \|\Phi_k - \hat{\Phi}\|_F^2 \\ &\leq \frac{1}{2(\sqrt{2} - 1)} \|P_{C_m}^k - \hat{P}\|_F^2 \\ &\leq \frac{1}{2(\sqrt{2} - 1)} \left(\frac{q_{\mathcal{A}}(1/\delta, \log(m), \log(n))}{n} \right) \end{aligned}$$

\square

where the second inequality holds from Lemma C.8 and the definition of \hat{P} , and the last inequality holds from Lemma C.9

Lemma C.11 (Estimation Error). *When the number of samples $n \geq \frac{2\lambda_1^2 q_{\mathcal{A}}(1/\delta, \log(m), \log(n))^2}{\lambda_k^2(\sqrt{2}-1)}$, then with probability at least $1 - \delta$, we have*

$$\epsilon_e \leq \frac{\lambda_1^2}{(\sqrt{2}-1)} \sqrt{\sum_{i=1}^k \left(\frac{2\lambda_i + 4\lambda_1}{\lambda_i^2} \right)^2 \frac{q_{\mathcal{A}}(1/\delta, \log(m), \log(n))^2}{n}}$$

Proof.

$$\begin{aligned} \epsilon_e &= \langle \mathfrak{A}P_{C_m}^k, \mathfrak{J}C \rangle_{HS(\rho)} - \langle \mathfrak{A}\hat{P}_{\mathcal{A}}, \mathfrak{J}C \rangle_{HS(\rho)} \\ &= \left\langle \mathfrak{A}P_{C_m}^k - \mathfrak{A}\hat{P}_{\mathcal{A}}, \mathfrak{J}C \right\rangle_{HS(\rho)} \\ &\leq \left\| \mathfrak{A}P_{C_m}^k - \mathfrak{A}\hat{P}_{\mathcal{A}} \right\|_{\mathcal{L}^1(\rho)} \|\mathfrak{J}C\|_2 \\ &\leq \lambda_1 \left\| \mathfrak{A}P_{C_m}^k - \mathfrak{A}\hat{P}_{\mathcal{A}} \right\|_{\mathcal{L}^1(\rho)} \\ &\leq \frac{\lambda_1^2}{(\sqrt{2}-1)} \sqrt{\sum_{i=1}^k \left(\frac{2\lambda_i + 4\lambda_1}{\lambda_i^2} \right)^2 \frac{q_{\mathcal{A}}(1/\delta, \log(m), \log(n))}{n}} \end{aligned}$$

The last inequality follows from Lemma C.5. \square

We now invoke the approximation and the estimation error bounds i.e. Lemma C.3 and Lemma C.11 with failure probabilities $\delta/2$ each. We then apply a union bound over them and get that with probability at least $1 - \delta$,

$$\langle \mathfrak{J}P_C^k, \mathfrak{J}C \rangle_{\rho} - \langle \mathfrak{A}\hat{P}_{\mathcal{A}}, \mathfrak{J}C \rangle_{\rho} \leq \frac{cB_k}{\sqrt{n}} + \frac{c'(k + \log(\delta/2) + 7B_k)}{\sqrt{n} \log(n)} + \sqrt{\frac{q_{\mathcal{A}}(2/\delta, \log(m), \log(n))}{n}},$$

This concludes the proof of the main theorem.

Also note that since $d(\mathfrak{A}\hat{P}_{\mathcal{A}}, \mathcal{P}_{HS(\rho)})$ decays as $O(1/\sqrt{n})$, we can bound the suboptimality of $\mathfrak{A}\hat{P}_{\mathcal{A}}$ projected onto the set of projection operators $\mathcal{P}_{HS(\rho)}^k$. It is now easy to give a bound on the objective with respect to the projection $\tilde{P}_{\mathcal{A}} \in \mathcal{P}_{HS(\rho)}$ of $\mathfrak{A}\hat{P}_{\mathcal{A}}$ onto the set of projection operators:

Corollary C.12. *Let $\tilde{P}_{\mathcal{A}}$ be the projection of $\mathfrak{A}\hat{P}_{\mathcal{A}}$ onto the set $\mathcal{P}_{HS(\rho)}$. Under the same conditions as in theorem 4.2, we have*

$$\langle \mathfrak{J}P_C^k, \mathfrak{J}C \rangle_{\rho} - \langle \tilde{P}_{\mathcal{A}}, \mathfrak{J}C \rangle_{\rho} \leq 24\kappa(B_k, k, m) + \frac{\log(\delta/2) + 7B_k}{m} + 2\sqrt{\frac{q_{\mathcal{A}}(2/\delta, \log(m), \log(n))}{n}}$$

Proof.

$$\begin{aligned} \langle \mathfrak{J}P_C^k, \mathfrak{J}C \rangle_{\rho} - \langle \tilde{P}_{\mathcal{A}}, \mathfrak{J}C \rangle_{\rho} &= \langle \mathfrak{J}P_C^k, \mathfrak{J}C \rangle_{\rho} - \langle \mathfrak{A}\hat{P}_{\mathcal{A}}, \mathfrak{J}C \rangle_{\rho} + \langle \mathfrak{A}\hat{P}_{\mathcal{A}} - \tilde{P}_{\mathcal{A}}, \mathfrak{J}C \rangle \\ &\leq 24\kappa(B_k, k, m) + \frac{\log(\delta/2) + 7B_k}{m} + \sqrt{\frac{q_{\mathcal{A}}(2/\delta, \log(m), \log(n))}{n}} \\ &\quad + d(\mathfrak{A}\hat{P}_{\mathcal{A}}, \mathcal{P}_{HS(\rho)}) \|\mathfrak{J}C\|_{HS(\rho)} \\ &\leq 24\kappa(B_k, k, m) + \frac{\log(\delta/2) + 7B_k}{m} + 2\sqrt{\frac{q_{\mathcal{A}}(2/\delta, \log(m), \log(n))}{n}}, \end{aligned}$$

where the second to last inequality follows from Cauchy-Schwartz in $HS(\rho)$. \square

We now give the proof of Corollary 4.3.

Proof of Corollary 4.3.

$$\begin{aligned}
\langle \mathfrak{P}_C^k, \mathfrak{J}C \rangle_\rho - \langle \hat{\mathfrak{P}}_{\mathcal{A}}, \mathfrak{J}C \rangle_\rho &= \langle \mathfrak{P}_C^k, \mathfrak{J}C \rangle_{HS(\rho)} - \langle \mathfrak{P}_{C_m}^k, \mathfrak{J}C \rangle_{HS(\rho)} \\
&\quad + \langle \mathfrak{P}_{C_m}^k, \mathfrak{J}C \rangle_{HS(\rho)} - \langle \hat{\mathfrak{P}}, \mathfrak{J}C \rangle_{HS(\rho)} \\
&\leq \frac{24B_k \log(m)}{\log(1/\alpha)m} + \frac{k + (1-\alpha)(11 \log(\delta/2)M + 7B_k)}{(1-\alpha)m} \\
&\quad + \lambda_1 \sqrt{\frac{q_{\mathcal{A}}(2/\delta, \log(m), \log(n))}{n}},
\end{aligned}$$

with probability at least $1 - \delta$. The last inequality follows from Lemma C.4 and Lemma C.11 with a union bound over them. \square

D Examples of ESL

In this section, we instantiate our framework with two popular learning algorithms, Empirical Risk Minimization (ERM) and Oja's Algorithm, and show that they satisfy the requirements of ESL.

D.1 Empirical Risk Minimizer

A natural candidate for an efficient subspace learner is the Empirical Risk Minimizer, which we call as \mathcal{A}_{ERM} . We first show that \mathcal{A}_{ERM} satisfies the sufficient condition of Definition 4.1 and then show that \mathcal{A}_{ERM} is an efficient subspace learner. We then discuss its computational aspects. Let $\{\mathbf{x}_i\}_{i=1}^n$ be n data samples and $\{\mathbf{z}(\mathbf{x})\}_{i=1}^n$ be the corresponding representations in \mathcal{F} . The empirical covariance matrix in \mathcal{F} is defined as

$$\hat{\mathbf{C}}_m = \frac{1}{n} \sum_{i=1}^n \mathbf{z}(\mathbf{x}_i) \mathbf{z}(\mathbf{x}_i)^\top$$

The algorithm \mathcal{A}_{ERM} computes the top k eigenvectors of $\hat{\mathbf{C}}_m$, and returns a rank k orthogonal matrix say $\hat{\Phi}$. Let the corresponding projection matrix be $\hat{\mathbf{P}}_{ERM}$. We first state the bound on covariance matrices \mathbf{C}_m and $\hat{\mathbf{C}}_m$.

Lemma D.1 (Covariance Estimation). *With probability at least $1 - \delta$,*

$$\|\hat{\mathbf{C}}_m - \mathbf{C}_m\|_2 \leq \frac{\kappa}{3n} \log\left(\frac{\delta}{2m}\right) + \sqrt{\frac{\kappa}{3n} \log\left(\frac{\delta}{2m}\right)^2 + \log\left(\frac{\delta}{2m}\right) \frac{\kappa \lambda_1}{n}}$$

Proof.

$$\begin{aligned} \|\hat{\mathbf{C}}_m - \mathbf{C}_m\|_2 &= \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{z}(\mathbf{x}_i) \mathbf{z}(\mathbf{x}_i)^\top - \mathbf{C}_m \right\|_2 \\ &= \left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{z}(\mathbf{x}_i) \mathbf{z}(\mathbf{x}_i)^\top - \mathbf{C}_m) \right\|_2 \\ &= \left\| \sum_{i=1}^n \Xi_i \right\|_2 \end{aligned}$$

where $\Xi_i = \frac{1}{n} (\mathbf{z}(\mathbf{x}_i) \mathbf{z}(\mathbf{x}_i)^\top - \mathbf{C}_m)$. Ξ_i 's are 0 mean random matrices, i.e. $\mathbb{E}[\Xi_i] = 0 \forall i \in [n]$.

Note that $\|\mathbf{z}(\mathbf{x})\|_2^2 = \int_{\mathcal{X}} z_\omega(\mathbf{x})^2 d\rho(\mathbf{x}) \leq \tau^2$, since $z_\omega(\mathbf{x}) \leq \tau \forall \omega \in \Omega, \mathbf{x} \in \mathcal{X}$ by Assumption 3.1. We have,

$$\begin{aligned} \|\Xi_i\|_2 &\leq \frac{1}{n} (\|\mathbf{z}(\mathbf{x}_i) \mathbf{z}(\mathbf{x}_i)^\top\|_2 + \|\mathbf{C}_m\|_2) \\ &= \frac{1}{n} (\text{Tr}(\mathbf{z}(\mathbf{x}_i) \mathbf{z}(\mathbf{x}_i)^\top) + \|\mathbb{E}_{\mathbf{x}} [\mathbf{z}(\mathbf{x}) \mathbf{z}(\mathbf{x})^\top]\|_2) \\ &\leq \frac{1}{n} (\|\mathbf{z}(\mathbf{x}_i)\|_2^2 + \mathbb{E}_{\mathbf{x}} [\|\mathbf{z}(\mathbf{x}) \mathbf{z}(\mathbf{x})^\top\|_2]) \\ &\leq \frac{1}{n} (\|\mathbf{z}(\mathbf{x}_i)\|_2^2 + \mathbb{E}_{\mathbf{x}} [\|\mathbf{z}(\mathbf{x})\|_2^2]) \\ &\leq \frac{2\tau^2}{n} \end{aligned}$$

so that $L(\Xi) := \max_i \{\|\Xi_i\|_2\} \leq \frac{2\tau^2}{n}$. where in the second inequality, we apply Jensen's inequality. Define $v(\Xi) := \left\| \sum_{i=1}^n \mathbb{E} [\Xi_i \Xi_i^\top] \right\|_2$. We have,

$$\begin{aligned}
v(\Xi) &= \left\| \sum_{i=1}^n \mathbb{E} [\Xi_i \Xi_i^\top] \right\|_2 \\
&= \left\| \sum_{i=1}^n \mathbb{E} \left[\frac{1}{n^2} (z(x_i)z(x_i)^\top - C_m) (z(x_i)z(x_i)^\top - C_m)^\top \right] \right\|_2 \\
&= \frac{1}{n^2} \left\| \sum_{i=1}^n \mathbb{E} \left[\|z(x_i)\|^2 z(x_i)^\top z(x_i) - z(x_i)^\top z(x_i) C_m - C_m z(x_i)^\top z(x_i) + C_m^2 \right] \right\|_2 \\
&\leq \frac{1}{n^2} \left\| \sum_{i=1}^n \mathbb{E} \left[\tau^2 z(x_i)^\top z(x_i) - z(x_i)^\top z(x_i) C_m - C_m z(x_i)^\top z(x_i) + C_m^2 \right] \right\|_2 \\
&= \frac{1}{n^2} \left\| \sum_{i=1}^n \tau^2 C_m - C_m^2 - C_m^2 + C_m^2 \right\|_2 \\
&= \frac{1}{n} \left\| \tau^2 C_m - C_m^2 \right\|_2 \\
&\leq \frac{\tau^2}{n} \|C_m\|_2 \\
&\leq \frac{\tau^2 \lambda_1}{n}
\end{aligned}$$

where the second last inequality holds because C_m is a positive semi-definite matrix. From matrix Bernstein concentration (Lemma F.2, restated from [Tropp et al., 2015]), we have, with probability at least $1 - \delta$

$$\begin{aligned}
\left\| \hat{C}_m - C_m \right\|_2 &= \left\| \sum_{i=1}^m \Xi_i \right\|_2 \leq \frac{L(\Xi)}{6} \log \left(\frac{\delta}{2m} \right) + \sqrt{\frac{L(\Xi)^2}{12} \log \left(\frac{\delta}{2m} \right)^2 + \log \left(\frac{\delta}{2m} \right) v(\Xi)} \\
&\leq \frac{\tau^2}{3n} \log \left(\frac{\delta}{2m} \right) + \sqrt{\frac{\tau^2}{3n} \log \left(\frac{\delta}{2m} \right)^2 + \log \left(\frac{\delta}{2m} \right) \frac{\tau^2 \lambda_1}{n}}
\end{aligned}$$

□

In the following lemma, we show that \mathcal{A}_{ERM} is an efficient subspace learner.

Lemma D.2. \mathcal{A}_{ERM} is an efficient subspace learner.

Proof. We invoke Theorem F.1 with the sub-multiplicative norm being the spectral norm. With $A = C_m, B = \hat{C}_m, U = \hat{\Phi}, V = \Phi_k^\perp$. Let $\epsilon = \left\| C_m - \hat{C}_m \right\|_2$. From Weyl's inequality, we have $\lambda_k(\hat{C}_m) \geq \lambda_k - \epsilon = \lambda_{k+1} + \text{gap} - \epsilon \geq \lambda_{k+1}$, if $\epsilon < \text{gap}$. Therefore, setting $\mu = \lambda_{k+1}$, and

$\alpha = \text{gap} = \lambda_k - \lambda_{k+1}$, then with probability $1 - \delta$, we get

$$\begin{aligned}
\left\| (\Phi_k^\perp)^\top \widehat{\Phi} \right\|_F^2 &\leq k \left\| (\Phi_k^\perp)^\top \widehat{\Phi} \right\|_2^2 \leq \frac{k\epsilon}{\alpha^2} \\
&\leq \frac{k}{\alpha^2} \left(\frac{\tau}{3n} \log \left(\frac{\delta}{2m} \right) + \sqrt{\frac{\tau}{3n} \log \left(\frac{\delta}{2m} \right)^2 + \log \left(\frac{\delta}{2m} \right) \frac{\tau \lambda_1}{n}} \right)^2 \\
&\leq \frac{k}{\alpha^2} \left(\frac{\tau}{3n} \log \left(\frac{\delta}{2m} \right) + \sqrt{\frac{2\lambda_1 \tau}{n} \log \left(\frac{\delta}{2m} \right)} \right)^2 \\
&\leq \frac{k}{\alpha^2} \left(\frac{\lambda_1 \tau^2}{n} \log \left(\frac{\delta}{2m} \right)^2 \right)
\end{aligned}$$

Setting $q_{ERM}(1/\delta, \log(m), \log(n)) = \frac{\lambda_1 \tau^2}{\alpha^2} \log \left(\frac{\delta}{2m} \right)^2 = \frac{k \lambda_1 \tau^2}{(\lambda_k - \lambda_{k+1})^2} \log \left(\frac{\delta}{2m} \right)^2$, we get,

$$\left\| (\Phi_k^\perp)^\top \widehat{\Phi} \right\|_F^2 \leq \frac{q_{ERM}(1/\delta, \log(m), \log(n))}{n}$$

□

Space and Computational Complexity of ERM: ERM requires computing and storing the empirical covariance matrix \widehat{C}_m , which takes $O(m^2)$ memory. A rank k SVD on \widehat{C}_m , generally, takes $O(m^2 k)$ computations. We note that there are methods to scale this up but it is out of the scope of this work.

D.2 Oja's Algorithm

Having shown that ERM achieves optimal statistical rates, we now discuss a (relatively) more efficient algorithm in terms of space and computational complexity. We leverage the recent analysis of the classical Oja's algorithm and show how the algorithmic parameters affect the main result. We first restate the theorem statement from the analysis of Oja in [Allen-Zhu and Li \[2016a\]](#).

Theorem D.3. Let $\text{gap} := \lambda_k - \lambda_{k+1} \in (0, \frac{1}{k}]$ and $\Lambda := \sum_{i=1}^k \lambda_i \in (0, 1]$, for every $\epsilon, \delta \in (0, 1)$ define learning rates

$$T_0 = \Theta \left(\frac{4k\Lambda}{\text{gap}^2 \delta^2} \right), T_1 = \Theta \left(\frac{\Lambda}{\text{gap}^2} \right), \eta_t = \begin{cases} \Theta \left(\frac{1}{\text{gap} T_0} \right) & 1 \leq t \leq T_0 \\ \Theta \left(\frac{1}{\text{gap}^2 T_1} \right) & T_0 < t \leq T_0 + T_1 \\ \Theta \left(\frac{1}{\text{gap}(t-T_0)} \right) & t > T_0 + T_1 \end{cases}$$

Let Z be the column orthonormal matrix consisting of all eigenvectors of C_m with values no more than λ_{k+1} . Then the output Q_T of the algorithm satisfies with at least $1 - \frac{\delta}{2}$,

$$\text{for every } T = T_0 + T_1 + \Theta \left(\frac{T_1}{\epsilon} \right), \text{ it satisfies } \|Z^\top Q_T\|_F^2 \leq \epsilon$$

The above theorem gives guarantees of the form required by the definition of efficient subspace learner. Therefore, implicitly, Oja is an efficient subspace learner. This is formally stated in the following lemma.

Lemma D.4. \mathcal{A}_{Oja} is an Efficient Subspace Learner.

Proof. From Theorem D.3, we have

$$\begin{aligned}
\|Z^\top Q_n\|_F^2 &\leq \epsilon \\
&= \tilde{\Theta}\left(\frac{T_1}{n - T_0 - T_1}\right) \\
&\leq \tilde{\Theta}\left(\frac{2\Lambda}{\text{gap}^2 n}\right) \quad (\text{for large } n)
\end{aligned}$$

Setting $q_{oja}(1/\delta, \log(m), \log(n)) = \tilde{\Theta}\left(\frac{\Lambda}{\text{gap}^2}\right)$, we get,

$$\|Z^\top Q_n\|_F^2 \leq \frac{q_{oja}(1/\delta, \log(m), \log(n))}{n}$$

□

Moreover the requirement of an initial constant number of samples as stated in Theorem 4.2 also appears in Theorem D.3 as warm-up phase. Therefore, the requirement of initial samples can be absorbed in the warm-up phase of Oja.

Space and Computational Complexity of Oja's Algorithm: Oja's algorithm takes $O(mk)$ memory. The per iteration computational cost is $O(mk)$. Therefore, for an ϵ -suboptimal solution, the total computational cost is $O\left(\frac{mk}{\epsilon^2}\right)$.

E Experiments

We now need some lemmas which gives us analytical forms which would be used to calculate the objective with respect to empirical measure in the experiments.

Let $\hat{P}_{\mathcal{A}}$ be the output of an efficient subspace learner \mathcal{A} . Let $\hat{P}_{\mathcal{A}} = \tilde{\Phi}\tilde{\Phi}^\top$ be its eigendecomposition. We define $\hat{\Phi} = \tilde{\Phi}R^*$, where let

$$R^* = \arg \min_{R^\top R = RR^\top = I} \left\| \tilde{\Phi}R - \Phi_k \right\|_{\mathcal{F}}^2$$

The following gives an explicit form for R^* .

Lemma E.1. *For any orthogonal matrix $\tilde{\Phi} \in \mathbb{R}^{m \times k}$ and $\Phi \in \mathbb{R}^{m \times k}$, the solution of the optimization problem*

$$\arg \min_{R^\top R = RR^\top = I} \left\| \tilde{\Phi}R - \Phi \right\|_{\mathcal{F}}^2$$

is $R^* = \tilde{\Phi}^\top \Phi$

Proof.

$$\arg \min_{R^\top R = RR^\top = I} \left\| \tilde{\Phi}R - \Phi_k \right\|_{\mathcal{F}}^2 = \arg \min_{R^\top R = RR^\top = I} -\text{Tr} \left(R^\top \tilde{\Phi}^\top \Phi_k \right)$$

$$\max_{R^\top R = RR^\top = I} \text{Tr} \left(R^\top \tilde{\Phi}^\top \Phi_k \right) \leq \|R\|_F \left\| \tilde{\Phi}^\top \Phi_k \right\|_F = k$$

Note that $\text{Tr} \left(R^* \tilde{\Phi}^\top \Phi_k \right) = k$. So, the maximum is achieved at $R = R^* = \tilde{\Phi}^\top \Phi_k$. \square

We have $\hat{\Phi} = \tilde{\Phi}R^*$. We now use this and apply Lemma E.2 to evaluate the objective.

Lemma E.2. *For a projection matrix $P = UU^\top = \sum_{i=1}^k u_i \otimes_{\mathcal{F}} u_i$, $\langle \mathfrak{A}P, \mathfrak{I}C \rangle_{\rho} = \frac{1}{n} \text{Tr} (V^\top KV)$, where $V = \Phi^\top U S^{-\frac{1}{2}}$ and $S = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_k)$, $\lambda_i = \langle C_m u_i, u_i \rangle_{\mathcal{F}}$*

Proof of Lemma E.2.

$$\begin{aligned} \langle \mathfrak{A}P, \mathfrak{I}C \rangle_{HS(\rho)} &= \left\langle \sum_{i=1}^k \frac{A u_i}{\sqrt{\lambda_i}} \otimes_{\rho} \frac{A u_i}{\sqrt{\lambda_i}}, \sum_{j=1}^n \bar{\phi}_j \otimes_{\rho} \bar{\phi}_j \right\rangle_{HS(\rho)} \\ &= \sum_{i,j=1}^{k,n} \frac{1}{\lambda_i} \langle A u_i, \bar{\phi}_j \rangle_{\rho}^2 \\ &= \sum_{i,j=1}^{k,n} \frac{1}{\lambda_i} \langle u_i, A^* \bar{\phi}_j \rangle_{\mathcal{F}}^2 \end{aligned}$$

where the second equality follows from bi-linearity of inner products, third from the definition of adjoints.

$$\begin{aligned} \langle u_i, A^* \bar{\phi}_j \rangle_{\mathcal{F}} &= \sum_{l=1}^m (u_i)_l (A^* \bar{\phi}_j)_l \\ &= \sum_{l=1}^m (u_i)_l \frac{1}{n} \sum_{q=1}^n \bar{\phi}_j(x_q) z(x_q)_l = \frac{1}{n} u_i^\top \Phi \bar{\Phi}_j \end{aligned}$$

where $\bar{\Phi}_j \in \mathbb{R}^n$ and $(\bar{\Phi}_j)_q = \bar{\phi}_j(x_q)$.

Note that

$$\begin{aligned}
\bar{\lambda}_j &= \langle C\bar{\phi}_j, \bar{\phi}_j \rangle_{\mathcal{H}} = \langle I^* I \bar{\phi}_j, \bar{\phi}_j \rangle_{\mathcal{H}} \\
&= \langle I \bar{\phi}_j, I \bar{\phi}_j \rangle_{\rho} = \langle \bar{\phi}_j, \bar{\phi}_j \rangle_{\rho} \\
&= \frac{1}{n} \sum_{q=1}^n \bar{\phi}_j(x_q)^2 = \frac{1}{n} \|\bar{\Phi}_j\|_2^2
\end{aligned}$$

where the third equality follows from the property of adjoints.

Moreover, $(V_j^*)^\top K V_j^* = \bar{\lambda}_j$. Therefore $\bar{\Phi}_j^\top \bar{\Phi}_j = n \bar{\lambda}_j = n (V_j^*)^\top K V_j^*$.

So we have $\bar{\Phi}_j = \sqrt{n} K^{1/2} V_j^*$. Hence,

$$\begin{aligned}
\langle \mathfrak{A}P, \mathfrak{I}C \rangle_{HS(\rho)} &= \sum_{i,j=1}^{k,n} \frac{1}{\lambda_i} \left(\frac{1}{n} u_i^\top \Phi \sqrt{n} K^{1/2} V_j^* \right)^2 \\
&= \frac{1}{n} \sum_{i,j=1}^{k,n} \frac{1}{\lambda_i} \left(u_i^\top \Phi K^{1/2} V_j^* \right)^2 \\
&= \frac{1}{n} \sum_{i,j=1}^{k,n} \left(\frac{1}{\sqrt{\lambda_i}} u_i^\top \Phi K^{1/2} V_j^* \right)^2 \\
&= \frac{1}{n} \sum_{i,j=1}^{k,n} \left(V_i^\top K^{1/2} V_j^* \right)^2
\end{aligned}$$

where $V = \Phi^\top U S^{-\frac{1}{2}}$. Therefore, we have,

$$\begin{aligned}
\langle \mathfrak{A}P, \mathfrak{I}C \rangle_{HS(\rho)} &= \frac{1}{n} \sum_{i,j=1}^{k,n} \text{Tr} \left(V_i^\top K^{1/2} V_j^* \right)^2 \\
&= \frac{1}{n} \sum_{i,j=1}^{k,n} \text{Tr} \left(V_i^\top K^{1/2} V_j^* (V_j^*)^\top K^{1/2} V_i \right) \\
&= \frac{1}{n} \sum_{i=1}^k \text{Tr} \left(V_i^\top K^{1/2} \sum_{j=1}^n V_j^* (V_j^*)^\top K^{1/2} V_i \right) \\
&= \frac{1}{n} \sum_{i=1}^k \text{Tr} \left(V_i^\top K^{1/2} V^* (V^*)^\top K^{1/2} V_i \right) \\
&= \frac{1}{n} \sum_{i=1}^k \text{Tr} (V_i^\top K V_i) \\
&= \frac{1}{n} \text{Tr} (V^\top K V)
\end{aligned}$$

□

F Auxillary Results

Here we state some Auxillary results used in the proofs.

Theorem F.1 (Generalized Gap free Wedin Theorem). *For $\epsilon > 0$, let A and B be two PSD matrices. For every $\mu > 0, \alpha > 0$, let U be column orthonormal matrix consisting of eigenvectors of A with eigenvalue $\leq \mu$, let V be column orthonormal matrix consisting of eigenvectors of B with eigenvalue $\geq \mu + \alpha$, then we have*

$$\|U^\top V\| \leq \frac{\|A - B\|}{\alpha}$$

where the norm $\|\cdot\|$ is any sub-multiplicative norm.

Proof. The above theorem is stated in [Allen-Zhu and Li, 2016b, Lemma B.3] in the sense of spectral norm. For the sake of completeness, we present the proof and show that it can easily be generalized to any sub-multiplicative norm.

Let the SVD of A and B be $A = U\Sigma U^\top + U'\Sigma'U'^\top$, $B = V\tilde{\Sigma}V^\top + V'\tilde{\Sigma}'V'^\top$, where Σ is a diagonal matrix which contains all eigenvalues of A which are $\leq \mu$. Similarly, $\tilde{\Sigma}$ contains all eigenvalues $\geq \mu + \alpha$. Let $E := A - B$.

$$\Sigma U^\top = U^\top A = U^\top (B + E)$$

where the first equality follows because U is orthogonal to U' . Multiply by V on the right on both sides, we get,

$$\Sigma U^\top V = U^\top B V + U^\top E V = U^\top V \tilde{\Sigma} + U^\top E V$$

where the second equality follows because V is orthogonal to V' . Multiplying by $\tilde{\Sigma}^{-1}$ on the right on both sides, we get,

$$\Sigma U^\top V \tilde{\Sigma}^{-1} = U^\top V + U^\top E V \tilde{\Sigma}^{-1}$$

Taking any sub-multiplicative norm on the left hand side, we obtain an upper bound on it as follows,

$$\begin{aligned} \|\Sigma U^\top V \tilde{\Sigma}^{-1}\| &\leq \|\Sigma\|_2 \|\tilde{\Sigma}^{-1}\|_2 \|U^\top V\| \\ &\leq \frac{\mu}{\mu + \alpha} \|U^\top V\| \end{aligned}$$

where the first inequality follows from the property of sub-multiplicative norms, and the second from the definition of Σ and $\tilde{\Sigma}$.

Similarly, taking any sub-multiplicative norm on the right hand side, we get a lower bound on it as follows,

$$\begin{aligned} \|U^\top V + U^\top E V \tilde{\Sigma}^{-1}\| &\geq \|U^\top V\| - \|U^\top E V \tilde{\Sigma}^{-1}\| \\ &\geq \|U^\top V\| - \|U^\top\|_2 \|E\| \|V\|_2 \|\tilde{\Sigma}^{-1}\|_2 \\ &\geq \|U^\top V\| - \frac{\|E\|}{\mu + \alpha} \end{aligned}$$

where the first inequality follows from (reverse) triangle inequality, the second from property of sub-multiplicative norms and third because U and V are orthonormal matrices and by definition of $\tilde{\Sigma}$. Combining both the bounds, we get,

$$\begin{aligned} \|U^\top V\| \left(1 - \frac{\mu}{\mu + \alpha}\right) &\leq \frac{\|E\|}{\mu + \alpha} \\ \implies \|U^\top V\| &\leq \frac{\|E\|}{\alpha} \end{aligned}$$

□

Theorem F.2 (Matrix Bernstein [Tropp et al., 2015]). Let S_1, S_2, \dots, S_n be n i.i.d $d_1 \times d_2$ random matrices such that $\mathbb{E}S_i = 0, \|S_i\| \leq L \forall i \in [n]$. Let $Z = \sum_{i=1}^n S_i$. Let $v(Z)$ denote the matrix variance statistic of the sum defined as,

$$v(Z) = \max\{\mathbb{E}ZZ^\top, \mathbb{E}Z^\top Z\}$$

Then, with probability at least $1 - \delta$, we have,

$$\mathbb{P}\{\|Z\| \geq t\} \leq (d_1 + d_2) \exp\left(\frac{-t^2/2}{v(Z) + Lt/3}\right) \forall t \geq 0$$

Theorem F.3 (Local Rademacher Complexity [Bartlett et al., 2002]). Let \mathcal{X} be a measurable space. Let \mathcal{P} be a probability distribution on \mathcal{X} and let x_1, x_2, \dots, x_n be i.i.d. samples drawn from \mathcal{P} . Let \mathcal{P}_n denote the empirical measure. Let \mathcal{F} be a class of functions on \mathcal{X} ranging from $[-1, 1]$ and assume that there exists some constant B such that for every $f \in \mathcal{F}, \mathcal{P}^2 f \leq B\mathcal{P}f$. Let ψ be a sub-root function and let r^* be the fixed point of ψ . If ψ satisfies

$$\psi(r) \geq B\mathbb{E}_{X,\sigma} [\mathcal{R}_n\{f \in \text{star}(\mathcal{F}) | \mathcal{P}f^2 \leq r\}]$$

where $\text{star}(\mathcal{F}) = \{\lambda f | f \in \mathcal{F}, \lambda \in [0, 1]\}$ is the star shaped hull of \mathcal{F} and $\mathcal{R}_n\mathcal{F} = \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i)$ is the empirical Rademacher complexity of \mathcal{F} given data points $\{x_i\}_{i=1}^n$; then for every $K > 0$ and $x > 0$, with probability at least $1 - e^{-\delta}$

$$\forall f \in \mathcal{F}, \mathcal{P}f \leq \frac{K}{K-1} \mathcal{P}_n f + \frac{6K}{B} r^* + \frac{\delta(11 + 5BK)}{n} \quad (6)$$

Also, with probability at least $1 - e^{-\delta}$

$$\forall f \in \mathcal{F}, \mathcal{P}_n f \leq \frac{K}{K+1} \mathcal{P}f + \frac{6K}{B} r^* + \frac{\delta(11 + 5BK)}{n} \quad (7)$$

Furthermore, if $\hat{\psi}_n$ is a data-dependent sub-root function with fixed point \hat{r}^* such that

$$\psi^*(r) > 2(10 \vee B)\mathbb{E}_\sigma [\mathcal{R}_n\{f \in \text{star}(\mathcal{F}) | \mathcal{P}^n f^2 \leq 2r\}] + \frac{2(10 \vee B + 11)\delta}{n}$$

then with probability at least $1 - 2e^{-\delta}$, it holds that $\hat{r}^* \geq r^*$; as a consequence, equations 6 and 7 holds with r^* replaced by \hat{r}^*