
Provably Strict Generalisation Benefit for Invariance in Kernel Methods

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

12 1 Introduction

13 Recently, there has been significant interest in models that are invariant to the action of a group
14 on their inputs. It is believed that engineering models in this way improves sample efficiency and
15 generalisation. Intuitively, if a task has an invariance, then a model that is constructed to be invariant
16 ahead of time should require fewer examples to generalise than one that must learn to be invariant.
17 Indeed, there are many application domains, such as fundamental physics or medical imaging, in
18 which the invariance is known a priori [28, 30]. Although this intuition is certainly not new (e.g. [31]),
19 it has inspired much recent work (for instance, see [34, 15]).

20 However, while implementations and practical applications abound, until very recently a rigorous
21 theoretical justification for invariance was missing. As pointed out in [8], many prior works such
22 as [27, 23] provide only worst-case guarantees on the performance of invariant algorithms. It follows
23 that these results do not rule out the possibility of modern training algorithms automatically favouring
24 invariant models, irrespective of the choice of architecture. Steps towards a more concrete theory of
25 the benefit of invariance have been taken by [8, 20] and our work is a continuation along the path set
26 by [8].

27 In this work we provide a precise characterisation of the generalisation benefit of invariance in
28 kernel ridge regression. In contrast to [27, 23], this proves a *provably strict* generalisation benefit for
29 invariant, feature-averaged models. In deriving this result, we provide insights into the structure of
30 reproducing kernel Hilbert spaces in relation to invariant functions that we believe will be useful for
31 analysing invariance in other kernel algorithms.

32 The use of feature averaging to produce invariant predictors enjoys both theoretical and practical
33 success [17, 9]. For the purposes of this work, feature averaging is defined as training a model
34 as normal (according to any algorithm) and then transforming the learned model to be invariant.
35 This transformation is done by *orbit-averaging*, which means projecting the model on the space of
36 invariant functions using the operator \mathcal{O} introduced in Section 2.3.

Kernel methods have a long been a mainstay of machine learning (see [29, 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 tangent kernel* [11], which states that the functions realised by an infinitely wide neural network belong to an RKHS with a kernel determined by the network architecture. This relation has led to many results on the theory of optimisation and generalisation of wide neural networks (e.g. [14, 3]).

1.1 Summary of Contributions

This paper builds towards Theorem 5 in Section 4, which gives a precise characterisation of the benefit of incorporating invariance in kernel methods by feature averaging. We find that for a predictor from a reproducing kernel Hilbert space (RKHS) \mathcal{H} with kernel k , the benefit is $O(\dim_{\text{eff}}(\mathcal{H}_A)/n)$ where $\dim_{\text{eff}}(\mathcal{H}_A)$ is a notion of effective dimension of an RKHS $\mathcal{H}_A \subset \mathcal{H}$ that arises from the interaction between the kernel k and the group. Lemma 3, given in Section 3, forms the basis of Theorem 5 and shows that \mathcal{H} decomposes into an orthogonal direct sum $\mathcal{H} = \mathcal{H}_S \oplus \mathcal{H}_A$, where \mathcal{H}_S is an RKHS consisting of all of the invariant functions in \mathcal{H} . We stress that while Theorem 5 is specialised to kernel ridge regression, Lemma 3 holds regardless of training algorithm and could be used to explore invariance in other kernel methods. In Section 2 we outline our assumptions and the ideas from [8] on which we build. We discuss related works in Section 5.

2 Background and Preliminaries

In this section we provide a brief introduction to reproducing kernel Hilbert spaces (RKHS) and the ideas we borrow from Elesedy and Zaidi [8]. Throughout this paper, \mathcal{H} will be an RKHS with kernel k . In Section 2.2 we state some topological and measurability assumptions that are needed for our proofs. These conditions are benign, and the reader not interested in technicalities need take from Section 2.2 only that μ is \mathcal{G} -invariant and that the kernel k is bounded and satisfies Eq. (1). We defer some background results to Appendix A of the Supplementary Material.

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 $f : \mathcal{X} \rightarrow \mathbb{R}$ on which the evaluation functional $\delta_x : \mathcal{H} \rightarrow \mathbb{R}$ with $\delta_x[f] = f(x)$ is continuous $\forall x \in \mathcal{X}$, or, equivalently is a bounded operator. The Riesz Representation Theorem tells us that there 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 $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ is the inner product on \mathcal{H} . We identify the function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ with $k(x, y) = \langle k_x, k_y \rangle_{\mathcal{H}}$ as the *reproducing kernel* of \mathcal{H} . Using the inner product representation, one can see that k is positive-definite and symmetric. Conversely, the Moore-Aronszajn Theorem shows that for any positive-definite and symmetric kernel k , there is a unique RKHS with reproducing kernel k . In addition, any Hilbert space admitting a reproducing kernel is an RKHS. Finally, another characterisation of \mathcal{H} is as the completion of linear combinations of the form $f_c(x) = \sum_{i=1}^n c_i k(x, x_i)$ for $c_1, \dots, c_n \in \mathbb{R}$ and $x_1, \dots, x_n \in \mathcal{X}$. For (many) more details, see [29, Chapter 4].

2.2 Technical Setup and Assumptions

Let \mathcal{G} be a compact group with Haar measure λ . Let \mathcal{X} be a non-empty Polish space admitting a finite, \mathcal{G} -invariant Borel measure μ , with $\text{supp } \mu = \mathcal{X}$. We normalise $\mu(\mathcal{X}) = \lambda(\mathcal{G}) = 1$, the latter is possible because λ is a Radon measure. We assume that \mathcal{G} has a measurable action on \mathcal{X} that we will write as gx for $g \in \mathcal{G}$, $x \in \mathcal{X}$. A function $f : \mathcal{X} \rightarrow \mathbb{R}$ is \mathcal{G} -invariant if $f(gx) = f(x) \forall x \in \mathcal{X} \forall g \in \mathcal{G}$. Similarly, a measure μ on \mathcal{X} is \mathcal{G} -invariant if $\forall g \in \mathcal{G}$ and any μ -measurable $B \subset \mathcal{X}$ the pushforward of μ by the action of \mathcal{G} equals μ , i.e. $(g_*\mu)(B) = \mu(B)$. This means that if $X \sim \mu$ then $gX \sim \mu \forall g \in \mathcal{G}$.

Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a measurable kernel with RKHS \mathcal{H} such that $k(\cdot, x) : \mathcal{X} \rightarrow \mathbb{R}$ 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

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

Every $f \in \mathcal{H}$ is μ -measurable, bounded and continuous by [29, Lemmas 4.24 and 4.28] and in addition \mathcal{H} is separable using [29, Lemma 4.33]. These conditions allow the application of [29,

89 Theorem 4.26] to relate \mathcal{H} to $L_2(\mathcal{X}, \mu)$ in the proofs building towards Lemma 3. We assume that the
 90 kernel satisfies, for all $x, y \in \mathcal{X}$,

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

91 For this it is sufficient to have $k(gx, y)$ equal to $k(x, gy)$ or $k(x, g^{-1}y)$ (the latter using unimodularity
 92 of \mathcal{G}). Highlighting two special cases: any inner product kernel $k(x, x') = \kappa(\langle x, x' \rangle)$ such that the
 93 action of \mathcal{G} is unitary with respect to $\langle \cdot, \cdot \rangle$ satisfies Eq. (1), as does any stationary kernel $k(x, x') =$
 94 $\kappa(\|x - x'\|)$ with norm that is preserved by \mathcal{G} in the sense that $\|gx - gx'\| = \|x - x'\|$ for any $g \in \mathcal{G}$,
 95 $x, x' \in \mathcal{X}$.

96 2.3 Invariance from a Function Space Perspective

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

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

99 $\mathcal{O}f$ will exist whenever f is μ -measurable. Note that \mathcal{O} is a linear operator and $\mathcal{O}f$ is always \mathcal{G} -
 100 invariant. Interestingly, f is \mathcal{G} -invariant *only* if $f = \mathcal{O}f$. Elesedy and Zaidi [8] use these observations
 101 to characterise invariant functions and study their generalisation properties. In short, this work
 102 extends these insights to kernel methods. Along the way, we will make frequent use of the following
 103 (well known) facts about \mathcal{O} .

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

106 **Lemma 2** ([8, Lemma 1]). $\mathcal{O} : L_2(\mathcal{X}, \mu) \rightarrow L_2(\mathcal{X}, \mu)$ is well-defined and self-adjoint. Hence,
 107 $L_2(\mathcal{X}, \mu)$ has the orthogonal decomposition

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

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

109 In [8], S and A are referred to as the symmetric and anti-symmetric parts of $L_2(\mathcal{X}, \mu)$. We will use
 110 the same terminology.

111 The meaning of Lemma 2 is that any $f \in L_2(\mathcal{X}, \mu)$ has a decomposition $f = \bar{f} + f^\perp$ where
 112 $\bar{f} = \mathcal{O}f$ is \mathcal{G} -invariant and $\mathcal{O}f^\perp = 0$. A noteworthy consequence of this setup, as discussed
 113 in [8], is a provably non-negative generalisation benefit for feature averaging. In particular, for
 114 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 $R(f) =$
 115 $\mathbb{E}_{X \sim \mu}[(f(X) - f^*(X))^2]$ satisfies

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

116 The same holds if the target is corrupted by independent, zero mean (additive) noise.

117 3 Induced Structure of \mathcal{H}

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

123 It may seem at first glance that Lemma 3 should follow immediately from Lemma 2, but this is not the
 124 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,
 125 in contrast with $L_2(\mathcal{X}, \mu)$, an explicit form for the inner product on \mathcal{H} is not immediate, which means
 126 that some work is needed to check that \mathcal{O} is self-adjoint on \mathcal{H} . These are important requirements for
 127 the proofs of both Lemmas 2 and 3 and we establish them, along with \mathcal{O} being continuous on \mathcal{H} , in
 128 the Supplementary Material. The assumption that the kernel satisfies Eq. (1) plays a central role.

129 **Lemma 3.** \mathcal{H} admits an orthogonal decomposition into symmetric and anti-symmetric parts

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

130 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
 131 with kernel

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

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

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

133 Finally, \bar{k} is \mathcal{G} -invariant in both arguments.

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

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

136 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
 137 and $f^\perp = f - \mathcal{O}f \in \mathcal{H}_A$ and these terms are mutually orthogonal.

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

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

141 and the uniqueness afforded by the Reisz representation theorem tells us that the reproducing kernel
 142 for \mathcal{H}_S is $k(x, y) = \int_{\mathcal{G}} k(x, gy) \, d\lambda(g)$. We have $\|\text{id} - \mathcal{O}\| \leq 2$ and we can do the same argument
 143 to show that \mathcal{H}_A is an RKHS with reproducing kernel k^\perp as claimed. Note that one can write
 144 $k^\perp(x, y) = \langle k_x^\perp, k_y^\perp \rangle_{\mathcal{H}}$ so it must be positive-definite. The \mathcal{G} -invariance of $\bar{k}(x, y)$ in both arguments
 145 is immediate from Eq. (1) and Lemma 1. \square

146 As stated earlier, the perspective provided by Lemma 3 will support our analysis of generalisation.
 147 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
 148 \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*
 149 *Kernel*, a concept introduced by Haasdonk, Vossen, and Burkhardt [10]. We will see that a notion
 150 of effective dimension of the RKHS \mathcal{H}_A with kernel k^\perp governs the generalisation gap between an
 151 arbitrary predictor f and its invariant version $\mathcal{O}f$. This effective dimension arises from the spectral
 152 theory of an integral operator related to k , which we develop in the next section.

153 3.1 Spectral Representation and Effective Dimension

154 In this section we consider the spectrum of an integral operator related to the kernel k . This analysis
 155 will ultimately allow us to define a notion of effective dimension of \mathcal{H}_A that we will later see is
 156 important to the generalisation of invariant predictors. While the integral operator setup is standard,
 157 the use of this technique to identify an effective dimension of \mathcal{H}_A is novel.

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

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

159 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
 160 the other hand, every $f \in \mathcal{H}$ is bounded so has $\|f\|_{L_2(\mathcal{X}, \mu)} < \infty$ and belongs to some element
 161 of $L_2(\mathcal{X}, \mu)$. We write $\iota : \mathcal{H} \rightarrow L_2(\mathcal{X}, \mu)$ for the *inclusion map* that sends f to the element of
 162 $L_2(\mathcal{X}, \mu)$ that contains f . In the Supplementary Material we show that ι is injective, so any element
 163 of $L_2(\mathcal{X}, \mu)$ contains at most one $f \in \mathcal{H}$.

164 One can define $T_k : L_2(\mathcal{X}, \mu) \rightarrow L_2(\mathcal{X}, \mu)$ by $T_k = \iota \circ S_k$, and [29, Theorem 4.27] says that T_k is
 165 compact, positive, self-adjoint and trace-class. In addition, $L_2(\mathcal{X}, \mu)$ is separable by [7, Proposition
 166 3.4.5], because \mathcal{X} is Polish and μ is a Borel measure, so has a countable orthonormal basis. Hence,
 167 by the Spectral Theorem, there exists a countable orthonormal basis $\{\tilde{e}_i\}$ for $L_2(\mathcal{X}, \mu)$ such that
 168 $T_k \tilde{e}_i = \lambda_i \tilde{e}_i$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$ are the eigenvalues of T_k . Moreover, since ι is injective, for
 169 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 $S_k \tilde{e}_i = \lambda_i e_i$.

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

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

171 From now on we permit ourself to drop the ι to reduce clutter. We use the above to define

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

172 These quantities will appear again in our analysis of the generalisation of invariant kernel methods.
 173 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 .
 174 Following Eq. (2), one finds the series representations given below in Lemma 4.

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

180 **Lemma 4.** We have

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

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

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

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

183 *Proof.* We show in the Supplementary Material that \mathcal{O} and S_k commute on $L_2(\mathcal{X}, \mu)$ and \mathcal{O} is
 184 self-adjoint on $L_2(\mathcal{X}, \mu)$ by Lemma 1, so \mathcal{O} and ι (the adjoint of S_k by [29, Theorem 4.26]) must
 185 also commute. The first comment is then immediate from the observation that if $a \in \mathcal{H}_S$ and $b \in \mathcal{H}_A$
 186 one has

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

187 We also have both of

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

188 and

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

189 Therefore $\bar{k}_x = \sum_i \lambda_i \bar{e}_i(x) \tilde{e}_i$ 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
 190 the remaining results. \square

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

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

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

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

194 The series converges by the comparison test because $\lambda_i \geq 0$ and $\sum_i \lambda_i = \text{Tr}(T_k) < \infty$. We have
 195 $\dim_{\text{eff}}(\mathcal{H}) = \text{Tr}(T_k^2)$ and we can think of this (very informally) as taking $L_2(\mathcal{X}, \mu)$, pushing it
 196 through \mathcal{H} twice using T_k and then measuring its size. Now because $j = \bar{j} + j^\perp$ we get

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

197 with

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

198 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
 199 twice and measuring the size of the output. In the next section we will consider the generalisation of
 200 kernel ridge regression and find that $\dim_{\text{eff}}(\mathcal{H}_A)$ plays a critical role.

201 **4 Generalisation**

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

205 **4.1 Kernel Ridge Regression**

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

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

208 and $\rho > 0$ is the regularisation parameter. KRR can be thought of as performing ridge regression with
 209 a possibly infinite dimensional feature space \mathcal{H} . The representer theorem tells us that the solution to
 210 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

$$\operatorname{argmin}_{\alpha \in \mathbb{R}^n} \{\|\mathbf{Y} - K\alpha\|_2^2 + \rho \alpha^\top K \alpha\} \quad (4)$$

211 and $\mathbf{Y} \in \mathbb{R}^n$ is the typical row-stacking of the training outputs with $\mathbf{Y}_i = y_i$. K is the typical kernel
 212 Gram matrix with $K_{ij} = k(x_i, x_j)$. We consider solutions of the form¹ $\alpha = (K + \rho I)^{-1} \mathbf{Y}$ which
 213 results in the predictor

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

214 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
 215 generalisation performance of this predictor with that of its averaged version

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

216 To do this we look at the generalisation gap.

217 **4.2 Generalisation Gap**

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

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

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

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

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

226 **4.3 Generalisation Benefit of Feature Averaging**

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

231 **Theorem 5.** Let the training data be $\{(X_i, Y_i) : i = 1, \dots, n\}$ with $Y_i = f^*(X_i) + \xi_i$ where $X \sim \mu$,
 232 $f^* \in L_2(\mathcal{X}, \mu)$ is \mathcal{G} -invariant and $\{\xi_i : i = 1, \dots, n\}$ are independent, with $\mathbb{E}[\xi_i] = 0$ and $\mathbb{E}[\xi_i^2] = 0$.
 233 Let $f = \operatorname{argmin}_{f \in \mathcal{H}} C(f)$ be the solution to Eq. (3) and let $\bar{f} = \mathcal{O}f \in \mathcal{H}_S$ be the result of applying
 234 feature averaging to f , then the generalisation gap with the squared-error loss satisfies

$$\mathbb{E}[\Delta(f, \bar{f})] \geq \frac{\sigma^2 \dim_{\text{eff}}(\mathcal{H}_A) + \sum_{\alpha} \lambda_{\alpha}^2 \langle (f^*)^2, (\tilde{c}_{\alpha}^{\perp})^2 \rangle_{L_2(\mathcal{X}, \mu)}}{(\sqrt{n} M_k + \rho / \sqrt{n})^2}$$

¹When K is a positive definite matrix this will be the *only* solution. If K is singular then $\exists c \in \mathbb{R}^n$ with $\sum_{ij} K_{ij} c_i c_j = \|\sum_i c_i k_{x_i}\|_{\mathcal{H}}^2 = 0$ so $\sum_i c_i k_{x_i}$ is identically 0 and $\forall f \in \mathcal{H}$ we get $\sum_i c_i f(x_i) = 0$ (see [18, Section 4.6.2]). Clearly, this can't happen if \mathcal{H} is sufficiently expressive. In any case, the chosen α is the minimum in Euclidean norm of all possible solutions.

235 where the squares are to be interpreted pointwise as $(h)^2(x) = h(x)^2$ and

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

236 is the *effective dimension* of \mathcal{H}_A .

237 *Proof.* Let J^+ be the Gram matrix with components $J_{ij}^+ = j^+(X_i, X_j)$ let $u \in \mathbb{R}^n$ have components
238 $u_i = f^*(x_i)$. We can use Lemma 2 to get

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

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

$$\begin{aligned} \mathbb{E}[\Delta(f, \bar{f}) | \mathbf{X}] &= \mathbb{E}[(k_{\bar{X}}^{\perp}(\mathbf{X})^{\top} (K + \rho I)^{-1} u)^2 | \mathbf{X}] + \mathbb{E}[(k_{\bar{X}}^{\perp}(\mathbf{X})^{\top} (K + \rho I)^{-1} \xi)^2 | \mathbf{X}] \\ &= u^{\top} (K + \rho I)^{-1} J^+ (K + \rho I)^{-1} u + \sigma^2 \text{Tr}(J^+ (K + \rho I)^{-2}) \end{aligned}$$

240 where the first equality follows because ξ has mean 0 and the second comes from the trace trick.

241 Let $\lambda_{\min}(A)$ be the smallest eigenvalue of a matrix A . Consider the first term

$$\begin{aligned} u^{\top} (K + \rho I)^{-1} J^+ (K + \rho I)^{-1} u &\geq \lambda_{\min}((K + \rho I)^{-1})^2 u^{\top} J^+ u \\ &\geq \frac{1}{(M_k n + \rho)^2} u^{\top} J^+ u \\ &= \frac{1}{(M_k n + \rho)^2} \sum_{ij} \sum_{\alpha} \lambda_{\alpha}^2 f^*(X_i) e_{\alpha}^{\perp}(X_i) e_{\alpha}^{\perp}(X_j) f^*(X_j). \end{aligned}$$

242 where the second inequality follows from $\|A\|_{\text{op}} \leq n \max_{ij} A_{ij}$ for $n \times n$ matrix A and the last
243 line comes from Lemma 4. Now $\mathbb{E}[f^*(X) e_{\alpha}^{\perp}(X)] = \langle \iota f^*, \tilde{e}_{\alpha}^{\perp} \rangle_{L_2(\mathcal{X}, \mu)}$ so the above terms vanish in
244 expectation when $i \neq j$ and (temporarily suspending suspicion regarding the infinite sum) it follows
245 that

$$\mathbb{E}[u^{\top} (K + \rho I)^{-1} J^+ (K + \rho I)^{-1} u] \geq \frac{n}{(nM_k + \rho)^2} \sum_{\alpha} \lambda_{\alpha}^2 \langle \iota(f^*)^2, (\tilde{e}_{\alpha}^{\perp})^2 \rangle_{L_2(\mathcal{X}, \mu)}$$

246 where the squares are interpreted pointwise $(h)^2(x) = h(x)^2$.

247 Interchanging the sum and the expectation above is valid by Fubini's theorem [12, Theorem 1.27]
248 (interpreting the index α as having the counting measure, which is σ -finite) as long as

$$\sum_{\alpha} \lambda_{\alpha}^2 \mathbb{E}[f^*(X_i) e_{\alpha}^{\perp}(X_i) e_{\alpha}^{\perp}(X_j) f^*(X_j)] < \infty \quad (5)$$

249 for all i, j . When $i \neq j$ we have established that this holds, while when $i = j$ we look at each term
250 individually. We have

$$\mathbb{E}[f^*(X)^2 e_{\alpha}^{\perp}(X)^2] \leq \|\tilde{e}_{\alpha}^{\perp}\|_{L_2(\mathcal{X}, \mu)}^2 \text{ess sup}_{x \in \mathcal{X}} f^*(x)^2 \leq \text{ess sup}_{x \in \mathcal{X}} f^*(x)^2 < \infty$$

251 since

$$\|\mathcal{O}\tilde{e}_{\alpha}\|_{L_2(\mathcal{X}, \mu)}^2 + \|\tilde{e}_{\alpha}^{\perp}\|_{L_2(\mathcal{X}, \mu)}^2 = \|\tilde{e}_{\alpha}\|_{L_2(\mathcal{X}, \mu)}^2 = 1$$

252 by Lemma 2. Hence the series in Eq. (5) converges if $\sum_{\alpha} \lambda_{\alpha}^2$ converges. Recall from Section 3.1 that
253 T_k is positive and trace-class, so $\sum_{\alpha} \lambda_{\alpha} < \infty$ and $\lambda_{\alpha} \geq 0 \forall \alpha$, so $\sum_{\alpha} \lambda_{\alpha}^2 < \infty$ by the comparison
254 test and we're safe.

255 Moving to the second term, we have

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

256 and then

$$\begin{aligned}
\frac{1}{n} \mathbb{E}[\text{Tr}(J^\perp)] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\sum_{\alpha} \lambda_{\alpha}^2 e_{\alpha}^{\perp}(X_i) e_{\alpha}^{\perp}(X_i) \right] \\
&= \sum_{\alpha} \lambda_{\alpha}^2 \|\tilde{e}_{\alpha}^{\perp}\|_{L_2(\mathcal{X}, \mu)}^2 \\
&= \sum_{\alpha} \lambda_{\alpha}^2 - \sum_{\alpha} \lambda_{\alpha}^2 \|\mathcal{O} \tilde{e}_{\alpha}\|_{L_2(\mathcal{X}, \mu)}^2 \\
&= \text{Tr}(T_k^2) - \text{Tr}(T_k^2 \mathcal{O})
\end{aligned}$$

257 with the sums converging and the expectations being okay by virtue of the above. \square

258 Theorem 5 shows that feature averaging is provably beneficial in terms of generalisation if the mean
259 of the target distribution is invariant. One might think that, given enough training examples, the
260 solution f to Eq. (3) would *learn* to be \mathcal{G} -invariant. Theorem 5 shows that this cannot happen unless
261 the number of examples dominates the effective dimension of \mathcal{H}_A .

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

266 Finally, we must be careful to state that our setting does not directly reduce to that of [8, Theorem 6]
267 for two reasons. First, [8, Theorem 6] considers \mathcal{G} invariant linear models without regularisation. This
268 may turn out to be accessible by a $\rho \rightarrow 0^+$ limit (the so called ridgeless limit) of Theorem 5. More
269 importantly, linear regression is equivalent to kernel regression with the linear kernel. However, the
270 linear kernel can be unbounded (e.g. on \mathbb{R}), so does not meet our technical conditions in Section 2.2.
271 We therefore conjecture that the boundedness assumption on k can be removed.

272 5 Related Work

273 Incorporating invariance into machine learning models is not a new idea. The majority of
274 modern applications concern neural networks, but earlier work has used kernels [10], support
275 vector machines [24] and polynomial feature spaces [25, 26]. Indeed, early work also considered
276 invariant neural networks [31], using methods that seem to have been rediscovered in [22]. Modern
277 implementations include invariant/equivariant convolutional architectures [4, 6] that are inspired by
278 concepts from mathematical physics and harmonic analysis [13, 5]. Some of these models even enjoy
279 universal approximation properties [19, 33].

280 The earliest attempt at theoretical justification for invariance of which we are aware is [1], which
281 roughly states that enforcing invariance cannot increase the VC dimension of a model. Anselmi
282 et al. [2] and Mroueh, Voinea, and Poggio [21] propose heuristic arguments for improved sample
283 complexity of invariant models. Sokolic et al. [27] build on the work of Xu and Mannor [32] to obtain
284 a generalisation bound for certain types of classifiers that are invariant to a finite set of transformations,
285 while Sannai and Imaizumi [23] obtain a bound for models that are invariant to finite permutation
286 groups. The PAC Bayes formulation is considered in [16, 17].

287 The above works guarantee only a worst-case improvement and it was not until very recently
288 that Elesedy and Zaidi [8] derived a strict benefit for invariant/equivariant models. Our work is similar
289 to [8] in that we provide a provably strict benefit, but differs in its application to kernels and RKHSs
290 as opposed to linear models. Also very recently, Mei, Misiakiewicz, and Montanari [20] analyse the
291 generalisation benefit of invariance in kernels and random feature models. Our results differ from [20]
292 in some key aspects. First, Mei, Misiakiewicz, and Montanari [20] focus kernel ridge regression with
293 an invariant inner product kernel whereas we study symmetrised predictors from general kernels.
294 Second, they obtain an expression for the generalisation error that is conditional on the training data
295 and in terms of the projection of the predictor onto a space of high degree polynomials, while we are
296 able to integrate against the training data and express the generalisation benefit directly in terms of
297 properties of the kernel and the RKHS.

298 **References**

- 299 [1] Yaser S Abu-Mostafa. “Hints and the VC dimension”. In: *Neural Computation* 5.2 (1993),
300 pp. 278–288 (page 8).
- 301 [2] Fabio Anselmi et al. *Unsupervised Learning of Invariant Representations in Hierarchical*
302 *Architectures*. 2014. arXiv: [1311.4158](https://arxiv.org/abs/1311.4158) [cs.CV] (page 8).
- 303 [3] Sanjeev Arora et al. “Fine-Grained Analysis of Optimization and Generalization for
304 Overparameterized Two-Layer Neural Networks”. In: *Proceedings of the 36th International*
305 *Conference on Machine Learning*. Ed. by Kamalika Chaudhuri and Ruslan Salakhutdinov.
306 Vol. 97. Proceedings of Machine Learning Research. PMLR, 2019, pp. 322–332. URL: <http://proceedings.mlr.press/v97/arora19a.html> (page 2).
- 307 [4] Taco Cohen and Max Welling. “Group equivariant convolutional networks”. In: *International*
308 *conference on machine learning*. 2016, pp. 2990–2999 (page 8).
- 309 [5] Taco S Cohen, Mario Geiger, and Maurice Weiler. “A general theory of equivariant cnns
310 on homogeneous spaces”. In: *Advances in Neural Information Processing Systems*. 2019,
311 pp. 9145–9156 (page 8).
- 312 [6] Taco S Cohen et al. “Spherical cnns”. In: *arXiv preprint arXiv:1801.10130* (2018) (page 8).
- 313 [7] Cohn, Donald L. *Measure Theory*. 2nd ed. Springer, 2013 (page 4).
- 314 [8] Bryn Elesedy and Sheheryar Zaidi. “Provably Strict Generalisation Benefit for Equivariant
315 Models”. In: (2021). arXiv: [2102.10333](https://arxiv.org/abs/2102.10333) [stat.ML] (pages 1–3, 6, 8).
- 316 [9] Adam Foster, Rattana Pukdee, and Tom Rainforth. “Improving Transformation Invariance in
317 Contrastive Representation Learning”. In: *arXiv preprint arXiv:2010.09515* (2020) (page 1).
- 318 [10] B. Haasdonk, A. Vossen, and H. Burkhardt. “Invariance in Kernel Methods by Haar Integration
319 Kernels”. In: *SCIA 2005, Scandinavian Conference on Image Analysis*. Springer-Verlag, 2005,
320 pp. 841–851 (pages 4, 8).
- 321 [11] Arthur Jacot, Franck Gabriel, and Clément Hongler. “Neural tangent kernel: Convergence and
322 generalization in neural networks”. In: *Advances in neural information processing systems*.
323 2018, pp. 8571–8580 (page 2).
- 324 [12] Olav Kallenberg. *Foundations of modern probability*. Springer Science & Business Media,
325 2006 (page 7).
- 326 [13] Risi Kondor and Shubhendu Trivedi. “On the Generalization of Equivariance and Convolution
327 in Neural Networks to the Action of Compact Groups”. In: *International Conference on*
328 *Machine Learning*. 2018, pp. 2747–2755 (page 8).
- 329 [14] Jaehoon Lee et al. “Wide neural networks of any depth evolve as linear models under gradient
330 descent”. In: *Advances in neural information processing systems*. 2019, pp. 8572–8583
331 (page 2).
- 332 [15] Juho Lee et al. “Set Transformer: A Framework for Attention-based Permutation-Invariant
333 Neural Networks”. In: *Proceedings of the 36th International Conference on Machine Learning*.
334 Ed. by Kamalika Chaudhuri and Ruslan Salakhutdinov. Vol. 97. Proceedings of Machine
335 Learning Research. PMLR, 2019, pp. 3744–3753. URL: <http://proceedings.mlr.press/v97/lee19d.html> (page 1).
- 336 [16] Clare Lyle, Marta Kwiatkowska, and Yarin Gal. “An analysis of the effect of invariance on
337 generalization in neural networks”. In: *International conference on machine learning Workshop*
338 *on Understanding and Improving Generalization in Deep Learning*. 2019 (page 8).
- 339 [17] Clare Lyle et al. *On the Benefits of Invariance in Neural Networks*. 2020. arXiv: [2005.00178](https://arxiv.org/abs/2005.00178)
340 [cs.LG] (pages 1, 8).
- 341 [18] Jonathan H Manton and Pierre-Olivier Amblard. “A primer on reproducing kernel hilbert
342 spaces”. In: *arXiv preprint arXiv:1408.0952* (2014) (page 6).
- 343 [19] Haggai Maron et al. “On the Universality of Invariant Networks”. In: *International Conference*
344 *on Machine Learning*. 2019, pp. 4363–4371 (page 8).
- 345 [20] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. “Learning with invariances in
346 random features and kernel models”. In: *arXiv preprint arXiv:2102.13219* (2021) (pages 1, 8).
- 347 [21] Youssef Mroueh, Stephen Voinea, and Tomaso A Poggio. “Learning with Group Invariant
348 Features: A Kernel Perspective.” In: *Advances in Neural Information Processing Systems*.
349 2015, pp. 1558–1566 (page 8).
- 350 [22] Siamak Ravanbakhsh, Jeff Schneider, and Barnabas Poczos. “Equivariance through parameter-
351 sharing”. In: *International Conference on Machine Learning*. PMLR. 2017, pp. 2892–2901
352 (page 8).

- 355 [23] Akiyoshi Sannai and Masaaki Imaizumi. *Improved Generalization Bound of Group Invariant*
356 */ Equivariant Deep Networks via Quotient Feature Space*. 2019. arXiv: [1910 . 06552](https://arxiv.org/abs/1910.06552)
357 [\[stat.ML\]](#) (pages 1, 8).
- 358 [24] Bernhard Schölkopf, Chris Burges, and Vladimir Vapnik. “Incorporating Invariances in Support
359 Vector Learning Machines”. In: Springer, 1996, pp. 47–52 (page 8).
- 360 [25] H. Schulz-Mirbach. “Constructing invariant features by averaging techniques”. In: *Proceedings*
361 *of the 12th IAPR International Conference on Pattern Recognition, Vol. 3 - Conference C:*
362 *Signal Processing (Cat. No.94CH3440-5)*. Vol. 2. 1994, 387–390 vol.2 (page 8).
- 363 [26] Hanns Schulz-Mirbach. “On the existence of complete invariant feature spaces in pattern
364 recognition”. In: *International Conference On Pattern Recognition*. Citeseer. 1992, pp. 178–
365 178 (page 8).
- 366 [27] Jure Sokolic et al. “Generalization error of invariant classifiers”. In: *Artificial Intelligence and*
367 *Statistics*. 2017, pp. 1094–1103 (pages 1, 8).
- 368 [28] James S Spencer et al. “Better, Faster Fermionic Neural Networks”. In: *arXiv preprint*
369 *arXiv:2011.07125* (2020) (page 1).
- 370 [29] Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Information science and
371 statistics. Springer, 2008. ISBN: 978-0-387-77241-7 (pages 2, 4, 5).
- 372 [30] Marysia Winkels and Taco S Cohen. “3D G-CNNs for pulmonary nodule detection”. In: *arXiv*
373 *preprint arXiv:1804.04656* (2018) (page 1).
- 374 [31] Jeffrey Wood and John Shawe-Taylor. “Representation theory and invariant neural networks”.
375 In: *Discrete applied mathematics* 69.1-2 (1996), pp. 33–60 (pages 1, 8).
- 376 [32] Huan Xu and Shie Mannor. “Robustness and generalization”. In: *Machine learning* 86.3
377 (2012), pp. 391–423 (page 8).
- 378 [33] Dmitry Yarotsky. *Universal approximations of invariant maps by neural networks*. 2018. arXiv:
379 [1804.10306](https://arxiv.org/abs/1804.10306) [\[cs.NE\]](#) (page 8).
- 380 [34] Manzil Zaheer et al. “Deep sets”. In: *Advances in neural information processing systems*. 2017,
381 pp. 3391–3401 (page 1).

382 Checklist

- 383 1. For all authors...
- 384 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
385 contributions and scope? [\[Yes\]](#)
- 386 (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 2.2 which describes
387 our assumptions.
- 388 (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
- 389 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
390 them? [\[Yes\]](#)
- 391 2. If you are including theoretical results...
- 392 (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See ?? for
393 general technical conditions not given in statements of results.
- 394 (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) Proofs are given in
395 the supplementary material
- 396 3. If you ran experiments...
- 397 (a) Did you include the code, data, and instructions needed to reproduce the main
398 experimental results (either in the supplemental material or as a URL)? [\[N/A\]](#)
- 399 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
400 were chosen)? [\[N/A\]](#)
- 401 (c) Did you report error bars (e.g., with respect to the random seed after running
402 experiments multiple times)? [\[N/A\]](#)
- 403 (d) Did you include the total amount of compute and the type of resources used (e.g., type
404 of GPUs, internal cluster, or cloud provider)? [\[N/A\]](#)
- 405 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 406 (a) If your work uses existing assets, did you cite the creators? [N/A]
407 (b) Did you mention the license of the assets? [N/A]
408 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
409
410 (d) Did you discuss whether and how consent was obtained from people whose data you're
411 using/curating? [N/A]
412 (e) Did you discuss whether the data you are using/curating contains personally identifiable
413 information or offensive content? [N/A]
414 5. If you used crowdsourcing or conducted research with human subjects...
- 415 (a) Did you include the full text of instructions given to participants and screenshots, if
416 applicable? [N/A]
417 (b) Did you describe any potential participant risks, with links to Institutional Review
418 Board (IRB) approvals, if applicable? [N/A]
419 (c) Did you include the estimated hourly wage paid to participants and the total amount
420 spent on participant compensation? [N/A]