

MIT Open Access Articles

Learning with group invariant features: A Kernel perspective

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Mroueh, Youssef, Stephen Voinea and Tomaso Poggio. "Learning with Group Invariant Features: A Kernel Perspective." Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1 (NIPS '15), December 7-12, 2015, Montreal, Canada, Association of Computing Machinery, December 2015. © 2015 Association of Computing Machinery ACM

As Published: <https://dl.acm.org/citation.cfm?id=2969413>

Publisher: Association for Computing Machinery

Persistent URL: <http://hdl.handle.net/1721.1/112309>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike



Learning with Group Invariant Features: A Kernel Perspective.

Youssef Mroueh
Multimodal Algorithms & Engines Group
IBM T.J Watson Reseach Center
mroueh@us.ibm.com

Stephen Voinea*
CBMM, MIT.
voinea@mit.edu
*Co-first author

Tomaso Poggio
CBMM, MIT.
tp@ai.mit.edu

Abstract

We analyze in this paper a random feature map based on a theory of invariance (*I-theory*) introduced in [1]. More specifically, a group invariant signal signature is obtained through cumulative distributions of group-transformed random projections. Our analysis bridges invariant feature learning with kernel methods, as we show that this feature map defines an expected Haar-integration kernel that is invariant to the specified group action. We show how this non-linear random feature map approximates this group invariant kernel uniformly on a set of N points. Moreover, we show that it defines a function space that is dense in the equivalent Invariant Reproducing Kernel Hilbert Space. Finally, we quantify error rates of the convergence of the empirical risk minimization, as well as the reduction in the sample complexity of a learning algorithm using such an invariant representation for signal classification, in a classical supervised learning setting.

1 Introduction

Encoding signals or building similarity kernels that are invariant to the action of a group is a key problem in unsupervised learning, as it reduces the complexity of the learning task and mimics how our brain represents information invariantly to symmetries and various nuisance factors (change in lighting in image classification and pitch variation in speech recognition) [1, 2, 3, 4]. Convolutional neural networks [5, 6] achieve state of the art performance in many computer vision and speech recognition tasks, but require a large amount of labeled examples as well as augmented data, where we reflect symmetries of the world through virtual examples [7, 8] obtained by applying identity-preserving transformations such as shearing, rotation, translation, etc., to the training data. In this work, we adopt the approach of [1], where the representation of the signal is designed to reflect the invariant properties and model the world symmetries with group actions. The ultimate aim is to bridge unsupervised learning of invariant representations with invariant kernel methods, where we can use tools from classical supervised learning to easily address the statistical consistency and sample complexity questions [9, 10]. Indeed, many invariant kernel methods and related invariant kernel networks have been proposed. We refer the reader to the related work section for a review (Section 5) and we start by showing how to accomplish this invariance through group-invariant Haar-integration kernels [11], and then show how random features derived from a memory-based theory of invariances introduced in [1] approximate such a kernel.

1.1 Group Invariant Kernels

We start by reviewing group-invariant Haar-integration kernels introduced in [11], and their use in a binary classification problem. This section highlights the conceptual advantages of such kernels as well as their practical inconvenience, putting into perspective the advantage of approximating them with explicit and invariant random feature maps.

Invariant Haar-Integration Kernels. We consider a subset \mathcal{X} of the hypersphere in d dimensions \mathbb{S}^{d-1} . Let $\rho_{\mathcal{X}}$ be a measure on \mathcal{X} . Consider a kernel k_0 on \mathcal{X} , such as a radial basis function kernel. Let G be a group acting on \mathcal{X} , with a normalized Haar measure μ . G is assumed to be a compact and unitary group. Define an invariant kernel \mathcal{K} between $x, z \in \mathcal{X}$ through Haar-integration [11] as follows:

$$\mathcal{K}(x, z) = \int_G \int_G k_0(gx, g'z) d\mu(g) d\mu(g'). \quad (1)$$

As we are integrating over the entire group, it is easy to see that: $\mathcal{K}(g'x, gz) = \mathcal{K}(x, z)$, $\forall g, g' \in G, \forall x, z \in \mathcal{X}$. Hence the Haar-integration kernel is invariant to the group action. The symmetry of \mathcal{K} is obvious. Moreover, if k_0 is a positive definite kernel, it follows that \mathcal{K} is positive definite as well [11]. One can see the Haar-integration kernel framework as another form of data augmentation, since we have to produce group-transformed points in order to compute the kernel.

Invariant Decision Boundary. Turning now to a binary classification problem, we assume that we are given a labeled training set: $S = \{(x_i, y_i) \mid x_i \in \mathcal{X}, y_i \in \mathcal{Y} = \{\pm 1\}\}_{i=1}^N$. In order to learn a decision function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we minimize the following empirical risk induced by an L -Lipschitz, convex loss function V , with $V'(0) < 0$ [12]: $\min_{f \in \mathcal{H}_{\mathcal{K}}} \hat{\mathcal{E}}_V(f) := \frac{1}{N} \sum_{i=1}^N V(y_i f(x_i))$, where we restrict f to belong to a hypothesis class induced by the invariant kernel \mathcal{K} , the so called Reproducing Kernel Hilbert Space $\mathcal{H}_{\mathcal{K}}$. The representer theorem [13] shows that the solution of such a problem, or the optimal decision boundary f_N^* has the following form: $f_N^*(x) = \sum_{i=1}^N \alpha_i^* \mathcal{K}(x, x_i)$. Since the kernel \mathcal{K} is group-invariant it follows that: $f_N^*(gx) = \sum_{i=1}^N \alpha_i^* \mathcal{K}(gx, x_i) = \sum_{i=1}^N \alpha_i^* \mathcal{K}(x, x_i) = f_N^*(x)$, $\forall g \in G$. Hence the decision boundary f^* is group-invariant as well, and we have: $f_N^*(gx) = f_N^*(x)$, $\forall g \in G, \forall x \in \mathcal{X}$.

Reduced Sample Complexity. We have shown that a group-invariant kernel induces a group-invariant decision boundary, but how does this translate to the sample complexity of the learning algorithm? To answer this question, we will assume that the input set \mathcal{X} has the following structure: $\mathcal{X} = \mathcal{X}_0 \cup \mathcal{G}\mathcal{X}_0$, $\mathcal{G}\mathcal{X}_0 = \{z \mid z = gx, x \in \mathcal{X}_0, g \in G / \{e\}\}$, where e is the identity group element. This structure implies that for a function f in the invariant RKHS $\mathcal{H}_{\mathcal{K}}$, we have:

$$\forall z \in \mathcal{G}\mathcal{X}_0, \exists x \in \mathcal{X}_0, \exists g \in G \text{ such that } z = gx, \text{ and } f(z) = f(x).$$

Let $\rho_y(x) = \mathbb{P}(Y = y \mid x)$ be the label posteriors. We assume that $\rho_y(gx) = \rho_y(x)$, $\forall g \in G$. This is a natural assumption since the label is unchanged given the group action. Assume that the set \mathcal{X} is endowed with a measure $\rho_{\mathcal{X}}$ that is also group-invariant. Let f be the group-invariant decision function and consider the expected risk induced by the loss V , $\mathcal{E}_V(f)$, defined as follows:

$$\mathcal{E}_V(f) = \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} V(yf(x)) \rho_y(x) \rho_{\mathcal{X}}(x) dx, \quad (2)$$

$\mathcal{E}_V(f)$ is a proxy to the misclassification risk [12]. Using the invariant properties of the function class and the data distribution we have by invariance of f , ρ_y , and ρ :

$$\begin{aligned} \mathcal{E}_V(f) &= \int_{\mathcal{X}_0} \sum_{y \in \mathcal{Y}} V(yf(x)) \rho_y(x) \rho_{\mathcal{X}}(x) dx + \int_{\mathcal{G}\mathcal{X}_0} \sum_{y \in \mathcal{Y}} V(yf(z)) \rho_y(z) \rho_{\mathcal{X}}(z) dz \\ &= \int_G d\mu(g) \int_{\mathcal{X}_0} \sum_{y \in \mathcal{Y}} V(yf(gx)) \rho_y(gx) \rho_{\mathcal{X}}(x) dx \\ &= \int_G d\mu(g) \int_{\mathcal{X}_0} \sum_{y \in \mathcal{Y}} V(yf(x)) \rho_y(x) \rho_{\mathcal{X}}(x) dx \quad (\text{By invariance of } f, \rho_y, \text{ and } \rho) \\ &= \int_{\mathcal{X}_0} \sum_{y \in \mathcal{Y}} V(yf(x)) \rho_y(x) \rho_{\mathcal{X}}(x) dx. \end{aligned}$$

Hence, given an invariant kernel to a group action that is identity preserving, it is sufficient to minimize the empirical risk on the core set \mathcal{X}_0 , and it generalizes to samples in $\mathcal{G}\mathcal{X}_0$.

Let us imagine that \mathcal{X} is finite with cardinality $|\mathcal{X}|$; the cardinality of the core set \mathcal{X}_0 is a small fraction of the cardinality of \mathcal{X} : $|\mathcal{X}_0| = \alpha |\mathcal{X}|$, where $0 < \alpha < 1$. Hence, when we sample training points from \mathcal{X}_0 , the maximum size of the training set is $N = \alpha |\mathcal{X}| \ll |\mathcal{X}|$, yielding a reduction in the sample complexity.

1.2 Contributions

We have just reviewed the group-invariant Haar-integration kernel. In summary, a group-invariant kernel implies the existence of a decision function that is invariant to the group action, as well as a reduction in the sample complexity due to sampling training points from a reduced set, a.k.a the core set \mathcal{X}_0 .

Kernel methods with Haar-integration kernels come at a very expensive computational price at both training and test time: computing the Kernel is computationally cumbersome as we have to integrate over the group and produce virtual examples by transforming points explicitly through the group action. Moreover, the training complexity of kernel methods scales cubically in the sample size. Those practical considerations make the usefulness of such kernels very limited.

The contributions of this paper are on three folds:

1. We first show that a non-linear random feature map $\Phi : \mathcal{X} \rightarrow \mathbb{R}^D$ derived from a memory-based theory of invariances introduced in [1] induces an expected group-invariant Haar-integration kernel K . For fixed points $x, z \in \mathcal{X}$, we have: $\mathbb{E} \langle \Phi(x), \Phi(z) \rangle = K(x, z)$, where K satisfies: $K(gx, g'z) = K(x, z), \forall g, g' \in G, x, z \in \mathcal{X}$.
2. We show a Johnson-Lindenstrauss type result that holds uniformly on a set of N points that assess the concentration of this random feature map around its expected induced kernel. For sufficiently large D , we have $\langle \Phi(x), \Phi(z) \rangle \approx K(x, z)$, uniformly on an N points set.
3. We show that, with a linear model, an invariant decision function can be learned in this random feature space by sampling points from the core set \mathcal{X}_0 i.e: $f_N^*(x) \approx \langle w^*, \Phi(x) \rangle$ and generalizes to unseen points in $\mathcal{G}\mathcal{X}_0$, reducing the sample complexity. Moreover, we show that those features define a function space that approximates a dense subset of the invariant RKHS, and assess the error rates of the empirical risk minimization using such random features.
4. We demonstrate the validity of these claims on three datasets: text (artificial), vision (MNIST), and speech (TIDIGITS).

2 From Group Invariant Kernels to Feature Maps

In this paper we show that a random feature map based on I-theory [1]: $\Phi : \mathcal{X} \rightarrow \mathbb{R}^D$ approximates a group-invariant Haar-integration kernel K having the form given in Equation (1):

$$\langle \Phi(x), \Phi(z) \rangle \approx K(x, z).$$

We start with some notation that will be useful for defining the feature map. Denote the cumulative distribution function of a random variable X by,

$$F_X(\tau) = \mathbb{P}(X \leq \tau),$$

Fix $x \in \mathcal{X}$, Let $g \in G$ be a random variable drawn according to the normalized Haar measure μ and let t be a random template whose distribution will be defined later. For $s > 0$, define the following truncated cumulative distribution function (CDF) of the dot product $\langle x, gt \rangle$:

$$\psi(x, t, \tau) = \mathbb{P}_g(\langle x, gt \rangle \leq \tau) = F_{\langle x, gt \rangle}(\tau), \tau \in [-s, s], x \in \mathcal{X},$$

Let $\varepsilon \in (0, 1)$. We consider the following Gaussian vectors (sampling with rejection) for the templates t :

$$t = n \sim \mathcal{N}\left(0, \frac{1}{d}I_d\right), \text{ if } \|n\|_2^2 < 1 + \varepsilon, t = \perp \text{ else.}$$

The reason behind this sampling is to keep the range of $\langle x, gt \rangle$ under control: The squared norm $\|n\|_2^2$ will be bounded by $1 + \varepsilon$ with high probability by a classical concentration result (See proof of Theorem 1 for more details). The group being unitary and $x \in \mathbb{S}^{d-1}$, we know that: $|\langle x, gt \rangle| \leq \|n\|_2 < \sqrt{1 + \varepsilon} \leq 1 + \varepsilon$, for $\varepsilon \in (0, 1)$.

Remark 1. We can also consider templates t , drawn uniformly on the unit sphere \mathbb{S}^{d-1} . Uniform templates on the sphere can be drawn as follows:

$$t = \frac{\nu}{\|\nu\|_2}, \nu \sim \mathcal{N}(0, I_d),$$

since the norm of a gaussian vector is highly concentrated around its mean \sqrt{d} , we can use the gaussian sampling with rejection. Results proved for gaussian templates (with rejection) will hold true for templates drawn at uniform on the sphere with different constants.

Define the following kernel function,

$$K_s(x, z) = \mathbb{E}_t \int_{-s}^s \psi(x, t, \tau) \psi(z, t, \tau) d\tau,$$

where s will be fixed throughout the paper to be $s = 1 + \varepsilon$ since the gaussian sampling with rejection controls the dot product to be in that range.

Let $\bar{g} \in G$. As the group is closed, we have $\psi(t, \bar{g}x, \tau) = \int_G \mathbb{1}_{\langle g\bar{g}x, t \rangle \leq \tau} d\mu(g) = \int_G \mathbb{1}_{\langle gx, t \rangle \leq \tau} d\mu(g) = \psi(t, x, \tau)$ and hence $K_s(gx, g'z) = K_s(x, z)$, for all $g, g' \in G$. It is clear now that K is a group-invariant kernel.

In order to approximate K_s , we sample $|G|$ elements uniformly and independently from the group G , i.e. $g_i, i = 1 \dots |G|$, and define the normalized empirical CDF :

$$\phi(x, t, \tau) = \frac{1}{|G|\sqrt{m}} \sum_{i=1}^{|G|} \mathbb{1}_{\langle g_i t, x \rangle \leq \tau}, \quad -s \leq \tau \leq s.$$

We discretize the continuous threshold τ as follows:

$$\phi\left(x, t, \frac{sk}{n}\right) = \frac{\sqrt{s}}{\sqrt{nm}|G|} \sum_{i=1}^{|G|} \mathbb{1}_{\langle g_i t, x \rangle \leq \frac{s}{n}k}, \quad -n \leq k \leq n.$$

We sample m templates independently according to the Gaussian sampling with rejection, $t_j, j = 1 \dots m$. We are now ready to define the random feature map Φ :

$$\Phi(x) = \left[\phi\left(x, t_j, \frac{sk}{n}\right) \right]_{j=1 \dots m, k=-n \dots n} \in \mathbb{R}^{(2n+1) \times m}.$$

It is easy to see that:

$$\lim_{n \rightarrow \infty} \mathbb{E}_{t, g} \langle \Phi(x), \Phi(z) \rangle_{\mathbb{R}^{(2n+1) \times m}} = \lim_{n \rightarrow \infty} \mathbb{E}_{t, g} \sum_{j=1}^m \sum_{k=-n}^n \phi\left(x, t_j, \frac{sk}{n}\right) \phi\left(z, t_j, \frac{sk}{n}\right) = K_s(x, z).$$

In Section 3 we study the geometric information captured by this kernel by stating explicitly the similarity it computes.

Remark 2 (Efficiency of the representation). 1) The main advantage of such a feature map, as outlined in [1], is that we store transformed templates in order to compute Φ , while if we wanted to compute an invariant kernel of type \mathcal{K} (Equation (1)), we would need to explicitly transform the points. The latter is computationally expensive. Storing transformed templates and computing the signature Φ is much more efficient. It falls in the category of memory-based learning, and is biologically plausible [1].

2) As $|G|, m, n$ get large enough, the feature map Φ approximates a group-invariant Kernel, as we will see in next section.

3 An Equivalent Expected Kernel and a Uniform Concentration Result

In this section we present our main results, with proofs given in the supplementary material. Theorem 1 shows that the random feature map Φ , defined in the previous section, corresponds in expectation to a group-invariant Haar-integration kernel $K_s(x, z)$. Moreover, $s - K_s(x, z)$ computes the average pairwise distance between all points in the orbits of x and z , where the orbit is defined as the collection of all group-transformations of a given point $x : \mathcal{O}_x = \{gx, g \in G\}$.

Theorem 1 (Expectation). Let $\varepsilon \in (0, 1)$ and $x, z \in \mathcal{X}$. Define the distance d_G between the orbits \mathcal{O}_x and \mathcal{O}_z :

$$d_G(x, z) = \frac{1}{\sqrt{2\pi d}} \int_G \int_G \|gx - g'z\|_2 d\mu(g) d\mu(g'),$$

and the group-invariant expected kernel

$$K_s(x, z) = \lim_{n \rightarrow \infty} \mathbb{E}_{t, g} \langle \Phi(x), \Phi(z) \rangle_{\mathbb{R}^{(2n+1) \times m}} = \mathbb{E}_t \int_{-s}^s \psi(x, t, \tau) \psi(z, t, \tau) d\tau, \quad s = 1 + \varepsilon.$$

1. The following inequality holds with probability 1:

$$\varepsilon - \delta_2(d, \varepsilon) \leq K_s(x, z) - (1 - d_G(x, z)) \leq \varepsilon + \delta_1(d, \varepsilon), \quad (3)$$

$$\text{where } \delta_1(\varepsilon, d) = \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}} - \frac{1}{2} \frac{e^{-\varepsilon d/2}(1+\varepsilon)^{\frac{d}{2}}}{\sqrt{d}} \text{ and } \delta_2(\varepsilon, \delta) = \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}} + (1 + \varepsilon)e^{-d\varepsilon^2/8}.$$

2. For any $\varepsilon \in (0, 1)$ as the dimension $d \rightarrow \infty$ we have $\delta_1(\varepsilon, d) \rightarrow 0$ and $\delta_2(\varepsilon, d) \rightarrow 0$, and we have asymptotically $K_s(x, z) \rightarrow 1 - d_G(x, z) + \varepsilon = s - d_G(x, z)$.

3. K_s is symmetric and K_s is positive semi-definite.

Remark 3. 1) $\varepsilon, \delta_1(d, \varepsilon)$, and $\delta_2(d, \varepsilon)$ are not errors due to results holding with high probability but are due to the truncation and are a technical artifact of the proof. 2) Local invariance can be defined by restricting the sampling of the group elements to a subset $\mathcal{G} \subset G$. Assuming that for each $g \in \mathcal{G}, g^{-1} \in \mathcal{G}$, the equivalent kernel has asymptotically the following form:

$$K_s(x, z) \approx s - \frac{1}{\sqrt{2\pi d}} \int_{\mathcal{G}} \int_{\mathcal{G}} \|gx - g'z\|_2 d\mu(g)d\mu(g').$$

3) The norm-one constraint can be relaxed, let $R = \sup_{x \in \mathcal{X}} \|x\|_2 < \infty$, hence we can set $s = R(1 + \varepsilon)$, and

$$-\delta_2(d, \varepsilon) \leq K_s(x, z) - (R(1 + \varepsilon) - d_G(x, z)) \leq \delta_1(d, \varepsilon), \quad (4)$$

$$\text{where } \delta_1(\varepsilon, d) = R \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}} - \frac{R}{2} \frac{e^{-\varepsilon d/2}(1+\varepsilon)^{\frac{d}{2}}}{\sqrt{d}} \text{ and } \delta_2(\varepsilon, \delta) = R \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}} + R(1 + \varepsilon)e^{-d\varepsilon^2/8}.$$

Theorem 2 is, in a sense, an invariant Johnson-Lindenstrauss [14] type result where we show that the dot product defined by the random feature map Φ , i.e. $\langle \Phi(x), \Phi(z) \rangle$, is concentrated around the invariant expected kernel uniformly on a data set of N points, given a sufficiently large number of templates m , a large number of sampled group elements $|G|$, and a large bin number n . The error naturally decomposes to a numerical error ε_0 and statistical errors $\varepsilon_1, \varepsilon_2$ due to the sampling of the templates and the group elements respectively.

Theorem 2. [Johnson-Lindenstrauss type Theorem- N point Set] Let $\mathcal{D} = \{x_i \mid x_i \in \mathcal{X}\}_{i=1}^N$ be a finite dataset. Fix $\varepsilon_0, \varepsilon_1, \varepsilon_2, \delta_1, \delta_2 \in (0, 1)$. For a number of bins $n \geq \frac{1}{\varepsilon_0}$, templates $m \geq \frac{C_1}{\varepsilon_1} \log(\frac{N}{\delta_1})$, and group elements $|G| \geq \frac{C_2}{\varepsilon_2} \log(\frac{Nm}{\delta_2})$, where C_1, C_2 are universal numeric constants, we have:

$$|\langle \Phi(x_i), \Phi(x_j) \rangle - K_s(x_i, x_j)| \leq \varepsilon_0 + \varepsilon_1 + \varepsilon_2, i = 1 \dots N, j = 1 \dots N, \quad (5)$$

with probability $1 - \delta_1 - \delta_2$.

Putting together Theorems 1 and 2, the following Corollary shows how the group-invariant random feature map Φ captures the invariant distance between points uniformly on a dataset of N points.

Corollary 1 (Invariant Features Maps and Distances between Orbits). Let $\mathcal{D} = \{x_i \mid x_i \in \mathcal{X}\}_{i=1}^N$ be a finite dataset. Fix $\varepsilon_0, \delta \in (0, 1)$. For a number of bins $n \geq \frac{3}{\varepsilon_0}$, templates $m \geq \frac{9C_1}{\varepsilon_0^2} \log(\frac{N}{\delta})$, and group elements $|G| \geq \frac{9C_2}{\varepsilon_0^2} \log(\frac{Nm}{\delta})$, where C_1, C_2 are universal numeric constants, we have:

$$\varepsilon - \delta_2(d, \varepsilon) - \varepsilon_0 \leq \langle \Phi(x_i), \Phi(x_j) \rangle - (1 - d_G(x_i, x_j)) \leq \varepsilon_0 + \varepsilon + \delta_1(d, \varepsilon), \quad (6)$$

$i = 1 \dots N, j = 1 \dots N$, with probability $1 - 2\delta$.

Remark 4. Assuming that the templates are unitary and drawn from a general distribution $p(t)$, the equivalent kernel has the following form:

$$K_s(x, z) = \int_{\mathcal{G}} \int_{\mathcal{G}} d\mu(g)d\mu(g') \left(\int s - \max(\langle x, gt \rangle, \langle z, g't \rangle) p(t) dt \right).$$

Indeed when we use the gaussian sampling with rejection for the templates, the integral $\int \max(\langle x, gt \rangle, \langle z, g't \rangle) p(t) dt$ is asymptotically proportional to $\|g^{-1}x - g'^{-1}z\|_2$. It is interesting to consider different distributions that are domain-specific for the templates and assess the number of the templates needed to approximate such kernels. It is also interesting to find the optimal templates that achieve the minimum distortion in equation 6, in a data dependent way, but we will address these points in future work.

4 Learning with Group Invariant Random Features

In this section, we show that learning a linear model in the invariant, random feature space, on a training set sampled from the reduced core set \mathcal{X}_0 , has a low expected risk, and generalizes to unseen test points generated from the distribution on $\mathcal{X} = \mathcal{X}_0 \cup \mathcal{G}\mathcal{X}_0$. The architecture of the proof follows ideas from [15] and [16]. Recall that given an L -Lipschitz convex loss function V , our aim is to minimize the expected risk given in Equation (2). Denote the CDF by $\psi(x, t, \tau) = \mathbb{P}(\langle gt, x \rangle \leq \tau)$, and the empirical CDF by $\hat{\psi}(x, t, \tau) = \frac{1}{|G|} \sum_{i=1}^{|G|} \mathbb{I}_{\langle g_i t, x \rangle \leq \tau}$. Let $p(t)$ be the distribution of templates t . The RKHS defined by the invariant kernel K_s , $K_s(x, z) = \int_{-s}^s \psi(x, t, \tau) \psi(z, t, \tau) p(t) dt d\tau$ denoted \mathcal{H}_{K_s} , is the completion of the set of all finite linear combinations of the form:

$$f(x) = \sum_i \alpha_i K_s(x, x_i), x_i \in \mathcal{X}, \alpha_i \in \mathbb{R}. \quad (7)$$

Similarly to [16], we define the following infinite-dimensional function space:

$$\mathcal{F}_p = \left\{ f(x) = \int \int_{-s}^s w(t, \tau) \psi(x, t, \tau) dt d\tau \mid \sup_{\tau, t} \frac{|w(t, \tau)|}{p(t)} \leq C \right\}.$$

Lemma 1. \mathcal{F}_p is dense in \mathcal{H}_{K_s} . For $f \in \mathcal{F}_p$ we have $\mathcal{E}_V(f) = \int_{\mathcal{X}_0} \sum_{y \in \mathcal{Y}} V(yf(x)) \rho_y(x) d\rho_{\mathcal{X}}(x)$, where \mathcal{X}_0 is the reduced core set.

Since \mathcal{F}_p is dense in \mathcal{H}_{K_s} , we can learn an invariant decision function in the space \mathcal{F}_p , instead of learning in \mathcal{H}_{K_s} . Let $\Psi(x) = \left[\hat{\psi}\left(x, t_j, \frac{sk}{n}\right) \right]_{j=1 \dots m, k=-n \dots n}$. Ψ , and Φ are equivalent up to constants. We will approximate the set \mathcal{F}_p as follows:

$$\tilde{\mathcal{F}} = \left\{ f(x) = \langle w, \Psi(x) \rangle = \frac{s}{n} \sum_{j=1}^m \sum_{k=-n}^n w_{j,k} \hat{\psi}\left(x, t_j, \frac{sk}{n}\right), t_j \sim p, j = 1 \dots m \mid \|w\|_{\infty} \leq \frac{C}{m} \right\}.$$

Hence, we learn the invariant decision function via empirical risk minimization where we restrict the function to belong to $\tilde{\mathcal{F}}$, and the sampling in the training set is restricted to the core set \mathcal{X}_0 . Note that with this function space we are regularizing for convenience the norm infinity of the weights but this can be relaxed in practice to a classical Tikhonov regularization.

Theorem 3 (Learning with Group invariant features). *Let $S = \{(x_i, y_i) \mid x_i \in \mathcal{X}_0, y_i \in \mathcal{Y}, i = 1 \dots N\}$, a training set sampled from the core set \mathcal{X}_0 . Let $f_N^* = \arg \min_{f \in \tilde{\mathcal{F}}} \hat{\mathcal{E}}_V(f) = \frac{1}{N} \sum_{i=1}^N V(y_i f(x_i))$. Fix $\delta > 0$, then*

$$\begin{aligned} \mathcal{E}_V(f_N^*) &\leq \min_{f \in \mathcal{F}_p} \mathcal{E}_V(f) + 2 \frac{1}{\sqrt{N}} \left(4LsC + 2V(0) + LC \sqrt{\frac{1}{2} \log\left(\frac{1}{\delta}\right)} \right) \\ &\quad + \frac{2sLC}{\sqrt{m}} \left(1 + \sqrt{2 \log\left(\frac{1}{\delta}\right)} \right) + L \left(\frac{2sC}{\sqrt{|G|}} \left(1 + \sqrt{2 \log\left(\frac{m}{\delta}\right)} \right) + \frac{2sC}{n} \right), \end{aligned}$$

with probability at least $1 - 3\delta$ on the training set and the choice of templates and group elements.

The proof of Theorem 3 is given in Appendix B. Theorem 3 shows that learning a linear model in the invariant random feature space defined by Φ (or equivalently Ψ), has a low expected risk. More importantly, this risk is arbitrarily close to the optimal risk achieved in an infinite-dimensional class of functions, namely \mathcal{F}_p . The training set is sampled from the reduced core set \mathcal{X}_0 , and invariant learning generalizes to unseen test points generated from the distribution on $\mathcal{X} = \mathcal{X}_0 \cup \mathcal{G}\mathcal{X}_0$, hence the reduction in the sample complexity. Recall that \mathcal{F}_p is dense in the RKHS of the Haar-integration invariant Kernel, and so the expected risk achieved by a linear model in the invariant random feature space is not far from the one attainable in the invariant RKHS. Note that the error decomposes into two terms. The first, $O(\frac{1}{\sqrt{N}})$, is statistical and it depends on the training sample complexity N . The other is governed by the approximation error of functions \mathcal{F}_p , with functions in $\tilde{\mathcal{F}}$, and depends on the number of templates m , number of group elements sampled $|G|$, the number of bins n , and has the following form $O(\frac{1}{\sqrt{m}}) + O\left(\sqrt{\frac{\log m}{|G|}}\right) + \frac{1}{n}$.

5 Relation to Previous Work

We now put our contributions in perspective by outlining some of the previous work on invariant kernels and approximating kernels with random features.

Approximating Kernels. Several schemes have been proposed for approximating a non-linear kernel with an explicit non-linear feature map in conjunction with linear methods, such as the Nyström method [17] or random sampling techniques in the Fourier domain for translation-invariant kernels [15]. Our features fall under the random sampling techniques where, unlike previous work, we sample both projections and group elements to induce invariance with an integral representation. We note that the relation between random features and quadrature rules has been thoroughly studied in [18], where sharper bounds and error rates are derived, and can apply to our setting.

Invariant Kernels. We focused in this paper on Haar-integration kernels [11], since they have an integral representation and hence can be represented with random features [18]. Other invariant kernels have been proposed: In [19] authors introduce transformation invariant kernels, but unlike our general setting, the analysis is concerned with dilation invariance. In [20], multilayer arccosine kernels are built by composing kernels that have an integral representation, but does not explicitly induce invariance. More closely related to our work is [21], where kernel descriptors are built for visual recognition by introducing a kernel view of histogram of gradients that corresponds in our case to the cumulative distribution on the group variable. Explicit feature maps are obtained via kernel PCA, while our features are obtained via random sampling. Finally the convolutional kernel network of [22] builds a sequence of multilayer kernels that have an integral representation, by convolution, considering spatial neighborhoods in an image. Our future work will consider the composition of Haar-integration kernels, where the convolution is applied not only to the spatial variable but to the group variable akin to [2].

6 Numerical Evaluation

In this paper, and specifically in Theorems 2 and 3, we showed that the random, group-invariant feature map Φ captures the invariant distance between points, and that learning a linear model trained in the invariant, random feature space will generalize well to unseen test points. In this section, we validate these claims through three experiments. For the claims of Theorem 2, we will use a nearest neighbor classifier, while for Theorem 3, we will rely on the regularized least squares (RLS) classifier, one of the simplest algorithms for supervised learning. While our proofs focus on norm-infinity regularization, RLS corresponds to Tikhonov regularization with square loss. Specifically, for performing T -way classification on a batch of N training points in \mathbb{R}^d , summarized in the data matrix $X \in \mathbb{R}^{N \times d}$ and label matrix $Y \in \mathbb{R}^{N \times T}$, RLS will perform the optimization, $\min_{W \in \mathbb{R}^{m \times T}} \left\{ \frac{1}{N} \|Y - \Phi(X)W\|_F^2 + \lambda \|W\|_F^2 \right\}$, where $\|\cdot\|_F$ is the Frobenius norm, λ is the regularization parameter, and Φ is the feature map, which for the representation described in this paper will be a CDF pooling of the data projected onto group-transformed random templates. All RLS experiments in this paper were completed with the GURLS toolbox [23]. The three datasets we explore are:

X_{perm} (Figure 1): An artificial dataset consisting of all sequences of length 5 whose elements come from an alphabet of 8 characters. We want to learn a function which assigns a positive value to any sequence that contains a target set of characters (in our case, two of them) regardless of their position. Thus, the function label is globally invariant to permutation, and so we project our data onto all permuted versions of our random template sequences.

MNIST (Figure 2): We seek local invariance to translation and rotation, and so all random templates are translated by up to 3 pixels in all directions and rotated between -20 and 20 degrees.

TIDIGITS (Figure 3): We use a subset of TIDIGITS consisting of 326 speakers (men, women, children) reading the digits 0-9 in isolation, and so each datapoint is a waveform of a single word. We seek local invariance to pitch and speaking rate [25], and so all random templates are pitch shifted up and down by 400 cents and warped to play at half and double speed. The task is 10-way classification with one class-per-digit. See [24] for more detail.

Acknowledgements: Stephen Voinea acknowledges the support of a Nuance Foundation Grant. This work was also supported in part by the Center for Brains, Minds and Machines (CBMM), funded by NSF STC award CCF 1231216.

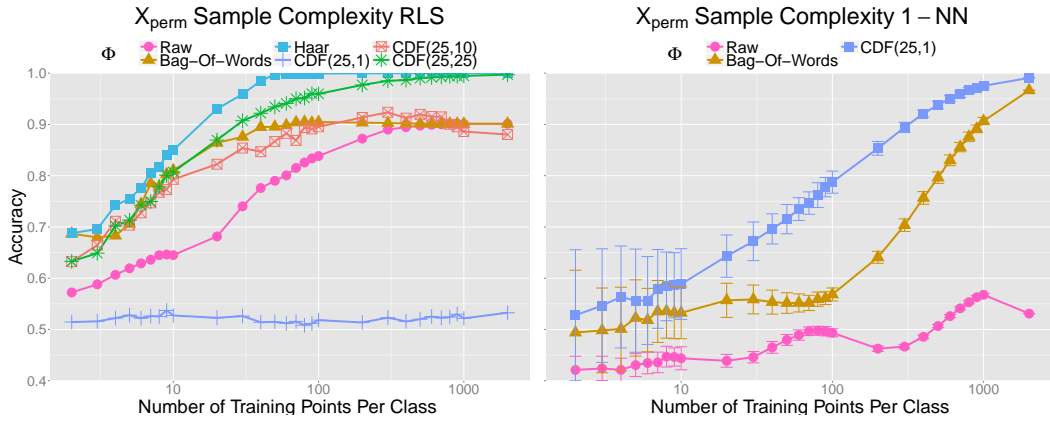


Figure 1: Classification accuracy as a function of training set size, averaged over 100 random training samples at each size. $\Phi = CDF(n, m)$ refers to a random feature map with n bins and m templates. With 25 templates, the random feature map outperforms the raw features and a bag-of-words representation (also invariant to permutation) and even approaches an RLS classifier with a Haar-integration kernel. Error bars were removed from the RLS plot for clarity. See supplement.

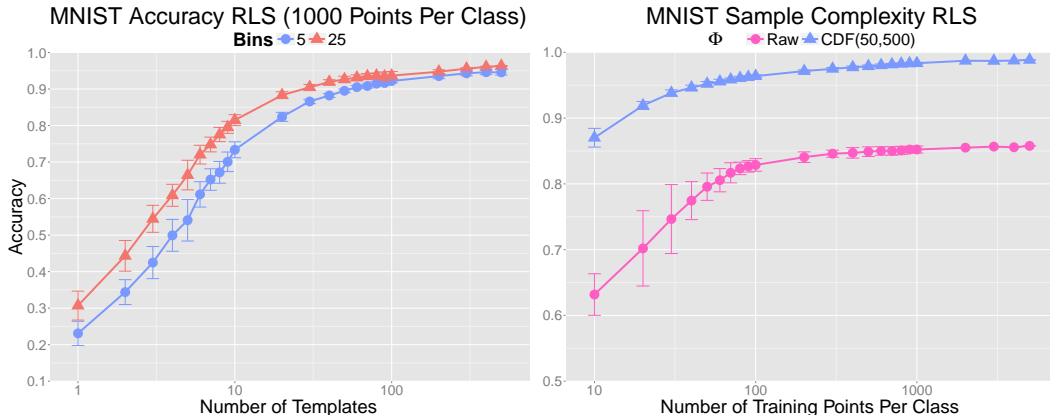


Figure 2: Left Plot) Mean classification accuracy as a function of number of bins and templates, averaged over 30 random sets of templates. Right Plot) Classification accuracy as a function of training set size, averaged over 100 random samples of the training set at each size. At 1000 examples per class, we achieve an accuracy of 98.97%.

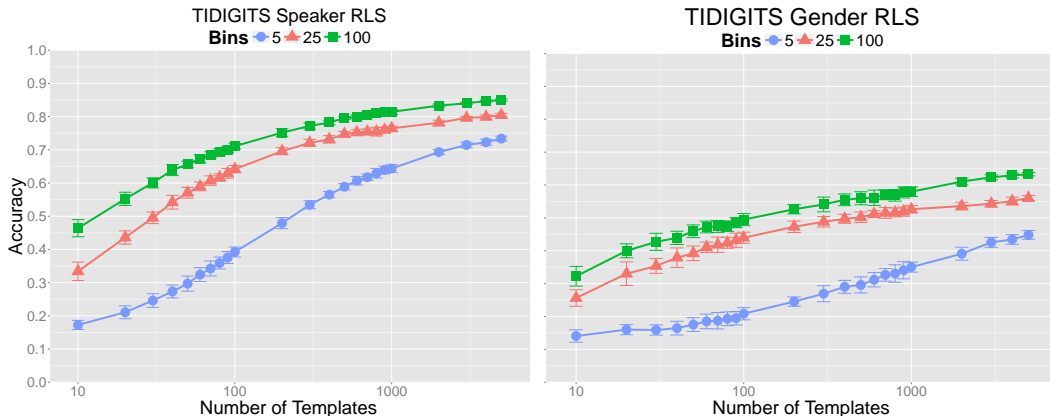


Figure 3: Mean classification accuracy as a function of number of bins and templates, averaged over 30 random sets of templates. In the “Speaker” dataset, we test on unseen speakers, and in the “Gender” dataset, we test on a new gender, giving us an extreme train/test mismatch. [25].

A Proofs of Theorems 1 and 2

Proof of Theorem 1. 1)

$$\begin{aligned}
K_s(x, z) &= \mathbb{E}_t \int_{-s}^s \mathbb{E}_g [\mathbb{1}_{\langle x, gt \rangle \leq \tau}] \mathbb{E}_{g'} [\mathbb{1}_{\langle z, g't \rangle \leq \tau}] d\tau \\
&= \mathbb{E}_t \int d\mu(g) d\mu(g') \int_{-s}^s \mathbb{1}_{\langle x, gt \rangle \leq \tau} \mathbb{1}_{\langle z, g't \rangle \leq \tau} d\tau \\
&= \int d\mu(g) d\mu(g') \mathbb{E}_t (s - \max(\langle x, gt \rangle, \langle z, g't \rangle)).
\end{aligned}$$

where the second equality is by Fubini theorem and the last one holds since for $a, b \in [-s, s]$:

$$\int_{-s}^s \mathbb{1}_{a \leq \tau} \mathbb{1}_{b \leq \tau} d\tau = s - \max(a, b).$$

Recall that the sampling of t is the following for $\varepsilon \in (0, 1)$ let :

$$t = n \sim \mathcal{N}\left(0, \frac{1}{d} I_d\right), \text{ if } \|n\|_2^2 < 1 + \varepsilon, t = \perp \text{ else ,}$$

since our group is unitary, x being norm one, and by virtue of this sampling the dot product $|\langle x, gt \rangle| \leq \|n\|_2 \leq \sqrt{1 + \varepsilon} \leq 1 + \varepsilon$. Hence $\langle x, gt \rangle \in [-(1 + \varepsilon), 1 + \varepsilon]$, and we can choose $s = 1 + \varepsilon$. Using again the fact the group is unitary and compact we have:

$$K_s(x, z) = \int d\mu(g) d\mu(g') \mathbb{E}_t (s - \max(\langle g^{-1}x, t \rangle, \langle g'^{-1}z, t \rangle)).$$

Now using this particular sampling of templates we have:

$$K_s(x, z) = \int_G \int_G d\mu(g) d\mu(g') \mathbb{E}_n \left(\mathbb{1}_{\|n\|_2^2 < 1 + \varepsilon} [1 + \varepsilon - \max(\langle g^{-1}x, n \rangle, \langle g'^{-1}z, n \rangle)] \right).$$

Let

$$Z_{x,z}(n, g, g') = \max(\langle g^{-1}x, n \rangle, \langle g'^{-1}z, n \rangle),$$

It follows that:

$$\begin{aligned}
K_s(x, z) &= \int_G \int_G d\mu(g) d\mu(g') \mathbb{E}_n \left(\mathbb{1}_{\|n\|_2^2 < 1 + \varepsilon} [1 + \varepsilon - Z_{x,z}(n, g, g')] \right) \\
&= (1 + \varepsilon) \mathbb{P}(\|n\|_2^2 < 1 + \varepsilon) - \int_G \int_G d\mu(g) d\mu(g') \mathbb{E}_n \left(\mathbb{1}_{\|n\|_2^2 < 1 + \varepsilon} Z_{x,z}(n, g, g') \right) \\
&= (1 + \varepsilon) \mathbb{P}(\|n\|_2^2 < 1 + \varepsilon) - \int_G \int_G d\mu(g) d\mu(g') \mathbb{E}_n \left((1 - \mathbb{1}_{\|n\|_2^2 \geq 1 + \varepsilon}) Z_{x,z}(n, g, g') \right) \\
&= (1 + \varepsilon) \mathbb{P}(\|n\|_2^2 < 1 + \varepsilon) - \int_G \int_G d\mu(g) d\mu(g') \mathbb{E}_n Z_{x,z}(n, g, g') \\
&\quad + \int_G \int_G d\mu(g) d\mu(g') \mathbb{E}_n \left(\mathbb{1}_{\|n\|_2^2 \geq 1 + \varepsilon} Z_{x,z}(n, g, g') \right) \tag{8}
\end{aligned}$$

We are left with evaluating or bounding two expectations: $I_1 = \mathbb{E}_n Z_{x,z}(n, g, g')$, and $I_2 = \mathbb{E}_n \left(\mathbb{1}_{\|n\|_2^2 \geq 1 + \varepsilon} Z_{x,z}(n, g, g') \right)$, that involve the maximum of correlated gaussian variables as we will see in the following.

By rotation invariance of Gaussians we have that $\langle g^{-1}x, n \rangle$, and $\langle g'^{-1}z, n \rangle$ are two correlated random gaussian variables with correlation coefficient that we note by $\cos(\theta_{g,g'}) = \langle g^{-1}x, g'^{-1}z \rangle$. Hence by a change of a basis we can write:

$$\langle g^{-1}x, n \rangle = \frac{1}{\sqrt{d}} u, \quad \langle g'^{-1}z, n \rangle = \frac{1}{\sqrt{d}} \cos(\theta_{g,g'}) u + \frac{1}{\sqrt{d}} \sqrt{1 - \cos^2(\theta_{g,g'})} v$$

where $\cos(\theta_{g,g'}) = \langle g^{-1}x, g'^{-1}z \rangle$, and $u, v \sim \mathcal{N}(0, 1)$ iids.

Hence,

$$I_1 = \frac{1}{\sqrt{d}} \mathbb{E}_{u,v} \max \left(u, \cos(\theta_{g,g'})u + \sqrt{1 - \cos^2(\theta_{g,g'})}v \right).$$

The following Lemma from [26] gives the expectation and the variance of the maximum of two gaussians with correlation coefficient ρ .

Lemma 2 (Mean and Variance of Maximum of Correlated Gaussians [26]). *Let $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$, two correlated gaussians with correlation coefficient ρ . Define $\phi_{\mathcal{N}}(x) = \frac{1}{\sqrt{2\pi}} \exp(-x^2/2)$, and $\Phi_{\mathcal{N}}(y) = \int_{-\infty}^y \phi_{\mathcal{N}}(x)dx$. Let $a = \sqrt{\sigma_X^2 + \sigma_Y^2 - 2\rho\sigma_X\sigma_Y}$, and $\alpha = \frac{\mu_X - \mu_Y}{a}$.*

The mean μ_Z and variance σ_Z^2 of $Z = \max(X, Y)$ are expressed analytically as follows:

$$\mu_Z = \mu_X \Phi_{\mathcal{N}}(\alpha) + \mu_Y \Phi_{\mathcal{N}}(-\alpha) + a \phi_{\mathcal{N}}(\alpha). \quad (9)$$

$$\sigma_Z^2 = \underbrace{(\sigma_X^2 + \mu_X^2) \Phi_{\mathcal{N}}(\alpha) + (\sigma_Y^2 + \mu_Y^2) \Phi_{\mathcal{N}}(-\alpha) + (\mu_X + \mu_Y) a \phi_{\mathcal{N}}(\alpha)}_{\mathbb{E}Z^2} - \mu_Z^2. \quad (10)$$

Applying Lemma 2 to our case ($\mu_X = \mu_Y = 0, \sigma_X = \sigma_Y = 1, \rho = \cos(\theta_{g,g'})$). We have: $a = \sqrt{2(1 - \cos(\theta_{g,g'}))}$ and $\alpha = 0$.

$$\begin{aligned} I_1 &= \frac{1}{\sqrt{d}} a \phi_{\mathcal{N}}(0) \\ &= \frac{1}{\sqrt{2\pi d}} \sqrt{2(1 - \cos(\theta_{g,g'}))} \\ &= \frac{1}{\sqrt{2\pi d}} \|g^{-1}x - g'^{-1}z\|_2. \end{aligned} \quad (11)$$

We turn now to I_2 that we bound using Cauchy-Schwarz inequality:

$$\begin{aligned} |I_2| &= \left| \mathbb{E}_n \left(\mathbf{1}_{\|n\|_2^2 \geq 1+\varepsilon} Z_{x,z}(n, g, g') \right) \right| \\ &\leq \sqrt{E(\mathbf{1}_{\|n\|_2^2 \geq 1+\varepsilon})} \sqrt{E(Z_{x,z}^2(n, g, g'))} \\ &= \sqrt{\mathbb{P}(\|n\|_2^2 \geq 1 + \varepsilon)} \sqrt{E(Z_{x,z}^2(n, g, g'))}. \end{aligned} \quad (12)$$

On the first hand, applying again Lemma 2 (for $\mathbb{E}Z^2$) we have:

$$\begin{aligned} E(Z_{x,z}^2(n, g, g')) &= \frac{1}{d} \mathbb{E}_{u,v} \left(\max \left(u, \cos(\theta_{g,g'})u + \sqrt{1 - \cos^2(\theta_{g,g'})}v \right) \right)^2 \\ &= \frac{1}{d} (2\Phi_{\mathcal{N}}(0)) \\ &= \frac{1}{d}. \end{aligned} \quad (13)$$

On the other hand, note that $\|n\|_2^2$ has a (normalized) chi squared distribution with d degree of freedom χ_d^2 , with mean 1. The following Lemma gives upper bounds for the upper and lower tails of a chi square distribution.

Lemma 3 (χ^2 tail bounds). *Let $X \sim \chi_k^2$, a chi squared random variable with k degree of freedom. The following hold true for any $\varepsilon \in (0, 1)$:*

- *Upper Bound for the upper tail [27]: $\mathbb{P}(\frac{1}{k}X \geq 1 + \varepsilon) \leq e^{-k\varepsilon^2/8}$.*
- *Upper Bound for the lower tail [28]: For all $k \geq 2, u \geq k - 1$ we have:*

$$\mathbb{P}(X < u) \leq 1 - \frac{1}{2} \exp \left(-\frac{1}{2} (u - k - (k - 2) \log(u/k) + \log(k)) \right).$$

More specifically for $u = k(1 + \varepsilon)$ we have:

$$\mathbb{P}\left(\frac{1}{k}X < 1 + \varepsilon\right) \leq 1 - \frac{1}{2} \frac{e^{-\varepsilon k/2} (1 + \varepsilon)^{\frac{k-2}{2}}}{\sqrt{k}}.$$

Applying Lemma 3, for $\|n\|_2^2$. We have $\|n\|_2^2 = \frac{1}{d}X$, where $X \sim \chi_d^2$, hence:

$$\mathbb{P}\left(\|n\|_2^2 \geq 1 + \varepsilon\right) \leq e^{-d\varepsilon^2/8}, \quad (14)$$

Putting together Equations (12),(14), (13) we have finally:

$$|I_2| \leq \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}}. \quad (15)$$

Putting together Equations (8), (11), and (15), and using upper and lower bounds for $\mathbb{P}(\|n\|_2^2 < 1 + \varepsilon)$ from Lemma 3:

$$\begin{aligned} K_s(x, z) &\leq (1 + \varepsilon) \mathbb{P}(\|n\|_2^2 < 1 + \varepsilon) - \frac{1}{\sqrt{2\pi d}} \int_G \int_G \|g^{-1}x - g'^{-1}z\|_2 d\mu(g)d\mu(g') + \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}} \\ &\leq (1 + \varepsilon) \left(1 - \frac{1}{2} \frac{e^{-\varepsilon d/2} (1 + \varepsilon)^{\frac{d-2}{2}}}{\sqrt{d}}\right) - \frac{1}{\sqrt{2\pi d}} \int_G \int_G \|g^{-1}x - g'^{-1}z\|_2 d\mu(g)d\mu(g') \\ &\quad + \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}}. \end{aligned}$$

$$\begin{aligned} K_s(x, z) &\geq (1 + \varepsilon) \mathbb{P}(\|n\|_2^2 < 1 + \varepsilon) - \frac{1}{\sqrt{2\pi d}} \int_G \int_G \|g^{-1}x - g'^{-1}z\|_2 d\mu(g)d\mu(g') - \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}} \\ &\geq (1 + \varepsilon) \left(1 - e^{-d\varepsilon^2/8}\right) - \frac{1}{\sqrt{2\pi d}} \int_G \int_G \|g^{-1}x - g'^{-1}z\|_2 d\mu(g)d\mu(g') - \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}}. \end{aligned}$$

Noting by d_G the integral and using that the group is compact and unitary:

$$\begin{aligned} d_G(x, z) &= \frac{1}{\sqrt{2\pi d}} \int_G \int_G \|g^{-1}x - g'^{-1}z\|_2 d\mu(g)d\mu(g') \\ &= \frac{1}{\sqrt{2\pi d}} \int_G \int_G \|gx - g'z\|_2 d\mu(g)d\mu(g'). \end{aligned}$$

We finally have:

$$-\frac{e^{-d\varepsilon^2/16}}{\sqrt{d}} - (1 + \varepsilon)e^{-d\varepsilon^2/8} + \varepsilon \leq K_s(x, z) - (1 - d_G(x, z)) \leq \frac{e^{-d\varepsilon^2/16}}{\sqrt{d}} - \frac{1}{2} \frac{e^{-\varepsilon d/2} (1 + \varepsilon)^{\frac{d}{2}}}{\sqrt{d}} + \varepsilon. \quad (16)$$

For any $\varepsilon \in (0, 1)$, as the dimension $d \rightarrow \infty$, we have asymptotically:

$$K_s(x, z) \rightarrow 1 - d_G(x, z) + \varepsilon = s - d_G(x, z).$$

2) The symmetry of K is obvious. Let $p(t)$ be the distribution of the templates t . Define the following weighted dot product: $\langle f(x, \cdot, \cdot), g(z, \cdot, \cdot) \rangle = \int_t p(t) \int_{-s}^s d\tau f(x, t, \tau)g(z, t, \tau)$. Recall that:

$$\begin{aligned} K_s(x, z) &= \int p(t)dt \int_{-s}^s \psi(x, t, \tau)\psi(z, t, \tau)d\tau \\ &= \langle \psi(x, \cdot, \cdot), \psi(z, \cdot, \cdot) \rangle. \end{aligned}$$

Hence K is symmetric and positive semidefinite. □

Proof of Theorem 2. In the following we fix two points x and z in \mathcal{X} and a random template t . Let $X_j = \int_{-s}^s \mathbb{P}(\langle gt_j, x \rangle \leq \tau) \mathbb{P}(\langle gt_j, z \rangle \leq \tau) d\tau$, we have $0 \leq X_j \leq 2s$, where $s = 1 + \varepsilon$. Recall that $K_s(x, z) = \frac{1}{m} \mathbb{E}_t(\sum_{j=1}^m X_j)$. By Hoeffding's inequality we have:

$$\mathbb{P}_t \left\{ \left| \frac{1}{m} \sum_{j=1}^m X_j - K_s(x, z) \right| > \epsilon \right\} \leq 2 \exp\left(\frac{-2m\epsilon^2}{(2s)^2}\right)$$

Turning now to the CDF $\psi(x, t, \tau) = \mathbb{P}(\langle gt, x \rangle \leq \tau)$, and the empirical CDF $\hat{\psi}(x, t, \tau) = \frac{1}{|G|} \sum_{i=1}^{|G|} \mathbb{1}_{\langle g_i t, x \rangle \leq \tau}$. By the theorem on convergence of the empirical CDF [29] (Theorem 4 given in Appendix D) we have, for $\gamma > 0$:

$$\mathbb{P}_g \left\{ \sup_{\tau} \left| \hat{\psi}(x, t, \tau) - \psi(x, t, \tau) \right| > \gamma \right\} \leq 2 \exp(-2|G|\gamma^2)$$

Hence we have $\forall \tau \in [-s, s]$:

$$\left| \hat{\psi}(x, t, \tau) - \psi(x, t, \tau) \right| \leq \gamma \text{ and } \left| \hat{\psi}(x, t, \tau) - \psi(z, t, \tau) \right| \leq \gamma$$

with a probability at least $1 - 4 \exp(-2|G|\gamma^2)$.

Define $X = \int_{-s}^s \psi(x, t, \tau) \psi(z, t, \tau) d\tau$, $\hat{X} = \int_{-s}^s \hat{\psi}(x, t, \tau) \hat{\psi}(z, t, \tau) d\tau$, and $\tilde{X} = \frac{(2s)}{n} \sum_{k=-n}^n \hat{\psi}(x, t, \frac{ks}{n}) \hat{\psi}(z, t, \frac{ks}{n})$, choose $0 < \gamma < 1$:

$$\begin{aligned} |\hat{X} - X| &= \left| \int_{-s}^s \left(\hat{\psi}(x, t, \tau) \hat{\psi}(z, t, \tau) - \psi(x, t, \tau) \psi(z, t, \tau) \right) d\tau \right| \\ &= \left| \int_{-s}^s \left(\hat{\psi}(x, t, \tau) - \psi(x, t, \tau) + \psi(x, t, \tau) \right) \left(\hat{\psi}(z, t, \tau) - \psi(z, t, \tau) + \psi(z, t, \tau) \right) - \psi(x, t, \tau) \psi(z, t, \tau) d\tau \right| \\ &\leq (2\gamma + \gamma^2) 2s \\ &\leq 6s\gamma, \end{aligned}$$

with probability $1 - 4 \exp(-2|G|\gamma^2)$. Define $X_j = \int_{-s}^s \psi(x, t_j, \tau) \psi(z, t_j, \tau) d\tau$, $\hat{X}_j = \int_{-s}^s \hat{\psi}(x, t_j, \tau) \hat{\psi}(z, t_j, \tau) d\tau$, and $\tilde{X}_j = \frac{(2s)}{n} \sum_{k=-n}^n \hat{\psi}(x, t_j, \frac{ks}{n}) \hat{\psi}(z, t_j, \frac{ks}{n})$, Then for all $j = 1 \dots m$, we have

$$|\hat{X}_j - X_j| \leq 6s\gamma$$

with probability $1 - 4m \exp(-2|G|\gamma^2) - 2 \exp\left(\frac{-2m\epsilon^2}{(2s)^2}\right)$.

Now we turn to the numerical approximation of the integra by a Riemann sum, we have for all $j = 1 \dots m$:

$$\left| \hat{X}_j - \tilde{X}_j \right| \leq \frac{s}{n}.$$

Hence the error decomposes in the following way:

$$\begin{aligned} |\langle \Phi(x), \Phi(z) \rangle - K_s(x, z)| &= \left| \frac{1}{m} \sum_{j=1}^m \tilde{X}_j - K_s(x, z) \right| \\ &= \left| \left(\frac{1}{m} \sum_{j=1}^m \tilde{X}_j - \frac{1}{m} \sum_{j=1}^m \hat{X}_j \right) + \left(\frac{1}{m} \sum_{j=1}^m \hat{X}_j - \frac{1}{m} \sum_{j=1}^m X_j \right) + \left(\frac{1}{m} \sum_{j=1}^m X_j - K_s(x, z) \right) \right| \\ &\leq \underbrace{\left| \frac{1}{m} \sum_{j=1}^m \tilde{X}_j - \frac{1}{m} \sum_{j=1}^m \hat{X}_j \right|}_{\text{Numerical Binning Error}} + \underbrace{\left| \frac{1}{m} \sum_{j=1}^m \hat{X}_j - \frac{1}{m} \sum_{j=1}^m X_j \right|}_{\text{Group CDF Approximation Error}} + \underbrace{\left| \frac{1}{m} \sum_{j=1}^m X_j - K_s(x, z) \right|}_{\text{Templates Concentration Error}} \\ &\leq \frac{s}{n} + 6s\gamma + \epsilon. \end{aligned}$$

with probability $1 - 4m \exp(-2|G|\gamma^2) - 2 \exp\left(\frac{-2m\epsilon^2}{(2s)^2}\right)$. For this to hold on all pairs of points in a set of cardinality N we have:

$$|\langle \Phi(x_i), \Phi(x_j) \rangle - K(x_i, x_j)| \leq \frac{s}{n} + 6s\gamma + \epsilon, i = 1 \dots N, j = 1 \dots N,$$

with probability $1 - 4mN(N-1) \exp(-2|G|\gamma^2) - 2N(N-1) \exp\left(\frac{-m\epsilon^2}{2(s)^2}\right)$.

Hence we have for numerical constants C_1 , and C_2 , $0 < \delta_1, \delta_2 < 1$, and $0 < \epsilon_0, \epsilon_1, \epsilon_2 < 1$, for $n \geq \frac{s}{\epsilon_0}$, $m \geq \frac{C_1}{\epsilon_1^2} \log\left(\frac{N}{\delta_1}\right)$, $|G| \geq \frac{C_2}{\epsilon_2^2} \log\left(\frac{Nm}{\delta_2}\right)$, :

$$|\langle \Phi(x_i), \Phi(x_j) \rangle - K_s(x_i, x_j)| \leq \epsilon_0 + \epsilon_1 + \epsilon_2, i = 1 \dots N, j = 1 \dots N,$$

with probability $1 - \delta_1 - \delta_2$.

□

B Proof of Theorem 3

Proof of Lemma 1. Our proof parallels similar proofs in [16]. Note that functions of the form (7) are dense in \mathcal{H}_K . $f(x) = \sum_i \alpha_i K_s(x, x_i) = \sum_i \alpha_i \int \int_{-s}^s \psi(x, t, \tau) \psi(x_i, t, \tau) p(t) dt d\tau$
 $= \int \int_{-s}^s (p(t) \sum_i \alpha_i \psi(x_i, t, \tau)) \psi(x, t, \tau) dt d\tau$. Let $\beta(t, \tau) = p(t) \sum_i \alpha_i \psi(x_i, t, \tau)$, since $0 \leq \psi(x, t, \tau) \leq 1, \forall x, t, \tau$, we have $\frac{|\beta(t, \tau)|}{p(t)} \leq \sum_i |\alpha_i| < \infty$, since α_i are finite. Hence f can be written in the form:

$$f(x) = \int \int_{-s}^s \beta(t, \tau) \psi(x, t, \tau) dt d\tau, \sup_{\tau, t} \frac{|\beta(t, \tau)|}{p(t)} < \infty,$$

and $f \in \mathcal{F}_p$.

□

In order to prove Theorem 3, we need some preliminary lemmas. The following Lemma assess the approximation of any function $f \in \mathcal{F}_p$, by a certain $\tilde{f} \in \tilde{\mathcal{F}}$.

Lemma 4 ($\tilde{\mathcal{F}}$ Approximation of \mathcal{F}_p). *Let f be a function in \mathcal{F}_p . Then for $\delta_1, \delta_2 > 0$, there exists a function $\tilde{f} \in \tilde{\mathcal{F}}$ such that:*

$$\|\tilde{f} - f\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} \leq \frac{2sC}{\sqrt{m}} \left(1 + \sqrt{2 \log\left(\frac{1}{\delta_1}\right)}\right) + \frac{2sC}{\sqrt{|G|}} \left(1 + \sqrt{2 \log\left(\frac{m}{\delta_2}\right)}\right) + \frac{2sC}{n},$$

with probability at least $1 - \delta_1 - \delta_2$.

Proof of Lemma 4. Let $f \in \mathcal{F}_p, f(x) = \int \int_{-s}^s w(t, \tau) \psi(x, t, \tau) d\tau dt$.

Let $f_j(x) = \int_{-s}^s \frac{w(t_j, \tau)}{p(t_j)} \psi(x, t_j, \tau) d\tau, \hat{f}_j(x) = \int_{-s}^s \frac{w(t_j, \tau)}{p(t_j)} \hat{\psi}(x, t_j, \tau) d\tau$, and $\tilde{f}_j(x) = \frac{s}{n} \sum_{k=-n}^n \frac{w(t_j, \frac{ks}{n})}{p(t_j)} \hat{\psi}(x, t_j, \frac{ks}{n})$. We have the following: $\mathbb{E}_t(f_j) = f$, and $\frac{1}{m} \mathbb{E}_t(\sum_{j=1}^m f_j) = f$. Consider the Hilbert space $\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})$, with dot product: $\langle f, g \rangle_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} = \int_{\mathcal{X}} f(x) g(x) d\rho_{\mathcal{X}}(x)$.

Note that : $\int_{-s}^s g(\tau) d\tau \leq \sqrt{2s} \sqrt{\int_{-s}^s g^2(\tau) d\tau}$

$$\|f_j\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} = \sqrt{\int_{\mathcal{X}} \left(\int_{-s}^s \frac{w(t_j, \tau)}{p(t_j)} \psi(x, t_j, \tau) d\tau \right)^2 d\rho_{\mathcal{X}}(x)} \leq (2sC),$$

Fix $\delta_1 > 0$, applying Lemma 7 we have therefore with probability $1 - \delta_1$:

$$\left\| \frac{1}{m} \sum_{j=1}^m f_j - f \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} \leq \frac{2sC}{\sqrt{m}} \left(1 + \sqrt{2 \log\left(\frac{1}{\delta_1}\right)}\right), \quad (17)$$

Now turn to:

$$\begin{aligned}
\left\| \frac{1}{m} \sum_{j=1}^m (\hat{f}_j - f_j) \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} &\leq \frac{1}{m} \sum_{j=1}^m \left\| \hat{f}_j - f_j \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})}, \\
\left\| \hat{f}_j - f_j \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})}^2 &= \int_{\mathcal{X}} \left(\int_{-s}^s \frac{w(t_j, \tau)}{p(t_j)} (\psi(x, t_j, \tau) - \hat{\psi}(x, t_j, \tau)) d\tau \right)^2 d\rho_{\mathcal{X}}(x) \\
&\leq 2s \int_{\mathcal{X}} \int_{-s}^s \frac{w^2(t_j, \tau)}{p^2(t_j)} (\psi(x, t_j, \tau) - \hat{\psi}(x, t_j, \tau))^2 d\tau d\rho_{\mathcal{X}}(x) \\
&\leq 2sC^2 \int_{\mathcal{X}} \int_{-s}^s (\hat{\psi}(x, t_j, \tau) - \psi(x, t_j, \tau))^2 d\tau d\rho_{\mathcal{X}}(x) \\
&= 2sC^2 \int_{-s}^s \int_{\mathcal{X}} (\hat{\psi}(x, t_j, \tau) - \psi(x, t_j, \tau))^2 d\rho_{\mathcal{X}}(x) d\tau \\
&= 2sC^2 \int_{-s}^s \left\| \hat{\psi}(\cdot, t_j, \tau) - \psi(\cdot, t_j, \tau) \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})}^2 d\tau \\
&\leq (2sC)^2 \sup_{\tau, j=1 \dots m} \left\| \hat{\psi}(\cdot, t_j, \tau) - \psi(\cdot, t_j, \tau) \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})}^2.
\end{aligned}$$

Recall that: $\hat{\psi}(x, t, \tau) = \frac{1}{|G|} \sum_{i=1}^{|G|} \mathbb{1}_{(g_i t, x) \leq \tau}$, and $\psi(x, t, \tau) = \mathbb{E}_g \hat{\psi}(x, t, \tau)$.

Clearly $\left\| \mathbb{1}_{(\cdot, gt) \leq \tau} \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} \leq 1$, hence applying again Lemma 7, for $\delta_2 > 0$ we have with probability $1 - \delta_2$:

$$\left\| \hat{\psi}(\cdot, t_j, \tau) - \psi(\cdot, t_j, \tau) \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})}^2 \leq \frac{1}{|G|} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_2} \right)} \right)^2,$$

It follows that: $\forall j = 1 \dots m, \left\| \hat{f}_j - f_j \right\| \leq \frac{2Cs}{\sqrt{|G|}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_2} \right)} \right)$, with probability $1 - m\delta_2$.

Hence with probability $1 - m\delta_2$, we have:

$$\left\| \frac{1}{m} \sum_{j=1}^m (\hat{f}_j - f_j) \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} \leq \frac{2Cs}{\sqrt{|G|}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_2} \right)} \right). \quad (18)$$

and by the approximation of a Riemann sum we have that:

$$\left\| \frac{1}{m} \sum_{j=1}^m (\hat{f}_j - \tilde{f}_j) \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} \leq \frac{2sC}{n}. \quad (19)$$

It is clear that $\tilde{f} = \frac{1}{m} \sum_{j=1}^m \tilde{f}_j \in \tilde{\mathcal{F}}$, hence, putting together equations (17), (18), and (19) we finally have:

$$\begin{aligned}
\left\| \frac{1}{m} \sum_{j=1}^m \tilde{f}_j - f \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} &\leq \left\| \frac{1}{m} \sum_{j=1}^m (\tilde{f}_j - \hat{f}_j) \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} + \left\| \frac{1}{m} \sum_{j=1}^m (\hat{f}_j - f_j) \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} + \left\| \frac{1}{m} \sum_{j=1}^m f_j - f \right\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})} \\
&\leq \frac{2sC}{n} + \frac{2Cs}{\sqrt{|G|}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_2} \right)} \right) + \frac{2sC}{\sqrt{m}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_1} \right)} \right)
\end{aligned}$$

with probability $1 - \delta_1 - m\delta_2$. \square

The following Lemma shows how the approximation of functions in \mathcal{F}_p , by functions in $\tilde{\mathcal{F}}$, translates to the expected Risk:

Lemma 5 (Bound on the Approximation Error). *Let $f \in \mathcal{F}_p$, fix $\delta_1, \delta_2 > 0$. There exists a function $\tilde{f} \in \tilde{\mathcal{F}}$, such that:*

$$\mathcal{E}_V(\tilde{f}) \leq \mathcal{E}_V(f) + \frac{2sLC}{\sqrt{m}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_1} \right)} \right) + L \left(\frac{2sC}{\sqrt{|G|}} \left(1 + \sqrt{2 \log \left(\frac{m}{\delta_2} \right)} \right) + \frac{2sC}{n} \right),$$

with probability at least $1 - \delta_1 - \delta_2$.

Proof of Lemma 5. $\mathcal{E}_V(\tilde{f}) - \mathcal{E}_V(f) \leq \int_{\mathcal{X}} |V(y\tilde{f}(x)) - V(yf(x))| d\rho_{\mathcal{X}}(x) \leq L \int_{\mathcal{X}} |\tilde{f}(x) - f(x)| d\rho_{\mathcal{X}}(x) \leq L \sqrt{\int_{\mathcal{X}} (\tilde{f}(x) - f(x))^2 d\rho_{\mathcal{X}}(x)} = L \|\tilde{f} - f\|_{\mathcal{L}^2(\mathcal{X}, \rho_{\mathcal{X}})}$, where we used the Lipschitz condition and Jensen inequality. The rest of the proof follows from Lemma 4. \square

The following Lemma gives a bound on the estimation of the expected Risk with finite training samples:

Lemma 6 (Bound on the Estimation Error). *Fix $\delta > 0$, then*

$$\sup_{f \in \tilde{\mathcal{F}}} |\mathcal{E}_V(f) - \hat{\mathcal{E}}_V(f)| \leq \frac{1}{\sqrt{N}} \left(4LsC + 2V(0) + LC \sqrt{\frac{1}{2} \log \left(\frac{1}{\delta} \right)} \right),$$

with probability $1 - \delta$.

Proof. The proof follows from Theorem 5 given in Appendix D. It is sufficient to bound the Rademacher complexity of the class $\tilde{\mathcal{F}}$:

$$\begin{aligned} \mathcal{R}_N(\tilde{\mathcal{F}}) &= \mathbb{E}_{x, \sigma} \left[\sup_{f \in \tilde{\mathcal{F}}} \left| \frac{1}{N} \sum_{i=1}^N \sigma_i f(x_i) \right| \right] = \mathbb{E}_{x, \sigma} \left[\sup_{f \in \tilde{\mathcal{F}}} \left| \frac{s}{Nn} \sum_{i=1}^N \sigma_i \left(\sum_{j=1}^m \sum_{k=-n}^n w_{j,k} \hat{\psi} \left(x_i, t_j, \frac{sk}{n} \right) \right) \right| \right] \\ &= \mathbb{E}_{x, \sigma} \left[\sup_{f \in \tilde{\mathcal{F}}} \left| \frac{s}{Nn} \sum_{j=1}^m \sum_{k=-n}^n w_{j,k} \sum_{i=1}^N \sigma_i \hat{\psi} \left(x_i, t_j, \frac{sk}{n} \right) \right| \right] \\ &\leq \mathbb{E}_{x, \sigma} \frac{sC}{mNn} \sum_{j=1}^m \sum_{k=-n}^n \left| \sum_{i=1}^N \sigma_i \hat{\psi} \left(x_i, t_j, \frac{sk}{n} \right) \right| \quad \text{By Holder inequality: } \langle a, b \rangle \leq \|a\|_{\infty} \|b\|_1 \\ &\leq \frac{sC}{mNn} \mathbb{E}_x \sum_{j=1}^m \sum_{k=-n}^n \sqrt{\mathbb{E}_{\sigma} \left(\sum_{i=1}^N \sigma_i \hat{\psi} \left(x_i, t_j, \frac{sk}{n} \right) \right)^2} \quad \text{Jensen inequality, concavity of square root} \end{aligned}$$

Note that $\mathbb{E}(\sigma_i \sigma_j) = 0$, for $i \neq j$ it follows that:

$$\begin{aligned} \mathbb{E}_{\sigma} \left(\sum_{i=1}^N \sigma_i \hat{\psi} \left(x_i, t_j, \frac{sk}{n} \right) \right)^2 &= \mathbb{E}_{\sigma} \sum_{i=1}^N \sum_{\ell=1}^N \sigma_i \sigma_{\ell} \hat{\psi} \left(x_i, t_j, \frac{sk}{n} \right) \hat{\psi} \left(x_{\ell}, t_j, \frac{sk}{n} \right) = \\ \sum_{i=1}^N \hat{\psi}^2 \left(x_i, t_j, \frac{sk}{n} \right) &\leq N, \text{ since } \hat{\psi}(\cdot, \cdot, \cdot) \leq 1. \text{ Finally:} \end{aligned}$$

$$\mathcal{R}_m(\tilde{\mathcal{F}}) \leq \frac{Cs}{\sqrt{N}}.$$

\square

We are now ready to prove Theorem 3:

Proof of Theorem 3. Let $f_N^* = \arg \min_{f \in \tilde{\mathcal{F}}} \hat{\mathcal{E}}_V(f)$, $\tilde{f} = \arg \min_{f \in \tilde{\mathcal{F}}} \mathcal{E}_V(f)$, $f_p = \arg \min_{f \in \mathcal{F}_p} \mathcal{E}_V(f)$.

$$\mathcal{E}_V(f_N^*) - \min_{f \in \mathcal{F}_p} \mathcal{E}_V(f) = \underbrace{\left(\mathcal{E}_V(f_N^*) - \mathcal{E}_V(\tilde{f}) \right)}_{\text{Statistical Error}} + \underbrace{\left(\mathcal{E}_V(\tilde{f}) - \mathcal{E}_V(f_p) \right)}_{\text{Approximation Error}}$$

The first term is the usual estimation or statistical error than we can bound using Lemma 6, we have:

$$\begin{aligned}\mathcal{E}_V(f_N^*) - \mathcal{E}_V(\tilde{f}) &= \left(\mathcal{E}_V(f_N^*) - \hat{\mathcal{E}}_V(f_N^*) \right) + \underbrace{\left(\hat{\mathcal{E}}_V(f_N^*) - \hat{\mathcal{E}}_V(\tilde{f}) \right)}_{\leq 0, \text{by optimality of } f_N^*} + \left(\hat{\mathcal{E}}_V(\tilde{f}) - \mathcal{E}_V(\tilde{f}) \right) \\ &\leq 2 \sup_{f \in \tilde{\mathcal{F}}} \left| \mathcal{E}_V(f) - \hat{\mathcal{E}}_V(f) \right| \\ &\leq 2 \frac{1}{\sqrt{N}} \left(4LsC + 2V(0) + LC \sqrt{\frac{1}{2} \log \left(\frac{1}{\delta} \right)} \right),\end{aligned}$$

with probability $1 - \delta$ over the training samples. Let \tilde{f}_p , the function defined in Lemma 4, that approximates f_p in $\tilde{\mathcal{F}}$. By Lemma 5 we know that:

$$\mathcal{E}_V(\tilde{f}_p) \leq \mathcal{E}_V(f_p) + \frac{2sLC}{\sqrt{m}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_1} \right)} \right) + L \left(\frac{2sC}{\sqrt{|G|}} \left(1 + \sqrt{2 \log \left(\frac{m}{\delta_2} \right)} \right) + \frac{2sC}{n} \right),$$

with probability $1 - \delta_1 - \delta_2$, on the choice of the templates and the sampled group elements. By optimality of $\tilde{f} \in \tilde{\mathcal{F}}$, we have

$$\mathcal{E}_V(\tilde{f}) \leq \mathcal{E}_V(\tilde{f}_p) \leq \mathcal{E}_V(f_p) + \frac{2sLC}{\sqrt{m}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_1} \right)} \right) + L \left(\frac{2sC}{\sqrt{|G|}} \left(1 + \sqrt{2 \log \left(\frac{m}{\delta_2} \right)} \right) + \frac{2sC}{n} \right)$$

Hence by a union bound with probability $1 - \delta - \delta_1 - \delta_2$, on the training set, the templates and the group elements we have:

$$\begin{aligned}\mathcal{E}_V(f_N^*) - \min_{f \in \tilde{\mathcal{F}}_p} \mathcal{E}_V(f) &\leq 2 \frac{1}{\sqrt{N}} \left(4LsC + 2V(0) + LC \sqrt{\frac{1}{2} \log \left(\frac{1}{\delta} \right)} \right) \\ &\quad + \frac{2sLC}{\sqrt{m}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta_1} \right)} \right) + L \left(\frac{2sC}{\sqrt{|G|}} \left(1 + \sqrt{2 \log \left(\frac{m}{\delta_2} \right)} \right) + \frac{2sC}{n} \right).\end{aligned}$$

□

C Technical tools

Theorem 4. [29] Let X_1, X_2, \dots, X_m be i.i.d. random variables with cumulative distribution function F , and let \hat{F}_m be the associated empirical cumulative density function $\hat{F}_m = \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{X_i \leq \tau}$. Then for any $\gamma > 0$

$$\mathbb{P} \left\{ \sup_{\tau} \left| \hat{F}_m(\tau) - F(\tau) \right| > \gamma \right\} \leq 2 \exp(-2m\gamma^2).$$

Lemma 7 ([15], Concentration of the mean of bounded random variables in a Hilbert Space). Let $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ be a Hilbert space. Let $X_j, j = 1 \dots K$, be iid random, such that $\|X_j\|_{\mathcal{H}} \leq M$. Then for any $\delta > 0$, with probability $1 - \delta$,

$$\left\| \frac{1}{K} \sum_{j=1}^K X_j - \frac{1}{K} \mathbb{E} \sum_{j=1}^K X_j \right\|_{\mathcal{H}} \leq \frac{M}{\sqrt{K}} \left(1 + \sqrt{2 \log \left(\frac{1}{\delta} \right)} \right).$$

Theorem 5 ([15]). Let \mathcal{F} be a bounded class of function, $\sup_{x \in \mathcal{X}} |f(x)| \leq C$ for all $f \in \mathcal{F}$. Let V be an L -Lipschitz loss. Then with probability $1 - \delta$, with respect to training samples $\{x_i, y_i\}_{i=1 \dots N}$, every f satisfies:

$$\mathcal{E}_V(f) \leq \hat{\mathcal{E}}_V(f) + 4L\mathcal{R}_N(\mathcal{F}) + \frac{2V(0)}{\sqrt{N}} + LC \sqrt{\frac{1}{2N} \log \frac{1}{\delta}},$$

where $\mathcal{R}_N(\mathcal{F})$ is the Rademacher complexity of the class \mathcal{F} :

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

the variables σ_i are iid symmetric Bernoulli random variables taking value in $\{-1, 1\}$, with equal probability and are independent from x_i .

D Numerical Evaluation

D.1 Permutation Invariance Experiment

For our first experiment, we created an artificial dataset which was designed to exploit permutation invariance, providing us with a finite group to which we had complete access. The dataset X_{perm} consists of all sequences of length $L = 5$, where each element of the sequence is taken from an alphabet A of 8 characters, giving us a total of 32,768 data points. Two characters $c_1, c_2 \in A$ were randomly chosen and designated as targets, so that a sequence $x \in X_{perm}$ is labeled positive if it contains both c_1 and c_2 , where the position of these characters in the sequence does not matter. Likewise, any sequence that does not contain both characters is labeled negative. This provides us with a binary classification problem (positive sequences vs. negative sequences), for which the label is preserved by permutations of the sequence indices, i.e. two sequences will belong to the same orbit if and only if they are permuted versions of one another.

The i^{th} character in A is encoded as an 8-dimensional vector which is 0 in every position but the i^{th} , where it is 1. Each sequence $x \in X_{perm}$ is formed by concatenating the 5 such vectors representing its characters, resulting in a binary vector of length 40. To build the permutation-invariant representation, we project a binary sequences onto an equal-length sequence consisting of standard-normal gaussian vectors, as well as all of its permutations, and then pool over the projections with a CDF.

As a baseline, we also used a bag-of-words representation, where each $x \in X_{perm}$ was encoded with an 8-dimensional vector with i^{th} element equal to the count of how many times character i appears in x . Note that this representation is also invariant to permutations, and so should share many of the benefits of our feature map.

For all classification results, 4000 points were randomly chosen from X_{perm} to form the training set, with an even split of 2000 positive points and 2000 negative points. The remaining 28,768 points formed the test set.

We know from Theorem 3 that the expected risk is dependent on the number of templates used to encode our data and on the number of bins used in the CDF-pooling step. The right panel of Figure 4 shows RLS classification accuracy on X_{perm} for different numbers of templates and bins. We see that, for a fixed number of templates, increasing the number of bins will improve accuracy, and for a fixed number of bins, adding more templates will improve accuracy. We also know there is a further dependence on the number of transformation samples from the group G . The left panel of Figure 4 shows how classification accuracy, for a fixed number of training points, bins, and templates, depends on the number of transformation we have access to. We see the curve is rather flat, and there is a very graceful degradation in performance.

In Figure 5, we include the sample complexity plot (for RLS) with the error bars added.

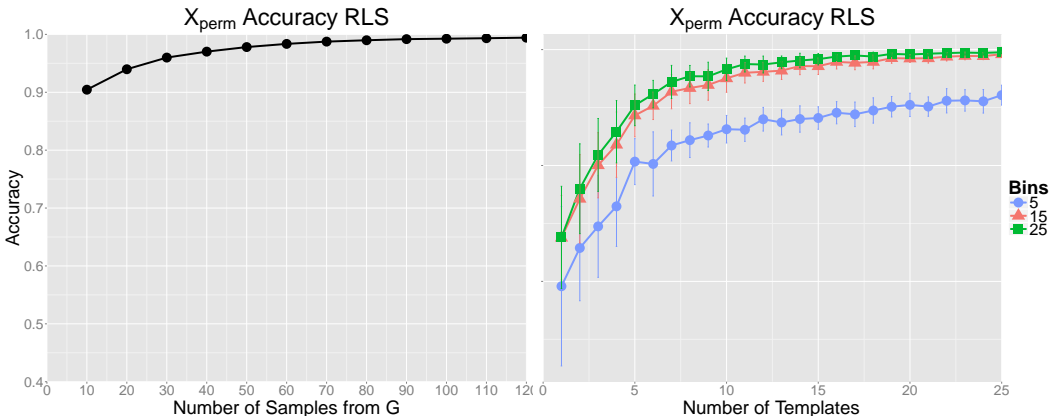


Figure 4: *Left) Classification accuracy of random invariant features as function of the number of sampled group elements on X_{perm} . Right) Classification accuracy of random invariant features as function of the number of templates and bin sizes on X_{perm} .*

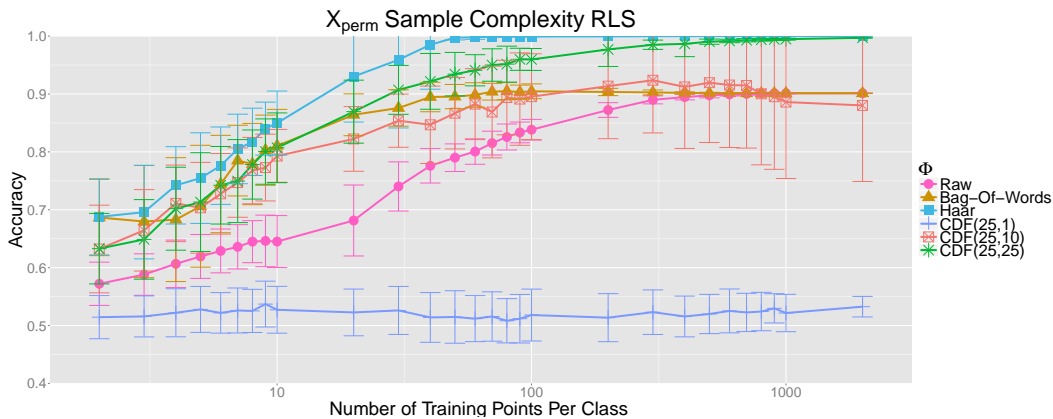


Figure 5: Classification accuracy as a function of training set size. $\Phi = CDF(n, m)$ refers to a random feature map with n bins and m templates. For each training set size, the accuracy is averaged over 100 random training samples. With enough templates/bins, the random feature map outperforms the raw features as well as a bag-of-words representation (also invariant to permutation). We also train an RLS classifier with a haar-invariant kernel, which naturally gives the best performance. However, by increasing the number of templates, we come close to matching this performance with random feature maps.

D.2 TIDIGITS Experiment

Here, we add plots (Figures 6,7 and 8) showing performance as a function of number of templates and bins for some other splits of the TIDIGITS data.

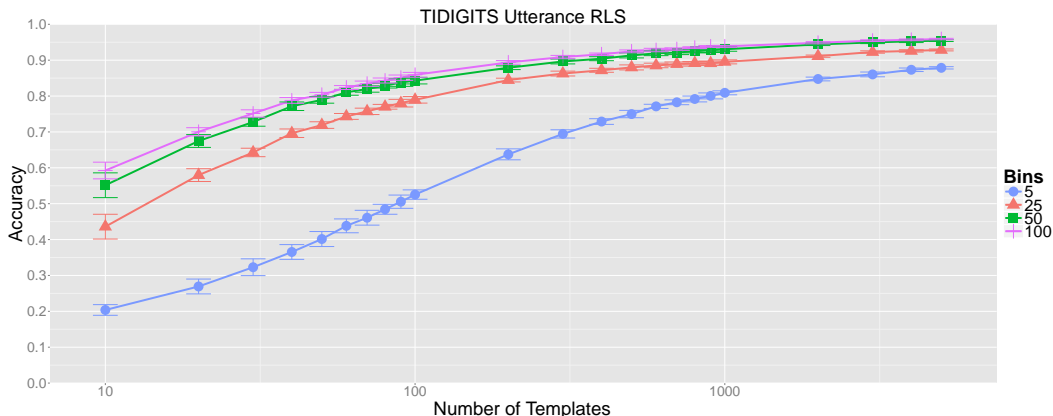


Figure 6: Mean classification accuracy as a function of number of templates, m , and bins, n . Accuracy is averaged over 30 random template samples for each m and error bars are displayed. In the “Utterance” dataset, we train and test on the same speakers, but the test set contains new utterances of each digit. This is the easiest dataset, representing only intraspeaker variability, and the performance is quite good even for a small number of bins.

References

- [1] F. Anselmi, J. Z. Leibo, L. Rosasco, J. Mutch, A. Tacchetti, and T. Poggio, “Unsupervised learning of invariant representations in hierarchical architectures,” *CoRR*, vol. abs/1311.4158, 2013.
- [2] J. Bruna and S. Mallat, “Invariant scattering convolution networks,” *CoRR*, vol. abs/1203.1513, 2012.
- [3] G. Hinton, A. Krizhevsky, and S. Wang, “Transforming auto encoders,” *ICANN-11*, 2011.

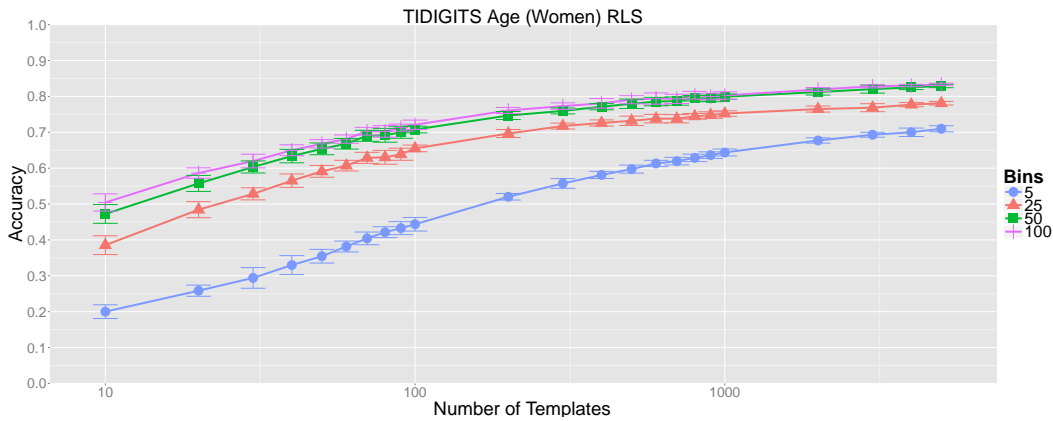


Figure 7: Mean classification accuracy as a function of number of templates, m , and bins, n . Accuracy is averaged over 30 random template samples for each m and error bars are displayed. In the “Age (Women)” dataset, we train on adult women and test on children, giving us an age mismatch. Despite this mismatch, performance remains strong.

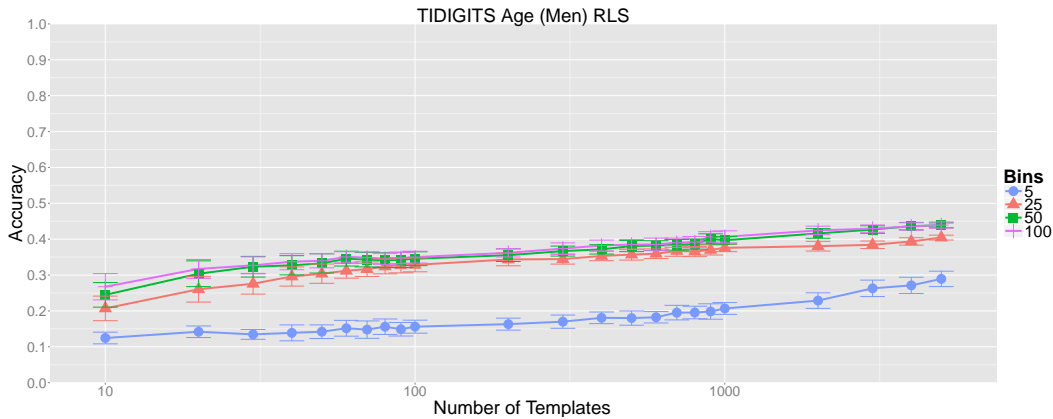


Figure 8: Mean classification accuracy as a function of number of templates, m , and bins, n . Accuracy is averaged over 30 random template samples for each m and error bars are displayed. In the “Age (Men)” dataset, we train on adult men and test on children, giving us an age mismatch. We see the weakest performance in this dataset, much worse than on the “Age (Women)” dataset. This is possibly due to the fact that women have higher pitched voices than men, creating less of a mismatch between women and children than men and children.

[4] Y. Bengio, A. C. Courville, and P. Vincent, “Representation learning: A review and new perspectives,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1798–1828, 2013.

[5] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” in *Proceedings of the IEEE*, vol. 86, pp. 2278–2324, 1998.

[6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *NIPS*, pp. 1106–1114, 2012.

[7] P. Niyogi, F. Girosi, and T. Poggio, “Incorporating prior information in machine learning by creating virtual examples,” in *Proceedings of the IEEE*, pp. 2196–2209, 1998.

[8] Y.-A. Mostafa, “Learning from hints in neural networks,” *Journal of complexity*, vol. 6, pp. 192–198, June 1990.

[9] V. N. Vapnik, *Statistical learning theory*. A Wiley-Interscience Publication 1998.

[10] I. Steinwart and A. Christmann, *Support vector machines*. Information Science and Statistics, New York: Springer, 2008.

- [11] B. Haasdonk, A. Vossen, and H. Burkhardt, “Invariance in kernel methods by haar-integration kernels,” in *SCIA*, Springer, 2005.
- [12] P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe, “Convexity, classification, and risk bounds,” *Journal of the American Statistical Association*, vol. 101, no. 473, pp. 138–156, 2006.
- [13] G. Wahba, *Spline models for observational data*, vol. 59 of *CBMS-NSF Regional Conference Series in Applied Mathematics*. Philadelphia, PA: SIAM, 1990.
- [14] W. B. Johnson and J. Lindenstrauss, “Extensions of lipschitz mappings into a hilbert space,” *Conference in modern analysis and probability*, 1984.
- [15] A. Rahimi and B. Recht, “Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning,” in *NIPS* 2008.
- [16] A. Rahimi and B. Recht, “Uniform approximation of functions with random bases,” in *Proceedings of the 46th Annual Allerton Conference*, 2008.
- [17] C. Williams and M. Seeger, “Using the nystrm method to speed up kernel machines,” in *NIPS*, 2001.
- [18] F. R. Bach, “On the equivalence between quadrature rules and random features,” *CoRR*, vol. abs/1502.06800, 2015.
- [19] C. Walder and O. Chapelle, “Learning with transformation invariant kernels,” in *NIPS*, 2007.
- [20] Y. Cho and L. K. Saul, “Kernel methods for deep learning,” in *NIPS*, pp. 342–350, 2009.
- [21] L. Bo, X. Ren, and D. Fox, “Kernel descriptors for visual recognition,” in *NIPS*, 2010.
- [22] J. Mairal, P. Koniusz, Z. Harchaoui, and C. Schmid, “Convolutional kernel networks,” in *NIPS*, 2014.
- [23] A. Tacchetti, P. K. Mallapragada, M. Santoro, and L. Rosasco, “Gurls: a least squares library for supervised learning,” *CoRR*, vol. abs/1303.0934, 2013.
- [24] S. Voinea, C. Zhang, G. Evangelopoulos, L. Rosasco, and T. Poggio, “Word-level invariant representations from acoustic waveforms,” vol. 14, pp. 3201–3205, September 2014.
- [25] M. Benzeghiba, R. De Mori, O. Deroo, S. Dupont, T. Erbes, D. Jouviet, L. Fissore, P. Laface, A. Mertins, C. Ris, R. Rose, V. Tyagi, and C. Wellekens, “Automatic speech recognition and speech variability: A review,” *Speech Communication*, vol. 49, pp. 763–786, 01 2007.
- [26] C. E. Clark, “The greatest of a finite set of random variables,” *Operations Research*, vol. 9, pp. 145–162, Mar-Apr 1961.
- [27] R. Vershynin, “Introduction to the non-asymptotic analysis of random matrices,” *Compressed Sensing: Theory and Applications*, Y. Eldar and G. Kutyniok, Eds. Cambridge University Press., 2011.
- [28] T. Inglot, “Inequalities for quantiles of the chi-square distribution,” *Probability and Mathematical Statistics*, vol. 30(2):339351, 2010.
- [29] A. Dvoretzky, J. Kiefer, and J. Wolfowitz, “Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator,” *Ann. Math. Statist.*, vol. 27, pp. 642–669, 09 1956.