# Provably Strict Generalisation Benefit for Invariance in Kernel Methods

Anonymous Author(s) Affiliation Address email

## Abstract

It is a commonly held belief that enforcing invariance improves generalisation. 1 Although this approach enjoys widespread popularity, it is only very recently that 2 a rigorous theoretical demonstration of this benefit has been established. In this 3 work we build on the function space perspective of Elesedy and Zaidi [8] to derive 4 a strictly non-zero generalisation benefit of incorporating invariance in kernel ridge 5 regression when the target is invariant to the action of a compact group. We study 6 invariance enforced by feature averaging and find that generalisation is governed 7 by a notion of effective dimension that arises from the interplay between the kernel 8 and the group. In building towards this result, we find that the action of the group 9 induces an orthogonal decomposition of both the reproducing kernel Hilbert space 10 and its kernel, which may be of interest in its own right. 11

## 12 **1** Introduction

Recently, there has been significant interest in models that are invariant to the action of a group on their inputs. It is believed that engineering models in this way improves sample efficiency and generalisation. Intuitively, if a task has an invariance, then a model that is constructed to be invariant ahead of time should require fewer examples to generalise than one that must learn to be invariant. Indeed, there are many application domains, such as fundamental physics or medical imaging, in which the invariance is known a priori [30, 32]. Although this intuition is certainly not new (e.g. [33]), it has inspired much recent work (for instance, see [36, 16]).

However, while implementations and practical applications abound, until very recently a rigorous theoretical justification for invariance was missing. As pointed out in [8], many prior works such as [29, 25] provide only worst-case guarantees on the performance of invariant algorithms. It follows that these results do not rule out the possibility of modern training algorithms automatically favouring invariant models, irrespective of the choice of architecture. Steps towards a more concrete theory of the benefit of invariance have been taken by [8, 21] and our work is a continuation along the path set by [8].

In this work we provide a precise characterisation of the generalisation benefit of invariance in kernel ridge regression. In contrast to [29, 25], this proves a *provably strict* generalisation benefit for invariant, feature-averaged models. In deriving this result, we provide insights into the structure of reproducing kernel Hilbert spaces in relation to invariant functions that we believe will be useful for analysing invariance in other kernel algorithms.

The use of feature averaging to produce invariant predictors enjoys both theoretical and practical success [18, 9]. For the purposes of this work, feature averaging is defined as training a model as normal (according to any algorithm) and then transforming the learned model to be invariant. This transformation is done by *orbit-averaging*, which means projecting the model on the space of

invariant functions using the operator O introduced in Section 2.3.

Kernel methods have a long been a mainstay of machine learning (see [31, Section 4.7] for a brief historical overview). Kernels can be viewed as mapping the input data into a potentially infinite

dimensional feature space, which allows for analytically tractable inference with non-linear predictors.

While modern machine learning practice is dominated by neural networks, kernels remain at the core

of much of modern theory. The most notable instance of this is the theory surrounding the *neural* 

42 *tangent kernel* [12], which states that the functions realised by an infinitely wide neural network

43 belong to a reproducing kernel Hilbert space (RKHS) with a kernel determined by the network

44 architecture. This relation has led to many results on the theory of optimisation and generalisation

45 of wide neural networks (e.g. [15, 3]). In the same vein, via the NTK, we believe the results of this

<sup>46</sup> paper can be extended to study wide, invariant neural networks.

## 47 **1.1 Summary of Contributions**

This paper builds towards a precise characterisation of the benefit of incorporating invariance in
 kernel ridge regression by feature averaging.

<sup>50</sup> Lemma 3, given in Section 3, forms the basis of our work, showing that the action of the group  $\mathcal{G}$  on <sup>51</sup> the input space induces an orthogonal decomposition of the RKHS  $\mathcal{H}$  as

$$\mathcal{H} = \mathcal{H}_S \oplus \mathcal{H}_A$$

where each term is an RKHS and  $\mathcal{H}_S$  consists of all of the invariant functions in  $\mathcal{H}$ . We stress that, while the main results of this paper concern kernel ridge regression, Lemma 3 holds regardless of

training algorithm and could be used to explore invariance in other kernel methods.

<sup>55</sup> Our main results are given in Section 4 and we outline them here. We define the generalisation gap

56  $\Delta(f, f')$  for two predictors f, f' as the difference in their test errors. If  $\Delta(f, f') > 0$  then f has

57 strictly better test performance than f'. Theorem 5 describes  $\Delta(f, f')$  for f being the solution to

kernel ridge regression and f' its invariant (feature averaged) version and shows that it is positive when the target is invariant.

More specifically, let  $X \sim \mu$  where  $\mu$  is  $\mathcal{G}$ -invariant and  $Y = f^*(X) + \xi$  with  $f^* \mathcal{G}$ -invariant and  $\mathbb{E}[\xi] = 0$ ,  $\mathbb{E}[\xi^2] = \sigma^2 < \infty$ . Let f be the solution to kernel ridge regression with kernel k and regularisation parameter  $\rho > 0$  on n i.i.d. training examples  $\{(X_i, Y_i) \sim (X, Y) : i = 1, ..., n\}$  and let f' be its feature averaged version. Our main result, Theorem 5, says that

$$\mathbb{E}[\Delta(f, f')] \ge \frac{\sigma^2 \dim_{\text{eff}}(\mathcal{H}_A) + \mathcal{E}}{(\sqrt{n}M_k + \rho/\sqrt{n})^2}$$

where  $M_k = \sup_x k(x, x) < \infty$ ,  $\mathcal{E} \ge 0$  describes the approximation errors and  $\dim_{\text{eff}}(\mathcal{H}_A)$  is the effective dimension of the RHKS  $\mathcal{H}_A$ . For an RKHS  $\mathcal{H}$  with kernel k the effective dimension is

66 defined by

$$\dim_{\text{eff}}(\mathcal{H}) = \int_{\mathcal{X}} k(x, y)^2 \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y).$$

<sup>67</sup> where  $\mathcal{X} = \sup \mu$ . We return to this quantity at various points in the paper. Finally, for intuition, <sup>68</sup> in Theorem 7 we specialise Theorem 5 to the linear setting and compute the bound exactly.

Assumptions and technical conditions are given in Section 2 along with an outline of the ideas of Elesedy and Zaidi [8] on which we build. Related works are discussed in Section 5.

## 71 **2** Background and Preliminaries

<sup>72</sup> In this section we provide a brief introduction to reproducing kernel Hilbert spaces (RKHS) and <sup>73</sup> the ideas we borrow from Elesedy and Zaidi [8]. Throughout this paper,  $\mathcal{H}$  with be an RKHS with <sup>74</sup> kernel *k*. In Section 2.2 we state some topological and measurability assumptions that are needed <sup>75</sup> for our proofs. These conditions are benign and the reader not interested in technicalities need take <sup>76</sup> from Section 2.2 only that  $\mu$  is  $\mathcal{G}$ -invariant and that the kernel *k* is bounded and satisfies Eq. (1). We <sup>77</sup> defer some background and technical results to Appendices **B** and **C** respectively.

#### 78 2.1 RKHS Basics

A Hilbert space is an inner product space that is complete with respect to the norm topology induced by the inner product. A reproducing kernel Hilbert space (RKHS)  $\mathcal{H}$  is Hilbert space of real functions  $f: \mathcal{X} \to \mathbb{R}$  on which the evaluation functional  $\delta_x : \mathcal{H} \to \mathbb{R}$  with  $\delta_x[f] = f(x)$  is continuous

 $\forall x \in \mathcal{X}$ , or, equivalently is a bounded operator. The Reisz Representation Theorem tells us that there 82 is a unique function  $k_x \in \mathcal{H}$  such that  $\delta_x[f] = \langle k_x, f \rangle_{\mathcal{H}}$  for any  $f \in \mathcal{H}$ , where  $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \to \mathbb{R}$ 83 is the inner product on  $\mathcal{H}$ . We identify the function  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  with  $k(x, y) = \langle k_x, k_y \rangle_{\mathcal{H}}$  as the 84 reproducing kernel of  $\mathcal{H}$ . Using the inner product representation, one can see that k is positive-definite 85 and symmetric. Conversely, the Moore-Aronszajn Theorem shows that for any positive-definite and 86 symmetric function k, there is a unique RKHS with reproducing kernel k. In addition, any Hilbert 87 space admitting a reproducing kernel is an RKHS. Finally, another characterisation of  $\mathcal{H}$  is as the 88 completion of the set of linear combinations of the form  $f_c(x) = \sum_{i=1}^n c_i k(x, x_i)$  for  $c_1, \ldots, c_n \in \mathbb{R}$ and  $x_1, \ldots, x_n \in \mathcal{X}$ . For (many) more details, see [31, Chapter 4]. 89 90 2.2 Technical Setup and Assumptions 91 **Input Space, Group and Measures** Let  $\mathcal{G}$  be a compact, second countable, Hausdorff topological 92 group with Haar measure  $\lambda$  (see [13, Theorem 2.27]). Let  $\mathcal{X}$  be a non-empty Polish space admitting 93 a finite,  $\mathcal{G}$ -invariant Borel measure  $\mu$ , with supp  $\mu = \mathcal{X}$ . We normalise  $\mu(\mathcal{X}) = \lambda(\mathcal{G}) = 1$ , the latter 94 is possible because  $\lambda$  is a Radon measure. We assume that  $\mathcal{G}$  has a measurable action on  $\mathcal{X}$  that we 95 will write as gx for  $g \in \mathcal{G}, x \in \mathcal{X}$ . A measurable action is one such that the map  $g: \mathcal{G} \times \mathcal{X} \to \mathcal{X}$ 96 is  $(\lambda \otimes \mu)$ -measurable. A function  $f : \mathcal{X} \to \mathbb{R}$  is  $\mathcal{G}$ -invariant if  $f(gx) = f(x) \ \forall x \in \mathcal{X} \ \forall g \in \mathcal{G}$ . 97 Similarly, a measure  $\mu$  on  $\mathcal{X}$  is  $\mathcal{G}$ -invariant if  $\forall g \in \mathcal{G}$  and any  $\mu$ -measurable  $B \subset \mathcal{X}$  the pushforward 98

<sup>95</sup> of  $\mu$  by the action of  $\mathcal{G}$  equals  $\mu$ , i.e.  $(g_*\mu)(B) = \mu(B)$ . This means that if  $X \sim \mu$  then  $gX \sim \mu$ <sup>100</sup>  $\forall g \in \mathcal{G}$ . We will make use of the fact that the Haar measure is  $\mathcal{G}$ -invariant when  $\mathcal{G}$  acts on itself by <sup>101</sup> either left or right multiplication, the latter holding because  $\mathcal{G}$  is compact. Up to normalisation,  $\lambda$  is <sup>102</sup> the unique measure on  $\mathcal{G}$  with this property.

The Kernel and the RKHS Let  $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  be a measurable kernel with RKHS  $\mathcal{H}$  such that  $k(\cdot, x) : \mathcal{X} \to \mathbb{R}$  is continuous for any  $x \in \mathcal{X}$ . Assume that  $\sup_{x \in \mathcal{X}} k(x, x) = M_k < \infty$  and note that this implies that k is bounded since

$$k(x,x') = \langle k_x, k_{x'} \rangle_{\mathcal{H}} \le ||k_x||_{\mathcal{H}} ||k_{x'}||_{\mathcal{H}} = \sqrt{k(x,x)}\sqrt{k(x',x')} \le M_k$$

Every  $f \in \mathcal{H}$  is  $\mu$ -measurable, bounded and continuous by [31, Lemmas 4.24 and 4.28] and in addition  $\mathcal{H}$  is separable using [31, Lemma 4.33]. These conditions allow the application of [31, Theorem 4.26] to relate  $\mathcal{H}$  to  $L_2(\mathcal{X}, \mu)$  in the proofs building towards Lemma 3, given in Appendix C. We assume that the kernel satisfies, for all  $x, y \in \mathcal{X}$ ,

$$\int_{\mathcal{G}} k(gx, y) \,\mathrm{d}\lambda(g) = \int_{\mathcal{G}} k(x, gy) \,\mathrm{d}\lambda(g). \tag{1}$$

For this it is sufficient to have k(gx, y) equal to k(x, gy) or  $k(x, g^{-1}y)$ , the latter uses compactness (hence unimodularity) of  $\mathcal{G}$  to change variables  $g \leftrightarrow g^{-1}$ . Highlighting two special cases: any inner product kernel  $k(x, x') = \kappa(\langle x, x' \rangle)$  such that the action of  $\mathcal{G}$  is unitary with respect to  $\langle \cdot, \cdot \rangle$ satisfies Eq. (1), as does any stationary kernel  $k(x, x') = \kappa(||x - x'||)$  with norm that is preserved by  $\mathcal{G}$  in the sense that ||gx - gx'|| = ||x - x'|| for any  $g \in \mathcal{G}, x, x' \in \mathcal{X}$ .

## 115 2.3 Invariance from a Function Space Perspective

Given a function  $f : \mathcal{X} \to \mathbb{R}$  we can define a corresponding orbit-averaged function  $\mathcal{O}f : \mathcal{X} \to \mathbb{R}$ with values

$$\mathcal{O}f(x) = \int_{\mathcal{G}} f(gx) \,\mathrm{d}\lambda(g).$$

<sup>118</sup>  $\mathcal{O}f$  will exist whenever f is  $\mu$ -measurable. Note that  $\mathcal{O}$  is a linear operator and, from the invariance <sup>119</sup> of  $\lambda$ ,  $\mathcal{O}f$  is always  $\mathcal{G}$ -invariant. Interestingly, f is  $\mathcal{G}$ -invariant *only* if  $f = \mathcal{O}f$ . Elesedy and Zaidi [8] <sup>120</sup> use these observations to characterise invariant functions and study their generalisation properties. In <sup>121</sup> short, this work extends these insights to kernel methods. Along the way, we will make frequent use <sup>122</sup> of the following (well known) facts about  $\mathcal{O}$ .

Lemma 1 ([8, Propositions 23 and 24]). A function f is  $\mathcal{G}$ -invariant if and only if  $\mathcal{O}f = f$ . This implies that  $\mathcal{O}$  is a projection operator, so can have only two eigenvalues 0 and 1.

Lemma 2 ([8, Lemma 1]).  $\mathcal{O}: L_2(\mathcal{X}, \mu) \to L_2(\mathcal{X}, \mu)$  is well-defined and self-adjoint. Hence, L<sub>26</sub>  $L_2(\mathcal{X}, \mu)$  has the orthogonal decomposition

$$L_2(\mathcal{X},\mu) = S \oplus A$$

where  $S = \{f \in L_2(\mathcal{X}, \mu) : f \text{ is } \mathcal{G} \text{ invariant}\}$  and  $A = \{f \in L_2(\mathcal{X}, \mu) : \mathcal{O}f = 0\}.$ 

The meaning of Lemma 2 is that any  $f \in L_2(\mathcal{X}, \mu)$  has a (unique) decomposition  $f = \overline{f} + f^{\perp}$ where  $\overline{f} = \mathcal{O}f$  is  $\mathcal{G}$ -invariant and  $\mathcal{O}f^{\perp} = 0$ . A noteworthy consequence of this setup, as discussed in [8], is a provably non-negative generalisation benefit for feature averaging. In particular, for any predictor  $f \in L_2(\mathcal{X}, \mu)$ , if the target  $f^* \in L_2(\mathcal{X}, \mu)$  is  $\mathcal{G}$ -invariant then the test error R(f) = $\mathbb{E}_{X \sim \mu}[(f(X) - f^*(X))^2]$  satisfies

$$R(f) - R(\bar{f}) = ||f^{\perp}||^2_{L_2(\mathcal{X},\mu)} \ge 0.$$

<sup>133</sup> The same holds if the target is corrupted by independent, zero mean (additive) noise.

## 134 **3 Induced Structure of** $\mathcal{H}$

In this section we present Lemma 3, which is an analog of Lemma 2 for RKHSs. Lemma 3 shows that for any compact group  $\mathcal{G}$  and RKHS  $\mathcal{H}$ , if the kernel for  $\mathcal{H}$  satisfies the assumptions in Section 2.2, then  $\mathcal{H}$  can be viewed as being built from two orthogonal RKHSs, one consisting of invariant functions and another of those that vanish when averaged over  $\mathcal{G}$ . Later in the paper, this decomposition will allow us to analyse the generalisation benefit of invariant predictors.

It may seem at first glance that Lemma 3 should follow immediately from Lemma 2, but this is not the case. First, it is not obvious that for any  $f \in \mathcal{H}$ , its orbit averaged version  $\mathcal{O}f$  is also in  $\mathcal{H}$ . Moreover, in contrast with  $L_2(\mathcal{X}, \mu)$ , an explicit form for the inner product on  $\mathcal{H}$  is not immediate, which means that some work is needed to check that  $\mathcal{O}$  is self-adjoint on  $\mathcal{H}$ . These are important requirements for the proofs of both Lemmas 2 and 3 and we establish them, along with  $\mathcal{O}$  being continuous on  $\mathcal{H}$ , in Lemmas C.6 and C.7 and Corollary C.8 respectively. The assumption that the kernel satisfies Eq. (1) plays a central role.

#### 147 **Lemma 3.** $\mathcal{H}$ admits the orthogonal decomposition

$$\mathcal{H} = \mathcal{H}_S \oplus \mathcal{H}_A$$

where  $\mathcal{H}_S = \{f \in \mathcal{H} : f \text{ is } \mathcal{G}\text{-invariant}\}$  and  $\mathcal{H}_A = \{f \in \mathcal{H} : \mathcal{O}f = 0\}$ . Moreover,  $\mathcal{H}_S$  is an RKHS with kernel

$$\bar{k}(x,y) = \int_{\mathcal{G}} k(x,gy) \,\mathrm{d}\lambda(g)$$

150 and  $\mathcal{H}_A$  is an RKHS with kernel

$$k^{\perp}(x,y) = k(x,y) - \bar{k}(x,y).$$

<sup>151</sup> Finally,  $\bar{k}$  is *G*-invariant in both arguments.

*Proof.* From Lemma 1 we know that  $\mathcal{O}$  is a projection operator. Since it is self-adjoint,  $\mathcal{O}$  is even an orthogonal projection on  $\mathcal{H}$ : let  $h_S$  have eigenvalue 1 and  $h_A$  have eigenvalue 0 under  $\mathcal{O}$ , then

$$\langle h_S, h_A \rangle_{\mathcal{H}} = \langle \mathcal{O}h_S, h_A \rangle_{\mathcal{H}} = \langle h_S, \mathcal{O}h_A \rangle_{\mathcal{H}} = 0.$$

Therefore, by linearity, for any  $f \in \mathcal{H}$  we can write  $f = \bar{f} + f^{\perp}$  where  $\bar{f} = \mathcal{O}f \in \mathcal{H}_S$  is  $\mathcal{G}$ -invariant and  $f^{\perp} = f - \mathcal{O}f \in \mathcal{H}_A$  and these terms are mutually orthogonal.

By the linearity of  $\mathcal{O}$ , it is clear that  $\mathcal{H}_S = \mathcal{OH}$  is an inner product space. It is easy to show that  $\mathcal{O}$  being continuous implies  $\mathcal{H}_S$  is complete. Thus  $\mathcal{H}_S$  is a Hilbert space, and an RKHS since the evaluation functional is clearly continuous on  $\mathcal{H}_S \subset \mathcal{H}$ . For any  $h_S \in \mathcal{H}_S$  we have

$$h_S(x) = \langle h_S, k_x \rangle_{\mathcal{H}} = \langle h_S, \mathcal{O}k_x \rangle_{\mathcal{H}} = \langle h_S, \bar{k}_x \rangle_{\mathcal{H}}$$

and the uniqueness afforded by the Reisz representation theorem tells us that the reproducing kernel for  $\mathcal{H}_S$  is  $\bar{k}(x,y) = \int_{\mathcal{G}} k(x,gy) d\lambda(g)$ . We have  $\| \text{id} - \mathcal{O} \| \le 2$  and we can do the same argument to show that  $\mathcal{H}_A$  is an RKHS with reproducing kernel  $k^{\perp}$  as claimed. Note that one can write  $k^{\perp}(x,y) = \langle k_x^{\perp}, k_y^{\perp} \rangle_{\mathcal{H}}$  so it must be positive-definite. The  $\mathcal{G}$ -invariance of  $\bar{k}(x,y)$  in both arguments is immediate from Eq. (1) and Lemma 1.

As stated earlier, the perspective provided by Lemma 3 will support our analysis of generalisation. Just as with Lemma 2, Lemma 3 says that any  $f \in \mathcal{H}$  can be written as  $f = \bar{f} + f^{\perp}$  where  $\bar{f}$  is  $\mathcal{G}$ -invariant and  $\mathcal{O}f^{\perp} = 0$  with  $\langle \bar{f}, f^{\perp} \rangle_{\mathcal{H}} = 0$ . As an aside,  $\bar{k}$  happens to qualify as a *Haar Integration Kernel*, a concept introduced by Haasdonk, Vossen, and Burkhardt [10]. We will see that a notion of effective dimension of the RKHS  $\mathcal{H}_A$  with kernel  $k^{\perp}$  governs the generalisation gap between an arbitrary predictor f and its invariant version  $\mathcal{O}f$ . This effective dimension arises from the spectral theory of an integral operator related to k, which we develop in the next section.

#### 171 3.1 Spectral Representation and Effective Dimension

In this section we consider the spectrum of an integral operator related to the kernel k. This analysis will ultimately allow us to define a notion of effective dimension of  $\mathcal{H}_A$  that we will later see is

important to the generalisation of invariant predictors. While the integral operator setup is standard, the use of this technique to identify an effective dimension of  $\mathcal{H}_A$  is novel.

176 Define the integral operator  $S_k : L_2(\mathcal{X}, \mu) \to \mathcal{H}$  by

$$S_k f(x) = \int_{\mathcal{X}} k(x, x') f(x') \,\mathrm{d}\mu(x').$$

One way of viewing things is that  $S_k$  assigns to every element in  $L_2(\mathcal{X}, \mu)$  a function in  $\mathcal{H}$ . On the other hand, every  $f \in \mathcal{H}$  is bounded so has  $||f||_{L_2(\mathcal{X},\mu)} < \infty$  and belongs to some element of  $L_2(\mathcal{X}, \mu)$ . We write  $\iota : \mathcal{H} \to L_2(\mathcal{X}, \mu)$  for the *inclusion map* that sends f to the element of  $L_2(\mathcal{X}, \mu)$ that contains f. In Lemma C.1 we show that  $\iota$  is injective, so any element of  $L_2(\mathcal{X}, \mu)$  contains at most one  $f \in \mathcal{H}$ .

One can define  $T_k : L_2(\mathcal{X}, \mu) \to L_2(\mathcal{X}, \mu)$  by  $T_k = \iota \circ S_k$ , and [31, Theorem 4.27] says that  $T_k$  is compact, positive, self-adjoint and trace-class. In addition,  $L_2(\mathcal{X}, \mu)$  is separable by [7, Proposition 3.4.5], because  $\mathcal{X}$  is Polish and  $\mu$  is a Borel measure, so has a countable orthonormal basis. Hence, by the Spectral Theorem, there exists a countable orthonormal basis  $\{\tilde{e}_i\}$  for  $L_2(\mathcal{X}, \mu)$  such that  $T_k \tilde{e}_i = \lambda_i \tilde{e}_i$  where  $\lambda_1 \ge \lambda_2 \ge \cdots \ge 0$  are the eigenvalues of  $T_k$ . Moreover, since  $\iota$  is injective, for each of the  $\tilde{e}_i$  for which  $\lambda_i > 0$  there is a unique  $e_i \in \mathcal{H}$  such that  $\iota e_i = \tilde{e}_i$  and  $S_k \tilde{e}_i = \lambda_i e_i$ .

188 Now, since  $\iota k_x \in L_2(\mathcal{X}, \mu)$  we have

$$\iota k_x = \sum_i \langle \iota k_x, \tilde{e}_i \rangle_{L_2(\mathcal{X},\mu)} \tilde{e}_i = \sum_i (S_k \tilde{e}_i)(x) \tilde{e}_i = \sum_i \lambda_i e_i(x) \tilde{e}_i.$$
(2)

From now on we permit ourself to drop the  $\iota$  to reduce clutter. We use the above to define

$$j(x,y) = \langle k_x, k_y \rangle_{L_2(\mathcal{X},\mu)}, \quad \overline{j}(x,y) = \langle \overline{k}_x, \overline{k}_y \rangle_{L_2(\mathcal{X},\mu)} \quad \text{and} \quad j^{\perp}(x,y) = \langle k_x^{\perp}, k_y^{\perp} \rangle_{L_2(\mathcal{X},\mu)}.$$

<sup>190</sup> These quantities will appear again in our analysis of the generalisation of invariant kernel methods.

Indeed, we will see later in this section that  $\mathbb{E}[j^{\perp}(X, X)]$  is a type of effective dimension of  $\mathcal{H}_A$ . Following Eq. (2), one finds the series representations given below in Lemma 4.

The reader may have noticed that our setup is very similar to the one provided by Mercer's theorem. However, we do not assume compactness of  $\mathcal{X}$  and so (the classical form of) Mercer's Theorem does not apply. In particular, the set  $\{e_i\}$  (even when scaled appropriately) need not form an orthonormal basis in  $\mathcal{H}$ . This aspect of our work is a feature, rather than a bug: the loosening of the compactness condition allows application to common settings such as  $\mathcal{X} = \mathbb{R}^n$ .

198 Lemma 4. We have

$$j = \bar{j} + j^{\perp}.$$

199 Furthermore, let  $\bar{e}_i = \mathcal{O}e_i$  and  $e_i^{\perp} = e_i - \bar{e}_i$  then

$$j(x,y) = \sum_i \lambda_i^2 e_i(x) e_i(y), \quad \bar{j}(x,y) = \sum_i \lambda_i^2 \bar{e}_i(x) \bar{e}_i(y), \quad \text{and} \quad j^{\scriptscriptstyle \perp}(x,y) = \sum_i \lambda_i^2 e_i^{\scriptscriptstyle \perp}(x) e_j^{\scriptscriptstyle \perp}(y).$$

Finally, the function  $\sum_i \lambda_i^2 \bar{e}_i \otimes e_i^{\perp} : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  vanishes everywhere.

*Proof.* We show in Lemma C.2 that  $\mathcal{O}$  and  $S_k$  commute on  $L_2(\mathcal{X}, \mu)$  and  $\mathcal{O}$  is self-adjoint on  $L_2(\mathcal{X}, \mu)$  by Lemma 1, so  $\mathcal{O}$  and  $\iota$  (the adjoint of  $S_k$  by [31, Theorem 4.26]) must also commute. The first comment is then immediate from the observation that if  $a \in \mathcal{H}_S$  and  $b \in \mathcal{H}_A$  one has

$$\langle \iota a, \iota b \rangle_{L_2(\mathcal{X},\mu)} = \langle \iota \mathcal{O}a, \iota b \rangle_{L_2(\mathcal{X},\mu)} = \langle \mathcal{O}\iota a, \iota b \rangle_{L_2(\mathcal{X},\mu)} = \langle \iota a, \iota \mathcal{O}b \rangle_{L_2(\mathcal{X},\mu)} = 0.$$

204 We also have both of

$$\langle \iota \bar{k}_x, \tilde{e}_i \rangle_{L_2(\mathcal{X}, \mu)} = \langle \iota k_x, \mathcal{O} \tilde{e}_i \rangle_{L_2(\mathcal{X}, \mu)} = S_k \mathcal{O} \tilde{e}_i = \mathcal{O} S_k \tilde{e}_i = \lambda_i \bar{e}_i$$

205 and

$$\langle \iota k_x^{\perp}, \tilde{e}_i \rangle_{L_2(\mathcal{X}, \mu)} = \langle \iota k_x, (\mathrm{id} - \mathcal{O}) \tilde{e}_i \rangle_{L_2(\mathcal{X}, \mu)} = S_k (\mathrm{id} - \mathcal{O}) \tilde{e}_i = (\mathrm{id} - \mathcal{O}) S_k \tilde{e}_i = \lambda_i e_i^{\perp}.$$

Therefore  $\iota \bar{k}_x = \sum_i \lambda_i \bar{e}_i(x) \tilde{e}_i$  and  $\iota k_x^{\perp} = \sum_i \lambda_i e_i^{\perp}(x) \tilde{e}_i$ . Taking inner products on  $L_2(\mathcal{X}, \mu)$  gives the remaining results.

Before turning to generalisation, we describe how the above quantities can be used to define a measure effective dimension. We define

$$\dim_{\text{eff}}(\mathcal{H}) = \mathbb{E}[j(X, X)]$$

where  $X \sim \mu$ . Applying Fubini's theorem, we find

$$\dim_{\text{eff}}(\mathcal{H}) = \sum_{i} \lambda_i^2 \mathbb{E}[e_i(X)^2] = \sum_{i} \lambda_i^2 \|\tilde{e}_i\|_{L_2(\mathcal{X},\mu)}^2 = \sum_{i} \lambda_i^2.$$

The series converges by the comparison test because  $\lambda_i \ge 0$  and  $\sum_i \lambda_i = \text{Tr}(T_k) < \infty$ . We have dim<sub>eff</sub>( $\mathcal{H}$ ) =  $\text{Tr}(T_k^2)$  and we can think of this (very informally) as taking  $L_2(\mathcal{X}, \mu)$ , pushing it through  $\mathcal{H}$  twice using  $T_k$  and then measuring its size. Now because  $j = \overline{j} + j^{\perp}$  we get

$$\dim_{\mathrm{eff}}(\mathcal{H}) = \dim_{\mathrm{eff}}(\mathcal{H}_S) + \dim_{\mathrm{eff}}(\mathcal{H}_A)$$

214 with

$$\dim_{\text{eff}}(\mathcal{H}_A) = \sum_i \lambda_i^2 \|\tilde{e}_i^{\perp}\|_{L_2(\mathcal{X},\mu)}^2 = \operatorname{Tr}(T_k^2) - \operatorname{Tr}((\mathcal{O}T_k)^2)$$

where  $\tilde{e}_i^{\perp} = \iota e_i^{\perp}$ . Again, very informally, this can be thought of as pushing  $L_2(\mathcal{X}, \mu)$  through  $\mathcal{H}_A$ twice and measuring the size of the output. In the next section we will consider the generalisation of kernel ridge regression and find that  $\dim_{\text{eff}}(\mathcal{H}_A)$  plays a critical role.

# 218 **4** Generalisation

In this section we apply the theory developed in Section 3 to study the impact of invariance on kernel ridge regression with an invariant target. We analyse the generalisation benefit of feature averaging, finding a strict benefit when the target is *G*-invariant.

#### 222 4.1 Kernel Ridge Regression

Given input/output pairs  $\{(x_i, y_i) : i = 1, ..., n\}$  where  $x_i \in \mathcal{X}$  and  $y_i \in \mathbb{R}$ , kernel ridge regression (KRR) returns a predictor that solves the optimisation problem

$$\underset{f \in \mathcal{H}}{\operatorname{argmin}} C(f) \quad \text{where} \quad C(f) = \sum_{i=1}^{n} (f(x_i) - y_i)^2 + \rho \|f\|_{\mathcal{H}}^2$$
(3)

and  $\rho > 0$  is the regularisation parameter. KRR can be thought of as performing ridge regression in a possibly infinite dimensional feature space  $\mathcal{H}$ . The representer theorem tells us that the solution to this problem is of the form  $f(x) = \sum_{i=1}^{n} \alpha_i k_{x_i}(x)$  where  $\alpha \in \mathbb{R}^n$  solves

$$\operatorname*{argmin}_{\alpha \in \mathbb{R}^{n}} \left\{ \| \boldsymbol{Y} - K\alpha \|_{2}^{2} + \rho \alpha^{\top} K\alpha \right\},$$
(4)

228  $\mathbf{Y} \in \mathbb{R}^n$  is the standard row-stacking of the training outputs with  $\mathbf{Y}_i = y_i$  and K is the kernel Gram 229 matrix with  $K_{ij} = k(x_i, x_j)$ . We consider solutions of the form<sup>1</sup>  $\alpha = (K + \rho I)^{-1} \mathbf{Y}$  which results 230 in the predictor

$$f(x) = k_x(\boldsymbol{X})^\top (K + \rho I)^{-1} \boldsymbol{Y}$$

where  $k_x(\mathbf{X}) \in \mathbb{R}^n$  is the vector with components  $k_x(\mathbf{X})_i = k_x(x_i)$ . We will compare the generalisation performance of this predictor with that of its averaged version

$$\bar{f} = \bar{k}_x(\boldsymbol{X})^\top (K + \rho I)^{-1} \boldsymbol{Y} \in \mathcal{H}_S.$$

<sup>233</sup> To do this we look at the generalisation gap.

<sup>&</sup>lt;sup>1</sup>When K is a positive definite matrix this will be the *only* solution. If K is singular then  $\exists c \in \mathbb{R}^n$  with  $\sum_{ij} K_{ij}c_ic_j = \|\sum_i c_ik_{x_i}\|_{\mathcal{H}}^2 = 0$  so  $\sum_i c_ik_{x_i}$  is identically 0 and  $\forall f \in \mathcal{H}$  we get  $\sum_i c_if(x_i) = 0$  (see [19, Section 4.6.2]). Clearly, this can't happen if  $\mathcal{H}$  is sufficiently expressive. In any case, the chosen  $\alpha$  is the minimum in Euclidean norm of all possible solutions.

#### 234 4.2 Generalisation Gap

The generalisation gap is a quantity that compares the expected test performances of two predictors on a given task. Given a probability distribution  $\mathbb{P}$ , data  $(X, Y) \sim \mathbb{P}$  and loss function l defining a supervised learning task, we define the generalisation gap between two predictors f and f' to be

$$\Delta(f, f') = \mathbb{E}[l(f(X), Y)] - \mathbb{E}[l(f'(X), Y)]$$

where the expectations are conditional on the given realisations of f, f' if the predictors are random. In this paper we consider  $l(a,b) = (a-b)^2$  the squared-error loss and we will assume  $Y = f^*(X) + \xi$ for some target function  $f^*$  where  $\xi$  is has mean 0, finite variance and is independent of X. In this case, the generalisation gap reduces to

$$\Delta(f, f') = \mathbb{E}[(f(X) - f^*(X))^2] - \mathbb{E}[(f'(X) - f^*(X))^2]$$

<sup>242</sup> Clearly, if  $\Delta(f, f') > 0$  then we expect strictly better test performance from f' than f.

#### 243 4.3 Generalisation Benefit of Feature Averaging

We are now in a position to give our main result, which is a characterisation of the generalisation benefit of invariance in kernel methods. This is in some sense a generalisation of [8, Theorem 6] and we will return to this comparison later. We emphasise that Theorem 5 holds under quite general conditions that cover many practical applications.

**Theorem 5.** Let the training data be  $\{(X_i, Y_i) : i = 1, ..., n\}$  i.i.d. with  $Y_i = f^*(X_i) + \xi_i$  where X<sub>i</sub> ~  $\mu$ ,  $f^* \in L_2(\mathcal{X}, \mu)$  is  $\mathcal{G}$ -invariant and  $\{\xi_i : i = 1, ..., n\}$  are independent of each other and the X<sub>i</sub> >  $\mu$ ,  $f^* \in L_2(\mathcal{X}, \mu)$  is  $\mathcal{G}$ -invariant and  $\{\xi_i : i = 1, ..., n\}$  are independent of each other and the X<sub>i</sub> >  $\lambda_i$ , with  $\mathbb{E}[\xi_i] = 0$  and  $\mathbb{E}[\xi_i^2] = \sigma^2 < \infty$ . Let  $f = \operatorname{argmin}_{f \in \mathcal{H}} C(f)$  be the solution to Eq. (3) and let  $\overline{f} = \mathcal{O}f \in \mathcal{H}_S$  be the result of applying feature averaging to f, then the generalisation gap with the squared-error loss satisfies

$$\mathbb{E}[\Delta(f,\bar{f})] \ge \frac{\sigma^2 \dim_{\mathrm{eff}}(\mathcal{H}_A) + \mathbb{E}[f^*(X)^2 j^{\perp}(X,X)]}{(\sqrt{n}M_k + \rho/\sqrt{n})^2}$$

<sup>253</sup> where each term is non-negative and

$$\dim_{\text{eff}}(\mathcal{H}_A) \coloneqq \operatorname{Tr}(T_k^2) - \operatorname{Tr}((\mathcal{O}T_k)^2) = \mathbb{E}[j^{\perp}(X,X)] = \sum_{\alpha} \lambda_{\alpha}^2 \|\tilde{e}_{\alpha}^{\perp}\|_{L_2(\mathcal{X},\mu)}^2 \ge 0$$

- is the *effective dimension* of  $\mathcal{H}_A$ .
- *Proof.* Let  $J^{\perp}$  be the Gram matrix with components  $J_{ij}^{\perp} = j^{\perp}(X_i, X_j)$  let  $u \in \mathbb{R}^n$  have components  $u_i = f^*(x_i)$ . We can use Lemma 2 to get

$$\Delta(f,\bar{f}) = \mathbb{E}[(k_X^{\perp}(\boldsymbol{X})^{\top}(K+\rho I)^{-1}\boldsymbol{Y})^2|\boldsymbol{X},\boldsymbol{Y}]$$

where  $k_x^{\perp}(X) \in \mathbb{R}^n$  with  $k_x^{\perp}(X)_i = k_x^{\perp}(X_i)$ . Let  $\boldsymbol{\xi} \in \mathbb{R}^n$  have components  $\boldsymbol{\xi}_i = \xi_i$  then one finds

$$\mathbb{E}[\Delta(f,\bar{f})|\boldsymbol{X}] = \mathbb{E}[(k_X^{\perp}(\boldsymbol{X})^{\top}(K+\rho I)^{-1}u)^2|\boldsymbol{X}] + \mathbb{E}[(k_X^{\perp}(\boldsymbol{X})^{\top}(K+\rho I)^{-1}\boldsymbol{\xi})^2|\boldsymbol{X}] \\ = u^{\top}(K+\rho I)^{-1}J^{\perp}(K+\rho I)^{-1}u + \sigma^2 \operatorname{Tr}\left(J^{\perp}(K+\rho I)^{-2}\right)$$

- where the first equality follows because  $\boldsymbol{\xi}$  has mean 0 and the second comes from the trace trick.
- The first term vanishes in expectation. To see this, first note that it is non-negative because  $J^{\perp}$  is positive semi-definite, while at the same time Consider the first term. We have

$$u^{\top}(K+\rho I)^{-1}J^{\perp}(K+\rho I)^{-1}u = \operatorname{Tr}((K+\rho I)^{-1}J^{\perp}(K+\rho I)^{-1}uu^{\top})$$

applying Corollary B.2 twice and using Lemma B.3 with boundedness of the kernel gives

$$u^{\top}(K+\rho I)^{-1}J^{\perp}(K+\rho I)^{-1}u \ge \lambda_{\min}((K+\rho I)^{-1})^{2}\operatorname{Tr}(J^{\perp}uu^{\top}) \ge \frac{u^{\top}J^{\perp}u}{(Mn+\rho)^{2}}$$

262 SO

$$\mathbb{E}[u^{\top}(K+\rho I)^{-1}J^{\perp}(K+\rho I)^{-1}u] \geq \frac{\mathbb{E}[u^{\top}J^{\perp}u]}{(Mn+\rho)^2} = \frac{\sum_{ij}\mathbb{E}[f^*(X_i)f^*(X_j)j^{\perp}(X_i,X_j)]}{(Mn+\rho)^2}.$$

It remains to show that the above vanishes when  $i \neq j$ . Using Lemma 4 we have

$$\mathbb{E}[f^*(X_i)f^*(X_j)j^{\perp}(X_i,X_j)] = \mathbb{E}[f^*(X_i)f^*(X_j)\sum_{\alpha}\lambda_{\alpha}^2 e_{\alpha}^{\perp}(X_i)e_{\alpha}^{\perp}(X_j)].$$

By Fubini's theorem, since the counting measure on  $\alpha$  is  $\sigma$ -finite, we can exchange the expectation and the sum as long as

$$\sum_{\alpha} \lambda_{\alpha}^{2} \mathbb{E}[f^{*}(X_{i})f^{*}(X_{j})e_{\alpha}^{\perp}(X_{i})e_{\alpha}^{\perp}(X_{j})] < \infty.$$

On the other hand, since  $i \neq j$  each term is just  $\mathbb{E}[f^*(X)e^{\perp}_{\alpha}(X)]^2 = \langle \iota f^*, \tilde{e}^{\perp}_{\alpha} \rangle^2_{L_2(\mathcal{X},\mu)} = 0$  by the *G*-invariance of  $f^*$  and the orthogonality in Lemma 2.

Moving to the second term, we have again by two applications of Corollary B.2 and then Lemma B.3 with boundedness of the kernel that

$$\operatorname{Tr}\left(J^{\perp}(K+\rho I)^{-2}\right) \geq \lambda_{\min}\left((K+\rho I)^{-2}\right)\operatorname{Tr}(J^{\perp}) \geq \frac{\operatorname{Tr}(J^{\perp})}{(M_k n+\rho)^2}$$

270 and then

$$\frac{1}{n} \mathbb{E}[\mathrm{Tr}(J^{\perp})] = \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}\left[\sum_{\alpha} \lambda_{\alpha}^{2} e_{\alpha}^{\perp}(X_{i}) e_{\alpha}^{\perp}(X_{i})\right]$$
$$= \sum_{\alpha} \lambda_{\alpha}^{2} \|\tilde{e}_{\alpha}^{\perp}\|_{L_{2}(\mathcal{X},\mu)}^{2}$$
$$= \sum_{\alpha} \lambda_{\alpha}^{2} - \sum_{\alpha} \lambda_{\alpha}^{2} \|\mathcal{O}\tilde{e}_{\alpha}\|_{L_{2}(\mathcal{X},\mu)}^{2}$$
$$= \mathrm{Tr}(T_{k}^{2}) - \mathrm{Tr}(T_{k}^{2}\mathcal{O})$$

where we exchange the expectation and sum again using Fubini's theorem as we will justify now. Considering the sum in the second line, note that  $\|\tilde{e}_{\alpha}\|_{L_{2}(\mathcal{X},\mu)}^{2} = 1 = \|\mathcal{O}\tilde{e}_{\alpha}\|_{L_{2}(\mathcal{X},\mu)}^{2} + \|\tilde{e}_{\alpha}^{\perp}\|_{L_{2}(\mathcal{X},\mu)}^{2}$ by Lemma 2 so the sum converges if  $\sum_{i} \lambda_{i}^{2}$  converges. However,  $T_{k}$  is both positive and trace-class from Section 3.1 so  $\lambda_{i} \geq 0$  and  $\sum_{i} \lambda_{i} < \infty$  (using Lidskii's theorem) so  $\sum_{i} \lambda_{i}^{2}$  converges by the comparison test.

Theorem 5 shows that feature averaging is provably beneficial in terms of generalisation if the mean of the target distribution is invariant. If  $\mathcal{H}$  contains any functions that are not  $\mathcal{G}$ -invariant then the lower bound is strictly positive. One might think that, given enough training examples, the solution fto Eq. (3) would *learn* to be  $\mathcal{G}$ -invariant. Theorem 5 shows that this cannot happen unless the number of examples dominates the effective dimension of  $\mathcal{H}_A$ .

Recall the subspace A in Lemma 2. The role of  $\dim_{\text{eff}}(\mathcal{H}_A)$  mirrors that of dim A in [8, Theorem 6] and in the context of the theorem (linear models) A can be thought of as  $\mathcal{H}_A$  when k is the linear kernel. In this sense Theorem 5 is a generalisation of [8, Theorem 6]. It is for this reason that we believe that, although the constant  $M_k$  in the denominator is likely not optimal, the O(1/n) rate that matches [8] is tight. We leave a more precise analysis of the constants to future work.

The second term in the numerator can be interpreted as quantifying the differences in bias. One has by the definition of  $j^{\perp}$ , that

$$\mathbb{E}[f^*(X)^2 j^{\perp}(X,X)] = \int_{\mathcal{X}} f^*(y)^2 k^{\perp}(x,y)^2 \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y)$$
(5)

and we also have the following.

**Proposition 6.** 

$$\int_{\mathcal{X}} f^*(y)^2 k^{\perp}(x,y)^2 \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y) = \int_{\mathcal{X}} f^*(y)^2 \left(k(x,y)^2 - \bar{k}(x,y)^2\right) \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y)$$

289 Proof. Using  $k^{\perp} = k - \bar{k}$ 

$$\begin{split} \int_{\mathcal{X}} f^*(y)^2 k^{\perp}(x,y)^2 \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) &= \int_{\mathcal{X}} f^*(y)^2 k(x,y)^2 \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) \\ &- 2 \int_{\mathcal{X}} f^*(y)^2 \bar{k}(x,y) k(x,y) \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) \\ &+ \int_{\mathcal{X}} f^*(y)^2 \bar{k}(x,y)^2 \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) \end{split}$$

while, since  $f^*$  is  $\mathcal{G}$ -invariant,  $\mu$  is  $\mathcal{G}$ -invariant and  $\mathcal{G}$  is unimodular (because it is compact),

$$\begin{split} \int_{\mathcal{X}} f^*(y)^2 \bar{k}(x,y) k(x,y) \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) &= \int_{\mathcal{X}} \int_{\mathcal{G}} f^*(gy)^2 \, \mathrm{d}\lambda(g) \bar{k}(x,y) k(x,y) \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) \\ &= \int_{\mathcal{X}} \int_{\mathcal{G}} f^*(gy)^2 \, \mathrm{d}\lambda(g) \bar{k}(x,y) k(x,y) \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) \\ &= \int_{\mathcal{X}} f^*(y)^2 \int_{\mathcal{G}} \bar{k}(x,gy) k(x,gy) \, \mathrm{d}\lambda(g) \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) \\ &= \int_{\mathcal{X}} f^*(y)^2 \bar{k}(x,y)^2 \, \mathrm{d}\mu(x) \, \mathrm{d}\mu(y) \end{split}$$

where the final line follows because  $\bar{k}$  is  $\mathcal{G}$ -invariant.

For intuition, we present a simple special case of Theorem 5. In particular, the next result shows that Eq. (5) reduces to an approximation error that is reminiscent of the one in [8, Theorem 6] in a linear setting. For the rest of this section we find it helpful to refer to the action  $\phi$  of  $\mathcal{G}$  explicitly, writing  $\phi(g)x$  instead of gx.

**Theorem 7.** Assume the setting and notation of Theorem 5. In addition, let  $\mathcal{X} = \mathbb{S}_{d-1}$  be the unit d - 1 sphere and let  $\mu = \text{Unif}(\mathcal{X})$ . Let  $\mathcal{G}$  act via an orthogonal representation  $\phi$  on  $\mathcal{X}$  and define the matrix  $\Phi = \int_{\mathcal{G}} \phi(g) \, d\lambda(g)$ . Let  $k(x, y) = x^{\top} y$  be the linear kernel and suppose  $f^*(x) = \theta^{\top} x$  for some  $\theta \in \mathbb{R}^d$ . Then the bound in Theorem 5 becomes

$$\mathbb{E}[\Delta(f,\bar{f})] \ge \frac{1}{(\sqrt{n}+\rho/\sqrt{n})^2} \left(\frac{d-\|\Phi\|_{\mathrm{F}}^2}{d^2} + \frac{(d-\|\Phi\|_{\mathrm{F}}^2)\|\theta\|_2^2}{d^2(d+2)}\right)$$

where  $\|\cdot\|_{\mathrm{F}}$  is the Frobenius norm. The first term in the parenthesis is exactly  $\dim_{\mathrm{eff}}(\mathcal{H}_A)$  and the second term is exactly  $\mathbb{E}[f^*(X)^2 j^{\perp}(X,X)]$ .

Proof. We will make use of the Einstein convention of summing repeated indices. Since  $\mu$  is finite, by Fubini's theorem we are free to integrate in any order throughout the proof. First of all notice that sup<sub>x</sub> k(x, x) = 1 so  $M_k = 1$ . Now observe that

$$\bar{k}(x,y) = x^{\top} \int_{\mathcal{G}} \phi(g) y \, \mathrm{d}\lambda(g) = x^{\top} \Phi y.$$

305 Then the first term in the numerator becomes

$$\begin{aligned} \dim_{\text{eff}}(\mathcal{H}_A) &= \mathbb{E}[j^{\perp}(X,X)] \\ &= \mathbb{E}[j(X,X)] - \mathbb{E}[\bar{j}(X,X)] \\ &= \int_X k(x,y)^2 \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y) - \int_X \bar{k}(x,y)^2 \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y) \\ &= \int_{\mathcal{X}} x_a x_b y_a y_b \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y) - \int_X x_a x_b y_c y_e \Phi_{ac} \Phi_{be} \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y) \\ &= \frac{1}{d} - \frac{1}{d^2} \|\Phi\|_{\text{F}}^2. \end{aligned}$$

Now for the second term. We calculate each term of the right hand side of Proposition 6 separately.
 We know that

$$f^*(x)^2 k(x,y)^2 = (\theta^\top x)^2 (x^\top y)^2 = \theta_a \theta_b y_c y_e x_a x_b x_c x_e.$$

308 Integrating y first, we get

$$\int_{\mathcal{X}} f^*(x)^2 k(x,y)^2 d\mu(x) d\mu(y) = \int_{\mathcal{X}} \theta_a \theta_b y_c y_e x_a x_b x_c x_e d\mu(x) d\mu(y)$$
$$= \frac{1}{d} \int_{\mathcal{X}} \theta_a \theta_b x_a x_b d\mu(x)$$
$$= \frac{1}{d^2} \|\theta\|_2^2$$

309 Similarly, we find

$$\int_{\mathcal{X}} f^*(x)^2 \bar{k}(x,y)^2 d\mu(x) d\mu(y) = \int_{\mathcal{X}} \theta_a \theta_b x_a x_b x_c x_e y_f y_h \Phi_{cf} \Phi_{eh} d\mu(x) d\mu(y)$$
$$= \frac{1}{d} \theta_a \theta_b \Phi_{cf} \Phi_{ef} \int_{\mathcal{X}} x_a x_b x_c x_e d\mu(x).$$

The 4-tensor  $\int_{\mathcal{X}} x_a x_b x_c x_e d\mu(x)$  is isotropic, so must have the form

$$\int_{\mathcal{X}} x_a x_b x_c x_e \, \mathrm{d}\mu(x) = \alpha \delta_{ab} \delta_{ce} + \beta \delta_{ac} \delta_{be} + \gamma \delta_{ae} \delta_{bc}$$

(see, e.g. Hodge [11]). By symmetry and exchangeability we have  $\alpha = \beta = \gamma$ . Then contracting the first two indices gives

$$\int_{\mathcal{X}} x_a x_a x_c x_e \, \mathrm{d}\mu(x) = \frac{1}{d} \delta_{ce} = \alpha (d+2) \delta_{ce}$$

so  $\alpha = \frac{1}{d(d+2)}$  and we end up with

$$\int_{\mathcal{X}} f^*(x)^2 \bar{k}(x,y)^2 \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y) = \frac{\|\theta\|_2^2 \|\Phi\|_F^2 + 2\|\Phi\theta\|_2^2}{d^2(d+2)} = \frac{\|\theta\|_2^2 (\|\Phi\|_F^2 + 2)}{d^2(d+2)}$$

314 where the second equality comes from

$$\theta^{\top} \Phi x = \int_{\mathcal{G}} \theta^{\top} \phi(g) x \, \mathrm{d}\lambda(g) = \int_{\mathcal{G}} f^*(\phi(g)x) \, \mathrm{d}\lambda(g) = f^*(x) = \theta^{\top} x$$

for any  $x \in \mathcal{X}$ . Putting everything together gives the result.

One can confirm that the generalisation gap cannot be negative in Theorem 7 using Jensen's inequality

$$\|\Phi\|_{\mathrm{F}}^{2} = \left\|\int_{\mathcal{G}} \phi(g) \,\mathrm{d}\lambda(g)\right\|_{\mathrm{F}}^{2} \leq \int_{\mathcal{G}} \|\phi(g)\|_{\mathrm{F}}^{2} \,\mathrm{d}\lambda(g) = \int_{\mathcal{G}} \mathrm{Tr}(\phi(g)^{\top} \phi(g)) \,\mathrm{d}\lambda(g) = \mathrm{Tr}(I) = d$$

because the representation  $\phi$  is orthgonal.

The matrix  $\Phi$  in Theorem 7 can be computed analytically for various  $\mathcal{G}$  and in the linear setting describes the importance of the symmetry to the task. For instance, in the simple case that  $\mathcal{G} = S_d$  the permutation group on d elements and  $\phi$  is the natural representation in terms of permutation matrices, we have  $\Phi = \frac{1}{d} \mathbf{1} \mathbf{1}^{\top}$  where  $\mathbf{1} \in \mathbb{R}^d$  is the vector of all 1s. In this case, since the target is assumed to be  $\mathcal{G}$ -invariant, we must have  $\theta = t\mathbf{1}$  for some  $t \in \mathbb{R}$ . Specifically, Theorem 7 then asserts

$$\mathbb{E}[\Delta(f,\bar{f})] \ge \frac{(d-1)(dt^2 + d + 2)}{d^2(d+2)(\sqrt{n} + \rho/\sqrt{n})^2}$$

# 323 5 Related Work

Incorporating invariance into machine learning models is not a new idea. The majority of modern applications concern neural networks, but earlier work has used kernels [10], support vector machines [26] and polynomial feature spaces [27, 28]. Indeed, early work also considered invariant neural networks [33], using methods that seem to have been rediscovered in [24]. Modern implementations include invariant/equivariant convolutional architectures [4, 6] that are inspired by concepts from mathematical physics and harmonic analysis [14, 5]. Some of these models even enjoy universal approximation properties [20, 35].

The earliest attempt at theoretical justification for invariance of which we are aware is [1], which 331 roughly states that enforcing invariance cannot increase the VC dimension of a model. Anselmi 332 et al. [2] and Mroueh, Voinea, and Poggio [23] propose heuristic arguments for improved sample 333 complexity of invariant models. Sokolic et al. [29] build on the work of Xu and Mannor [34] to obtain 334 335 a generalisation bound for certain types of classifiers that are invariant to a finite set of transformations, while Sannai and Imaizumi [25] obtain a bound for models that are invariant to finite permutation 336 337 groups. The PAC Bayes formulation is considered in [17, 18]. The above works guarantee only a worst-case improvement and it was not until very recently 338 that Elesedy and Zaidi [8] derived a strict benefit for invariant/equivariant models. Our work is similar 339

to [8] in that we provide a provably strict benefit, but differs in its application to kernels and RKHSs 340 as opposed to linear models. We are careful to state that our setting does not directly reduce to that 341 of [8, Theorem 6] for two reasons. First, [8, Theorem 6] considers  $\mathcal{G}$  invariant linear models without 342 regularisation. This may turn out to be accessible by a  $\rho \rightarrow 0^+$  limit (the so called ridgeless limit) 343 of Theorem 5. More importantly, linear regression is equivalent to kernel regression with the linear 344 kernel. However, the linear kernel can be unbounded (e.g. on  $\mathbb{R}$ ), so does not meet our technical 345 conditions in Section 2.2. We conjecture that the boundedness assumption on k can be removed, or at 346 least with mild care weakened to hold  $\mu$ -almost-surely. 347

Also very recently, Mei, Misiakiewicz, and Montanari [21] analyse the generalisation benefit of 348 invariance in kernels and random feature models. Our results differ from [21] in some key aspects. 349 First, Mei, Misiakiewicz, and Montanari [21] focus kernel ridge regression with an invariant inner 350 product kernel whereas we study symmetrised predictors from more general kernels. Second, they 351 obtain an expression for the generalisation error that is conditional on the training data and in terms of 352 the projection of the predictor onto a space of high degree polynomials, while we are able to integrate 353 against the training data and express the generalisation benefit directly in terms of properties of the 354 kernel and the group. 355

## 356 6 Discussion

We have demonstrated a provably strict generalisation benefit for feature averaging in kernel ridge regression. In doing this we have leveraged an observation on the structure of RKHSs under the action of compact groups. We believe that this observation is applicable to other kernel methods too.

There are many possibilities for future work. As we remarked in the introduction, there is an established connection between kernels and wide neural networks via the neural tangent kernel. Using this connection, generalisation properties of wide, invariant neural networks might be accessible through the techniques of this paper. Another natural extension of this paper is to equivariant (sometimes called *steerable*) matrix valued kernels. Finally, the ideas of this paper should also be applicable to Gaussian processes.

## **366** A Notation and Definitions

Trace of a linear operator  $A: V \to V$  on a inner product space V is defined by

$$\operatorname{Tr}(A) = \sum_{i} \langle Av_i, v_i \rangle$$

where the collection  $\{v_i\}$  forms an orthonormal basis of V. In this paper we will only encounter situations in which the basis is countable. This expressions is independent of the basis. We say A is *trace-class* if  $Tr(A) < \infty$ .

For any matrix  $A \in \mathbb{R}^{n \times n}$ , we define  $||A||_2 = \sup_{x \in \mathbb{R}^n} \frac{||Ax||_2}{||x||}$  which is the operator norm induced by the Euclidean norm. For any symmetric matrix A, we denote by  $\lambda_{\max}(A)$  and  $\lambda_{\min}(A)$  the largest and smallest eigenvalues of A respectively.

# 374 **B** Useful Results

<sup>375</sup> This section contains some results that are relied upon elsewhere in the paper.

**Lemma B.1** (Mori [22]). Let  $A, B \in \mathbb{R}^{n \times n}$  and suppose B is symmetric. Define  $A' = \frac{1}{2}(A + A^{\top})$ , then

$$\lambda_{\min}(A')\operatorname{Tr}(B) \leq \operatorname{Tr}(AB) \leq \lambda_{\max}(A')\operatorname{Tr}(B),$$

- where  $\lambda_{\min}$  and  $\lambda_{\max}$  denote the smallest and largest eigenvalues respectively.
- 379 **Corollary B.2.** Let  $A, B \in \mathbb{R}^{n \times n}$  and suppose A is symmetric, then

 $\lambda_{\min}(A)\operatorname{Tr}(B) \leq \operatorname{Tr}(AB) \leq \lambda_{\max}(A)\operatorname{Tr}(B).$ 

<sup>380</sup> Proof. Let  $B' = \frac{1}{2}(B + B^{\top})$ , then using Lemma B.1 we have  $\lambda_{\min}(A) \operatorname{Tr}(B') \leq \operatorname{Tr}(AB') \leq \lambda_{\max}(A) \operatorname{Tr}(B').$ <sup>381</sup> On the other hand,  $\operatorname{Tr}(B') = \operatorname{Tr}(B)$  and  $2 \operatorname{Tr}(AB') = \operatorname{Tr}(AB) + \operatorname{Tr}(AB^{\top}) = \operatorname{Tr}(AB) + \operatorname{Tr}(BA).$ 

382

383 **Lemma B.3.** Let  $A \in \mathbb{R}^{n \times n}$ , then

$$\|A\|_{2} \le n \max_{ij} |A_{ij}|.$$

384 *Proof.* Let  $a_i \in \mathbb{R}^n$  be the *i*<sup>th</sup> column of A, then

$$\sup_{\|x\|_{2}=1} \|Ax\|_{2} = \sup_{\|x\|_{2}=1} \sqrt{\sum_{i} (a_{i}^{\top}x)^{2}} \le \sup_{\|x\|_{2}=1} \sqrt{\sum_{i} \|a_{i}\|_{2}^{2}} \|x\|_{2}^{2} \le \sqrt{\sum_{i} \|a_{i}\|_{2}^{2}} \le \sqrt{n^{2} \max_{ij} A_{ij}^{2}}.$$

385

# **386 C Results leading to Lemma 3**

Recall from Section 3 the integral operator  $S_k : L_2(\mathcal{X}, \mu) \to \mathcal{H}$  defined by

$$S_k f(x) = \int_{\mathcal{X}} k(x, y) f(y) \, \mathrm{d}\mu(y)$$

- with adjoint  $\iota : L_2(\mathcal{X}, \mu) \to \mathcal{H}$ .
- **Lemma C.1.** The image of  $L_2(\mathcal{X}, \mu)$  under  $S_k$  is dense in  $\mathcal{H}$  and  $\iota$  is injective.

Proof. By [31, Theorem 4.26]  $||f||_{L_2(\mathcal{X},\mu)} < \infty \forall f \in \mathcal{H}$  and  $S_k(L_2(\mathcal{X},\mu))$  is dense in  $\mathcal{H}$  if and only if the inclusion  $\iota : \mathcal{H} \to L_2(\mathcal{X},\mu)$  is injective. Injectivity of the inclusion is equivalent to the statement that for any  $f, f' \in \mathcal{H}$  the set

$$A(f, f') = \{x \in \mathcal{X} : f(x) \neq f'(x)\}$$

has  $A \neq \emptyset \implies \mu(A) > 0$ . Continuity implies that for any  $f, f' \in \mathcal{H}$ , either f = f' pointwise or A(f, f') contains an open set. By the support of  $\mu$  this implies  $\mu(A) > 0$ . Thus,  $\iota$  is injective.  $\Box$ 

From [8, Proposition 22] we know that  $\mathcal{O}: L_2(\mathcal{X}, \mu) \to L_2(\mathcal{X}, \mu)$  is well-defined and that  $\|\mathcal{O}\| \leq 1$ . Let the image of  $L_2(\mathcal{X}, \mu)$  under  $S_k$  be  $\mathcal{H}_2$ , then Lemma C.1 states that  $\overline{\mathcal{H}_2} = \mathcal{H}$ .

Lemma C.2. For any  $f \in L_2(\mathcal{X}, \mu)$ ,  $\mathcal{O}S_k f = S_k \mathcal{O}f \in \mathcal{H}_2$ . This implies  $\mathcal{O} : \mathcal{H}_2 \to \mathcal{H}_2$  is well defined.

*Proof.*  $\lambda$  is a Radon measure [13, Theorem 2.27] so is finite because  $\mathcal{G}$  is compact and all  $f \in \mathcal{H}$  are bounded so we can apply Fubini's theorem [13, Theorem 1.27] as follows: taking  $f \in L_2(\mathcal{X}, \mu)$ 

$$\begin{split} S_k \mathcal{O}f(x) &= \int_{\mathcal{X}} \int_{\mathcal{G}} k(x,y) f(gy) \, \mathrm{d}\lambda(g) \, \mathrm{d}\mu(y) \\ &= \int_{\mathcal{X}} \int_{\mathcal{G}} k(x,g^{-1}y) f(y) \, \mathrm{d}\lambda(g) \, \mathrm{d}\mu(y) \quad \text{invariance of } \mu \\ &= \int_{\mathcal{X}} \int_{\mathcal{G}} k(gx,y) \, \mathrm{d}\lambda(g) f(y) \, \mathrm{d}\mu(y) \quad \text{Eq. (1) then unimodularity of } \mathcal{G} \\ &= \int_{\mathcal{G}} \int_{\mathcal{X}} k(gx,y) f(y) \, \mathrm{d}\mu(y) \, \mathrm{d}\lambda(g) \quad \text{Fubini} \\ &= \mathcal{O}S_k f(x). \end{split}$$

401

Lemma C.3. Let  $a, b \in \mathcal{H}_2$  with preimages  $a', b' \in L_2(\mathcal{X}, \mu)$  such that  $a = S_k a'$  and  $b = S_k b'$ , then

$$\langle a, b \rangle_{\mathcal{H}} = \int_{\mathcal{X}} a'(x)b'(y)k(x,y) \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y).$$

404 *Proof.* The inner product on  $\mathcal{H}$  is a bounded linear functional, hence commutes with integration. We 405 can thus calculate

$$\begin{aligned} \langle a,b \rangle_{\mathcal{H}} &= \langle \int_{\mathcal{X}} a'(x)k(x,\cdot) \,\mathrm{d}\mu(x), \int_{\mathcal{X}} b'(y)k(y,\cdot) \,\mathrm{d}\mu(y) \rangle_{\mathcal{H}} \\ &= \int_{\mathcal{X}} a'(x)b'(y)\langle k_x, k_y \rangle_{\mathcal{H}} \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y) \\ &= \int_{\mathcal{X}} a'(x)b'(y)k(x,y) \,\mathrm{d}\mu(x) \,\mathrm{d}\mu(y). \end{aligned}$$

406

407 **Lemma C.4.** For any  $f, h \in \mathcal{H}_2$ ,

$$\langle \mathcal{O}f,h\rangle_{\mathcal{H}} = \langle f,\mathcal{O}h\rangle_{\mathcal{H}}.$$

<sup>408</sup> *Proof.* Let f' and h' be the pre-images of f and h respectively under  $S_k$ . Using Lemma C.3, Fubini's <sup>409</sup> theorem [13, Theorem 1.27], the  $\mathcal{G}$ -invariance of  $\mu$  and Eq. (1) we can calculate

$$\begin{split} \langle \mathcal{O}f,h\rangle_{\mathcal{H}} &= \int_{\mathcal{X}} \int_{\mathcal{G}} f'(gx)h'(y)k(x,y)\,\mathrm{d}\lambda(g)\,\mathrm{d}\mu(x)\,\mathrm{d}\mu(y) \\ &= \int_{\mathcal{G}} \int_{\mathcal{X}} f'(x)h'(y)k(g^{-1}x,y)\,\mathrm{d}\mu(x)\,\mathrm{d}\mu(y)\,\mathrm{d}\lambda(g) \quad \mathcal{G}\text{-invariance of }\mu \\ &= \int_{\mathcal{G}} \int_{\mathcal{X}} f'(x)h'(y)k(x,g^{-1}y)\,\mathrm{d}\mu(x)\,\mathrm{d}\mu(y)\,\mathrm{d}\lambda(g) \quad \text{Eq. (1)} \\ &= \int_{\mathcal{G}} \int_{\mathcal{X}} f'(x)h'(gy)k(x,y)\,\mathrm{d}\mu(x)\,\mathrm{d}\mu(y)\,\mathrm{d}\lambda(g) \quad \mathcal{G}\text{-invariance of }\mu \\ &= \langle f,\mathcal{O}h\rangle_{\mathcal{H}}. \end{split}$$

410

411 Lemma C.5.  $\mathcal{O}: \mathcal{H}_2 \to \mathcal{H}_2$  is bounded and  $\|\mathcal{O}\| \leq 1$ .

412 *Proof.* Let  $f \in \mathcal{H}_2$ , then using Lemmas 1 and C.4 along with Cauchy-Schwarz

$$\|\mathcal{O}f\|_{\mathcal{H}}^2 = \langle \mathcal{O}f, \mathcal{O}f \rangle_{\mathcal{H}} = \langle f, \mathcal{O}f \rangle_{\mathcal{H}} \le \|f\|_{\mathcal{H}} \|\mathcal{O}f\|_{\mathcal{H}}.$$

413

414 **Lemma C.6.**  $f \in \mathcal{H} \implies \mathcal{O}f \in \mathcal{H}$  so  $\mathcal{O} : \mathcal{H} \to \mathcal{H}$  is well defined.

415 *Proof.* By Lemma C.1, for any  $f \in \mathcal{H}$  there is a sequence  $\{f_n\} \subset \mathcal{H}_2$  converging to f in  $\|\cdot\|_{\mathcal{H}}$ . 416 Lemma C.2 shows that  $\mathcal{O} : \mathcal{H}_2 \to \mathcal{H}_2$  is well defined, so the sequence  $\{\mathcal{O}f_n\} \subset \mathcal{H}_2$ . By Lemma C.5 417 we have  $\|\mathcal{O}f_n - \mathcal{O}f_m\|_{\mathcal{H}} \leq \|f_n - f_m\|_{\mathcal{H}}$  and so  $\{\mathcal{O}f_n\}$  is Cauchy. By completeness of  $\mathcal{H}, \bar{f} :=$ 418  $\lim_{n\to\infty} \mathcal{O}f_n \in \mathcal{H}$ . Moreover,  $\mathcal{O}$  bounded so is also continuous and we get  $\bar{f} = \lim_{n\to\infty} \mathcal{O}f_n =$ 419  $\mathcal{O}\lim_{n\to\infty} f_n = \mathcal{O}f$ . □

420 **Lemma C.7.**  $\mathcal{O}$  is self-adjoint with respect to the inner product on  $\mathcal{H}$ .

*Proof.* We will make use of the continuity of the inner product on  $\mathcal{H}$ . First let  $h \in \mathcal{H}$ ,  $f \in \mathcal{H}_2$ . We saw from the proof of Lemma C.6 that  $\exists$  sequence  $\{h_n\} \subset \mathcal{H}_2$  with limit h and  $\{\mathcal{O}h_n\} \subset \mathcal{H}_2$  with limit  $\mathcal{O}h$ . Then  $\langle \mathcal{O}h_n, f \rangle_{\mathcal{H}} \to \langle \mathcal{O}h, f \rangle_{\mathcal{H}}$  and simultaneously, applying Lemma C.4,  $\langle \mathcal{O}h_n, f \rangle_{\mathcal{H}} =$  $\langle h_n, \mathcal{O}f \rangle_{\mathcal{H}} \to \langle h, \mathcal{O}f \rangle_{\mathcal{H}}$  so the two limits must be equal. Then assuming instead that  $f \in \mathcal{H}$  one can do the same calculation again arrive at the conclusion. 426 **Corollary C.8.**  $\mathcal{O} : \mathcal{H} \to \mathcal{H}$  is bounded with  $\|\mathcal{O}\| \leq 1$ . Indeed, if  $\mathcal{H}$  contains any  $\mathcal{G}$ -invariant 427 functions then  $\|\mathcal{O}\| = 1$  and if not then  $\|\mathcal{O}\| = 0$ .

Proof. Using Lemma C.7 we can repeat the calculation in Lemma C.4. The second claim follows from Lemma 1 and the variational representation of the operator norm.  $\Box$ 

## 430 **References**

- [1] Yaser S Abu-Mostafa. "Hints and the VC dimension". In: *Neural Computation* 5.2 (1993),
   pp. 278–288 (page 11).
- Fabio Anselmi et al. Unsupervised Learning of Invariant Representations in Hierarchical
   Architectures. 2014. arXiv: 1311.4158 [cs.CV] (page 11).
- [3] Sanjeev Arora et al. "Fine-Grained Analysis of Optimization and Generalization for
   Overparameterized Two-Layer Neural Networks". In: *Proceedings of the 36th International Conference on Machine Learning*. Ed. by Kamalika Chaudhuri and Ruslan Salakhutdinov.
   Vol. 97. Proceedings of Machine Learning Research. PMLR, 2019, pp. 322–332. URL: http:
   //proceedings.mlr.press/v97/arora19a.html (page 2).
- [4] Taco Cohen and Max Welling. "Group equivariant convolutional networks". In: *International conference on machine learning*. 2016, pp. 2990–2999 (page 10).
- Taco S Cohen, Mario Geiger, and Maurice Weiler. "A general theory of equivariant cnns
   on homogeneous spaces". In: *Advances in Neural Information Processing Systems*. 2019,
   pp. 9145–9156 (page 10).
- [6] Taco S Cohen et al. "Spherical cnns". In: *arXiv preprint arXiv:1801.10130* (2018) (page 10).
- [7] Cohn, Donald L. *Measure Theory*. 2nd ed. Springer, 2013 (page 5).
- [8] Bryn Elesedy and Sheheryar Zaidi. "Provably Strict Generalisation Benefit for Equivariant Models". In: (2021). arXiv: 2102.10333 [stat.ML] (pages 1–4, 7–9, 11, 12).
- [9] Adam Foster, Rattana Pukdee, and Tom Rainforth. "Improving Transformation Invariance in Contrastive Representation Learning". In: *arXiv preprint arXiv:2010.09515* (2020) (page 1).
- [10] B. Haasdonk, A. Vossen, and H. Burkhardt. "Invariance in Kernel Methods by Haar Integration Kernels". In: *SCIA 2005, Scandinavian Conference on Image Analysis*. Springer-Verlag, 2005, pp. 841–851 (pages 4, 10).
- [11] Philip G. Hodge. "On Isotropic Cartesian Tensors". eng. In: *The American Mathematical Monthly* 68.8 (1961), pp. 793–795. ISSN: 00029890 (page 10).
- [12] Arthur Jacot, Franck Gabriel, and Clément Hongler. "Neural tangent kernel: Convergence and
   generalization in neural networks". In: *Advances in neural information processing systems*.
   2018, pp. 8571–8580 (page 2).
- [13] Olav Kallenberg. *Foundations of modern probability*. Springer Science & Business Media,
   2006 (pages 3, 12, 13).
- [14] Risi Kondor and Shubhendu Trivedi. "On the Generalization of Equivariance and Convolution in Neural Networks to the Action of Compact Groups". In: *International Conference on Machine Learning*. 2018, pp. 2747–2755 (page 10).
- 464 [15] Jaehoon Lee et al. "Wide neural networks of any depth evolve as linear models under gradient
  465 descent". In: Advances in neural information processing systems. 2019, pp. 8572–8583
  466 (page 2).
- International Conference on Machine Learning.
   Kamalika Chaudhuri and Ruslan Salakhutdinov. Vol. 97. Proceedings of Machine Learning Research. PMLR, 2019, pp. 3744–3753. URL: http://proceedings.mlr.press/
- 471 v97/lee19d.html (page 1).
- [17] Clare Lyle, Marta Kwiatkowksa, and Yarin Gal. "An analysis of the effect of invariance on generalization in neural networks". In: *International conference on machine learning Workshop* on Understanding and Improving Generalization in Deep Learning. 2019 (page 11).
- [18] Clare Lyle et al. On the Benefits of Invariance in Neural Networks. 2020. arXiv: 2005.00178
   [cs.LG] (pages 1, 11).
- In: *arXiv preprint arXiv:1408.0952* (2014) (page 6).
- Haggai Maron et al. "On the Universality of Invariant Networks". In: *International Conference* on Machine Learning. 2019, pp. 4363–4371 (page 10).

- 481 [21] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. "Learning with invariances in 482 random features and kernel models". In: *arXiv preprint arXiv:2102.13219* (2021) (pages 1, 483 11).
- T. Mori. "Comments on "A matrix inequality associated with bounds on solutions of algebraic
   Riccati and Lyapunov equation" by J.M. Saniuk and I.B. Rhodes". In: *IEEE Transactions on Automatic Control* 33.11 (1988), pp. 1088–. DOI: 10.1109/9.14428 (page 11).
- Youssef Mroueh, Stephen Voinea, and Tomaso A Poggio. "Learning with Group Invariant
   Features: A Kernel Perspective." In: *Advances in Neural Information Processing Systems*.
   2015, pp. 1558–1566 (page 11).
- [24] Siamak Ravanbakhsh, Jeff Schneider, and Barnabas Poczos. "Equivariance through parametersharing". In: *International Conference on Machine Learning*. PMLR. 2017, pp. 2892–2901 (page 10).
- 493 [25] Akiyoshi Sannai and Masaaki Imaizumi. Improved Generalization Bound of Group Invariant
   494 / Equivariant Deep Networks via Quotient Feature Space. 2019. arXiv: 1910.06552
   495 [stat.ML] (pages 1, 11).
- <sup>496</sup> [26] Bernhard Schölkopf, Chris Burges, and Vladimir Vapnik. "Incorporating Invariances in Support Vector Learning Machines". In: Springer, 1996, pp. 47–52 (page 10).
- 498 [27] H. Schulz-Mirbach. "Constructing invariant features by averaging techniques". In: *Proceedings* 499 of the 12th IAPR International Conference on Pattern Recognition, Vol. 3 - Conference C: 500 Signal Processing (Cat. No.94CH3440-5). Vol. 2. 1994, 387–390 vol.2 (page 10).
- [28] Hanns Schulz-Mirbach. "On the existence of complete invariant feature spaces in pattern recognition". In: *International Conference On Pattern Recognition*. Citeseer. 1992, pp. 178– 178 (page 10).
- <sup>504</sup> [29] Jure Sokolic et al. "Generalization error of invariant classifiers". In: *Artificial Intelligence and* 505 *Statistics*. 2017, pp. 1094–1103 (pages 1, 11).
- <sup>506</sup> [30] James S Spencer et al. "Better, Faster Fermionic Neural Networks". In: *arXiv preprint* <sup>507</sup> *arXiv:2011.07125* (2020) (page 1).
- <sup>508</sup> [31] Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Information science and statistics. Springer, 2008. ISBN: 978-0-387-77241-7 (pages 2, 3, 5, 12).
- [32] Marysia Winkels and Taco S Cohen. "3D G-CNNs for pulmonary nodule detection". In: *arXiv preprint arXiv:1804.04656* (2018) (page 1).
- Jeffrey Wood and John Shawe-Taylor. "Representation theory and invariant neural networks".
   In: *Discrete applied mathematics* 69.1-2 (1996), pp. 33–60 (pages 1, 10).
- <sup>514</sup> [34] Huan Xu and Shie Mannor. "Robustness and generalization". In: *Machine learning* 86.3 (2012), pp. 391–423 (page 11).
- [35] Dmitry Yarotsky. Universal approximations of invariant maps by neural networks. 2018. arXiv: 1804.10306 [cs.NE] (page 10).
- [36] Manzil Zaheer et al. "Deep sets". In: *Advances in neural information processing systems*. 2017, pp. 3391–3401 (page 1).

## 520 Checklist

522

523

527

528

- 521 1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
- (b) Did you describe the limitations of your work? [Yes] See Section 2.2 which describes
   our assumptions.
- 526 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 529 2. If you are including theoretical results...
- (a) Did you state the full set of assumptions of all theoretical results? [Yes] See ?? for
   general technical conditions not given in statements of results.
- (b) Did you include complete proofs of all theoretical results? [Yes] Proofs are given in
   the supplementary material

534	3. If you ran experiments
535 536	(a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]
537 538	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
539 540	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
541 542	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
543	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
544	(a) If your work uses existing assets, did you cite the creators? $[N/A]$
545	(b) Did you mention the license of the assets? [N/A]
546 547	(c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
548 549	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
550 551	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
552	5. If you used crowdsourcing or conducted research with human subjects
553 554	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
555 556	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
557 558	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]