

Generalization bounds for Kernel Canonical Correlation Analysis

Enayat Ullah

*Department of Computer Science
Johns Hopkins University*

enayat@jhu.edu

Raman Arora

*Department of Computer Science
Johns Hopkins University*

arora@cs.jhu.edu

Reviewed on OpenReview: <https://openreview.net/forum?id=KwWKB9Bqam>

Abstract

We study the problem of multiview representation learning using kernel canonical correlation analysis (KCCA) and establish non-asymptotic bounds on generalization error for regularized empirical risk minimization. In particular, we give fine-grained high-probability bounds on generalization error ranging from $O(n^{-1/6})$ to $O(n^{-1/5})$ depending on underlying distributional properties, where n is the number of data samples. For the special case of finite-dimensional Hilbert spaces (such as linear CCA), our rates improve, ranging from $O(n^{-1/2})$ to $O(n^{-1})$. Finally, our results generalize to the problem of functional canonical correlation analysis over abstract Hilbert spaces.

1 Introduction

Canonical correlation analysis (CCA) is a popular technique for multiview representation learning and statistical data analysis. Given a pair of random vectors, CCA finds maximally correlated linear components of the two vectors (Hotelling, 1936). CCA-based methods have recently been shown to improve unsupervised learning of low-dimensional representations of data when multiple “views” of data are available (Vinokourov et al., 2003; Hardoon et al., 2004; Arora and Livescu, 2013). The different views often contain complementary information, and CCA-based multiview representation learning methods can take advantage of this information to learn features that are useful for understanding the structure of the data and that is beneficial for downstream tasks.

Various nonlinear extensions of these multiview learning techniques have also been proposed including kernel CCA (Lai and Fyfe, 2000; Akaho, 2001; Hardoon et al., 2004; Fukumizu et al., 2007) based on positive definite kernels wherein data are represented as functions in associated reproducing kernel Hilbert spaces (RKHS), and deep neural network based extensions, e.g., deep CCA (Andrew et al., 2013)

While CCA and its nonlinear extensions have enjoyed tremendous empirical success, the theoretical understanding of the approaches to solving these problems has been somewhat limited. For example, only recently, were we (as a community) able to give statistical and computational complexity bounds for CCA as a stochastic optimization problem (aka a learning problem) (Allen-Zhu and Li, 2016; Ge et al., 2016; Arora et al., 2017). In a similar spirit of understanding the data analysis techniques as learning problems, in this paper, we look at Kernel CCA and focus on understanding the generalization properties.

However, moving from subspace learning (i.e., linear representations, e.g., using CCA) to learning representations in an RKHS has additional theoretical challenges associated with it. It is then natural to rely on kernel duality, i.e., the representer theorem to reduce the empirical risk minimization (ERM) problem to a finite dimensional optimization problem. Using kernel duality to formulate Kernel CCA was first studied by Lai and Fyfe (2000), Akaho (2001), Melzer et al. (2001) and Bach and Jordan (2002).

In this work, we are interested in understanding the statistical properties of the regularized empirical risk minimizer (defined formally in the subsequent sections) using excess *generalization error* as the error criterion. Informally, excess generalization error of an estimator is the excess error, incurred in objective (or cost), compared to the best, with respect to the underlying data distribution (see Section 3 for a precise formula). This problem has been studied in prior works of (Fukumizu et al., 2007; Fan and Lian, 2016), however their results are asymptotic (see paragraph “Relation to prior work” for more details). Further, these works have studied Kernel CCA in terms of estimation error (or convergence in parameters), and we emphasize that studying the problem in terms of generalization error (or convergence in objective) is important for the following reasons. (Modern) machine learning is typically posed as *risk minimization* problem where the goal is to find parameters that are good in terms of the objective (aka generalization error or population risk) rather than finding the *true* parameters (under some statistical model). Taking a learning view of the KCCA problem, we therefore measure the quality of the solution in terms of the objective rather than distance from a ground truth (which may or may not be unique). This error criterion has been used in the prior works, such as Arora et al. (2017) and Wang et al. (2016), for (linear) CCA. Hence, the main goal in our work, is to give “fine-grained” non-asymptotic guarantees on excess generalization error of regularized empirical risk minimizer (a widely used estimator) for kernel CCA.

1.1 Our Contributions

Our main contributions are as follows.

1. We pose kernel CCA as a learning problem and give upper bounds on excess generalization error of the regularized Empirical Risk Minimizer (ERM). Our results hold for the more general problem of functional CCA in abstract Hilbert spaces. To the best of our knowledge, this is the first work which establishes statistical rates of a finite sample estimator for functional CCA. As special cases, our results give generalization bounds for Kernel CCA and linear CCA, and for both of these special cases, we establish novel results compared to previous work (see below).
2. Under standard assumptions (see Assumption 1), we obtain non-asymptotic bounds on excess generalization error of regularized ERM for kernel CCA, which are between $O(n^{-1/6})$ to $O(n^{-1/5})$ depending on properties of the underlying distribution, where n is the number of data points (see Theorem 1). In contrast, previous works only yield asymptotic guarantees. In the setting when the Reproducing Kernel Hilbert Spaces (RKHS) are finite dimensional, we obtain faster rates ranging from $O(n^{-1/2})$ to $O(n^{-1})$ (see Corollary 3). In the special case of linear CCA, our optimistic rate (i.e. $O(n^{-1})$) is better than the previous result of Gao et al. (2017) and the worse case rate is better in the regime where eigengap of covariance matrix at k is $o(1/\sqrt{n})$ (see Section 4 for details).
3. Our analysis provides insights on the role that regularization parameter plays towards trading off approximation error (bias) and estimation error (variance) and in ensuring statistical consistency of the estimator. In particular, in our bounds, the regularization parameter can decay as $\omega(n^{-1/2})$ and ensure statistical consistency of the estimator – see paragraph “Regularization parameter” in Section 4 for details.

Our proof strategy is to decouple the estimation and approximation errors in the learning problem, bound them separately and balance the tradeoff (between them). To bound the estimation error, the primary tool we use is local Rademacher complexity analysis (Bartlett et al., 2002), which allows us to get a spectrum of rates, from worst case to optimistic (depending on how “easy” the problem is). In the context of kernel methods, these techniques have been applied to give improved rates for kernel principal component analysis (Blanchard et al., 2007), support vector machines (SVMs) with random Fourier features (Gilbert et al., 2018) and other kernel learning problems (Mendelson, 2003; Cortes et al., 2013; Ullah et al., 2018). Please see Section 5 for a detailed proof sketch.

Relation to prior work. Herein, we informally discuss how our work compares with prior results. We refer the reader to paragraph “Comparison with prior works”, in Section 4 for more details. Previous work has studied the the statistical properties of Kernel CCA through the lens of statistical estimation of

Problem	Error criterion	Convergence rate	Reference
Kernel CCA	Parameter	$o_p(1)$	Fukumizu et al. (2007)
Kernel CCA	Parameter	$O_p(n^{-\alpha/(\alpha+1)})^\dagger$	Fan and Lian (2016)
Kernel CCA	Objective	$O(n^{-1/6})$ to $O(n^{-1/5})$	Ours (Corollary 2)
(Linear) CCA	Objective	$O((\text{gap}^2 n)^{-1})^\ddagger$	Gao et al. (2017)
(Linear) CCA	Objective	$O((\text{gap}\sqrt{n})^{-1})$ to $O((\text{gap } n)^{-1})^\ddagger$	Ours (Corollary 3).

Table 1: Summarizing our results in context of relevant prior works. In the table, ‘‘Parameter’’ and ‘‘Objective’’ stand for convergence in parameter and objective respectively (see Section 4 for details). \dagger : obtained under additional assumption on eigenvalues of covariance operators – see Eqn. (3), and $\alpha = \min(\alpha_X, \alpha_Y)$ therein. \ddagger : $\text{gap} = \lambda_1(C) - \lambda_2(C)$ is the eigengap.

parameters. (i.e by bounding *estimation error*). In particular, the works that are most related to ours are that of Fukumizu et al. (2007) and Fan and Lian (2016). Under standard assumptions, Fukumizu et al. (2007) established statistical consistency of the regularized ERM solution if the regularization parameter $\lambda = \omega(n^{-1/3})$. More recently, Fan and Lian (2016) established minimax statistical rates for Kernel CCA under additional assumptions on the problem. The guarantees in both Fukumizu et al. (2007) and Fan and Lian (2016) are asymptotic. Fukumizu et al. (2007) show that as number of samples $n \rightarrow \infty$, the *estimation error* goes to 0, in probability. The work of Fan and Lian (2016) gave rates but these are also in the *convergence in probability* sense. To elaborate, they consider the event that the *estimation error* random variable grows faster than certain sequence in n , and show that probability of this event is limiting to 0. These notions do not give any quantitative finite sample guarantees, and are even weaker than convergence in mean. On the other hand, our guarantees are non-asymptotic - the bounds hold with probability, say at least $1 - \delta$, over the randomness in data for any sample size; and importantly the sample complexity bounds only has $\text{poly}(\log(1/\delta))$ dependence in the failure parameter δ - what are known as ‘‘high confidence’’ guarantees.

In Table 1, we give a summary of our results in context of prior works.

Organization. The rest of the paper is organized as follows. We give mathematical preliminaries in Section 2. In Section 3, we present functional and kernel CCA as learning problems, emphasizing the role of kernel duality and regularization. In Section 4, we present our main result and discuss various implications. Finally, in Section 5, we conclude a brief sketch of the proof.

2 Preliminaries

In this section, we quickly review some mathematical preliminaries in functional analysis; a didactic treatment of random variables in Hilbert spaces, reproducing kernel Hilbert spaces and Local Rademacher complexity is presented in Appendix A.

Let $(\mathcal{H}_X, \mathcal{F}_X, \rho_X)$ and $(\mathcal{H}_Y, \mathcal{F}_Y, \rho_Y)$ be two measurable Hilbert spaces where $\mathcal{H}_X, \mathcal{H}_Y$ are separable spaces, $\mathcal{F}_X, \mathcal{F}_Y$ are σ -fields and ρ_X, ρ_Y are probability measures. Let $\{e_i^X\}_{i \in \mathbb{N}}$ and $\{e_j^Y\}_{j \in \mathbb{N}}$ be an orthonormal basis for \mathcal{H}_X and \mathcal{H}_Y respectively. Let $h_1, h'_1 \in \mathcal{H}_X$ and $h_2, h'_2 \in \mathcal{H}_Y$. We use $\langle h_1, h'_1 \rangle_{\rho_X}$ or $\langle h_1, h'_1 \rangle_{\mathcal{H}_X}$, as per convenience, to denote the inner product between two elements. Similarly we use $\|h_1\|_{\rho_X}$ or $\|h_1\|_{\mathcal{H}_X}$ for norms.

An operator $D : \mathcal{H}_Y \rightarrow \mathcal{H}_X$ is bounded if its operator norm $\|D\|$, defined as $\|D\| := \sup\{\|Dh\|_{\mathcal{H}_X}, h \in \mathcal{H}_Y, \|h\|_{\mathcal{H}_Y} \leq 1\} < \infty$. The outer product $h_1 \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} h_2$ is an operator from \mathcal{H}_Y to \mathcal{H}_X , which acts as $(h_1 \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} h_2)h = \langle h_2, h \rangle_{\mathcal{H}_Y} h_1$ for $h \in \mathcal{H}_Y$. The adjoint operator of D , denoted as $D^* : \mathcal{H}_X \rightarrow \mathcal{H}_Y$, is defined as $\langle h_1, Dh_2 \rangle_{\mathcal{H}_X} = \langle D^*h_1, h_2 \rangle_{\mathcal{H}_Y}$.

A bounded operator is self-adjoint if $D^* = D$. An operator $D : \mathcal{H}_Y \rightarrow \mathcal{H}_X$ is Hilbert-Schmidt if its Hilbert-Schmidt norm, denoted as $\|D\|_{L(\mathcal{H}_Y, \mathcal{H}_X)} := \sum_{i \in \mathbb{N}} \|\text{De}_i^Y\|_{\mathcal{H}_X}^2 = \sum_{i, j \in \mathbb{N}} \langle \text{De}_i^Y, e_j^X \rangle_{\mathcal{H}_X} < \infty$. We use $L(\mathcal{H}_Y, \mathcal{H}_X)$ to denote all Hilbert-Schmidt operators from \mathcal{H}_Y to \mathcal{H}_X . For the sake of brevity, we use $L(\mathcal{H}_X)$ to denote Hilbert-Schmidt operators from \mathcal{H}_X to \mathcal{H}_X . An operator $\mathcal{D} : \mathcal{H}_X \rightarrow \mathcal{H}_Y$ is compact if the image of

every bounded subset of \mathcal{H}_X is a relatively compact subset of \mathcal{H}_Y . A compact operator $D : \mathcal{H}_X \rightarrow \mathcal{H}_X$ is trace-class if $\|D\|_{L^1(\mathcal{H})} := \sum_{i \geq 1} \langle (DD^*)^{1/2} e_i^1, e_i^1 \rangle_{\mathcal{H}_X} < \infty$, where $\|D\|_{L^1(\mathcal{H}_X)}$ denotes the nuclear norm of D . An operator $D : \mathcal{H}_X \rightarrow \mathcal{H}_X$ is positive if $\forall f \in \mathcal{H}_X, \langle f, Df \rangle_{\mathcal{H}_X} \geq 0$. The identity operator $I_X : \mathcal{H}_X \rightarrow \mathcal{H}_X$ is defined as $I_X f = f \forall f \in \mathcal{H}_X$.

Notation. We use capital Roman letters (e.g., A) to denote matrices and operators, small Roman letters (e.g., a) for vectors and small letters (e.g., a) for scalars. Operators over the space of Hilbert-Schmidt operators are represented using capital Fraktur letters, e.g., \mathfrak{A} . For a Hilbert-Schmidt operator D , $\lambda_i(D)$ denotes its i^{th} eigenvalue. Similarly, $\sigma_i(D)$ denotes the i^{th} singular value of D . We use P_D^k to denote the top rank k projection of D ; for example, if the Singular Value Decomposition (SVD) of D is $D = \sum_{i \in \mathbb{N}} \lambda_i u_i \otimes v_i$, then $P_D^k = \sum_{i=1}^k u_i \otimes v_i$. We use $I_k \in \mathbb{R}^{k \times k}$ to denote a $k \times k$ identity matrix. Natural numbers are denoted by \mathbb{N} ; $[n]$ denotes natural numbers from 1 to n .

3 Problem Setup and Background

We begin by recalling the finite dimensional CCA problem. For paired random vectors, $x \in \mathcal{X} \subseteq \mathbb{R}^{d_x}$ and $y \in \mathcal{Y} \subseteq \mathbb{R}^{d_y}$, with some unknown joint distribution ρ , Canonical Correlation Analysis (CCA) can be posed as the following problem.

$$\underset{U \in \mathbb{R}^{d_x \times k}, V \in \mathbb{R}^{d_y \times k}}{\text{maximize}} \quad \langle UV^\top, \mathbb{E}_{x,y} [xy^\top] \rangle \quad \text{such that } U^\top \mathbb{E}_x [xx^\top] U = I_k, V^\top \mathbb{E}_y [yy^\top] V = I_k,$$

where $\langle A, B \rangle = \text{Trace}(A^\top B)$ is the standard inner product on matrices.

Functional CCA. We can generalize the above formulation to abstract Hilbert spaces. In particular, when x and y are random variables in Hilbert spaces \mathcal{H}_X and \mathcal{H}_Y , respectively, the functional CCA problem can be formulated as,

$$\underset{U \in L(\mathcal{H}_X, \mathbb{R}^k), V \in L(\mathcal{H}_Y, \mathbb{R}^k)}{\text{maximize}} \quad \langle UV^*, C_{XY} \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \quad \text{such that } U^* C_X U = I_k, V^* C_Y V = I_k,$$

where $C_X = \mathbb{E} [x \otimes_{L(\mathcal{H}_X)} x]$ and $C_Y = \mathbb{E} [y \otimes_{L(\mathcal{H}_Y)} y]$ are auto-covariance operators, and $C_{XY} = \mathbb{E} [x \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} y]$ is a cross-covariance operator (we refer the reader to Appendix A for a definition).

Kernel CCA. Nonlinear CCA extends the problem of CCA in to a high dimensional feature space using nonlinear feature maps. Kernel CCA is one popular variant of nonlinear CCA where the feature maps are implicit in the kernel functions. It can be viewed as a special case of functional CCA with RKHS \mathcal{H}_X and \mathcal{H}_Y over \mathcal{X} and \mathcal{Y} associated with kernel functions k_X and k_Y , respectively:

$$\underset{U \in L(\mathcal{H}_X, \mathbb{R}^k), V \in L(\mathcal{H}_Y, \mathbb{R}^k)}{\text{maximize}} \quad \langle UV^*, C_{XY} \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \quad \text{such that } U^* C_X U = I_k, V^* C_Y V = I_k,$$

where $C_X = \mathbb{E} [k_X(x, \cdot) \otimes_{L(\mathcal{H}_X)} k_X(x, \cdot)]$, $C_Y = \mathbb{E} [k_Y(y, \cdot) \otimes_{L(\mathcal{H}_Y)} k_Y(y, \cdot)]$ are auto-covariance operators, and $C_{XY} = \mathbb{E} [k_X(x, \cdot) \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} k_Y(y, \cdot)]$ is the cross-covariance operator.

Basis transformation. An alternative equivalent formulation of CCA is obtained by rotating the components in the canonical basis. We then get the following formulation.

$$\underset{U \in L(\mathcal{H}_X, \mathbb{R}^k), V \in L(\mathcal{H}_Y, \mathbb{R}^k)}{\text{maximize}} \quad \left\langle UV^*, C_X^{-1/2} C_{XY} C_Y^{-1/2} \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \quad \text{such that } U^* U = I_k, V^* V = I_k.$$

At first, the reformulation above seems to require that the auto-covariance operators are invertible. In the context of kernel CCA, if the feature space corresponding to the kernel function is finite dimensional (for example, when using polynomial kernels), then invertibility can hold. However, it does not hold in general, for example when using a Gaussian kernel, the auto-covariance operators can no longer be trace-class. This is a problem since the standard assumption of random variables being bounded implies the corresponding

auto-covariance operators are trace class and therefore we have a contradiction. However, note that we can instead write the above equivalently as,

$$\underset{U \in L(\mathcal{H}_X, \mathbb{R}^k), V \in L(\mathcal{H}_Y, \mathbb{R}^k)}{\text{maximize}} \quad \langle UV^*, C \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \quad \text{s.t.} \quad C_X^{1/2} C C_Y^{1/2} = C_{XY}, U^*U = I_k, V^*V = I_k. \quad (1)$$

The existence and uniqueness of such an operator C , bounded as $\|C\| \leq 1$, is established in Baker (1973, Theorem 1) and also discussed in Section 2.2 of Fukumizu et al. (2007). As in Fukumizu et al. (2007), we assume that C is compact and abuse the notation to write $C = C_X^{-1/2} C_{XY} C_Y^{-1/2}$ even when these are not invertible.

We now discuss a final assumption. Note that since C_X and C_Y are self-adjoint, a spectral decomposition of C_X and C_Y exists. Let $C_X := \sum_{i=1}^{\infty} \lambda_i(C_X) \phi_i^X \otimes_{L(\mathcal{H}_X)} \phi_i^X$ and $C_Y := \sum_{i=1}^{\infty} \lambda_i(C_Y) \phi_i^Y \otimes_{L(\mathcal{H}_Y)} \phi_i^Y$ where $\{\phi_i^X\}_i$ and $\{\phi_i^Y\}_i$ are the eigenfunctions of C_X and C_Y respectively. We define $\gamma_{ij} := \mathbb{E}_{x,y} \left[\langle \phi_i^X \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} \phi_j^Y, C_{XY} \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \right]$. We assume that $\max_{i,j} \left| \frac{\gamma_{ij}}{\lambda_i(C_X) \sqrt{\lambda_j(C_Y)}} \right| \leq 1$ and $\max_{i,j} \left| \frac{\gamma_{ij}}{\sqrt{\lambda_i(C_X) \lambda_j(C_Y)}} \right| \leq 1$. With abuse of notation, this means that the operators $C_X^{-1} C_{XY} C_Y^{-1/2}$ and $C_X^{-1/2} C_{XY} C_Y^{-1}$ are bounded, and their operator norms bounded by 1. These assumptions have appeared in previous works Fukumizu et al. (2007) and Fan and Lian (2016) and ensures existence of a solution to the kernel CCA problem.

Observe that the solution to the CCA problem in Eqn. (1) is given by the singular value decomposition (SVD) of C . In particular, let $u_1^C, u_2^C, \dots, u_k^C$ and $v_1^C, v_2^C, \dots, v_k^C$ be the top- k left and right singular functions of C , respectively. We define $U_C : \mathbb{R}^k \rightarrow \mathcal{H}_X$ such that $U_C b = \sum_{i=1}^k b_i u_i^C$, where $b \in \mathbb{R}^k$; and similarly V_C . The solution to the CCA problem in Eqn. (1) is U_C and V_C .

Empirical Risk Minimization (ERM). We now discuss the learning problem and a finite sample estimator. Let $\{(x_i, y_i)\}_{i=1}^n$ be n data points drawn i.i.d. from ρ . We first define empirical counterparts of auto covariance and cross-covariance operators. We define $C_X^n : \mathcal{H}_X \rightarrow \mathcal{H}_X$ and $C_Y^n : \mathcal{H}_Y \rightarrow \mathcal{H}_Y$ as $C_X^n = \frac{1}{n} \sum_{i=1}^n x_i \otimes_{L(\mathcal{H}_X)} x_i$ and $C_Y^n = \frac{1}{n} \sum_{i=1}^n y_i \otimes_{L(\mathcal{H}_Y)} y_i$, respectively. Similarly, the empirical cross-covariance operator $C_{XY}^n = \frac{1}{n} \sum_{i=1}^n x_i \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} y_i$. We also define $C^n := (C_X^n)^{-1/2} C_{XY}^n (C_Y^n)^{-1/2}$. ERM is formulated as,

$$\underset{U \in L(\mathcal{H}_X, \mathbb{R}^k), V \in L(\mathcal{H}_Y, \mathbb{R}^k)}{\text{maximize}} \quad \langle UV^*, C^n \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \quad \text{such that} \quad U^*U = I_k, V^*V = I_k.$$

The solution to the ERM problem above, analogously, are the singular functions of C^n . Let $u_1^n, u_2^n, \dots, u_k^n$ and $v_1^n, v_2^n, \dots, v_k^n$ be the top- k left and right singular functions of C^n respectively. We define $U_n : \mathbb{R}^k \rightarrow \mathcal{H}_X$ such that $U_n b = \sum_{i=1}^k b_i u_i^n$, where $b \in \mathbb{R}^k$, and similarly V_n .

Excess Generalization error. As motivated in Section 1, we are interested in studying the kernel CCA problem from the point of view of learning, i.e., generalization error. The excess generalization error of an estimator is the excess error incurred, in objective, compared to the best, with respect to the underlying distribution. The excess generalization error of ERM solutions (i.e. (U_n, V_n)) is the following,

$$\mathcal{E}(U_n, V_n) = \langle U_C V_C^* - U_n V_n^*, C \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)}.$$

Kernel duality. Note that even if we establish statistical convergence of the estimator, we cannot talk about computational aspects as these are infinite dimensional objects. However, in the special case of RKHS, we can appeal to kernel duality to guarantee efficient computation. In the context of CCA, by observing that the solution should lie in the span of data, Akaho (2001) and Bach and Jordan (2002) show that the empirical problem is equivalent to solving the following finite dimensional optimization problem.

$$\underset{U \in \mathbb{R}^{n \times k}, V \in \mathbb{R}^{n \times k}}{\text{maximize}} \quad \left\langle UV^T, \frac{1}{n} K_X K_Y \right\rangle \quad \text{such that} \quad \frac{1}{n} U^T K_X^2 U = I_k, \frac{1}{n} V^T K_Y^2 V = I_k,$$

where $K_{\mathcal{X}}$ and $K_{\mathcal{Y}}$ are $n \times n$ kernel matrices with $(K_{\mathcal{X}})_{ij} = k_{\mathcal{X}}(x_i, x_j)$ and $(K_{\mathcal{Y}})_{ij} = k_{\mathcal{Y}}(y_i, y_j)$. Solving this typically takes $O(n^2k)$ time and $O(n^2)$ space. However, there are faster approximate alternatives, for example, those based on approximate feature maps (Rahimi and Recht, 2007; Kar and Karnick, 2012), Nyström’s method (Drineas and Mahoney, 2005) and sketching (Yang et al., 2015).

3.1 Regularization

Empirical risk minimization (ERM) is one of the most popular learning rules. However, without an appropriate inductive bias, e.g., in the form of a regularizer, the ERM solution may fail to generalize. For the ERM problem corresponding to canonical correlation analysis, a natural regularizer is a variant of Tikhonov regularization (Groetsch, 1984), which can be described as follows. Define $C_{\mathcal{X}}^{n, \lambda_{\mathcal{X}}} := C_{\mathcal{X}}^n + \lambda_{\mathcal{X}} I_{\mathcal{X}}$, $C_{\mathcal{Y}}^{n, \lambda_{\mathcal{Y}}} := C_{\mathcal{Y}}^n + \lambda_{\mathcal{Y}} I_{\mathcal{Y}}$ and $C^{n, \lambda} := (C_{\mathcal{X}}^{n, \lambda_{\mathcal{X}}})^{-\frac{1}{2}} C_{\mathcal{X}\mathcal{Y}}^n (C_{\mathcal{Y}}^{n, \lambda_{\mathcal{Y}}})^{-1/2}$, where $\lambda_{\mathcal{X}}$ and $\lambda_{\mathcal{Y}}$ are regularization parameters, and $I_{\mathcal{X}}$ and $I_{\mathcal{Y}}$ are the identity operators over $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$, respectively. The regularized ERM problem for CCA (Hardoon et al., 2004; Vinokourov et al., 2002) is then given as:

$$\underset{U \in L(\mathcal{H}_{\mathcal{X}}, \mathbb{R}^k), V \in L(\mathcal{H}_{\mathcal{Y}}, \mathbb{R}^k)}{\text{maximize}} \quad \langle UV^*, C^{n, \lambda} \rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} \quad \text{such that} \quad U^*U = I_k, V^*V = I_k.$$

The regularization above corresponds to shifting the spectrum of the auto-covariance operators for $\lambda_{\mathcal{X}}, \lambda_{\mathcal{Y}} > 0$, so that all eigenvalues are positive. Intuitively, we can see that this regularizer biases the problem to tradeoff solutions which maximize correlation and are not along directions with small (non-zero) variance. Specifically, since for all practical purposes, we only observe finitely many samples, and so the presence of very small eigenvalues magnifies the spurious correlations observed. It is therefore imperative to minimize the effect of such small eigenvalues.

We emphasize that regularization parameters $\lambda_{\mathcal{X}}$ and $\lambda_{\mathcal{Y}}$ should be viewed as parameters dependent on n , as is standard in statistical machine learning. Moreover, in such problems, it is usually expected that the regularization parameter decays *fast* with the number of samples, e.g., as $\text{poly}(1/n)$. In Section 4, we discuss the approximation and estimation error tradeoff and how to set the regularization parameters appropriately to optimize it, to have a small overall excess generalization error. In the case where the auto-covariance operators are already positive definite, i.e., their eigenvalues are lower bounded away from zero, the regularization parameters can be viewed as minimum eigenvalues of the corresponding operators. This gives us the rates for kernel CCA in finite dimensional Hilbert spaces, as presented in Corollary 3.

4 Main Results

We first collectively state all necessary and simplifying assumptions, and then state our main theorem. All these assumptions are standard and have appeared in the previous works, and we discussed some of these in the prior section.

Assumption 1. *We assume that the Hilbert spaces $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$ are separable, random variables x and y are centered and bounded as, $\|x\|_{\mathcal{H}_{\mathcal{X}}} \leq \beta$ and $\|y\|_{\mathcal{H}_{\mathcal{Y}}} \leq \beta$ almost surely, that C , defined in Eqn. (1), is compact and the singular values of C are distinct. We assume that $\|C_{\mathcal{X}}^{-1} C_{\mathcal{X}\mathcal{Y}} C_{\mathcal{Y}}^{-1/2}\| \leq 1$ and $\|C_{\mathcal{X}}^{-1/2} C_{\mathcal{X}\mathcal{Y}} C_{\mathcal{Y}}^{-1}\| \leq 1$.*

Theorem 1 (Functional CCA). *Let x and y be random variables in Hilbert spaces $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$, and let ρ be the joint distribution over $\mathcal{H}_{\mathcal{X}} \times \mathcal{H}_{\mathcal{Y}}$, and $\rho_{\mathcal{X}}$ and $\rho_{\mathcal{Y}}$ denote its marginals. Under the conditions of Assumption 1, given data samples $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ drawn i.i.d. from ρ , with probability at least $1 - \delta$, the excess generalization error of regularized ERM, $U_{n, \lambda}, V_{n, \lambda}$, is bounded as,*

$$\begin{aligned} \langle U_C V_C^* - U_{n, \lambda} V_{n, \lambda}^*, C \rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} &\leq 4k \left(\sqrt{\lambda_{\mathcal{X}}} + \sqrt{\lambda_{\mathcal{Y}}} \right) \left(1 + \frac{2}{\sigma_k(C)} \right) + \frac{16k \left(\sqrt{\lambda_{\mathcal{X}}} + \sqrt{\lambda_{\mathcal{Y}}} \right)^2}{\sigma_k(C)} \\ &+ \inf_{h \geq 0} \left\{ \frac{12\alpha_{\rho} h}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{24\alpha_{\rho} \sqrt{h}}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{24}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j > h} \lambda_j(\mathfrak{C}')} \right\} + \frac{12\alpha_{\rho}}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{22\beta \sqrt{k} \log(1/\delta)}{\left(\sqrt{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \right) n} + \frac{10\alpha_{\rho} \log(1/\delta)}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} \end{aligned}$$

where $\alpha_\rho = \mathbb{E}_{(x,y) \sim \rho, (x',y') \sim \rho} \langle x, x' \rangle_{\mathcal{H}_X}^2 \langle y, y' \rangle_{\mathcal{H}_Y}^2 / (\lambda_k(C) - \lambda_{k+1}(C))$ and $\mathcal{C}' \in L(L(\mathcal{H}_Y, \mathcal{H}_X))$, defined as $\mathcal{C}' := \mathbb{E}_{x,y} \left[(x \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} y) \otimes_{L(L(\mathcal{H}_Y, \mathcal{H}_X))} (x \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} y) \right] - C_{\mathcal{X}Y} \otimes_{L(L(\mathcal{H}_Y, \mathcal{H}_X))} C_{\mathcal{X}Y}$.

A few remarks are in order. First, for simplicity, consider the case when $\lambda_X = \lambda_Y = \lambda$. If we only consider the dependence on λ (which should be set as $1/\text{poly}(n)$), the operator \mathcal{C}' and n , with probability at least $1 - \delta$, we get the following bound,

$$\langle U_C V_C^* - U_{n,\lambda} V_{n,\lambda}^*, C \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \leq O \left(\sqrt{\lambda} + \frac{\log(1/\delta)}{\text{gap}_k(C) \lambda^2 n} + \inf_{h \geq 0} \left\{ \frac{h}{\text{gap}_k(C) \lambda^2 n} + \frac{1}{\lambda \sqrt{n}} \sqrt{\sum_{j>h} \lambda_j(\mathcal{C}')} \right\} \right) \quad (2)$$

where $\text{gap}_k(C) = \lambda_k(C) - \lambda_{k+1}(C)$ is the eigengap of C at k , which shows up in our generalization bound. We emphasize that the existence of eigengap at k as well as that the eigenvalues of C are distinct (in Assumption 1) are simplifying assumptions. The analysis goes through even otherwise; however then the dependence on gap at k is replaced by gap at p^{th} singular value where $p > k$. Moreover this dependence on gap, in general, is unavoidable, as evidenced by existing lower bounds in the special case of linear CCA (Gao et al., 2017, Lemma 20).

Spectrum decay of \mathcal{C}' . We see that the convergence rate is crucially controlled by the decay of the spectrum of \mathcal{C}' , i.e. the term $\frac{1}{\lambda \sqrt{n}} \sum_{j>h} \lambda_j(\mathcal{C}')$. In the worst case, it behaves as $O\left(\frac{1}{\lambda \sqrt{n}}\right)$; the best is when \mathcal{C}' is of finite rank, where it behaves as $O(1/\lambda n)$. Furthermore, if the spectrum has an exponential decay, it is $O(\log n/n)$. We emphasize that this term is the result of the local Rademacher complexity analysis which manifests as a spectrum of convergence rates based on higher-order distributional properties.

Optimizing the Approximation-Estimation error tradeoff. We call the first term on the right hand side of the inequality (2) as the *bias* and the second and third term, collectively, as the *variance*. We have that bias is in $O(\sqrt{\lambda})$ but the variance behaves differently depending on the spectrum decay of \mathcal{C}' . In the worst case, the variance is in $O(\lambda^{-1} n^{-1/2})$. Optimizing the tradeoff, we get a rate of $O(n^{-1/6})$ when $\lambda = \Theta(n^{-1/3})$. In the best case, the variance decays as $O(\lambda^{-2} n^{-1})$, which yields an optimistic rate of $O(n^{-1/5})$ when $\lambda = \Theta(n^{-5/2})$.

Kernel CCA. When Hilbert spaces \mathcal{H}_X and \mathcal{H}_Y are RKHS associated with kernel functions k_X and k_Y respectively, then Theorem 1 gives statistical rates of convergence of empirical kernel CCA in terms of the objective. In particular, we have the following corollary.

Corollary 2 (Kernel CCA). *Along with the notations and assumptions of Theorem 1, suppose that \mathcal{H}_X and \mathcal{H}_Y are reproducing kernel Hilbert spaces associated with kernel functions k_X and k_Y , respectively. Then, regularized ERM on n data points outputs $U_{n,\lambda} V_{n,\lambda}^*$ which satisfies the following with probability at least $1 - \delta$,*

$$\langle U_C V_C^* - U_{n,\lambda} V_{n,\lambda}^*, C \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \leq O \left(\sqrt{\lambda} + \frac{\log(1/\delta) \alpha_\rho}{\lambda^2 n} + \inf_{h \geq 0} \left\{ \frac{h \alpha_\rho}{\lambda^2 n} + \frac{1}{\lambda \sqrt{n}} \sqrt{\sum_{j>h} \lambda_j(\mathcal{C}')} \right\} \right)$$

where $\alpha_\rho = \mathbb{E}_{x,y,x',y'} [k_X(x, x')^2 k_Y(y, y')^2] / (\lambda_k(C) - \lambda_{k+1}(C))$.

Comparison with prior works. Previous work of Fukumizu et al. (2007) and Fan and Lian (2016) study Kernel CCA, for $k = 1$, in the sense of convergence of parameters, i.e. distance between the true and the estimated solution. In particular, let (u_1, v_1) be the true solution of Eqn. (1) and let (\hat{u}_1, \hat{v}_1) be a candidate solution. The error is then given by $\|C_X^{1/2}(u_1 - \hat{u}_1)\|_{\mathcal{H}_X}^2$ and similarly, $\|C_Y^{1/2}(v_1 - \hat{v}_1)\|_{\mathcal{H}_Y}^2$. We remark that this is a *stronger* notion of convergence and implies convergence in objective. However, (as we will discuss), their implied results are *weaker* than ours. Importantly, both these works only give asymptotic results, in sense of *convergence in probability*; this means that the success probability over the draw of samples is not quantified, but only limiting to 1 as $n \rightarrow \infty$. In contrast, our results being non-asymptotic hold for all sample sizes. Finally, our results hold for general k as opposed to the above works which are limited to $k = 1$.

We now discuss both of these works one by one. The goal in [Fukumizu et al. \(2007\)](#) is to establish statistical consistency of the (regularized) ERM. They show that as $n \rightarrow \infty$, the error $\|C_{\mathcal{X}}^{1/2}(u_1 - \hat{u}_1)\|_{\mathcal{H}_{\mathcal{X}}}^2 = o_P(1)$,¹ and similarly for error $\|C_{\mathcal{Y}}^{1/2}(v_1 - \hat{v}_1)\|_{\mathcal{H}_{\mathcal{Y}}}^2$. On the other hand, as remarked earlier, we give high-confidence upper bounds on error in objective, as a function of sample size n , for any n .

The other related work of [Fan and Lian \(2016\)](#) establishes minimax rates for Kernel CCA in the sense of convergence in parameters. However, their results hold under a strict assumption on the decay of eigenvalues of $C_{\mathcal{X}}$ and $C_{\mathcal{Y}}$. In particular, there require existence of constants $C, \alpha_{\mathcal{X}}$ and $\alpha_{\mathcal{Y}}$, such that all eigenvalues of $C_{\mathcal{X}}$ and $C_{\mathcal{Y}}$ are upper and lower bounded as,

$$\begin{aligned} (1/C)j^{-\alpha_{\mathcal{X}}} &\leq \lambda_j(C_{\mathcal{X}}) \leq Cj^{-\alpha_{\mathcal{X}}} \\ (1/C)j^{-\alpha_{\mathcal{Y}}} &\leq \lambda_j(C_{\mathcal{Y}}) \leq Cj^{-\alpha_{\mathcal{Y}}} \end{aligned} \quad (3)$$

The statistical rate then, in the sense of convergence of parameters, $\|C_{\mathcal{X}}^{1/2}(u_1 - \hat{u}_1)\|_{\mathcal{H}_{\mathcal{X}}}^2 = O_p(\max(n^{-\alpha_{\mathcal{X}}/(\alpha_{\mathcal{X}}+1)}, n^{-\alpha_{\mathcal{Y}}/(\alpha_{\mathcal{Y}}+1)}))$. As remarked before, these asymptotic guarantees do not quantify the number of samples required in terms of the confidence parameter δ , and are even weaker than convergence in mean. Moreover, there is a *large* regime where their additional assumptions do not hold, but our results still apply. In particular, suppose that eigenvalues decay exponentially, i.e., $\lambda_j(C_{\mathcal{X}}) = e^{-j}$. Since exponential decays faster than any inverse polynomial, so it goes below the inverse polynomial for large enough j (and we are in infinite dimensions) - this violates the condition of [Fan and Lian \(2016\)](#). In general, any sub-polynomial or super-polynomial eigenvalue value behavior of $C_{\mathcal{X}}$ or $C_{\mathcal{Y}}$ violates their condition; but our results are agnostic to it. Finally, their result doesn't provide any insights on how the regularization parameter be set to control and bias and variance, which is important from a practical perspective.

Regularization parameter. [Fukumizu et al. \(2007\)](#) suggest setting the regularization parameter as $\omega(n^{-1/3})$ to ensure statistical consistency of the ERM. In contrast, when we set $\lambda_{\mathcal{X}} = \lambda_{\mathcal{Y}} = \lambda$, our results showcase an improvement of the same estimator to $\omega(n^{-1/2})$, together with high-confidence guarantees. This follows because bias is in $O(\sqrt{\lambda})$ and variance is of the order $O(\lambda^{-1}n^{-1/2})$, so both decrease with n when $\lambda = \omega(n^{-1/2})$.

Finite dimensional Hilbert spaces. As a special case, our result can be applied to give guarantees for (unregularized) ERM wherein we assume that the smallest eigenvalues of auto-covariance operators $C_{\mathcal{X}}$ and $C_{\mathcal{Y}}$ are lower bounded away from 0 by $\lambda_{\mathcal{X}}$ and $\lambda_{\mathcal{Y}}$, respectively. This subsumes standard CCA problem as well as kernel CCA where the RKHS corresponding to the kernel map is finite dimensional. We obtain rates ranging from $n^{-1/2}$ to n^{-1} depending on the spectrum decay of \mathcal{C}' .

Corollary 3 (Finite dimensional kernel CCA). *Along with the notations and assumptions of Corollary 2, suppose that RKHS $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$ are finite dimensional with the eigenvalues of $C_{\mathcal{X}}$ and $C_{\mathcal{Y}}$ lower bounded by $\lambda_{\mathcal{X}} > 0$ and $\lambda_{\mathcal{Y}} > 0$ respectively. Then, ERM on n training data outputs $U_n V_n^*$ which satisfies the following with probability at least $1 - \delta$,*

$$\langle U_n V_n^* - U_{n,\lambda} V_{n,\lambda}^*, C \rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} \leq O\left(\frac{\log(1/\delta)}{\text{gap}_k(C)\lambda^2 n} + \inf_{h \geq 0} \left\{ \frac{h}{\text{gap}_k(C)\lambda^2 n} + \frac{1}{\lambda\sqrt{n}} \sqrt{\sum_{j>h} \lambda_j(\mathcal{C}')}\right\}\right)$$

where $\text{gap}_k(C) = \lambda_k(C) - \lambda_{k+1}(C)$ is the eigengap of C at k .

Linear CCA. We compare our results with that of [Gao et al. \(2017\)](#) for the special case of the linear CCA. For $k=1$, [Gao et al. \(2017\)](#) showed that ERM achieves ϵ sub-optimality with $O(1/\text{gap}^2 \lambda^2 \epsilon)$ sample complexity, where $\text{gap} = \lambda_1(C) - \lambda_2(C)$, ignoring log factors. In comparison, our sample complexity is $O(1/\text{gap} \lambda^2 \epsilon)$ in the optimistic case and $O(1/\lambda^2 \epsilon \min\{\text{gap}, \epsilon\})$ in the worst case; therefore, our optimistic case is *always* better than [Gao et al. \(2017\)](#), whereas the general/worst-case rate is better whenever $\text{gap} = \omega(\sqrt{\epsilon})$.

See Table 1 for a summary of our results and comparison with prior works.

¹ O_p and o_P are standard asymptotic notation for rates of convergence in probability - for a pair of random sequences, f_n and g_n , we write $f_n = O_p(g_n)$ if $f_n/|g_n|$ is bounded in probability, and $f_n = o_P(g_n)$ if $f_n/|g_n|$ converges to 0, in probability.

5 Proof Sketch

In this section, we sketch the main proof ideas underlying our analysis; full details are deferred to Appendix B. Recall that the optimal solution to the population objective are given by operators U_C and V_C that correspond to the first k eigenfunctions of C^*C and CC^* , where $C = C_{\mathcal{X}}^{-1/2}C_{\mathcal{X}\mathcal{Y}}C_{\mathcal{Y}}^{-1/2}$. For the sake of analysis, we introduce regularized whitened population cross-covariance operator $C_\lambda = (C_{\mathcal{X}} + \lambda_{\mathcal{X}})^{-1/2}C_{\mathcal{X}\mathcal{Y}}(C_{\mathcal{Y}} + \lambda_{\mathcal{Y}})^{-1/2}$, whose top k left and right singular functions are denoted by $U_{C,\lambda}$ and $V_{C,\lambda}$. The empirical counterparts are denoted by $C^n = (C_{\mathcal{X}}^n)^{-1/2}C_{\mathcal{X}\mathcal{Y}}^n(C_{\mathcal{Y}}^n)^{-1/2}$, U_n and V_n , and the regularized empirical counterparts by $C^{n,\lambda} = (C_{\mathcal{X}}^{n,\lambda_{\mathcal{X}}})^{-1/2}C_{\mathcal{X}\mathcal{Y}}^n(C_{\mathcal{Y}})^{-1/2}$, $U_{n,\lambda}$ and $V_{n,\lambda}$. We begin with the following decomposition of the excess error into approximation error, and estimation error,

$$\begin{aligned} \mathcal{E}(U_{n,\lambda}, V_{n,\lambda}) &= \left\langle U_C V_C^* - U_{n,\lambda} V_{n,\lambda}^*, C \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} = \underbrace{\left\langle U_C V_C^* - U_{C,\lambda} V_{C,\lambda}^*, C \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)}}_{\text{Approximation error}} \\ &\quad + \underbrace{\left\langle U_{C,\lambda} V_{C,\lambda}^* - U_{n,\lambda} V_{n,\lambda}^*, C \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)}}_{\text{Estimation error}} = \underbrace{\left\langle U_C V_C^* - U_{C,\lambda} V_{C,\lambda}^*, C \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)}}_{\text{Approximation error}} \\ &\quad + \underbrace{\left\langle U_{C,\lambda} V_{C,\lambda}^* - U_{n,\lambda} V_{n,\lambda}^*, C - C_\lambda \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)}}_{\text{Estimation error 2}} + \underbrace{\left\langle U_{C,\lambda} V_{C,\lambda}^* - U_{n,\lambda} V_{n,\lambda}^*, C_\lambda \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)}}_{\text{Error 3}} \end{aligned}$$

We can write the sum of the first two terms, i.e., the approximation error and the estimation error 2 as

$$\underbrace{\left\langle U_C V_C^*, C \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} - \left\langle U_{C,\lambda} V_{C,\lambda}^*, C_\lambda \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)}}_{\text{Error 1}} + \underbrace{\left\langle U_{n,\lambda} V_{n,\lambda}^*, C_\lambda - C \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)}}_{\text{Error 2}}$$

We bound each of the three error terms separately. The bound on Error 1 and Error 2 (see Appendix B) holds due to the following lemma which bounds the distance between the whitened covariance operator and its regularized counterpart.

Lemma 4. *Let $C = C_{\mathcal{X}}^{-1/2}C_{\mathcal{X}\mathcal{Y}}C_{\mathcal{Y}}^{-1/2}$ and $C_\lambda = (C_{\mathcal{X}} + \lambda_{\mathcal{X}}I_{\mathcal{X}})^{-1/2}C_{\mathcal{X}\mathcal{Y}}(C_{\mathcal{Y}} + \lambda_{\mathcal{Y}}I_{\mathcal{Y}})^{-1/2}$ be its regularized counterpart. Suppose that C is compact, then, there difference is bounded as, $\|C - C_\lambda\| \leq 4(\sqrt{\lambda_{\mathcal{X}}} + \sqrt{\lambda_{\mathcal{Y}}})$.*

It remains to bound Error 3; we prove the following result which follows using the local Rademacher complexity analysis of Bartlett et al. (2002).

Lemma 5 (Error 3). *With probability at least $1 - \delta$, the Error 3 is bounded as,*

$$\begin{aligned} \left\langle U_{C,\lambda} V_{C,\lambda}^* - U_{n,\lambda} V_{n,\lambda}^*, C_\lambda \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} &\leq \inf_{h \geq 0} \left\{ \frac{12\alpha_\rho h}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{24\alpha_\rho \sqrt{h}}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{24}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathcal{E}')} \right\} \\ &\quad + \frac{12\alpha_\rho}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{22\beta \sqrt{k} \log(1/\delta)}{\sqrt{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n}} + \frac{10\alpha_\rho \log(1/\delta)}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} \end{aligned}$$

where $\alpha_\rho = \mathbb{E}_{x,y,x',y'} \left[\langle x, x' \rangle_{\mathcal{H}_X}^2 \langle y, y' \rangle_{\mathcal{H}_Y}^2 \right] / (\lambda_k(C) - \lambda_{k+1}(C))$.

We sketch a proof of the key result in the next section.

5.1 Proof Sketch of Lemma 5

Our main tool is local Rademacher complexity analysis (Bartlett et al., 2002) of the problem, following its application in Kernel PCA (Blanchard et al., 2007). The exact statement of local Rademacher complexity result of (Bartlett et al., 2002) is stated as Theorem 10 in Appendix B. We now discuss how we apply it to our end. We first define the excess risk function class.

$$\mathcal{F} = \left\{ f_{U,V} : (x, y) \rightarrow \left\langle \bar{U}_{C,\lambda} \bar{V}_{C,\lambda}^* - UV^*, C_{x,y} \right\rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \mid U^* C_{\mathcal{X}}^{\lambda_{\mathcal{X}}} U = I_k, V^* C_{\mathcal{Y}}^{\lambda_{\mathcal{Y}}} V = I_k \right\},$$

where $\bar{U}_{C,\lambda} = \left(C_{\mathcal{X}}^{\lambda_x}\right)^{1/2} U_{C,\lambda}$, $\bar{V}_{C,\lambda} = \left(C_{\mathcal{Y}}^{\lambda_y}\right)^{1/2} V_{C,\lambda}$ and $C_{x,y} = x \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} y$.

We remark that we cannot directly apply the local Rademacher complexity technique to this class as Theorem 10 requires that the range of functions be contained in $[-1, 1]$. We therefore look at $\mathcal{G} = \tau\mathcal{F}$, where $\tau = \sqrt{\lambda_x \lambda_y} / (2\beta\sqrt{k})$ and show that it satisfies this requirement.

Lemma 6. *For any $f \in \mathcal{G} = \tau\mathcal{F}$ the range of f is contained in $[-1, 1]$, where $\tau = \frac{\sqrt{\lambda_x \lambda_y}}{2\beta\sqrt{k}}$.*

We also need to show that this function class \mathcal{G} satisfies the ‘‘local’’ assumption of Theorem 10, i.e., the second moment is at most a multiple of mean.

Lemma 7. *For any $f \in \mathcal{G}$, $\mathbb{E}[f^2] \leq \mu \mathbb{E}[f]$ where $\mu = 2\alpha_\rho \tau / (\lambda_x \lambda_y)$, and $\alpha_\rho = \mathbb{E}_{x,y,x',y'} [\langle x, x' \rangle_{\mathcal{H}_x}^2 \langle y, y' \rangle_{\mathcal{H}_y}^2] / (\lambda_k(C) - \lambda_{k+1}(C))$.*

Next, we bound the Rademacher complexity of the star shaped hull of the function class conditioned on bounded second moment. We simplify this by considering a bigger class which contains $\text{star}(\mathcal{G})$. Formally, we show that the family $\{g \in \text{star}(\mathcal{G}) \mid \mathbb{E}[g^2] \leq r\}$ is contained in

$$\mathcal{S}_r := \tau \left\{ (x, y) \rightarrow \langle \Gamma, C_{x,y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \Gamma \in L(\mathcal{H}_y, \mathcal{H}_x), \|\Gamma\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_x \lambda_y}, \langle \Gamma, \mathfrak{C}\Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \frac{r}{\tau^2} \right\},$$

where $\mathfrak{C}' : L(\mathcal{H}_y, \mathcal{H}_x) \rightarrow L(\mathcal{H}_y, \mathcal{H}_x)$ is a fourth moment operator, defined as $\mathfrak{C}' := \mathbb{E}_{x,y} [\langle x \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} y \rangle_{L(L(\mathcal{H}_y, \mathcal{H}_x))} \langle x \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} y \rangle_{L(L(\mathcal{H}_y, \mathcal{H}_x))}]$. A crucial technical result follows.

Lemma 8. *The Rademacher complexity of \mathcal{S}_r is bounded as follows,*

$$\mathfrak{R}_n(\mathcal{S}_r) \leq \sqrt{\frac{r}{n}} + \frac{1}{\sqrt{n}} \inf_{h \geq 0} \left\{ \sqrt{rh} + 2\tau \sqrt{\frac{k}{\lambda_x \lambda_y} \sum_{j>h} \lambda_j(\mathfrak{C}')} \right\}$$

where $\mathfrak{C}' = \mathfrak{C} - C_{\mathcal{X}\mathcal{Y}} \otimes_{L(L(\mathcal{H}_y, \mathcal{H}_x))} C_{\mathcal{X}\mathcal{Y}}$.

A final requirement to apply Theorem 10 is that $\psi(r) := \mu \mathfrak{R}_n(\mathcal{S}_r)$ is a sub-root function in r . This simply follows from the above lemma and the fact that the pointwise infimum of sub-root functions is a sub-root function. Furthermore, the existence of an upper bound on the fixed point is again guaranteed because it is a sub-root function. In particular, we show the following.

Lemma 9. *The fixed point of $\psi(r)$, denoted as r^* , is bounded as,*

$$r^* \leq \tau^2 \left(\inf_{h \geq 0} \left\{ \frac{\xi^2 h}{n} + \frac{2\xi^2 \sqrt{h}}{n} + \frac{4\xi}{\sqrt{m}} \sqrt{\frac{k}{\lambda_x \lambda_y} \sum_{j>h} \lambda_j(\mathfrak{C}')} \right\} + \frac{\xi^2}{n} \right)$$

where $\xi = \frac{2\alpha_\rho}{\lambda_x \lambda_y}$, and $\alpha_\rho = \mathbb{E}_{(x,y) \sim \rho, (x',y') \sim \rho} [\langle x, x' \rangle_{\mathcal{H}_x}^2 \langle y, y' \rangle_{\mathcal{H}_y}^2] / (\lambda_k(C) - \lambda_{k+1}(C))$.

To finish the proof, we set $U = (C_{\mathcal{X}}^{\lambda_x})^{-1/2} U_{n,\lambda}$ and $V = (C_{\mathcal{Y}}^{\lambda_y})^{-1/2} V_{n,\lambda}$, where $U_{n,\lambda}$ and $V_{n,\lambda}$ are the solutions to the regularized ERM problem. Note that $\mathbb{E}_n[f_{U,V}] \leq 0$ where \mathbb{E}_n is expectation with respect to the empirical measure. Applying Theorem 10, stated in Appendix B, and letting $K \rightarrow 1$ yields Theorem 1.

6 Conclusion

In this work, we established finite sample generalization bounds on regularized ERM for functional and kernel CCA. The focus here was on understanding the statistical complexity of regularized ERM for kernel CCA. A promising future direction is to study the problem from an algorithmic and computational perspective, in particular when using approximate feature maps based on random Fourier features (Rahimi and Recht, 2007), Nyström’s method (Drineas and Mahoney, 2005) or with randomized sketching (Yang et al., 2015).

References

- Shotaro Akaho. A kernel method for canonical correlation analysis. In *Proceedings of International Meeting on Psychometric Society (IMPS 2001)*, 2001.
- Zeyuan Allen-Zhu and Yuanzhi Li. Doubly accelerated methods for faster cca and generalized eigendecomposition. *arXiv preprint arXiv:1607.06017*, 2016.
- Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1247–1255, 2013.
- Raman Arora and Karen Livescu. Multi-view CCA-based acoustic features for phonetic recognition across speakers and domains. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 7135–7139. IEEE, 2013.
- Raman Arora, Teodor Vanislavov Marinov, Poorya Mianjy, and Nati Srebro. Stochastic approximation for canonical correlation analysis. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4775–4784. Curran Associates, Inc., 2017.
- Francis R Bach and Michael I Jordan. Kernel independent component analysis. *Journal of machine learning research*, 3(Jul):1–48, 2002.
- Charles R Baker. Joint measures and cross-covariance operators. *Transactions of the American Mathematical Society*, 186:273–289, 1973.
- Peter L Bartlett, Olivier Bousquet, and Shahar Mendelson. Localized rademacher complexities. In *International Conference on Computational Learning Theory*, pages 44–58. Springer, 2002.
- Gilles Blanchard, Olivier Bousquet, and Laurent Zwald. Statistical properties of kernel principal component analysis. *Machine Learning*, 66(2-3):259–294, 2007.
- Corinna Cortes, Marius Kloft, and Mehryar Mohri. Learning kernels using local rademacher complexity. In *Advances in neural information processing systems*, pages 2760–2768, 2013.
- Petros Drineas and Michael W Mahoney. On the nyström method for approximating a gram matrix for improved kernel-based learning. *journal of machine learning research*, 6(Dec):2153–2175, 2005.
- Zengyan Fan and Heng Lian. Minimax convergence rates for kernel cca. *Journal of Multivariate Analysis*, 150:183–190, 2016.
- Kenji Fukumizu, Francis R Bach, and Arthur Gretton. Statistical consistency of kernel canonical correlation analysis. *Journal of Machine Learning Research*, 8(Feb):361–383, 2007.
- Chao Gao, Dan Garber, Nathan Srebro, Jialei Wang, and Weiran Wang. Stochastic canonical correlation analysis. *arXiv preprint arXiv:1702.06533*, 2017.
- Rong Ge, Chi Jin, Praneeth Netrapalli, Aaron Sidford, et al. Efficient algorithms for large-scale generalized eigenvector computation and canonical correlation analysis. In *International Conference on Machine Learning*, pages 2741–2750, 2016.
- Anna Gilbert, Ambuj Tewari, and Yitong Sun. But how does it work in theory? linear svm with random features. *arXiv preprint arXiv:1809.04481*, 2018.
- CW Groetsch. The theory of tikhonov regularization for fredholm equations. *104p, Boston Pitman Publication*, 1984.
- David R Hardoon, Sandor Szedmak, and John Shawe-Taylor. Canonical correlation analysis: An overview with application to learning methods. *Neural computation*, 16(12):2639–2664, 2004.
- Harold Hotelling. Relations between two sets of variates. *Biometrika*, 28(3/4):321–377, 1936.

- Purushottam Kar and Harish Karnick. Random feature maps for dot product kernels. In *Artificial Intelligence and Statistics*, pages 583–591, 2012.
- Pei Ling Lai and Colin Fyfe. Kernel and nonlinear canonical correlation analysis. *International Journal of Neural Systems*, 10(05):365–377, 2000.
- Thomas Melzer, Michael Reiter, and Horst Bischof. Nonlinear feature extraction using generalized canonical correlation analysis. In *International Conference on Artificial Neural Networks*, pages 353–360. Springer, 2001.
- Shahar Mendelson. On the performance of kernel classes. *Journal of Machine Learning Research*, 4(Oct): 759–771, 2003.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2007.
- M Reed and B Simon. Methods of modern mathematical physics i: Functional analysis (academic, new york). *Google Scholar*, page 151, 1972.
- Dino Sejdinovic and Arthur Gretton. What is an rkhs?, 2012.
- Enayat Ullah, Poorya Mianjy, Teodor Vanislavov Marinov, and Raman Arora. Streaming kernel pca with $\tilde{o}(\sqrt{n})$ random features. *Advances in Neural Information Processing Systems*, 31, 2018.
- Alexei Vinokourov, Nello Cristianini, and John Shawe-Taylor. Inferring a semantic representation of text via cross-language correlation analysis. *Advances in neural information processing systems*, 15, 2002.
- Alexei Vinokourov, Nello Cristianini, and John Shawe-Taylor. Inferring a semantic representation of text via cross-language correlation analysis. In *Advances in neural information processing systems*, pages 1497–1504, 2003.
- Weiran Wang, Jialei Wang, Dan Garber, and Nati Srebro. Efficient globally convergent stochastic optimization for canonical correlation analysis. In *Advances in Neural Information Processing Systems*, pages 766–774, 2016.
- Yun Yang, Mert Pilanci, and Martin J Wainwright. Randomized sketches for kernels: Fast and optimal non-parametric regression. *arXiv preprint arXiv:1501.06195*, 2015.

Supplementary Material

A Preliminaries

A.1 Random variables in Hilbert spaces

We start with briefly discussing probability theory in Hilbert spaces. Let $(\mathcal{H}_X, \mathcal{F}_X, \rho_X)$ and $(\mathcal{H}_Y, \mathcal{F}_Y, \rho_Y)$ be measurable Hilbert spaces. A random variable x in \mathcal{H}_X is well-defined if and only if every continuous form $\langle f, x \rangle_{\mathcal{H}_X}$ is measurable for all $f \in \mathcal{H}_X$. Its expectation is $\mu_X \in \mathcal{H}_X$ if $\langle \mu_X, f \rangle_{\mathcal{H}_X} = \mathbb{E}_{\rho_X} [\langle x, f \rangle_{\mathcal{H}_X}] \forall f \in \mathcal{H}_X$. The existence and uniqueness of μ_X is guaranteed if $\|x\|_{\mathcal{H}_X} < \infty$. We similarly consider the space $(\mathcal{H}_Y, \mathcal{F}_Y, \rho_Y)$ with random variable y . Throughout the paper, we assume that the random variable x has mean $\mathbf{0}_X \in \mathcal{H}_X$, formally defined as $\mathbb{E}_x [\langle x, f \rangle_{\mathcal{H}_X}] = \langle \mathbf{0}_X, f \rangle_{\mathcal{H}_X} = 0 \forall f \in \mathcal{H}_X$, to simplify the presentation. Similarly, the mean of y is assumed to be $\mathbf{0}_Y$.

Definition 1 (Auto-covariance operator). *Let x be a random variable in a Hilbert space \mathcal{H}_X such that $\|x\|_{\mathcal{H}_X} < \infty$, with distribution ρ_X . Then, the auto-covariance operator of x , denoted as $C_X : \mathcal{H}_X \rightarrow \mathcal{H}_X$ is defined as, $\forall f, g \in \mathcal{H}_X$,*

$$\langle f, C_X g \rangle = \int_{\mathcal{X}} \langle f, x \rangle_{\mathcal{H}_X} \langle g, x \rangle_{\mathcal{H}_X} d\rho_X(x).$$

Equivalently, $C_X := \mathbb{E}_{\rho_X} [x \otimes_{L(\mathcal{H}_X)} x]$. Furthermore, C_X is self-adjoint, positive and trace class.

We similarly use C_Y to denote the auto-covariance operator of random variable y in \mathcal{H}_Y .

Let $\mathcal{H}_X \times \mathcal{H}_Y$ be the product space of \mathcal{H}_X and \mathcal{H}_Y , and let $\mathcal{F}_X \times \mathcal{F}_Y$ be the σ -field generated by the product of elements of \mathcal{F}_X and \mathcal{F}_Y . Let ρ be the joint probability measure on $(\mathcal{H}_X \times \mathcal{H}_Y, \mathcal{F}_X \times \mathcal{F}_Y)$ with its marginal/projection on $(\mathcal{H}_X, \mathcal{F}_X)$ and $(\mathcal{H}_Y, \mathcal{F}_Y)$ being ρ_X and ρ_Y respectively. We now define a cross-covariance operator (Baker, 1973).

Definition 2 (Cross-covariance operator). *Let (x, y) be paired random variables in $\mathcal{H}_X \times \mathcal{H}_Y$ distributed as ρ , such that $\|x\|_{\mathcal{H}_X} < \infty$ and $\|y\|_{\mathcal{H}_Y} < \infty$. The cross-covariance operator $C_{XY} : \mathcal{H}_Y \rightarrow \mathcal{H}_X$ is defined as*

$$\langle f, C_{XY} g \rangle_{\mathcal{H}_X} = \mathbb{E}_{\rho} [\langle f, x \rangle_{\mathcal{H}_X} \langle g, y \rangle_{\mathcal{H}_Y}].$$

Equivalently, $C_{XY} := \mathbb{E}_{\rho} [x \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} y]$. The adjoint of C_{XY} , denoted as $C_{XY}^ = C_{YX}$, which is defined similarly.*

The existence and uniqueness of bounded auto-covariance and cross-covariance operators is guaranteed by Riesz's representation theorem (Reed and Simon, 1972).

A.2 Reproducing Kernel Hilbert Spaces

Let $\mathcal{X} \subseteq \mathbb{R}^{d_X}$ and $\mathcal{Y} \subseteq \mathbb{R}^{d_Y}$ be finite dimensional data domains corresponding to the two views. Let ρ be the joint distribution over $\mathcal{X} \times \mathcal{Y}$ and ρ_X and ρ_Y denote the marginals. We start with a definition of a reproducing kernel. Our discussion is in the context of k_X but it analogously holds for k_Y as well.

Definition 3 (Reproducing kernel (Sejdicinovic and Gretton, 2012)). *Let \mathcal{H}_X be a Hilbert space of real valued functions over \mathcal{X} . A function $k_X : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a reproducing kernel of \mathcal{H}_X , if the following hold,*

- $\forall x \in \mathcal{X}, k_X(x, \cdot) \in \mathcal{H}_X$
- $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}_X, \langle f, k_X(x, \cdot) \rangle_{\mathcal{H}_X} = f(x)$ (Reproducing property)

The function $k_X(x, \cdot)$ is the *representer* of the evaluation functional at x . The inner product between two elements can therefore be expressed as $\langle k_X(x, \cdot), k_X(x', \cdot) \rangle_{\mathcal{H}_X} = k_X(x, x')$. Moreover, reproducing kernels are always positive definite.

A Reproducing Kernel Hilbert Space, conventionally abbreviated as RKHS, is defined as the completion of the linear span of $\{k_{\mathcal{X}}(x, \cdot), x \in \mathcal{X}\}$ with respect to the inner product $\langle k_{\mathcal{X}}(x, \cdot), k_{\mathcal{X}}(x', \cdot) \rangle_{\mathcal{H}_{\mathcal{X}}} = k_{\mathcal{X}}(x, x')$. The existence and uniqueness of a reproducing kernel for an RKHS is guaranteed by Riesz-representation theorem. We point the interested reader to [Sejdicinovic and Gretton \(2012\)](#) for more discussion on RKHS.

We now briefly discuss auto-covariance and cross-covariance operators in the context of RKHS. As before, we assume that $\sup_x \|k_{\mathcal{X}}(x, \cdot)\|_{\mathcal{H}_{\mathcal{X}}} < \infty$ and $\sup_y \|k_{\mathcal{Y}}(y, \cdot)\|_{\mathcal{H}_{\mathcal{Y}}} < \infty$. The auto covariance operator over $\mathcal{H}_{\mathcal{X}}$, $C_{\mathcal{X}} : \mathcal{H}_{\mathcal{X}} \rightarrow \mathcal{H}_{\mathcal{X}}$ is defined as, $\langle f_1, C_{\mathcal{X}} f_2 \rangle_{\mathcal{H}_{\mathcal{X}}} = \mathbb{E}_{\rho_{\mathcal{X}}} [f_1(x) f_2(x)] \quad \forall f_1, f_2 \in \mathcal{H}_{\mathcal{X}}$ or equivalently, $C_{\mathcal{X}} := \mathbb{E}_{\rho_{\mathcal{X}}} [k_{\mathcal{X}}(x, \cdot) \otimes_{L(\mathcal{H}_{\mathcal{X}})} k_{\mathcal{X}}(x, \cdot)]$. Similarly, $C_{\mathcal{Y}} : \mathcal{H}_{\mathcal{Y}} \rightarrow \mathcal{H}_{\mathcal{Y}}$ is defined as, $\langle g_1, C_{\mathcal{Y}} g_2 \rangle_{\mathcal{H}_{\mathcal{Y}}} = \mathbb{E}_{\rho_{\mathcal{Y}}} [g_1(y) g_2(y)] \quad \forall g_1, g_2 \in \mathcal{H}_{\mathcal{Y}}$ or equivalently, $C_{\mathcal{Y}} := \mathbb{E}_{\rho_{\mathcal{Y}}} [k_{\mathcal{Y}}(y, \cdot) \otimes_{L(\mathcal{H}_{\mathcal{Y}})} k_{\mathcal{Y}}(y, \cdot)]$. Furthermore, the cross-covariance operator $C_{\mathcal{X}\mathcal{Y}} : \mathcal{H}_{\mathcal{Y}} \rightarrow \mathcal{H}_{\mathcal{X}}$ is defined as, $\langle f, C_{\mathcal{X}\mathcal{Y}} g \rangle_{\mathcal{H}_{\mathcal{X}}} = \mathbb{E}_{\rho} [f(x) g(y)]$ or equivalently, $C_{\mathcal{X}\mathcal{Y}} := \mathbb{E}_{\rho} [k_{\mathcal{X}}(x, \cdot) \otimes_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} k_{\mathcal{Y}}(y, \cdot)]$. These are easy to see from definitions 1 and 2 and applying reproducing property from definition 3.

Kernel CCA notation. In the context of kernel CCA, we setup the following notation. Let $\mathcal{X} \subseteq \mathbb{R}^{d_x}$ and $\mathcal{Y} \subseteq \mathbb{R}^{d_y}$ be finite dimensional data domains corresponding to the two views. Let ρ be the joint distribution over $\mathcal{X} \times \mathcal{Y}$ and $\rho_{\mathcal{X}}$ and $\rho_{\mathcal{Y}}$ denote the marginals. Let $\{(x_i, y_i)\}_{i=1}^n$ be n data points drawn from ρ . Let $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ be two kernel functions and $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$ be the RKHS associated with $(\mathcal{X}, k_{\mathcal{X}})$ and $(\mathcal{Y}, k_{\mathcal{Y}})$ respectively.

A.3 Local Rademacher Complexity Technique

Rademacher complexity is a data-dependent notion of complexity which captures the richness of a class of real valued functions with respect to the data distribution. It is a fundamental concept in statistical learning theory and empirical process theory, and enables uniform convergence guarantees for the learning problem. We define it formally below.

Definition 4 (Rademacher Complexity). *Let \mathcal{F} be a set of functions over \mathcal{X} and let $\rho_{\mathcal{X}}$ be a probability distribution over \mathcal{X} . For a positive integer n , let x_1, x_2, \dots, x_n be i.i.d. samples drawn from $\rho_{\mathcal{X}}$ and let $\sigma_1, \sigma_2, \dots, \sigma_n$ be i.i.d. samples drawn from Rademacher distribution. The Rademacher complexity of \mathcal{F} , denoted by $\mathfrak{R}_n(\mathcal{F})$, is defined as,*

$$\mathfrak{R}_n(\mathcal{F}) = \mathbb{E}_{\sigma, x} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right].$$

Local Rademacher complexity technique is a concentration of measure tool, introduced in [Bartlett et al. \(2002\)](#), which aims to provide a *finer* analysis of the problem. It is motivated from the claim that the hypothesis selected by a learning algorithm usually has a low empirical error, and hence looking at the Rademacher complexity of the whole class is wasteful. [Bartlett et al. \(2002\)](#) use the variance of the empirical process to constrain the hypothesis class. In particular, if the variance is upper bounded by a constant multiple of the mean, then the empirical process is well-behaved and admits faster rates. This is formalized in Theorem 10 (in Section B) which is restated from [Bartlett et al. \(2002\)](#).

Local Rademacher complexity technique, therefore, essentially looks at the Rademacher complexity of a *smaller* subset of functions from the original hypothesis class. Formally, we have $\mathfrak{R}_n(\mathcal{F}, r) = \mathfrak{R}_n(\{f \in \text{star}(\mathcal{F}) \mid \mathbb{E}[f^2] \leq r\})$ where $\text{star}(\mathcal{F}) = \{\alpha f \mid f \in \mathcal{F}, \alpha \in [0, 1]\}$ is the star-hull of \mathcal{F} . This technique has been used to derive *sharper* generalization bounds for various kernel learning problems ([Blanchard et al., 2007](#); [Ullah et al., 2018](#); [Cortes et al., 2013](#); [Gilbert et al., 2018](#)). Informally, the rate obtained is usually controlled by the tail decay of the spectrum of a *higher order moment*. We discuss this in more detail in subsequent sections.

B Proof of the main theorem

We first restate the local Rademacher complexity result from [Bartlett et al. \(2002\)](#), which is a key tool in our analysis.

Theorem 10. (*Bartlett et al., 2002*) Let \mathcal{X} be a measurable space and let \mathcal{P} be a probability distribution on \mathcal{X} . Let $x_1, x_2 \dots x_n$ be i.i.d. samples drawn from \mathcal{P} , let \mathcal{P}_n denote its empirical measure, and let $\mathcal{P}f := \mathbb{E}_{x \sim \mathcal{P}}[f(x)]$ and $\mathcal{P}_n f := \frac{1}{n} \sum_{i=1}^n f(x_i)$ for a measurable function f . 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 . Suppose that ψ satisfies

$$\psi(r) \geq B \mathfrak{R}_n \{f \in \text{star}(\mathcal{F}) | \mathcal{P}f^2 \leq r\}$$

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

$$\forall f \in \mathcal{F}, \mathcal{P}f \leq \frac{K}{K-1} \mathcal{P}_n f + \frac{6K}{B} r^* + \frac{\delta(11+5BK)}{n} \quad (4)$$

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

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

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

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

then with probability at least $1 - 2e^{-\delta}$, it holds that $\hat{r}^* \geq r^*$; as a consequence, Equations (4) and (5) hold with r^* replaced by \hat{r}^*

We now start the proof of Theorem 1. As discussed in the proof sketch in Section 5, we decompose the error into three terms, and bound each of them one by one.

$$\begin{aligned} \mathcal{E}(U_{n,\lambda}, V_{n,\lambda}) &= \underbrace{\langle U_C V_C^*, C \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} - \langle U_{C,\lambda} V_{C,\lambda}^*, C_\lambda \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}}_{\text{Error 1}} \\ &\quad + \underbrace{\langle U_{n,\lambda} V_{n,\lambda}^*, C_\lambda - C \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}}_{\text{Error 2}} \\ &\quad + \underbrace{\langle U_{C,\lambda} V_{C,\lambda}^* - U_{n,\lambda} V_{n,\lambda}^*, C_\lambda \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}}_{\text{Error 3}} \end{aligned}$$

The three error terms are bound in Lemmas 11, 14, and 15 respectively. Combining them gives the guarantee stated in Theorem 1. In the following subsections, we give the proofs of the afore-mentioned lemmas.

B.1 Bounding the Error 1

Lemma 11 (Error 1). *Error 1 is bounded as,*

$$\langle U_C V_C^*, C \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} - \langle U_{C,\lambda} V_{C,\lambda}^*, C_\lambda \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \frac{8k}{\sigma_k(C)} \left((\sqrt{\lambda_x} + \sqrt{\lambda_y}) + 2(\sqrt{\lambda_x} + \sqrt{\lambda_y})^2 \right)$$

Proof. Note that U_C and V_C are operators corresponding to the first k left and right singular functions of C . Therefore $\langle U_C V_C^*, C \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} - \langle U_{C,\lambda} V_{C,\lambda}^*, C_\lambda \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}$ is difference of the sum of first k singular values of C and C_λ . That is, $\langle U_C V_C^*, C \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} = \sum_{i=1}^k \sigma_i(C)$ and $\langle U_{C,\lambda} V_{C,\lambda}^*, C_\lambda \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} = \sum_{i=1}^k \sigma_i(C_\lambda)$. We

therefore have,

$$\begin{aligned}
\langle U_C V_C^*, C \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} - \langle U_{C, \lambda} V_{C, \lambda}^*, C_\lambda \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} &= \sum_{i=1}^k (\sigma_i(C) - \sigma_i(C_\lambda)) \\
&= \sum_{i=1}^k \left(\sqrt{\lambda_i(CC^*)} - \sqrt{\lambda_i(C_\lambda C_\lambda^*)} \right) \\
&= \sum_{i=1}^k \frac{(\lambda_i(CC^*) - \lambda_i(C_\lambda C_\lambda^*))}{\sqrt{\lambda_i(CC^*)} + \sqrt{\lambda_i(C_\lambda C_\lambda^*)}} \\
&\leq \sum_{i=1}^k \frac{(\lambda_i(CC^*) - \lambda_i(C_\lambda C_\lambda^*))}{\sqrt{\lambda_i(CC^*)}}
\end{aligned}$$

where the first inequality follows because $\lambda(C_\lambda C_\lambda^*) \geq 0$. We now apply perturbation bounds, in particular Weyl's inequality to bound the difference of eigenvalues of operators using the norm of their difference. We therefore get

$$\begin{aligned}
&\langle U_C V_C^*, C \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} - \langle U_{C, \lambda} V_{C, \lambda}^*, C_\lambda \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \\
&\leq \sum_{i=1}^k \frac{1}{\sigma_i(C)} \|CC^* - C_\lambda C_\lambda^*\| \\
&= \sum_{i=1}^k \frac{1}{\sigma_i(C)} \|CC^* - CC_\lambda^* + CC_\lambda^* - C_\lambda C_\lambda^*\| \\
&\leq \frac{k}{\sigma_k(C)} (\|C\| + \|C_\lambda\|) \|C - C_\lambda\| \\
&\leq \frac{k}{\sigma_k(C)} (2\|C\| + \|C - C_\lambda\|) \|C - C_\lambda\| \\
&\leq \frac{k}{\sigma_k(C)} \left(2\|C - C_\lambda\| + \|C - C_\lambda\|^2 \right)
\end{aligned}$$

where in the second inequality we used the fact that singular values are ordered, in the last inequality that $\|C\| \leq 1$ and triangle inequality in others. We now appeal to Lemma 12 that bounds $\|C - C_\lambda\|$ which completes the proof. \square

Lemma 12. *Given $C = C_X^{-1/2} C_{XY} C_Y^{-1/2}$ and its regularized counterpart $C_\lambda = (C_X + \lambda_X I_X)^{-1/2} C_{XY} (C_Y + \lambda_Y I_Y)^{-1/2}$. Assuming C is compact, their difference is bounded as,*

$$\|C - C_\lambda\| \leq 2(\sqrt{\lambda_X} + \sqrt{\lambda_Y})$$

Proof. The proof follows the proof of Lemma 7 in [Fukumizu et al. \(2007\)](#).

$$\begin{aligned}
& \|C_\lambda - C\| \\
&= \left\| (C_X + \lambda_X I_X)^{-1/2} C_{XY} (C_Y + \lambda_Y I_Y)^{-1/2} - C_X^{-1/2} C_{XY} C_Y^{-1/2} \right\| \\
&\leq \left\| (C_X + \lambda_X I_X)^{-1/2} C_{XY} (C_Y + \lambda_Y I_Y)^{-1/2} - C_X^{-1/2} C_{XY} (C_Y + \lambda_Y I_Y)^{-1/2} \right\| \\
&\quad + \left\| (C_X)^{-1/2} C_{XY} (C_Y + \lambda_Y I_Y)^{-1/2} - C_X^{-1/2} C_{XY} C_Y^{-1/2} \right\| \\
&= \left\| ((C_X + \lambda_X I_X)^{-1/2} - C_X^{-1/2}) C_{XY} (C_Y + \lambda_Y I_Y)^{-1/2} \right\| \\
&\quad + \left\| C_X^{-1/2} C_{XY} ((C_Y + \lambda_Y I_Y)^{-1/2} - C_Y^{-1/2}) \right\| \\
&= \left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) C_X^{-1/2} C_{XY} (C_Y + \lambda_Y I_Y)^{-1/2} \right\| \\
&\quad + \left\| C (C_Y^{1/2} (C_Y + \lambda_Y I_Y)^{-1/2} - I_Y) \right\| \\
&= \left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) C C_Y^{1/2} (C_Y + \lambda_Y I_Y)^{-1/2} \right\| \\
&\quad + \left\| C (C_Y^{1/2} (C_Y + \lambda_Y I_Y)^{-1/2} - I_Y) \right\| \\
&\leq \left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) C \right\| + \left\| C (C_Y^{1/2} (C_Y + \lambda_Y I_Y)^{-1/2} - I_Y) \right\|
\end{aligned}$$

where the last inequality follows since $\left\| C_Y^{1/2} (C_Y + \lambda_Y I_Y)^{-1/2} \right\|_{L(\mathcal{H}_Y)} \leq 1$ for positive λ_Y . We will now bound the two terms separately. Suppose for the first term, the operator norm is realized for $w \in \mathcal{H}_Y$, $\|w\| = 1$ and $v = Cw$. [Fukumizu et al. \(2007\)](#) remarks that the range of C is contained in the closure of the range of C_X . That is, $v \in \{C\tilde{w} : \tilde{w} \in \mathcal{H}_Y, \|\tilde{w}\| \leq 1\} \cap \text{Closure}(\text{Range}(C_X))$. Since C_X is continuous, we have that for any $\epsilon > 0$, there exists $\tilde{v} \in \{C\tilde{w} : \tilde{w} \in \mathcal{H}_Y, \|\tilde{w}\| \leq 1\} \cap \text{Range}(C_X)$ such that $\|v - \tilde{v}\|_{\mathcal{H}_X} \leq \epsilon$. Further, there exists $\tilde{w} \in \mathcal{H}_Y, \tilde{u} \in \mathcal{H}_X$ such that $\tilde{v} = C\tilde{w} = C_X\tilde{u}$. Plugging this in the first term, we get,

$$\begin{aligned}
& \left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) v \right\|_{\mathcal{H}_X} \\
&= \left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) (\tilde{v} + v - \tilde{v}) \right\|_{\mathcal{H}_X} \\
&\leq \left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) C_X \tilde{u} \right\|_{\mathcal{H}_X} + \left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) (v - \tilde{v}) \right\|_{\mathcal{H}_X} \\
&\leq \left\| (C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} ((C_X^{1/2} - (C_X + \lambda_X I_X)^{1/2}) C_X^{1/2} \tilde{u} \right\|_{\mathcal{H}_X} + \left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) \right\| \|v - \tilde{v}\|_{\mathcal{H}_X} \\
&\leq \left\| (C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} \right\| \left\| C_X^{1/2} \tilde{u} \right\|_{\mathcal{H}_X} \left\| C_X^{1/2} - (C_X + \lambda_X I_X)^{1/2} \right\| + \epsilon \\
&\leq 2 \left\| C_X^{1/2} \tilde{u} \right\| \cdot \left\| C_X^{1/2} - (C_X + \lambda_X I_X)^{1/2} \right\|
\end{aligned}$$

where in the first inequality, we used the triangle inequality, where the second equality follows because $(C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} (C_X + \lambda_X I_X)^{1/2} C_X^{1/2} = C_X$ since these operators are commutative i.e. have the same eigenspaces. In the second to last inequality, we use that $\left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) \right\| \leq 1$ which follows because $\left\| ((C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} - I_X) \right\| = \max_i \left| \frac{\sqrt{\lambda_i(C_X)}}{\sqrt{\lambda_i(C_X) + \lambda_X}} - 1 \right| = \max_i \left| \frac{\sqrt{\lambda_i(C_X)} - \sqrt{\lambda_i(C_X) + \lambda_X}}{\sqrt{\lambda_i(C_X) + \lambda_X}} \right| \leq 1$ as C_X is trace-class so its eigenvalues go to 0. Finally, the last step follows from $\left\| C_X^{1/2} (C_X + \lambda_X I_X)^{-1/2} \right\| \leq 1$ for positive λ_X and by choosing $\epsilon = \left\| (C_X + \lambda_X I_X)^{-1/2} C_X^{1/2} \right\| \left\| C_X^{1/2} \tilde{u} \right\|_{\mathcal{H}_X} \left\| C_X^{1/2} - (C_X + \lambda_X I_X)^{1/2} \right\|$. We will now argue that $\left\| C_X^{1/2} \tilde{u} \right\| \leq 1$. This follows by our assumption 1 that the operator $C_X^{-1} C_{XY} C_Y^{-1/2}$

is Hilbert-Schmidt. With some abuse of notation, we have that $\|C_{\mathcal{X}}^{1/2}\tilde{u}\| = \|C_{\mathcal{X}}^{-1/2}\tilde{v}\| = \|C_{\mathcal{X}}^{-1/2}C\tilde{w}\| = \|C_{\mathcal{X}}^{-1}C_{\mathcal{X}Y}C_Y^{-1/2}\tilde{w}\| \leq \|C_{\mathcal{X}}^{-1}C_{\mathcal{X}Y}C_Y^{-1/2}\| \|\tilde{w}\| \leq 1$ where the first term is bounded by 1 by assumption, and the second because $\|\tilde{w}\| \leq 1$. Similarly, using the same argument for the other term, we get,

$$\begin{aligned} \left\|C(C_Y^{1/2}(C_Y + \lambda_Y I_Y)^{-1/2} - I_Y)\right\| &\leq 2\left\|C_Y^{1/2}\tilde{u}\right\| \left\|C_Y^{1/2} - (C_Y + \lambda_Y I_Y)^{1/2}\right\| \\ &\leq 2\left\|C_Y^{1/2} - (C_Y + \lambda_Y I_Y)^{1/2}\right\| \end{aligned}$$

Finally using Lemma 13, we get the bound $2(\sqrt{\lambda_{\mathcal{X}}} + \sqrt{\lambda_{\mathcal{Y}}})$. \square

Lemma 13. For self-adjoint trace-class operator $C_{\mathcal{X}}$ and positive $\lambda_{\mathcal{X}}$

$$\left\|C_{\mathcal{X}}^{1/2} - (C_{\mathcal{X}} + \lambda_{\mathcal{X}} I_{\mathcal{X}})^{1/2}\right\| \leq \sqrt{\lambda_{\mathcal{X}}}$$

Proof. Since these operators are commutative,

$$\left\|C_{\mathcal{X}}^{1/2} - (C_{\mathcal{X}} + \lambda_{\mathcal{X}} I_{\mathcal{X}})^{1/2}\right\| = \max_i |\sqrt{\lambda_i(C_{\mathcal{X}})} - \sqrt{\lambda_i(C_{\mathcal{X}} + \lambda_{\mathcal{X}} I_{\mathcal{X}})}| \leq \sqrt{\lambda_{\mathcal{X}}}$$

since the operator being trace class implies that the eigenvalues go to 0. \square

B.2 Bounding the Error 2

Lemma 14 (Error 2). Error 2 is bounded as

$$\langle U_{n,\lambda} V_{n,\lambda}^*, C - C_{\lambda} \rangle \leq k \|C_{\lambda} - C\| \leq 4k (\sqrt{\lambda_{\mathcal{X}}} + \sqrt{\lambda_{\mathcal{Y}}})$$

Proof. The first inequality simply follows from Holder's inequality with conjugates Schatten 1 and ∞ norms, and the second using Lemma 12. \square

B.3 Bounding the Error 3

Lemma 15 (Error 3). With probability at least $1 - \delta$, Error 3 is bounded as,

$$\begin{aligned} \langle U_{C,\lambda} V_{C,\lambda}^* - U_{n,\lambda} V_{n,\lambda}^*, C_{\lambda} \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} &\leq \inf_{h \geq 0} \left\{ \frac{12\alpha_{\rho} h}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{24\alpha_{\rho} \sqrt{h}}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{24}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathcal{C}')} \right\} \\ &\quad + \frac{12\alpha_{\rho}}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{22\beta \sqrt{k} \log(1/\delta)}{\sqrt{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n}} + \frac{10\alpha_{\rho} \log(1/\delta)}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} \end{aligned}$$

Proof. The proof follows the application of local Rademacher complexity analysis technique for Kernel PCA Blanchard et al. (2007), with modifications arising from differences in the problems. We start with the function class

$$\mathcal{F} = \{f_{U,V} : (x, y) \rightarrow \langle \bar{U}_{C,\lambda} \bar{V}_{C,\lambda}^* - UV^*, C_{x,y} \rangle_{L(\mathcal{H}_Y, \mathcal{H}_X)} \mid U^* C_{\mathcal{X}}^{\lambda_{\mathcal{X}}} U = I_k, V^* C_{\mathcal{Y}}^{\lambda_{\mathcal{Y}}} V = I_k\}$$

where $\bar{U}_{C,\lambda} = (C_{\mathcal{X}}^{\lambda_{\mathcal{X}}})^{1/2} U_C$, $\bar{V}_{C,\lambda} = (C_{\mathcal{Y}}^{\lambda_{\mathcal{Y}}})^{1/2} V_C$ and $C_{x,y} = x \otimes_{L(\mathcal{H}_Y, \mathcal{H}_X)} y$

We look at the function class $\mathcal{G} = \tau \mathcal{F}$, where $\tau = \frac{\sqrt{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}}}{2\beta \sqrt{k}}$. From Lemma 17, we get that for $f \in \mathcal{G}$ the range of f is contained in $[-1, 1]$. We then show in Lemma 22 that $\mathbb{E}[f^2] \leq \mu \mathbb{E}[f]$ where $\mu = \frac{2\alpha_{\rho} \tau}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}}$ and

$$\alpha_{\rho} = \frac{\mathbb{E}_{x,y,x',y'} \left[\langle x, x' \rangle_{\mathcal{H}_{\mathcal{X}}}^2 \langle y, y' \rangle_{\mathcal{H}_{\mathcal{Y}}}^2 \right]}{\lambda_k(C) - \lambda_{k+1}(C)} \text{ for } f \in \mathcal{G}.$$

From Lemma 16, we have that $\|UV^*\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{k}{\lambda_x \lambda_y}$. Similarly, we also have $\|\bar{U}_C, \lambda \bar{V}_C, \lambda\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{k}{\lambda_x \lambda_y}$. Therefore, we have

$$\begin{aligned} & \left\| \bar{U}_C \bar{V}_C^* - UV^* \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \\ & \leq 2 \left(\left\| \bar{U}_C \bar{V}_C^* \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 + \|UV^*\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right) \\ & \leq \frac{4k}{\lambda_x \lambda_y} \end{aligned}$$

Therefore, we can write,

$$\mathcal{F} \subseteq \left\{ (x, y) \rightarrow \langle \Gamma, C_{x,y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \Gamma \in L(\mathcal{H}_y, \mathcal{H}_x), \|\Gamma\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_x \lambda_y} \right\} =: \mathcal{H}$$

We will now concern ourselves with the set \mathcal{H} . We have

$$\mathbb{E}[f^2] = \mathbb{E} \left[\langle \Gamma, C_{x,y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] = \mathbb{E}_{x,y} \left[\langle \Gamma, C_{x,y} \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} C_{x,y} \Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \right] = \langle \Gamma, \mathfrak{C}\Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}$$

where $\mathfrak{C} \in L(L(\mathcal{H}_y, \mathcal{H}_x))$ is defined as $\mathfrak{C} := \mathbb{E}_{x,y} [C_{x,y} \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} C_{x,y}]$. Since the set \mathcal{F} is contained in \mathcal{H} , which is a convex set and contains origin, $\text{star}(\mathcal{F}_k)$ is also contained in \mathcal{H} .

$$\text{star}(\mathcal{F}) \subseteq \left\{ (x, y) \rightarrow \langle \Gamma, C_{x,y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \Gamma \in L(\mathcal{H}_y, \mathcal{H}_x), \|\Gamma\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_x \lambda_y} \right\}$$

Moreover,

$$\begin{aligned} \{g \in \text{star}(\mathcal{G}) \mid \mathbb{E}[g^2] \leq r\} &= \tau \{g \in \text{star}(\mathcal{F}) \mid \mathbb{E}[g^2] \leq \tau^{-2}r\} \\ &\subset \tau \{ (x, y) \rightarrow \langle \Gamma, C_{x,y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \Gamma \in L(\mathcal{H}_y, \mathcal{H}_x), \\ &\quad \|\Gamma\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_x \lambda_y}, \langle \Gamma, \mathfrak{C}\Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \tau^{-2}r \} \\ &=: \mathcal{S}_r \end{aligned}$$

We now want to bound the Rademacher Complexity of \mathcal{S}_r which is $\mathfrak{R}_n(\mathcal{S}_r) = \mathbb{E}_{x,y} \mathbb{E}_\sigma \left[\frac{1}{n} \sup_{f \in \mathcal{S}_r} \sum_{i=1}^n \sigma_i f(x_i, y_i) \right]$.

From Lemma 18, we get that the Rademacher complexity of \mathcal{S}_r is bounded as follows,

$$\mathfrak{R}_n(\mathcal{S}_r) \leq \sqrt{\frac{r}{n}} + \frac{1}{\sqrt{n}} \inf_{h \geq 0} \left(\sqrt{rh} + 2\tau \sqrt{\frac{k}{\lambda_x \lambda_y} \sum_{j>h} \lambda_j(\mathfrak{C}')} \right) =: \frac{\psi(r)}{\mu}$$

Note that this is a sub-root function in r , as the infimum for sub-root functions is sub-root. We now need to upper bound the fixed point of $\psi(r)$. Define $\xi := \frac{\mu}{\tau} = \frac{2\alpha_p}{\lambda_x \lambda_y}$. We want,

$$r^* = \psi(r^*) = \frac{\xi\tau}{\sqrt{n}} \left(\sqrt{r^*} (\sqrt{h} + 1) + 2\tau \sqrt{\frac{k}{\lambda_x \lambda_y} \sum_{j>h} \lambda_j(\mathfrak{C}')} \right)$$

From Lemma 21, we have that fixed point r^* is bounded as

$$r^* \leq \tau^2 \left(\inf_{h \geq 0} \left\{ \frac{\xi^2 h}{n} + \frac{2\xi^2 \sqrt{h}}{n} + \frac{4\xi}{\sqrt{n}} \sqrt{\frac{k}{\lambda_x \lambda_y} \sum_{j>h} \lambda_j(\mathfrak{C}')} \right\} + \frac{\xi^2}{n} \right)$$

Let

$$\kappa(\xi, n) = \inf_{h \geq 0} \left\{ \frac{\xi^2 h}{n} + \frac{2\xi^2 \sqrt{h}}{n} + \frac{4\xi}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right\}$$

Having shown that it satisfies all the conditions of Theorem 10, we now apply the theorem. We get that for any $K > 1$, with probability at least $1 - \delta$, $\forall U \in L(\mathcal{H}_{\mathcal{X}}, \mathbb{R}^k)$, $V \in L(\mathcal{H}_{\mathcal{Y}}, \mathbb{R}^k)$ in the feasible set, we have,

$$\begin{aligned} \mathbb{E}[\tau f_{U,V}] &\leq \frac{K\mathbb{E}_n f_{U,V}}{K-1} + \frac{6Kr^*}{\tau\xi} + \frac{(11+5\tau\xi K)\log\delta}{n} \\ &\leq \frac{K\mathbb{E}_n f_{U,V}}{K-1} + \frac{6K\kappa(\xi, k, n)\tau}{\xi} + \frac{6K\xi\tau}{n} + \frac{11\log\delta}{n} + \frac{5\xi K\log\delta\tau}{n} \end{aligned}$$

Therefore,

$$\mathbb{E}[f_{U,V}] \leq \frac{K\mathbb{E}_n f_{U,V}}{K-1} + \frac{6K\kappa(\xi, n)}{\xi} + \frac{6K\xi}{n} + \frac{11\log\delta}{\tau n} + \frac{5\xi K\log\delta}{n}$$

We set $U = C_{\mathcal{X}}^{-1/2}U_n$ and $V = C_{\mathcal{Y}}^{-1/2}V_n$, therefore we get $\mathbb{E}_n f_{U,V} \leq 0$. where \mathbb{E}_n is the expectation with respect to the empirical measure. Letting $K \rightarrow 1$, we get, with probability at least $1 - \delta$,

$$\begin{aligned} &\langle U_{C,\lambda} V_{C,\lambda}^* - U_{n,\lambda} V_{n,\lambda}^*, C_{\lambda} \rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} \\ &\leq \frac{6\kappa(\xi, n)}{\xi} + \frac{6\xi}{n} + \frac{11\log\delta}{\tau n} + \frac{5\xi\log\delta}{n} \end{aligned}$$

Note that $\xi = \frac{2\alpha_{\rho}}{\lambda_{\mathcal{X}}\lambda_{\mathcal{Y}}}$, where $\alpha_{\rho} = \frac{\mathbb{E}_{x,y,x',y'}[\langle x,x' \rangle_{\mathcal{H}_{\mathcal{X}}}^2 \langle y,y' \rangle_{\mathcal{H}_{\mathcal{Y}}}^2]}{(\lambda_k(C) - \lambda_{k+1}(C))}$, and $\tau = \frac{\sqrt{\lambda_{\mathcal{X}}\lambda_{\mathcal{Y}}}}{2\beta\sqrt{k}}$. We now substitute these to obtain the final bound. We have,

$$\begin{aligned} \frac{\kappa(\xi, k, n)}{\xi} &= \frac{1}{\xi} \inf_{h \geq 0} \left\{ \frac{\xi^2 h}{n} + \frac{2\xi^2 \sqrt{h}}{n} + \frac{4\xi}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right\} \\ &= \inf_{h \geq 0} \left\{ \frac{\xi h}{n} + \frac{2\xi \sqrt{h}}{n} + \frac{4}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right\} \\ &= \inf_{h \geq 0} \left\{ \frac{2\alpha_{\rho} h}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{4\alpha_{\rho} \sqrt{h}}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{4}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right\} \end{aligned}$$

Therefore, we get,

$$\begin{aligned} &\langle U_{C,\lambda} V_{C,\lambda}^* - U_{n,\lambda} V_{n,\lambda}^*, C_{\lambda} \rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} \\ &\leq \inf_{h \geq 0} \left\{ \frac{12\alpha_{\rho} h}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{24\alpha_{\rho} \sqrt{h}}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{24}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right\} + \frac{12\alpha_{\rho}}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} + \frac{22\beta\sqrt{k}\log\delta}{\sqrt{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n}} + \frac{10\alpha_{\rho}\log\delta}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}} n} \end{aligned}$$

□

Lemma 16. If U and V satisfies $U^* C_{\mathcal{X}}^{\lambda_{\mathcal{X}}} U = I_k$ and $V^* C_{\mathcal{Y}}^{\lambda_{\mathcal{Y}}} V = I_k$, where $C_{\mathcal{X}}^{\lambda_{\mathcal{X}}} = C_{\mathcal{X}} + \lambda_{\mathcal{X}} I_{\mathcal{X}}$ and $C_{\mathcal{Y}}^{\lambda_{\mathcal{Y}}} = C_{\mathcal{Y}} + \lambda_{\mathcal{Y}} I_{\mathcal{Y}}$, then,

$$\|UV^*\|_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})}^2 \leq \frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}}$$

Proof. We have,

$$\begin{aligned} k &= \left\| \mathbf{U}^* \mathbf{C}_x^{\lambda_x} \mathbf{U} \right\|_{L(L(\mathcal{H}_x, \mathbb{R}^k))}^2 \\ &\geq \lambda_x^2 \left\| \mathbf{U}^* \mathbf{U} \right\|_{L(L(\mathcal{H}_x, \mathbb{R}^k))}^2 \end{aligned}$$

Therefore, we get $\left\| \mathbf{U}^* \mathbf{U} \right\|_{L(L(\mathcal{H}_x, \mathbb{R}^k))} \leq \frac{\sqrt{k}}{\lambda_x}$. Similarly, we can show that $\left\| \mathbf{V}^* \mathbf{V} \right\|_{L(L(\mathcal{H}_y, \mathbb{R}^k))} \leq \frac{\sqrt{k}}{\lambda_y}$.

$$\begin{aligned} \left\| \mathbf{U} \mathbf{V}^* \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 &= \langle \mathbf{U} \mathbf{V}^*, \mathbf{U} \mathbf{V}^* \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\ &= \langle \mathbf{U}^* \mathbf{U}, \mathbf{V}^* \mathbf{V} \rangle_{L(\mathbb{R}^k)} \\ &\leq \left\| \mathbf{U}^* \mathbf{U} \right\|_{L(\mathbb{R}^k)} \left\| \mathbf{V}^* \mathbf{V} \right\|_{L(\mathbb{R}^k)} \\ &\leq \frac{k}{\lambda_x \lambda_y} \end{aligned}$$

where we use the definition of adjoints in the second step and Cauchy-Schwartz inequality in the third step. \square

Lemma 17. For any $f \in \mathcal{G} = \tau \mathcal{F}$ the range of f is contained in $[-1, 1]$, where $\tau = \frac{\sqrt{\lambda_x \lambda_y}}{2\beta\sqrt{k}}$

Proof. From Lemma 16, we get that $\left\| \mathbf{U} \mathbf{V}^* \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{k}{\lambda_x \lambda_y}$. Similarly, we also get $\left\| \mathbf{U}_{C, \lambda} \mathbf{V}_{C, \lambda}^* \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{k}{\lambda_x \lambda_y}$.

We now note that,

$$\begin{aligned} \left\langle \bar{\mathbf{U}}_{C, \lambda} \bar{\mathbf{V}}_{C, \lambda}^* - \mathbf{U} \mathbf{V}^*, \mathbf{C}_{x, y} \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 &\leq \left\| \bar{\mathbf{U}}_{C, \lambda} \bar{\mathbf{V}}_{C, \lambda}^* - \mathbf{U} \mathbf{V}^* \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \left\| \mathbf{C}_{x, y} \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \\ &\leq 2\beta^2 \left(\left\| \bar{\mathbf{U}}_{C, \lambda} \bar{\mathbf{V}}_{C, \lambda}^* \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 + \left\| \mathbf{U} \mathbf{V}^* \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right) \\ &\leq \frac{4\beta^2 k}{\lambda_x \lambda_y} \end{aligned}$$

where in the first step we applied Cauchy-Schwartz inequality and in the second step, we used that $\left\| \mathbf{C}_{x, y} \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)} = \left\| x \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} y \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \|x\|_{\mathcal{H}_x} \|y\|_{\mathcal{H}_y} \leq \beta^2$. Therefore, we have, $\left\langle \bar{\mathbf{U}}_{C, \lambda} \bar{\mathbf{V}}_{C, \lambda}^* - \mathbf{U} \mathbf{V}^*, \mathbf{C}_{x, y} \right\rangle \leq \frac{2\beta\sqrt{k}}{\sqrt{\lambda_x \lambda_y}}$.

Therefore, since any $\bar{f} \in \mathcal{F}$ range of $\bar{f} \leq \tau^{-1}$, we have for $f = \tau \bar{f} \in \tau \mathcal{F} = \mathcal{G}$ its range ≤ 1 . The lower bound holds because $\bar{\mathbf{U}}_C$ and $\bar{\mathbf{V}}_C$ correspond to the optimal solution, therefore for any function $f \in \mathcal{F}$, $\bar{f} \geq 0$. So, any $f \in \tau \mathcal{F} = \mathcal{G}$, $f \geq 0 \geq -1$. \square

Lemma 18. The Rademacher complexity of \mathcal{S}_r , defined as,

$$\mathcal{S}_r = \tau \left\{ (x, y) \rightarrow \langle \Gamma, \mathbf{C}_{x, y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \Gamma \in L(\mathcal{H}_y, \mathcal{H}_x), \left\| \Gamma \right\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_x \lambda_y}, \langle \Gamma, \mathbf{e} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \tau^{-2} r \right\}$$

is bounded as follows,

$$\mathfrak{R}_n(\mathcal{S}_r) \leq \sqrt{\frac{r}{n}} + \frac{1}{\sqrt{n}} \inf_{h \geq 0} \left(\sqrt{r} h + 2\tau \sqrt{\frac{k}{\lambda_x \lambda_y} \sum_{j > h} \lambda_j(\mathbf{e}')} \right)$$

Proof. Note that we can write $\langle \Gamma, C_{x,y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} = \langle \Gamma, C_{x,y} - C_{\mathcal{X}Y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} + \langle \Gamma, C_{\mathcal{X}Y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}$. Equivalently, we can decompose the function class \mathcal{S}_r into two classes \mathcal{P}_r and \mathcal{Q}_r , defined as,

$$\mathcal{P}_r = \tau \left\{ (x, y) \rightarrow \langle \Gamma, C_{\mathcal{X}Y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \langle \Gamma, \mathbf{C}\Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \tau^{-2}r \right\}$$

and

$$\begin{aligned} \mathcal{Q}_r = \tau \left\{ (x, y) \rightarrow \langle \Gamma, C_{x,y} - C_{\mathcal{X}Y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \|\Gamma\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_{\mathcal{X}}\lambda_{\mathcal{Y}}}, \right. \\ \left. \langle \Gamma, (\mathbf{C} - C_{\mathcal{X}Y} \otimes_{L(L(\mathcal{H}_y, \mathcal{H}_x))} C_{\mathcal{X}Y}) \Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \tau^{-2}r \right\} \end{aligned}$$

By a simple application of triangle inequality, we have,

$$\mathfrak{R}_n(\mathcal{S}_r) \leq \mathfrak{R}_n(\mathcal{P}_r) + \mathfrak{R}_n(\mathcal{Q}_r)$$

We bound the Rademacher complexities of the sets in Lemma 19 and 20 respectively. From them, we get,

$$\mathfrak{R}_n(\mathcal{P}_r) \leq \sqrt{\frac{r}{n}}$$

and

$$\mathfrak{R}_n(\mathcal{Q}_r) \leq \frac{1}{\sqrt{n}} \inf_{h \geq 0} \left(\sqrt{rh} + 2\tau \sqrt{\frac{k}{\lambda_{\mathcal{X}}\lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{C}')} \right)$$

Combining these, we have,

$$\mathfrak{R}_n(\mathcal{S}_r) \leq \sqrt{\frac{r}{n}} + \frac{1}{\sqrt{n}} \inf_{h \geq 0} \left(\sqrt{rh} + 2\tau \sqrt{\frac{k}{\lambda_{\mathcal{X}}\lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{C}')} \right)$$

□

Lemma 19. *The Rademacher complexity of the set \mathcal{P}_r , defined as,*

$$\mathcal{P}_r = \tau \left\{ (x, y) \rightarrow \langle \Gamma, C_{\mathcal{X}Y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \langle \Gamma, \mathbf{C}\Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \tau^{-2}r \right\}$$

is bounded as,

$$\mathfrak{R}_n(\mathcal{P}_r) \leq \sqrt{\frac{r}{n}}$$

Proof. Since \mathcal{P}_r contains only constant functions, we can easily bound its Rademacher complexity. In particular, for a set of scalars $Z \subset \mathbb{R}$, we have

$$\begin{aligned} \mathbb{E} \left[\sup_{z \in Z} \left(z \sum_{i=1}^n \sigma_i \right) \right] &= \left(\frac{\sup Z - \inf Z}{2} \right) \mathbb{E} \left[\left| \sum_{i=1}^n \sigma_i \right| \right] \\ &\leq \left(\frac{\sup Z - \inf Z}{2} \right) \sum_{i=1}^n \mathbb{E} [|\sigma_i|] \\ &= \frac{(\sup Z - \inf Z) \sqrt{n}}{2} \end{aligned}$$

where the second step follows from Jensen's inequality, and last step from the fact that $\mathbb{E}[|\sigma|] = 1$ for a Rademacher random variable σ . Let $\bar{f} \in \mathbb{R}^m$ where each co-ordinate is the value of the constant function. Therefore, we have

$$\begin{aligned} \mathfrak{R}_n(\mathcal{P}_r) &= \frac{1}{n} \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{P}_r} \sum_{i=1}^m \sigma_i \bar{f} \right] \\ &\leq \frac{1}{n} \cdot \frac{2\sqrt{n}}{2} \cdot \sup \left\{ \langle \Gamma, \mathbf{C}\mathcal{X}\mathcal{Y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \langle \Gamma, \mathbf{C}\Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \tau^{-2} r \right\} \\ &\leq \sqrt{\frac{r}{n}} \end{aligned}$$

This follows because

$$\langle \Gamma, \mathbf{C}\Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} = \mathbb{E} \left[\langle \mathbf{C}_{x,y}, \Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] \geq \left(\mathbb{E} \left[\langle \mathbf{C}_{x,y}, \Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \right] \right)^2$$

from Jensen's inequality. □

Lemma 20. *The Rademacher complexity of the set \mathcal{Q}_r , defined as*

$$\begin{aligned} \mathcal{Q}_r = &\tau \left\{ (x, y) \rightarrow \langle \Gamma, \mathbf{C}_{x,y} - \mathbf{C}\mathcal{X}\mathcal{Y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \mid \|\Gamma\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_x \lambda_y}, \right. \\ &\left. \langle \Gamma, (\mathbf{C} - \mathbf{C}\mathcal{X}\mathcal{Y} \otimes_{L(L(\mathcal{H}_y, \mathcal{H}_x))} \mathbf{C}\mathcal{X}\mathcal{Y}) \Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \tau^{-2} r \right\} \end{aligned}$$

is bounded as,

$$\mathfrak{R}_n(\mathcal{Q}_r) \leq \frac{1}{\sqrt{n}} \inf_{h \geq 0} \left(\sqrt{r} h + 2\tau \sqrt{\frac{k}{\lambda_x \lambda_y} \sum_{j>h} \lambda_j(\mathbf{C}')} \right)$$

Proof. Let ϕ_i 's be eigenfunctions of $\mathbf{C}' = \mathbf{C} - \mathbf{C}\mathcal{X}\mathcal{Y} \otimes_{L(L(\mathcal{H}_y, \mathcal{H}_x))} \mathbf{C}\mathcal{X}\mathcal{Y}$ which form an orthonormal basis. For $h \leq \text{rank}(\mathbf{C}')$, we have

$$\begin{aligned} \sum_{i=1}^n \sigma_i \langle \Gamma, \mathbf{C}_{x_i, y_i} - \mathbf{C}\mathcal{X}\mathcal{Y} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} &= \sum_{j \geq 1} \langle \Gamma, \phi_j \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \left\langle \phi_j, \sum_{i=1}^n \sigma_i (\mathbf{C}_{x_i, y_i} - \mathbf{C}\mathcal{X}\mathcal{Y}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\ &= \sum_{j=1}^h \langle \Gamma, \phi_j \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \sqrt{\lambda_j(\mathbf{C}')} \left\langle \phi_j, \sum_{i=1}^n \sigma_i (\mathbf{C}_{x_i, y_i} - \mathbf{C}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \cdot \frac{1}{\sqrt{\lambda_j(\mathbf{C}')}} \\ &+ \sum_{j>h} \langle \Gamma, \phi_j \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \left\langle \phi_j, \sum_{i=1}^n \sigma_i (\mathbf{C}_{x_i, y_i} - \mathbf{C}\mathcal{X}\mathcal{Y}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\ &\leq \left(\sum_{j=1}^h \langle \Gamma, \phi_j \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \lambda_j(\mathbf{C}') \right)^{1/2} \left(\sum_{j=1}^h \left\langle \phi_j, \sum_{i=1}^n \sigma_i (\mathbf{C}_{x_i, y_i} - \mathbf{C}\mathcal{X}\mathcal{Y}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \frac{1}{\lambda_j(\mathbf{C}')} \right)^{1/2} \\ &+ \left(\sum_{j>h} \langle \Gamma, \phi_j \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right)^{1/2} \left(\sum_{j>h} \left\langle \phi_j, \sum_{i=1}^n \sigma_i (\mathbf{C}_{x_i, y_i} - \mathbf{C}\mathcal{X}\mathcal{Y}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right)^{1/2} \\ &\leq \frac{\sqrt{r}}{\tau} \left(\sum_{j=1}^h \frac{1}{\lambda_j(\mathbf{C}')} \left\langle \phi_j, \sum_{i=1}^n \sigma_i (\mathbf{C}_{x_i, y_i} - \mathbf{C}\mathcal{X}\mathcal{Y}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right)^{1/2} \\ &+ \left(\frac{4k}{\lambda_x \lambda_y} \right)^{1/2} \left(\sum_{j>h} \left\langle \phi_j, \sum_{i=1}^n \sigma_i (\mathbf{C}_{x_i, y_i} - \mathbf{C}\mathcal{X}\mathcal{Y}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right)^{1/2} \end{aligned}$$

where in the third step, we used Cauchy Schwartz inequality. In the fourth step, we use that

$$\|\Gamma\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 = \sum_i \langle \Gamma, \phi_i \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_x \lambda_y}$$

and

$$\begin{aligned} \langle \Gamma, \mathbf{e}' \Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} &= \left\langle \Gamma, \left(\sum_i \lambda_i(\mathbf{e}') \phi_i \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} \phi_i \right) \Gamma \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\ &= \sum_i \lambda_i(\mathbf{e}') \langle \Gamma, (\phi_i \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} \phi_i) \Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\ &= \sum_i \lambda_i(\mathbf{e}') \left\langle \Gamma, \left(\langle \Gamma, \phi_i \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \phi_i \right) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\ &= \sum_i \lambda_i(\mathbf{e}') \langle \Gamma, \phi_i \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \tau^{-2} r \end{aligned}$$

where the third equality follows from the definition of the outer product.

We now look at,

$$\begin{aligned} &\mathbb{E}_{x,y,\sigma} \left[\left\langle \sum_{j=1}^n \sigma_j (C_{x_j, y_j} - C_{\mathcal{X}\mathcal{Y}}), \phi_i \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] = \mathbb{E}_{x_j, y_j, \sigma} \left[\sum_{j=1}^n \sigma_j^2 \langle C_{x_j, y_j} - C_{\mathcal{X}\mathcal{Y}}, \phi_i \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] \\ &= \mathbb{E}_{x,y} \left[\sum_{j=1}^n \langle C_{x_j, y_j} - C_{\mathcal{X}\mathcal{Y}}, \phi_i \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] \\ &= \mathbb{E}_{x,y} \left[\left\langle \phi_i, \left(\sum_{j=1}^n (C_{x_j, y_j} - C_{\mathcal{X}\mathcal{Y}}) \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} (C_{x_j, y_j} - C_{\mathcal{X}\mathcal{Y}}) \right) \phi_i \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \right] \\ &= \left\langle \phi_i, \mathbb{E}_{x,y} \left[\sum_{j=1}^n (C_{x_j, y_j} - C) \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} (C_{x_j, y_j} - C) \right] \phi_i \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\ &= n \langle \phi_i, \mathbf{e}' \phi_i \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\ &= n \lambda_i(\mathbf{e}') \end{aligned}$$

where in the first step, we use Pythagoras theorem by observing that ϕ_i 's form an orthonormal basis; and in the second step, we use that fact that for a Rademacher variable σ , $\mathbb{E}[\sigma^2] = 1$. In the fourth step, we use that $\mathbf{e}' = \mathbb{E}_{x,y} \left[\frac{1}{n} \sum_{j=1}^n (C_{x_j, y_j} - C_{\mathcal{X}\mathcal{Y}}) \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} (C_{x_j, y_j} - C_{\mathcal{X}\mathcal{Y}}) \right]$, and in the fifth step, we use the fact that ϕ_i is an eigenfunction of \mathbf{e}' .

Let \mathcal{G}_r denote the feasible set of Γ defined as

$$\mathcal{G}_r := \left\{ \Gamma \in L(\mathcal{H}_y, \mathcal{H}_x) \mid \|\Gamma\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \leq \frac{4k}{\lambda_x \lambda_y}, \langle \Gamma, \mathbf{e}' \Gamma \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \leq \tau^{-2} r \right\}$$

We have,

$$\begin{aligned}
\mathfrak{R}_m(\mathcal{Q}_r) &= \frac{\tau}{n} \mathbb{E}_{x,y,\sigma} \left[\sup_{f \in \mathcal{Q}_r} \sum_{i=1}^n \sigma_i f(x_i, y_i) \right] \\
&= \frac{\tau}{n} \mathbb{E}_{x,y,\sigma} \left[\sup_{\Gamma \in \mathcal{G}_r} \sum_{i=1}^n \sigma_i \langle \Gamma, C_{x_i, y_i} - C_{\mathcal{X}\mathcal{Y}} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \right] \\
&\leq \frac{\tau}{n} \mathbb{E} \left[\sup_{\Gamma \in \mathcal{G}_r} \left(\frac{r}{\tau^2} \sum_{j=1}^h \frac{1}{\lambda_j(\mathbf{e}')} \right)^{1/2} \left(\sum_{j=1}^h \left\langle \phi_j, \sum_{i=1}^n \sigma_i (C_{x_i, y_i} - C_{\mathcal{X}\mathcal{Y}}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right)^{1/2} \right] \\
&\quad + \left(\frac{4k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \right)^{1/2} \left(\sum_{j>h} \left\langle \phi_j, \sum_{i=1}^n \sigma_i (C_{x_i, y_i} - C_{\mathcal{X}\mathcal{Y}}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right)^{1/2} \\
&\leq \frac{\tau}{n} \left(\frac{\sqrt{r}}{\tau} \left(\sum_{j=1}^h \frac{1}{\lambda_j(\mathbf{e}')} \mathbb{E}_{x,y,\sigma} \left[\left\langle \phi_j, \sum_{i=1}^n \sigma_i (C_{x_i, y_i} - C_{\mathcal{X}\mathcal{Y}}) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] \right)^{1/2} \right) \\
&\quad + \left(\frac{4k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \right)^{1/2} \left(\sum_{j>h} \mathbb{E}_{x,y,\sigma} \left[\left\langle \phi_j, \sum_{i=1}^n \sigma_i (C_{x_i, y_i} - C) \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] \right)^{1/2} \\
&\leq \frac{1}{\sqrt{n}} \left(\sqrt{r} \sum_{j=1}^h \frac{\lambda_j(\mathbf{e}')}{\lambda_j(\mathbf{e}')} + \sqrt{\sum_{j>h} \lambda_j(\mathbf{e}') \frac{2\sqrt{k}}{\sqrt{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}}}} \right) \\
&= \frac{1}{\sqrt{n}} \left(\sqrt{r}h + 2\tau \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right)
\end{aligned}$$

Since the above holds for all $h \leq \text{rank}(\mathbf{e}')$ and can be trivially extended to $h \geq \text{rank}(\mathbf{e}')$ as $\lambda_j(\mathbf{e}') = 0$ for $j > \text{rank}(\mathbf{e}')$, it therefore holds for the infimum over h . We therefore have,

$$\mathfrak{R}_n(\mathcal{Q}_r) \leq \frac{1}{\sqrt{n}} \inf_{h \geq 0} \left(\sqrt{r}h + 2\tau \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right)$$

□

Lemma 21. *The fixed point of $\psi(r)$, i.e.*

$$r^* = \psi(r^*) = \frac{\xi\tau}{\sqrt{n}} \left(\sqrt{r^*} (\sqrt{h} + 1) + 2\tau \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right)$$

is bounded as,

$$r^* \leq \tau^2 \left(\inf_{h \geq 0} \left\{ \frac{\xi^2 h}{n} + \frac{2\xi^2 \sqrt{h}}{n} + \frac{4\xi}{\sqrt{m}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right\} + \frac{\xi^2}{n} \right)$$

Proof. We have,

$$r^* = \psi(r^*) = \frac{\xi\tau}{\sqrt{n}} \left(\sqrt{r^*} (\sqrt{h} + 1) + 2\tau \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathbf{e}')} \right)$$

Consider the quadratic equation $x - a\sqrt{x} - b \leq 0$, we have,

$$\begin{aligned} x &\leq \left(\frac{a + \sqrt{a^2 + 4b}}{2} \right)^2 \\ &\leq \frac{2a^2 + 4b + 2\sqrt{a^2(a^2 + 4b)}}{4} \\ &\leq a^2 + 2b \end{aligned}$$

where in the last step, we use that geometric mean \leq arithmetic mean. Plugging it here, we get

$$r^* \leq \frac{\xi^2 \tau^2}{n} (h + 1 + 2\sqrt{h}) + \frac{4\xi \tau^2}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathfrak{C}')}$$

Taking infimum over h , we get,

$$r^* \leq \tau^2 \left(\inf_{h \geq 0} \left\{ \frac{\xi^2 h}{n} + \frac{2\xi^2 \sqrt{h}}{n} + \frac{4\xi}{\sqrt{n}} \sqrt{\frac{k}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \sum_{j>h} \lambda_j(\mathfrak{C}')} \right\} + \frac{\xi^2}{n} \right)$$

□

Lemma 22. For any $f \in \mathcal{G}$, $\mathbb{E}[f^2] \leq \mu \mathbb{E}[f]$ where $\mu = \frac{2\alpha_\rho \tau}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}}$ and $\alpha_\rho = \frac{\mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}'} [\langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{H}_{\mathcal{X}}}^2 \langle \mathbf{y}, \mathbf{y}' \rangle_{\mathcal{H}_{\mathcal{Y}}}^2]}{(\lambda_k(\mathcal{C}) - \lambda_{k+1}(\mathcal{C}))}$.

Proof. Given U, V , define $\bar{U} := (C_{\mathcal{X}}^{\lambda_{\mathcal{X}}})^{1/2} U$, $\bar{V} := (C_{\mathcal{Y}}^{\lambda_{\mathcal{Y}}})^{1/2} V$. We remind that $U_{\mathcal{C}, \lambda} = (C_{\mathcal{X}}^{\lambda_{\mathcal{X}}})^{1/2} \bar{U}_{\mathcal{C}, \lambda}$ and $V_{\mathcal{C}, \lambda} = (C_{\mathcal{X}}^{\lambda_{\mathcal{X}}})^{1/2} \bar{V}_{\mathcal{C}, \lambda}$, so we get $U_{\mathcal{C}, \lambda}^* U_{\mathcal{C}, \lambda} = I$ and $V_{\mathcal{C}, \lambda}^* V_{\mathcal{C}, \lambda} = I$. Define the projection $P_{\mathcal{C}, \lambda} := U_{\mathcal{C}, \lambda} V_{\mathcal{C}, \lambda}^* = \sum_{i=1}^k u_i^{C_{\lambda}} \otimes_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} v_i^{C_{\lambda}}$ and $P := \bar{U} \bar{V}^* = \sum_{i=1}^k u_i \otimes_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} v_i$ using their singular value decomposition respectively. Let $f \in \mathcal{G}$. We first look at $\mathbb{E}[f^2]$.

$$\begin{aligned} \mathbb{E}[f^2] &= \mathbb{E} \left[\left\langle \bar{U}_{\mathcal{C}, \lambda} \bar{V}_{\mathcal{C}, \lambda}^* - UV^*, C_{x, y} \right\rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})}^2 \right] \\ &= \langle \bar{U}_{\mathcal{C}, \lambda} \bar{V}_{\mathcal{C}, \lambda}^* - UV^*, \mathbb{E}[(C_{x, y} \otimes_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} C_{x, y})] (\bar{U}_{\mathcal{C}, \lambda} \bar{V}_{\mathcal{C}, \lambda}^* - UV^*) \rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} \\ &= \langle \bar{U}_{\mathcal{C}, \lambda} \bar{V}_{\mathcal{C}, \lambda}^* - UV^*, \mathfrak{C}(\bar{U}_{\mathcal{C}, \lambda} \bar{V}_{\mathcal{C}, \lambda}^* - UV^*) \rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} \\ &= \|\mathfrak{C}\|_{L(L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}}))} \left\| \bar{U}_{\mathcal{C}, \lambda} \bar{V}_{\mathcal{C}, \lambda}^* - UV^* \right\|_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})}^2 \\ &= \|\mathfrak{C}\|_{L(L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}}))} \left\| (C_{\mathcal{X}}^{\lambda_{\mathcal{X}}})^{-1/2} (U_{\mathcal{C}, \lambda} V_{\mathcal{C}, \lambda}^* - \bar{U} \bar{V}^*) C_{\mathcal{Y}}^{\lambda_{\mathcal{Y}} - 1/2} \right\|_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})}^2 \\ &= \|\mathfrak{C}\|_{L(L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}}))} \left\| (C_{\mathcal{X}}^{\lambda_{\mathcal{X}}})^{-1/2} (P_{\mathcal{C}, \lambda} - P) (C_{\mathcal{Y}}^{\lambda_{\mathcal{Y}}})^{-1/2} \right\|_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})}^2 \\ &\leq \frac{1}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \|\mathfrak{C}\|_{L(L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}}))} \|P_{\mathcal{C}, \lambda} - P\|_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})}^2 \\ &= \frac{2}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \|\mathfrak{C}\|_{L(L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}}))} \left(k - \langle P_{\mathcal{C}, \lambda}, P \rangle_{L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}})} \right) \\ &= \frac{2}{\lambda_{\mathcal{X}} \lambda_{\mathcal{Y}}} \|\mathfrak{C}\|_{L(L(\mathcal{H}_{\mathcal{Y}}, \mathcal{H}_{\mathcal{X}}))} \left(k - \sum_{i, j=1}^k \langle u_i^{C_{\lambda}}, u_j \rangle_{\mathcal{H}_{\mathcal{X}}} \langle v_i^{C_{\lambda}}, v_j \rangle_{\mathcal{H}_{\mathcal{Y}}} \right) \end{aligned}$$

where in the seventh inequality we just expanded $\|P_{C,\lambda} - P\|_{L(\mathcal{H}_y, \mathcal{H}_x)}^2$. Now, note that

$$\begin{aligned}
\|\mathfrak{C}\|_{L(L(\mathcal{H}_y, \mathcal{H}_x))}^2 &= \langle \mathfrak{C}, \mathfrak{C} \rangle_{L(L(\mathcal{H}_y, \mathcal{H}_x))} \\
&= \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}'} \left[\langle C_{\mathbf{x}, \mathbf{y}} \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} C_{\mathbf{x}, \mathbf{y}}, C_{\mathbf{x}', \mathbf{y}'} \otimes_{L(\mathcal{H}_y, \mathcal{H}_x)} C_{\mathbf{x}', \mathbf{y}'} \rangle_{L(L(\mathcal{H}_y, \mathcal{H}_x))} \right] \\
&= \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}'} \left[\langle C_{\mathbf{x}, \mathbf{y}}, C_{\mathbf{x}', \mathbf{y}'} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] \\
&= \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}'} \left[\langle \mathbf{x} \otimes \mathbf{y}, \mathbf{x}' \otimes \mathbf{y}' \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)}^2 \right] \\
&= \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}'} \left[\langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{H}_x}^2 \langle \mathbf{y}, \mathbf{y}' \rangle_{\mathcal{H}_y}^2 \right] =: \alpha_\rho
\end{aligned}$$

Let us now look at $\mathbb{E}[f]$.

$$\begin{aligned}
\mathbb{E}[f] &= \mathbb{E} \left[\langle \bar{U}_{C,\lambda} \bar{V}_{C,\lambda}^* - UV^*, C_{\mathbf{x}, \mathbf{y}} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \right] \\
&= \langle \bar{U}_{C,\lambda} \bar{V}_{C,\lambda}^* - UV^*, C_{\mathcal{X}\mathcal{Y}} \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\
&= \left\langle P_{C,\lambda} - P, \left(C_{\mathcal{X}^x}^{\lambda_x} \right)^{-1/2} C_{\mathcal{X}\mathcal{Y}} \left(C_{\mathcal{Y}^y}^{\lambda_y} \right)^{-1/2} \right\rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\
&= \langle P_{C,\lambda} - P, C_\lambda \rangle_{L(\mathcal{H}_y, \mathcal{H}_x)} \\
&= \sum_{i=1}^k (\sigma_i(C_\lambda) - \langle u_i, C_\lambda v_i \rangle_{\mathcal{H}_x})
\end{aligned}$$

Let $u_i = \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_x} u_j^{C_\lambda} + r_i$, where r_i is orthogonal to $u_j^{C_\lambda}, j \in [k]$ and $v_i = \sum_{j=1}^k \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_y} v_j^{C_\lambda} + s_i$, where s_i is orthogonal to $v_j^{C_\lambda}, j \in [k]$. Then

$$\begin{aligned}
\langle u_i, C_\lambda v_i \rangle_{\mathcal{H}_x} &= \left\langle \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_x} u_j^{C_\lambda} + r_i, C_\lambda \left(\sum_{j=1}^k \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_y} v_j^{C_\lambda} + s_i \right) \right\rangle \\
&= \left\langle \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_x} u_j^{C_\lambda} + r_i, \sum_{j=1}^k \lambda_j(C_\lambda) \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_y} u_j^{C_\lambda} + C_\lambda s_i \right\rangle \\
&= \sum_{j=1}^k \lambda_j(C_\lambda) \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_x} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_y} + \langle r_i, C_\lambda s_i \rangle_{\mathcal{H}_x}
\end{aligned}$$

The cross terms are zero because $C_\lambda s_i$'s will be a linear combination of $u_i, i > k$ and so orthogonal to $u_j, j \in [k]$. Note that

$$\begin{aligned}
\langle r_i, C_\lambda s_i \rangle_{\mathcal{H}_X} &\leq \lambda_{k+1}(C_\lambda) \|r_i\|_{\mathcal{H}_X} \|s_i\|_{\mathcal{H}_Y} \\
&= \lambda_{k+1}(C_\lambda) \left\| u_i - \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle u_j^{C_\lambda} \right\|_{\mathcal{H}_X} \cdot \left\| v_i - \sum_{j=1}^k \langle v_i, v_j^{C_\lambda} \rangle v_j^{C_\lambda} \right\|_{\mathcal{H}_Y} \\
&= \lambda_{k+1}(C_\lambda) \sqrt{1 - \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X}^2} \sqrt{1 - \sum_{j=1}^k \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y}^2} \\
&\leq \lambda_{k+1}(C_\lambda) \left(1 - \frac{\sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X}^2 + \sum_{j=1}^k \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y}^2}{2} \right) \\
&\leq \lambda_{k+1}(C_\lambda) \left(1 - \sqrt{\left(\sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X}^2 \right) \left(\sum_{j=1}^k \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y}^2 \right)} \right) \\
&\leq \lambda_{k+1}(C_\lambda) \left(1 - \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \right)
\end{aligned}$$

where the first step follows because r_i and s_i don't include components along the first k u_i 's and v_i respectively. In the fourth and fifth steps, we use that arithmetic mean \geq geometric mean, and in the last step, we use Cauchy-Schwartz inequality. Plugging this in the previous bound, we get,

$$\langle u_i, C_\lambda v_i \rangle_{\mathcal{H}_X} \leq \sum_{j=1}^k \lambda_j(C_\lambda) \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} + \lambda_{k+1}(C_\lambda) \left(1 - \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \right)$$

Moreover,

$$\begin{aligned}
\mathbb{E} [f] &\geq \sum_{i=1}^k \left(\lambda_i(C_\lambda) - \left(\sum_{j=1}^k \lambda_j(C_\lambda) \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} + \lambda_{k+1}(C_\lambda) \left(1 - \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \right) \right) \right) \\
&= \sum_{i=1}^k \lambda_i(C_\lambda) \left(1 - \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \right) - \sum_{i=1}^k \lambda_{k+1}(C_\lambda) \left(1 - \sum_{j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \right) \\
&\geq k(\lambda_k(C_\lambda) - \lambda_{k+1}(C_\lambda)) - (\lambda_k(C_\lambda) - \lambda_{k+1}(C_\lambda)) \sum_{i,j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \\
&= (\lambda_k(C_\lambda) - \lambda_{k+1}(C_\lambda)) \left(k - \sum_{i,j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \right)
\end{aligned}$$

where in the second last step, we used that $\lambda_i \geq \lambda_k, i \in [k]$ and $\left(1 - \sum_{i,j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \right) \geq \|r_i\|_{\mathcal{H}_X} \|s_i\|_{\mathcal{H}_Y} \geq 0$ (see above). Therefore, we get,

$$\begin{aligned}
\mathbb{E} [f] &\geq (\lambda_k(C_\lambda) - \lambda_{k+1}(C_\lambda)) \left(k - \sum_{i,j=1}^k \langle u_i, u_j^{C_\lambda} \rangle_{\mathcal{H}_X} \langle v_i, v_j^{C_\lambda} \rangle_{\mathcal{H}_Y} \right) \\
&\geq \frac{(\lambda_k(C_\lambda) - \lambda_{k+1}(C_\lambda)) \lambda_X \lambda_Y}{2\alpha_\rho} \mathbb{E} [f^2]
\end{aligned}$$

Let $\xi = \frac{2\alpha_\rho}{\lambda_x \lambda_y}$. For $f \in \mathcal{G}$, let $\bar{f} = \tau^{-1}f$ where $\bar{f} \in \mathcal{F}$. Therefore,

$$\begin{aligned}\mathbb{E}[f^2] &= \tau^2 \mathbb{E}[\bar{f}^2] \\ &\leq \xi \tau^2 \mathbb{E}[f] \\ &= \xi \tau \mathbb{E}[f] = \mu \mathbb{E}[f]\end{aligned}$$

where $\mu = \frac{2\alpha_\rho \tau}{\lambda_x \lambda_y}$.

□