
Dimension and Margin Bounds for Reflection-invariant Kernels *

Thorsten Doliwa, Michael Kallweit, and Hans Ulrich Simon

Fakultät für Mathematik, Ruhr-Universität Bochum, D-44780 Bochum, Germany

{thorsten.doliwa,michael.kallweit,hans.simon}@rub.de

Abstract

A kernel over the Boolean domain is said to be reflection-invariant, if its value does not change when we flip the same bit in both arguments. (Many popular kernels have this property.) We study the geometric margins that can be achieved when we represent a specific Boolean function f by a classifier that employs a reflection-invariant kernel. It turns out $\|\hat{f}\|_\infty$ is an upper bound on the average margin. Furthermore, $\|\hat{f}\|_\infty^{-1}$ is a lower bound on the smallest dimension of a feature space associated with a reflection-invariant kernel that allows for a correct representation of f . This is, to the best of our knowledge, the first paper that exhibits margin and dimension bounds for *specific functions* (as opposed to *function families*). Several generalizations are considered as well. The main mathematical results are presented in a setting with arbitrary finite domains and a quite general notion of invariance.

1 Introduction

There has been much interest in margin and dimension bounds during the last decade. The simplest way to cast (most of) the existing results in this direction is offered by the notion of margin and dimension complexity associated with a given sign matrix $A \in \{-1, 1\}^{m \times n}$. A linear arrangement, given by unit vectors $u_1, \dots, u_m; v_1, \dots, v_n$ (taken from an inner product space), is said to represent A if, for all $i = 1, \dots, m$ and $j = 1, \dots, n$, $A_{i,j} = \text{sign}(\langle u_i, v_j \rangle)$. The dimension complexity of

A is the smallest dimension of an inner product space that allows for such a representation. The margin complexity is obtained similarly by looking for the linear arrangement that leads to the maximum average margin (or, alternatively, to the maximum margin that can be guaranteed for all choices of i and j). Applying counting arguments, Ben-David, Eiron, and Simon [1] have shown that, loosely speaking, an overwhelming majority of sign matrices of small VC-dimension do not allow for a linear arrangement whose margin or dimension is significantly better than what can be guaranteed in a trivial fashion. Starting with Forster's celebrated exponential lower bound on the dimension complexity of the Walsh-Hadamard matrix [4], there has been a series of papers [5, 6, 10, 7, 13, 15] presenting (increasingly powerful) techniques for deriving upper margin bounds or lower dimension bounds on the complexity of sign matrices.

Note that a sign matrix represents a *family* of Boolean functions, one Boolean function per column say. The lack of non-trivial margin or dimension bounds for a *specific* Boolean function has a simple explanation: a specific function $f(x)$ can always trivially be represented in a 1-dimensional space with geometric margin 1 by mapping an instance $x \in \{-1, 1\}^n$ to $f(x) \in \{-1, 1\}$. The corresponding kernel would map a pair (x, x') of instances to 1 if $f(x) = f(x')$, and to -1 otherwise. Clearly, the 1-dimensional "linear arrangement" for f does not say much about the ability of kernel-based large margin classifier systems to "learn" f because we would need to know f perfectly prior to the choice of the kernel. (If we had this knowledge, there would be nothing to learn anymore.) Nevertheless, this discussion shows that one cannot expect non-trivial margin or dimension bounds for *specific functions* that hold *uniformly for all kernels*.

In this paper, we introduce the concept of distributed functions that are invariant under a group \mathcal{G} of transformations. We present the mathematical results about invariant distributed functions in a quite general setting (because it does not make

*This work was supported in part by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778. This publication only reflects the authors' views. This work was furthermore supported by the Deutsche Forschungsgemeinschaft Grant SI 498/8-1.

sense to impose unnecessary restrictions). In particular, we derive non-trivial margin and dimension bounds for specific Boolean functions that are valid for all linear arrangements resulting from \mathcal{G} -invariant kernels. If the domain of the distributed function can be cast as a finite Abelian group, the margin and dimension bounds for a function f can be nicely expressed in terms of f 's Fourier-spectrum. As always, $\|\hat{f}\|_\infty$ denotes the largest absolute value found in the spectrum of f 's Fourier-coefficients. We show that $\|\hat{f}\|_\infty$ is an upper bound on the largest possible average margin, and $\|\hat{f}\|_\infty^{-1}$ is a lower bound on the smallest possible dimension. Our general results easily apply to a special case of high learning-theoretic relevance, namely the reflection-invariant kernels. Their relevance comes from the fact that, as demonstrated in the paper, many popular kernels actually happen to be reflection-invariant.

The remainder of the paper is structured as follows. In Section 2, we fix some notation and recall some facts about Fourier-expansions over finite Abelian groups and kernel-based classification. In Section 3, we present our results for arbitrary finite domains and a quite general notion of invariance. In Section 4, we introduce the concept of rotation-invariance and mention some connections between the Fourier-expansion over an arbitrary finite Abelian group and the spectral decomposition of such functions. In Section 5, we consider distributed functions over the Boolean domain and the concept of reflection-invariance, which is simply rotation-invariance over a Boolean domain. Section 6 presents the margin and dimension bounds that are valid for reflection-invariant kernels. Section 7 offers a possible interpretation of our results, and mentions a connection to a recent paper by Haasdonk and Burkhardt [8] along with some open problems.

2 Definitions and Notations

We assume familiarity with basics in matrix and learning theory. For example, notions like

- singular values, eigenvalues, spectral norm
- kernels, feature map, Reproducing Kernel Hilbert Space

are assumed as known (although we shall occasionally refresh the readers memory). Some central definitions and facts concerning

- linear arrangements representing a given sign matrix,
- margin and dimension associated with such a linear arrangement,

will be given later in the paper at the place where it is required. In the following we fix some notation

and recall the Fourier-expansion over finite Abelian groups as well as the notion of margin in kernel-based classification.

2.1 Preliminaries

Throughout the paper, δ denotes the Kronecker-symbol, i.e., $\delta(a, b) = 1$ if $a = b$ and $\delta(a, b) = 0$ otherwise. For two n -dimensional vectors x, y , we define $x \circ y$ to be the vector obtained by multiplying x and y componentwise, i.e., $(x \circ y)_i := x_i y_i$ for $i = 1, \dots, n$. The n -dimensional ‘‘all-ones vector’’ is given by

$$\vec{e} = (1, \dots, 1) .$$

The vector with 1 in component k and zeros elsewhere is denoted as \vec{e}_k . We consider functions over a finite domain D with values in \mathbb{R} (or in \mathbb{C} , resp.). These functions form a $|D|$ -dimensional vector space. A *distributed function over D* is a function over the domain $D \times D$. We will occasionally identify a distributed function f over D with the $(D \times D)$ -matrix F given by $F_{x,y} = f(x, y)$.

2.2 Fourier-expansions over Finite Abelian Groups

Let $(D, +)$ be a finite Abelian group of size $d = |D|$. A function $\chi : D \rightarrow \mathbb{C}$ is called a *character* over D if, for every $x, y \in D$,

$$\chi(x + y) = \chi(x) \cdot \chi(y) .$$

It is well-known that there are exactly d characters, and they form an orthonormal basis of the vector space \mathbb{C}^D with respect to the inner product

$$\langle f, g \rangle := \frac{1}{d} \cdot \sum_{x \in D} f(x) \cdot \overline{g(x)} . \quad (1)$$

We may fix a bijection between D and the set of characters and write χ_z for the character that corresponds to $z \in D$. Every function $f : D \rightarrow \mathbb{C}$ can be written in the form

$$f(x) = \sum_{z \in D} \hat{f}(z) \cdot \chi_z(x) \quad (2)$$

where

$$\hat{f}(z) := \langle f, \chi_z \rangle = \frac{1}{d} \cdot \sum_{y \in D} f(y) \cdot \overline{\chi_z(y)} .$$

Equation (2) is referred to as the *Fourier expansion of f* , and $\hat{f}(z)$ is called the *Fourier-coefficient of f at z* .

According to the ‘‘Fundamental Theorem for Finitely Generated Abelian Groups’’, every finite Abelian group is, up to isomorphism, of the form

$$D = \mathbb{Z}_{q_1} \times \dots \times \mathbb{Z}_{q_n} \quad (3)$$

for some sequence q_1, \dots, q_n of prime powers. Equation (3) is assumed henceforth so that

$$d = |D| = \prod_{k=1}^n q_k .$$

It is well-known that the characters over \mathbb{Z}_m are given by

$$\chi_k^{(m)}(j) = \omega_m^{jk} ,$$

where

$$\omega_m = \exp\left(\frac{2\pi i}{m}\right)$$

is a primitive root of unity of order m . The characters over D are then given by

$$\chi_z(x) = \prod_{k=1}^n \chi_{z_k}^{(q_k)}(x_k) .$$

Consider now the matrix $H = (H_{x,z})_{x,z \in D}$ given by

$$H_{x,z} = \chi_z(x) . \quad (4)$$

It is obvious that H is symmetric. By the orthonormality of the characters with respect to the inner product in (1), it follows that

$$H^* \cdot H = H \cdot H^* = d \cdot I ,$$

where I denotes the identity matrix.

2.3 Kernel-based Classification

Let $K : D \times D \rightarrow \mathbb{R}$ be a valid kernel over a finite domain D . In other words, $K(x, y)$ is a real-valued distributed function over D which, considered as matrix, is symmetric and positive semidefinite. Let Φ_K be the feature map and $\langle \cdot, \cdot \rangle_K$ the inner product that represent K in the Reproducing Kernel Hilbert Space, and let $\| \cdot \|_K$ be the norm induced by $\langle \cdot, \cdot \rangle_K$.¹ Then Φ satisfies

$$\forall x, y \in D : K(x, y) = \langle \Phi(x), \Phi(y) \rangle .$$

With every ‘‘dual vector’’ $\alpha : D \rightarrow \mathbb{R}$, we associate the ‘‘weight vector’’

$$w(\alpha) := \sum_{x \in D} \alpha(x) \Phi(x) . \quad (5)$$

In the context of ‘‘large margin classification’’, α is considered as a classifier that assigns the label $\text{sign}(\langle w(\alpha), \Phi(x) \rangle)$ to input x . Consider a target function $f : D \rightarrow \{-1, 1\}$ for a binary classification task. Then, a negative sign of $f(x) \cdot \langle w(\alpha), \Phi(x) \rangle$ indicates a ‘‘classification error’’ on x . So this expression should be *positive* and it is intuitively even better when it leads to a *large* positive value. Thus, the following number, called the (*geometric*) *margin achieved by α on x w.r.t. target function f and kernel K* , is of interest:

$$\mu_K(f|\alpha, x) := \frac{f(x) \cdot \langle w(\alpha), \Phi(x) \rangle}{\|w(\alpha)\| \cdot \|\Phi(x)\|} \quad (6)$$

By averaging over all $x \in D$, we obtain the function

$$\bar{\mu}_K(f|\alpha) := 2^{-n} \sum_{x \in D} \mu_K(f|\alpha, x) .$$

¹In the sequel, we drop index K unless we would like to stress the dependence on K .

Focusing on the margin that is guaranteed for every $x \in D$, we should consider the function

$$\mu_K(f|\alpha) := \min_{x \in D} \mu_K(f|\alpha, x) .$$

By taking the supremum over all $\alpha : D \rightarrow \mathbb{R}$, we get the respective parameters of a large margin classifier employing kernel function K :

$$\begin{aligned} \bar{\mu}_K(f) &:= \sup_{\alpha: D \rightarrow \mathbb{R}} \bar{\mu}_K(f|\alpha) \\ \mu_K(f) &:= \sup_{\alpha: D \rightarrow \mathbb{R}} \mu_K(f|\alpha) \end{aligned}$$

Finally, taking the supremum ranging over all K from a given kernel class \mathcal{C} , we get the respective parameters of a best possible large margin classifier among those that employ a kernel from \mathcal{C} :

$$\begin{aligned} \bar{\mu}_{\mathcal{C}}(f) &:= \sup_{K \in \mathcal{C}} \bar{\mu}_K(f) \\ \mu_{\mathcal{C}}(f) &:= \sup_{K \in \mathcal{C}} \mu_K(f) \end{aligned}$$

We briefly note that, obviously, the guaranteed margin is upper bounded by the average margin:

$$\begin{aligned} \mu_K(f|\alpha) &\leq \bar{\mu}_K(f|\alpha) \\ \mu_K(f) &\leq \bar{\mu}_K(f) \\ \mu_{\mathcal{C}}(f) &\leq \bar{\mu}_{\mathcal{C}}(f) \end{aligned}$$

3 A General Notion of Invariance

Throughout this section, D denotes an arbitrary finite domain, $\mathcal{S}(D)$ is the group of permutations over D , and $\mathcal{G} \leq \mathcal{S}(D)$ is an arbitrary but fixed subgroup. A distributed function over D with values in $V \subseteq \mathbb{C}$ is said to be \mathcal{G} -invariant if, for all $x, y \in D$ and every $\sigma \in \mathcal{G}$, the following holds:

$$f(\sigma(x), \sigma(y)) = f(x, y)$$

We clearly have the

Pointwise Closure Property: The pointwise limit of \mathcal{G} -invariant functions is a \mathcal{G} -invariant function. Furthermore, if f_1, \dots, f_d are \mathcal{G} -invariant functions and $g : V^d \rightarrow W$ is an arbitrary function with values in $W \subseteq \mathbb{C}$, then

$$g(f_1(x, y), \dots, f_d(x, y))$$

is \mathcal{G} -invariant too.

More interesting is the the following result:

Lemma 1 *\mathcal{G} -invariant distributed functions over a finite domain D are closed under the usual matrix product and under the tensor-product of matrices. More precisely, let $F(x, y)$ and $G(x, y)$ be two \mathcal{G} -invariant distributed functions (here viewed as matrices). Then, the functions $(F \cdot G)(x, y)$ is \mathcal{G} -invariant and the function $(F \otimes G)[(u, x), (v, y)]$ is invariant over $\mathcal{G} \times \mathcal{G}$ (as subgroup of $\mathcal{S}(D) \times \mathcal{S}(D)$).*

Proof: Consider first the function $(F \cdot G)(x, y)$. Let $x, y \in D$ and $\sigma \in \mathcal{G}$ be arbitrary but fixed. The following calculation shows that it is \mathcal{G} -invariant:

$$\begin{aligned} (F \cdot G)_{\sigma(x), \sigma(y)} &= \sum_{z \in D} F_{\sigma(x), z} \cdot G_{z, \sigma(y)} \\ &= \sum_{z \in D} F_{x, \sigma^{-1}(z)} \cdot G_{\sigma^{-1}(z), y} \\ &= \sum_{z \in D} F_{x, z} \cdot G_{z, y} \\ &= (F \cdot G)_{x, y} \end{aligned}$$

Now consider the tensor-product $(F \otimes G)[(u, x), (v, y)]$, which is a distributed function over $D \times D$, i.e., a function over domain $(D \times D) \times (D \times D)$. The following calculation shows that it is $(\mathcal{G} \times \mathcal{G})$ -invariant:

$$\begin{aligned} (F \otimes G)[(\sigma(u), \tau(x)), (\sigma(v), \tau(y))] &= \\ F(\sigma(u), \sigma(v)) \cdot G(\tau(x), \tau(y)) &= \\ F(u, v) \cdot G(x, y) &= \\ (F \otimes G)[(u, x), (v, y)] & \end{aligned}$$

■

In this section, we shall show the following. If $f : D \rightarrow \{-1, 1\}$ is a function on domain D and \mathcal{G} is a subgroup of $\mathcal{S}(D)$, then the largest average (or largest guaranteed, resp.) margin that can be obtained when f is represented by a \mathcal{G} -invariant kernel is upper-bounded by the largest average (or largest guaranteed, resp.) margin that can be obtained for the family

$$\mathcal{G}_f := \{f_\sigma : \sigma \in \mathcal{G}\}$$

where

$$f_\sigma(x) := f(\sigma(x)) .$$

Since there are classical margin bounds that apply to the family \mathcal{G}_f , we obtain corresponding bounds that apply to the single function f . An analogous remark holds for dimension bounds. Details follow.

Assume that $K(x, y)$ is a \mathcal{G} -invariant kernel and consider the feature map $\Phi = \Phi_K$ that represents K in the Reproducing Kernel Hilbert Space. Then, for all $x, y \in D$ and every $\sigma \in \mathcal{G}$, Φ satisfies

$$\langle \Phi(\sigma(x)), \Phi(\sigma(y)) \rangle = \langle \Phi(x), \Phi(y) \rangle . \quad (7)$$

Lemma 2 *If kernel K is \mathcal{G} -invariant, then the following holds for every $x \in D$ and every $\sigma \in \mathcal{G}$:*

$$\begin{aligned} \|\Phi_K(\sigma(x))\|_K &= \|\Phi_K(x)\|_K \\ \|w(\alpha)\|_K &= \|w(\alpha_\sigma)\|_K \end{aligned}$$

In other words, the norm $\|\cdot\|_K$ is constant on feature vectors of instances taken from the same orbit

$$x^{\mathcal{G}} := \{\sigma(x) : \sigma \in \mathcal{G}\}$$

and it assigns the same value to all dual vectors from the set

$$\{w(\alpha_\sigma) : \sigma \in \mathcal{G}\} .$$

Proof: Let $\Phi = \Phi_K$, $\|\cdot\| = \|\cdot\|_K$, and $\langle \cdot, \cdot \rangle = \langle \cdot, \cdot \rangle_K$. Clearly, $\|\Phi(\sigma(x))\| = \|\Phi(x)\|$ because of

$$\begin{aligned} \|\Phi(\sigma(x))\|^2 &= \langle \Phi(\sigma(x)), \Phi(\sigma(x)) \rangle \\ &\stackrel{(7)}{=} \langle \Phi(x), \Phi(x) \rangle \\ &= \|\Phi(x)\|^2 . \end{aligned}$$

As for the second statement, see the following calculation:

$$\begin{aligned} \|w(\alpha_\sigma)\|^2 &= \langle w(\alpha_\sigma), w(\alpha_\sigma) \rangle \\ &\stackrel{(5)}{=} \left\langle \sum_{x \in D} \alpha_\sigma(x) \Phi(x), \sum_{y \in D} \alpha_\sigma(y) \Phi(y) \right\rangle \\ &= \sum_{x, y \in D} \alpha(\sigma(x)) \alpha(\sigma(y)) \langle \Phi(x), \Phi(y) \rangle \\ &= \sum_{x, y \in D} \alpha(x) \alpha(y) \langle \Phi(\sigma^{-1}(x)), \Phi(\sigma^{-1}(y)) \rangle \\ &\stackrel{(7)}{=} \sum_{x, y \in D} \alpha(x) \alpha(y) \langle \Phi(x), \Phi(y) \rangle \\ &= \|w(\alpha)\|^2 \end{aligned}$$

■

Lemma 3 *For every \mathcal{G} -invariant kernel K , and every choice of $f : D \rightarrow \{-1, 1\}$, $x \in D$, $\sigma \in \mathcal{G}$, and $\alpha : D \rightarrow \mathbb{R}$, the following holds:*

$$\mu_K(f_\sigma | \alpha_\sigma, x) = \mu_K(f | \alpha, \sigma(x))$$

Proof: The proof starts as follows:

$$\begin{aligned} f_\sigma(x) \cdot \langle w(\alpha_\sigma), \Phi(x) \rangle &\stackrel{(5)}{=} \\ f_\sigma(x) \left\langle \sum_{y \in D} \alpha_\sigma(y) \Phi(y), \Phi(x) \right\rangle &= \\ f(\sigma(x)) \sum_{y \in D} \alpha(\sigma(y)) \langle \Phi(y), \Phi(x) \rangle &\stackrel{(7)}{=} \\ f(\sigma(x)) \sum_{y \in D} \alpha(\sigma(y)) \langle \Phi(\sigma(y)), \Phi(\sigma(x)) \rangle &= \\ f(\sigma(x)) \left\langle \sum_{y \in D} \alpha(\sigma(y)) \Phi(\sigma(y)), \Phi(\sigma(x)) \right\rangle &= \\ f(\sigma(x)) \left\langle \sum_{y \in D} \alpha(y) \Phi(y), \Phi(\sigma(x)) \right\rangle &= \\ f(\sigma(x)) \langle w(\alpha), \Phi(\sigma(x)) \rangle & \end{aligned}$$

Using this calculation in combination with Lemma 2, the proof is easy to accomplish:

$$\begin{aligned} \mu_K(f_\sigma | \alpha_\sigma, x) &\stackrel{(6)}{=} \frac{f_\sigma(x) \cdot \langle w(\alpha_\sigma), \Phi(x) \rangle}{\|w(\alpha_\sigma)\| \cdot \|\Phi(x)\|} \\ &= \frac{f(\sigma(x)) \cdot \langle w(\alpha), \Phi(\sigma(x)) \rangle}{\|w(\alpha)\| \cdot \|\Phi(\sigma(x))\|} \\ &\stackrel{(6)}{=} \mu_K(f | \alpha, \sigma(x)) \end{aligned}$$

■

Corollary 4 For every \mathcal{G} -invariant kernel K , and every choice of $f : D \rightarrow \{-1, 1\}$, $\sigma \in \mathcal{G}$, and $\alpha : D \rightarrow \mathbb{R}$, the following holds:

$$\begin{aligned}\bar{\mu}_K(f_\sigma | \alpha_\sigma) &= \bar{\mu}_K(f | \alpha) \\ \mu_K(f_\sigma | \alpha_\sigma) &= \mu_K(f | \alpha) \\ \bar{\mu}_K(f_\sigma) &= \bar{\mu}_K(f) \\ \mu_K(f_\sigma) &= \mu_K(f) \\ \bar{\mu}_{\mathcal{G}}(f_\sigma) &= \bar{\mu}_{\mathcal{G}}(f) \\ \mu_{\mathcal{G}}(f_\sigma) &= \mu_{\mathcal{G}}(f)\end{aligned}$$

Note that the last two equations in Corollary 4 basically say that the largest (average or guaranteed) margin that can be achieved for a function f by a large margin classifier is invariant under \mathcal{G} (provided that the underlying kernel is \mathcal{G} -invariant).

Let $M \in \{-1, 1\}^{r \times s}$ be a sign matrix. Consider a linear arrangement \mathcal{A} given by unit vectors $u_1, \dots, u_r; v_1, \dots, v_s \in \mathbb{R}^d$. The *average margin achieved by this arrangement for sign matrix M* is defined as follows:

$$\bar{\mu}(M | \mathcal{A}) := \frac{1}{rs} \cdot \sum_{i=1}^r \sum_{j=1}^s M_{i,j} \langle u_i, v_j \rangle$$

The *largest average margin that can be achieved for sign matrix M by any linear arrangement* is then given by

$$\bar{\mu}(M) := \sup_{\mathcal{A}} \bar{\mu}(M | \mathcal{A}) ,$$

where the supremum ranges over all linear arrangements \mathcal{A} for M . Forster and Simon [7] have shown that, for every $M \in \mathbb{R}^{r \times s}$, every $d \geq 1$, and every choice of unit vectors $u_1, \dots, u_r; v_1, \dots, v_s$ in a real inner-product space, the following holds:

$$\sum_{i=1}^r \sum_{j=1}^s M_{i,j} \langle u_i, v_j \rangle \leq \sqrt{rs} \|M\| .$$

From that, we conclude that

$$\bar{\mu}(M) \leq \frac{\|M\|}{\sqrt{rs}} .$$

Consider the sign matrix $M^{f, \mathcal{G}}$ given by

$$M_{x, \sigma}^{f, \mathcal{G}} := f_\sigma(x) . \quad (8)$$

In combination with Corollary 4, we arrive at the following

Theorem 5 Let D be a finite domain, and let \mathcal{G} be a subgroup of $\mathcal{S}(D)$. Then, every function $f : D \rightarrow \{-1, 1\}$ satisfies

$$\bar{\mu}_{\mathcal{G}}(f) \leq \frac{\|M^{f, \mathcal{G}}\|}{\sqrt{|D| \cdot |\mathcal{G}|}} .$$

In other words, no large margin classifier that employs a \mathcal{G} -invariant kernel can achieve an average margin for f which exceeds $\frac{\|M^{f, \mathcal{G}}\|}{\sqrt{|D| \cdot |\mathcal{G}|}}$.

As our input space D is finite, we can assume without loss of generality that the Reproducing Kernel Hilbert Space for a kernel K on D coincides with $\mathbb{R}^{d(K)}$ for some suitable $1 \leq d(K) \leq |D|$. We say that $\alpha : D \rightarrow \mathbb{R}$ represents target function f correctly w.r.t. kernel K if

$$\forall x \in D : \mu_K(f | \alpha, x) > 0 .$$

Corollary 6 Let $d_{\mathcal{G}}(f)$ denote the smallest dimension of a feature space associated with a \mathcal{G} -invariant kernel K that allows for a correct representation of f . Then,

$$d_{\mathcal{G}}(f) \geq \frac{\sqrt{|D| \cdot |\mathcal{G}|}}{\|M^{f, \mathcal{G}}\|} .$$

Proof: According to Lemma 3, a kernel that allows for a correct representation of f allows also for a correct representation of all f_σ . According to a result by Forster [4], the corresponding feature space must have dimension at least $\sqrt{|D| \cdot |\mathcal{G}|} / \|M^{f, \mathcal{G}}\|$. ■

Corollary 6 can be strengthened slightly:

Corollary 7 Let σ_i denote the i -th singular value of $M^{f, \mathcal{G}}$, where $\sigma_1, \sigma_2, \dots$ are in decreasing order. Then, $d_{\mathcal{G}}(f)$ satisfies the following lower bound:

$$d_{\mathcal{G}}(f) \cdot \sum_{i=1}^{d_{\mathcal{G}}(f)} \sigma_i^2 \geq 1 \quad (9)$$

Proof: Let $A \in \{-1, 1\}^{r \times s}$ be a matrix whose columns are viewed as binary functions f_1, \dots, f_s . It has been shown by Forster and Simon [7] that the dimension d of a feature space which allows for a correct representation of f_1, \dots, f_s satisfies

$$d \cdot \sum_{i=1}^d \sigma_i^2(A) \geq rs .$$

This trivially implies (9). ■

4 Rotation-invariant Functions

In Section 4.1 we will derive some facts about distributed functions over a finite Abelian group via the Fourier-expansion. Section 4.2 ties everything together and presents the resulting margin and dimension bounds obtained in this restricted setting.

4.1 Distributed Functions over Finite Abelian Groups

We apply the results of the preceding section to the case where D is a Abelian group of finite size d , and \mathcal{G}_{rot} is the subgroup of $\mathcal{S}(D)$ consisting of all permutations of the form $x \mapsto x + a$. Note that $d = |D| = |\mathcal{G}_{rot}|$.

We are interested in distributed functions $f : D \times D \rightarrow \mathbb{C}$ and arrange the d^2 Fourier-coefficients of such a function as a matrix as follows:

$$\widehat{F}_{a,b} = \widehat{f}(a, -b) \quad (10)$$

$$= d^{-2} \sum_{(x,y) \in D \times D} f(x,y) \overline{\chi_{(a,-b)}(x,y)} \quad (11)$$

$$= d^{-2} \cdot \sum_{x \in D} \sum_{y \in D} f(x,y) \overline{\chi_a(x)} \chi_b(y) \quad (12)$$

In matrix notation, this reads as

$$\widehat{F} = d^{-2} \cdot H^* \cdot F \cdot H \quad , \quad (13)$$

where H is the matrix from (4).

A distributed function $f(x,y)$ over D is said to be *rotation-invariant* if, for all $x,y,a \in D$, the following holds:

$$f(x+a, y+a) = f(x,y)$$

In the sense of the previous section, f is meant to be \mathcal{G}_{rot} -invariant.

Here are some examples for rotation-invariant functions:

- A distributed function of the form $f(x,y) = g(x-y)$ is obviously rotation-invariant. Conversely, any rotation-invariant function $f(x,y)$ can be written in this form by setting $g(x) := f(x,0)$ because rotation-invariance implies that

$$f(x,y) = f(x-y,0) = g(x-y) \quad .$$

- Because of the obvious identity

$$\chi_z(x-y) = \chi_z(x) \cdot \overline{\chi_z(y)} \quad ,$$

the distributed function $\chi_z(x) \cdot \overline{\chi_z(y)}$ is rotation-invariant too.

The fact that $f(x,y) = g(x-y)$ is a rotation-invariant function can be restated as follows: any function $f(x,y)$ that can be cast as a function in $x_1 - y_1 \bmod q_1, \dots, x_n - y_n \bmod q_n$ is rotation-invariant.

In terms of the matrix of Fourier-coefficients, \widehat{F} , rotation-invariant functions over D can be characterized as follows:

Lemma 8 *A distributed function $f(x,y)$ over D is rotation-invariant iff \widehat{F} is a diagonal matrix.*

Proof: Assume first that $f(x,y)$ is rotation-invariant. Consider a Fourier-coefficient in \widehat{F} outside the main diagonal, say $\widehat{F}_{a,b}$ so that $a_k \neq b_k$. Every pair (x,y) can be put into the equivalence class

$$\{(x + j\vec{e}_k, y + j\vec{e}_k) : j = 0, \dots, q_k - 1\} \quad .$$

We show that every equivalence class contributes 0 to (12):

$$\begin{aligned} & \sum_{j=0}^{q_k-1} f(x + j\vec{e}_k, y + j\vec{e}_k) \overline{\chi_a(x + j\vec{e}_k)} \cdot \chi_b(y + j\vec{e}_k) = \\ & f(x,y) \overline{\chi_a(x)} \cdot \chi_b(y) \sum_{j=0}^{q_k-1} \overline{\chi_{a_k}^{(q_k)}(j)} \chi_{b_k}^{(q_k)}(j) \end{aligned}$$

The latter sum vanishes because it equals

$$\sum_{j=0}^{q_k-1} \omega_{q_k}^{(b_k - a_k)j} \quad .$$

Recall that δ denotes the Kronecker symbol and it is well-known that

$$\sum_{j=0}^{m-1} \omega_m^{(l'-l)j} = m \cdot \delta_{l,l'} \quad .$$

This shows that $\widehat{F}_{a,b} = 0$.

Now assume that \widehat{F} is a diagonal matrix. We conclude from (13) that

$$F = H \cdot \widehat{F} \cdot H^* \quad , \quad (14)$$

which implies that

$$F_{x,y} = \sum_{z \in D} \widehat{F}_{z,z} \cdot \chi_x(z) \cdot \overline{\chi_y(z)} \quad .$$

Rotation-invariance is now easily obtained:

$$\begin{aligned} f(x+a, y+a) &= \sum_{z \in D} \widehat{F}_{z,z} \cdot \chi_{x+a}(z) \cdot \overline{\chi_{y+a}(z)} \\ &= \sum_{z \in D} \widehat{F}_{z,z} \cdot \chi_z(x+a) \cdot \overline{\chi_z(y+a)} \\ &= \sum_{z \in D} \widehat{F}_{z,z} \cdot \chi_z(x) \cdot \overline{\chi_z(y)} \\ &= f(x,y) \end{aligned}$$

In the second-last equation, we used the rotation-invariance of $\chi_z(x) \cdot \overline{\chi_z(y)}$. ■

Corollary 9 *Assume that $f(x,y)$ is a rotation-invariant distributed function over D and let $F_{x,y} = f(x,y)$ denote the corresponding matrix. Then the (complex) eigenvalues of $d^{-1} \cdot F$ are found on the main diagonal of \widehat{F} .*

Proof: Rewrite (14) as

$$d^{-1}F = (d^{-1/2}H) \cdot \widehat{F} \cdot (d^{-1/2}H^*)$$

and observe that this is nothing but the spectral decomposition of $d^{-1}F$ (since \widehat{F} is a diagonal matrix and $d^{-1/2}H$ is unitary). ■

We briefly note the following result:

Lemma 10 Let \hat{F} be the (diagonal) matrix that contains the Fourier-coefficients of the (rotation-invariant) distributed function $f(x - y)$. Then, for every $z \in D$, $\hat{f}(z) = \hat{F}_{z,z}$.

Proof: Consider the function $f_y(x) := f(x - y)$. We shall show below that the Fourier coefficients of f and f_y are related as follows:

$$\hat{f}_y(z) = \hat{f}(z) \cdot \overline{\chi_y(z)}. \quad (15)$$

The proof is now obtained by the following calculation:

$$\begin{aligned} \hat{F}_{z,z} &= d^{-2} \cdot \sum_{x,y \in D} f(x - y) \cdot \overline{\chi_x(x)} \cdot \chi_z(y) \\ &= d^{-1} \cdot \sum_{y \in D} \left(d^{-1} \cdot \sum_{x \in D} f_y(x) \overline{\chi_x(x)} \right) \chi_z(y) \\ &= d^{-1} \cdot \sum_{y \in D} \hat{f}_y(z) \cdot \chi_z(y) \\ &\stackrel{(15)}{=} \hat{f}(z) \cdot d^{-1} \cdot \sum_{y \in D} \underbrace{\overline{\chi_y(z)} \chi_z(y)}_{=1} \\ &= \hat{f}(z) \end{aligned}$$

The following calculation verifies (15):

$$\begin{aligned} \hat{f}_y(z) &= d^{-1} \cdot \sum_{x \in D} f(x - y) \cdot \overline{\chi_x(x)} \\ &= d^{-1} \cdot \sum_{x \in D} \sum_{w \in D} \hat{f}(w) \cdot \chi_w(x - y) \cdot \overline{\chi_x(x)} \\ &= d^{-1} \cdot \sum_{x \in D} \sum_{w \in D} \hat{f}(w) \cdot \chi_w(x) \cdot \overline{\chi_w(y)} \cdot \overline{\chi_x(x)} \\ &= d^{-1} \cdot \sum_{w \in D} \underbrace{\left(\sum_{x \in D} \chi_w(x) \cdot \overline{\chi_x(x)} \right)}_{=d \cdot \delta_{w,z}} \hat{f}(w) \cdot \overline{\chi_w(y)} \\ &= \hat{f}(z) \cdot \overline{\chi_z(y)} \end{aligned}$$

■

Corollary 9 and Lemma 10 yield the following.²

Corollary 11 Let F denote the matrix with entries $F_{x,y} = f(x - y)$. Then the spectrum of (complex) eigenvalues of $d^{-1} \cdot F$ coincides with the spectrum of (complex) Fourier-coefficients of f .

Consider the sign matrix $M^{f, \mathcal{G}_{rot}}$. From (8) and the definition of G_{rot} , we conclude that

$$M_{x,y}^{f, \mathcal{G}_{rot}} = f(x + y).$$

It follows that $M^{f, \mathcal{G}_{rot}}$ is a symmetric matrix. If f is real-valued, then $M^{f, \mathcal{G}_{rot}}$ has real eigenvalues. Note

²This result might be known, but we are not aware of an appropriate pointer to the literature.

that $M^{f, \mathcal{G}_{rot}}$ coincides with matrix $F_{x,y} = f(x - y)$ up to a permutation of columns (where the column indexed y is exchanged with the column indexed $-y$). Since the spectrum of eigenvalues (or singular values, resp.) of a matrix is left invariant under a permutation of columns, we obtain the following

Corollary 12 Let $f(x - y)$ be real-valued, and let F be the matrix with entries $F_{x,y} = f(x - y)$. Then, the following holds:

1. F coincides with the symmetric matrix $M^{f, \mathcal{G}_{rot}}$ up to a permutation of columns.
2. The spectrum of eigenvalues of $d^{-1} \cdot F$ coincides with the spectrum of (real) eigenvalues of $d^{-1} \cdot M^{f, \mathcal{G}_{rot}}$ and with the spectrum of Fourier-coefficients of f .

4.2 Margin and Dimension Bounds for Rotation-invariant Kernels

For every function $f : D \rightarrow \{-1, 1\}$,

$$\bar{\mu}_{rot}(f) := \bar{\mu}_{\mathcal{G}_{rot}}(f)$$

denotes the largest possible average margin that can be achieved by a linear arrangement for f resulting from a rotation-invariant kernel. As for the smallest possible dimension, parameter $d_{rot}(f)$ is understood analogously.

Corollary 13 Let D be a finite Abelian group of size d . Every function $f : D \rightarrow \{-1, 1\}$ satisfies

$$\bar{\mu}_{rot}(f) \leq \|\hat{f}\|_{\infty}. \quad (16)$$

In other words, no large margin classifier that employs a rotation-invariant kernel can achieve an average margin for f which exceeds $\|\hat{f}\|_{\infty}$.

Proof: According to Theorem 5,

$$\bar{\mu}_{rot}(f) \leq \frac{\|M^{f, \mathcal{G}_{rot}}\|}{\sqrt{|D| \cdot |\mathcal{G}_{rot}|}} = \frac{\|M^{f, \mathcal{G}_{rot}}\|}{d}.$$

We conclude from Corollary 12 that

$$\|M^{f, \mathcal{G}_{rot}}\| = \|F\| = d \cdot \|\hat{f}\|_{\infty},$$

which leads us to inequality (16). ■

Corollary 6 and 7 combined with Corollary 11 lead us to the following results:

Corollary 14 Let $d_{rot}(f)$ denote the smallest dimension of a feature space associated with a rotation-invariant kernel K that allows for a correct representation of f . Then, $d_{rot}(f) \geq \|\hat{f}\|_{\infty}^{-1}$.

Proof: According to Corollary 6, the corresponding feature space for the kernel must have dimension at least $\sqrt{|D| \cdot |\mathcal{G}_{rot}|} / \|M^{f, \mathcal{G}_{rot}}\| = d / \|M^{f, \mathcal{G}_{rot}}\|$. According to Corollary 12, the latter expression evaluates to $\|\hat{f}\|_{\infty}^{-1}$. ■

Corollary 15 Let \widehat{f}_i denote the i -th Fourier-coefficient of f , where $|\widehat{f}_1|, \dots, |\widehat{f}_d|$ are in decreasing order. Then,

$$d_{rot}(f) \cdot \sum_{i=1}^{d_{rot}(f)} |\widehat{f}_i|^2 \geq 1$$

Proof: From (9), we obtain

$$d_{rot}(f) \cdot \sum_{i=1}^{d_{rot}(f)} \sigma_i^2 \geq 1$$

where σ_i denotes the i -th largest singular value of $M^{f, \mathcal{G}_{rot}}$. We conclude from Corollary 12, that σ_i coincides with $|\widehat{f}_i|$. ■

5 Reflection-invariant Functions

In this section, we consider real-valued functions only. A distributed function $f(x, y)$ over $\{-1, 1\}^n$ is said to be *reflection-invariant* if, for all $x, y, a \in \{-1, 1\}^n$, the following holds:

$$f(x \circ a, y \circ a) = f(x, y) \quad (17)$$

Note that reflection-invariance corresponds to rotation-invariance with $(\mathbb{Z}_2^n, +)$ as the underlying (additive) Abelian group is or, equivalently, with $(\{-1, 1\}^n, \cdot)$ as the underlying (multiplicative) Abelian group. This is because the subgroup \mathcal{G}_{rot} of $\mathcal{S}(D)$ that we have used for rotation-invariant distributed functions collapses for $D = \{-1, 1\}^n$ (with a multiplicative group structure) to the following subgroup of $\mathcal{S}(\{-1, 1\}^n)$:

$$\mathcal{G}_{ref} = \{x \mapsto x \circ a : a \in \{-1, 1\}^n\}$$

Thus, reflection-invariant functions inherit all closure properties that hold, in general, for \mathcal{G} -invariant distributed functions (see the Pointwise Closure Property and Lemma 1 in Section 3):

Corollary 16 1. *The pointwise limit of reflection-invariant functions is a reflection-invariant function. Furthermore, if f_1, \dots, f_d are reflection-invariant functions and $g : \mathbb{R}^d \rightarrow \mathbb{R}$ is an arbitrary function, then*

$$g(f_1(x, y), \dots, f_d(x, y))$$

is reflection-invariant too.

2. *Reflection-invariant distributed functions over $\{-1, 1\}^n$ are closed under the usual matrix product and under the tensor-product of matrices.*

Furthermore, reflection-invariant functions inherit all properties that hold, in general, for distributed functions over a finite Abelian group:

- A reflection-invariant function $f(x, y)$ can be decomposed according to (2). Since $D = \{-1, 1\}^n$, the character χ_z coincides with the parity function induced by z , i.e., $\chi_z(x) = \prod_{z_i=-1} x_i$.

- The matrix \widehat{F} whose entries are the Fourier coefficients of f satisfies (13) where H is the matrix from (4). Since $D = \{-1, 1\}^n$, H equals the well-known $(2^n \times 2^n)$ -Walsh-Hadamard matrix.

Distributed functions $f(x, y)$ over \mathbb{R}^n that satisfy (17) for all $x, y \in \mathbb{R}^n$ and every $a \in \{-1, 1\}^n$ are said to be *reflection-invariant in the Euclidean space*. Here are some examples (with some overlap to our exemplification of rotation-invariant functions in Section 4):

- A distributed function of the form $f(x, y) = g(x \circ y)$ is reflection-invariant (in the Euclidean space provided that the domain is \mathbb{R}^n):

$$g((x \circ a) \circ (y \circ a)) = g(x \circ y \circ (a \circ a)) = g(x \circ y)$$

Conversely, any reflection-invariant function $f(x, y)$ (over domain $\{-1, 1\}^n$) can be written in this form by setting $g(x) := f(x, \vec{e})$ because reflection-invariance implies that

$$f(x, y) = f(x \circ y, y \circ y) = f(x \circ y, \vec{e}) = g(x \circ y) .$$

- Because of the obvious identity

$$\chi_z(x \circ y) = \chi_z(x) \cdot \chi_z(y) ,$$

the distributed function $\chi_z(x) \cdot \chi_z(y)$ is reflection-invariant too.

- The metric

$$L_p(x - y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

induced by the L_p -norm is clearly reflection-invariant in the Euclidean space.

In Section 6, we shall see that many popular kernel functions happen to be reflection-invariant.

The fact that $f(x, y) = g(x \circ y)$ is a reflection-invariant function can be restated as follows: any function $f(x, y)$ that can be cast as a function in $x_1 \cdot y_1, \dots, x_n \cdot y_n$ is reflection-invariant. Similarly, any function $f(x, y)$ that can be cast as a function in $L_p(x - y)$ (or, more generally, in $|x_1 - y_1|, \dots, |x_n - y_n|$) is reflection-invariant.

6 Reflection-invariant Kernels

In this section, we consider kernel functions $K(x, y)$ over the Boolean or over the Euclidean domain. In other words, $K(x, y)$ is a distributed function over $\{-1, 1\}^n$ or over \mathbb{R}^n with the additional property that every finite principal sub-matrix of K is symmetric and positive semidefinite. In Section 6.1, we demonstrate that the family of reflection-invariant kernels is quite rich and contains many popular kernels. In Section 6.2, we derive margin and dimension bounds for reflection-invariant kernels.

6.1 Examples and Closure Properties

Let us start with some examples. The following (quite popular) kernels (over \mathbb{R}^n except for the DNF-Kernel that has a Boolean domain) can be cast as functions in $x_1 \cdot y_1, \dots, x_n \cdot y_n$ or as functions in $\|x - y\|_2$ and are therefore reflection-invariant:

Polynomial Kernels: $K(x, y) = p(x^\top y)$ for an arbitrary polynomial p with positive coefficients.

All-subsets Kernel: $K(x, y) = \prod_{i=1}^n (1 + x_i y_i)$.

ANOVA Kernel: Let $1 \leq s \leq n$ and define

$$K_s(x, y) = \sum_{1 \leq i_1 < \dots < i_s \leq n} \prod_{j=1}^s x_{i_j} y_{i_j} .$$

DNF-Kernel: $K(x, y) = -1 + 2^{-n} \prod_{i=1}^n (x_i y_i + 3)$.

Exponential Kernels: $K(x, y) = e^{p(x^\top y)}$ for an arbitrary polynomial p with positive coefficients.

Gaussian Kernel: $K(x, y) = e^{-\|x-y\|_2^2/\sigma^2}$ for an arbitrary $\sigma > 0$.

These kernels have the usual nice properties like being efficiently evaluable although the number of (implicitly represented) features is exponentially large (or even infinite). Polynomial, Exponential, and Gaussian Kernels (first used in [2]) are found in almost any basic text-book that is relevant to the subject (e.g. [3]). The All-subsets Kernel is found in [18], and the ANOVA Kernel is found in [19]. As for the latter two kernels, see also [17]. The DNF-Kernel has been proposed in [16].³ The reader interested in more information about these (and other) kernels may consult the relevant literature. Here, we simply point to the fact that all kernels mentioned above are reflection-invariant.

We move on and consider the possibility of making new reflection-invariant kernels from kernels that are already known to be reflection-invariant. To this end, we briefly call into mind some basic closure properties of kernels:

Lemma 17 *Let K, K_1, K_2 be kernels, and let $c > 0$ be a positive constant. Then, the distributed functions*

$$\begin{aligned} K_1(x, y) + K_2(x, y) & , \quad c \cdot K(x, y) \\ K_1(x, y) \cdot K_2(x, y) & , \quad (K_1 \otimes K_2)[(u, x), (v, y)] \end{aligned}$$

are kernels too. Moreover, the pointwise limit of kernels yields a kernel.

³In [16], the kernel is defined over the Boolean domain $\{0, 1\}^n$. Our formula above is obtained from the formula in [16] by plugging in the affine transformation that identifies 1 with -1 and 0 with 1. A similar remark applies to the Monotone DNF-Kernel discussed at the end of this section.

The proof of Lemma 17 can be looked-up in [3], for example.

Corollary 18 *If K_1, \dots, K_d are kernels and $P : \mathbb{R}^d \rightarrow \mathbb{R}$ is a polynomial (or a converging power series) with positive coefficients, then*

$$P(K_1(x, y), \dots, K_d(x, y))$$

is a kernel too.

Note that closure properties of reflection-invariant functions (see Corollary 16) are comparably strong so that Lemma 17 and Corollary 18 remain valid (mutatis mutandis) for reflection-invariant kernels.

The following kernels (proposed in [11] and [9], respectively) define a new kernel-matrix K in terms of a given symmetric matrix B (called ‘‘similarity matrix’’ in this context):

Exponential Diffusion Kernel: For $\lambda \in \mathbb{R}$, define

$$K = e^{\lambda \cdot B} = \sum_{k \geq 0} \frac{\lambda^k}{k!} \cdot B^k .$$

von Neumann Diffusion Kernel: For $0 \leq \lambda < \|B\|^{-1}$, define

$$K = (I - \lambda \cdot B)^{-1} = \sum_{k \geq 0} \lambda^k \cdot B^k .$$

It follows from the closure properties of reflection-invariant functions that both diffusion kernels would inherit reflection-invariance from the underlying similarity matrix B .

The family of reflection-invariant kernels is quite rich. But here are two kernels (the first-one from [16], and the second-one from [12]) which are counterexamples:

Monotone DNF-Kernel:

$$K(x, y) = -1 + 2^{-2n} \prod_{i=1}^n (x_i y_i - x_j - y_j + 5) .$$

Spectrum Kernel: Here, $x, y \in \{-1, 1\}^n$ are considered as binary strings. For $1 \leq p \leq n$ and for every substring $u \in \{-1, 1\}^p$,

$$\Phi_v^p(x) = |\{(u, w) : x = uvw\}|$$

counts how often v occurs as a substring of x . The p -Spectrum Kernel is then given by

$$K(x, y) = \sum_{v \in \{-1, 1\}^p} \Phi_v^p(x) \cdot \Phi_v^p(y) .$$

It is easy to see that both kernels are *not* reflection-invariant. More generally, string kernels (measuring similarity between strings) often violate reflection-invariance.

6.2 Margin and Dimension Bounds for Reflection-invariant Kernels

For every function $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$,

$$\bar{\mu}_{ref}(f) := \bar{\mu}_{\mathcal{G}_{ref}}(f)$$

denotes the largest possible average margin that can be achieved by a linear arrangement for f resulting from a reflection-invariant kernel. Because reflection-invariance is a special case of rotation-invariance, the following result immediately follows from Corollaries 13, 14, and 15:

Corollary 19 *1. Every Boolean function f satisfies*

$$\bar{\mu}_{ref}(f) \leq \|\hat{f}\|_{\infty}.$$

In other words, no large margin classifier that employs a reflection-invariant kernel can achieve an average margin for f which exceeds $\|\hat{f}\|_{\infty}$.

- 2. Let $d_{ref}(f)$ denote the smallest dimension of a feature space associated with a reflection-invariant kernel K that allows for a correct representation of f . Then, $d_{ref}(f) \geq \|\hat{f}\|_{\infty}^{-1}$.*
- 3. Let \hat{f}_i denote the i -th Fourier-coefficient of f , where $|\hat{f}_1|, \dots, |\hat{f}_{2^n}|$ are in decreasing order. Then, $d_{ref}(f)$ satisfies the following lower bound:*

$$d_{ref}(f) \cdot \sum_{i=1}^{d_{ref}(f)} |\hat{f}_i|^2 \geq 1$$

7 Conclusions and Open Problems

We start with some remarks which offer a possible interpretation of our results. Finally, some open problems are mentioned.

7.1 Discussion of our Results

Ideally the invariance-properties of a kernel reflect symmetries in the data. For example, assume that there exists a set of transformations, say T , so that, for every instance $x \in D$ and every transformation $t \in T$, the label assigned to x by target function f equals the label assigned to $t(x)$ by f . Then, it looks desirable to apply a kernel that is invariant under the transformations from T . It would be surprising if our results implied that such kernels (that sort of perfectly model the symmetries in the data) would inherently lead to small margins or high-dimensional feature spaces. It is, however, easy to argue that (as expected) the contrary is true and our margin and dimension bounds trivialize whenever the invariance of the kernel perfectly matches with symmetries in the data. To see this, consider again (compare with the introduction) the “super-kernel”

$$K(x, y) = \begin{cases} +1 & \text{if } f(x) = f(y) \\ -1 & \text{otherwise} \end{cases}$$

that allows for a 1-dimensional halfspace representation of f with margin 1, and note that K actually *is* invariant under all transformations from T . Thus, no upper margin bound that holds uniformly for all T -invariant kernels can be smaller than 1. Similarly, no lower dimension bound can be larger than 1. Note that this is no contradiction to the main results in this paper because the family $\{f_t : t \in T\}$ of functions $f_t(x) = f(t(x))$ collapses to the singleton $\{f\}$. Thus Forster’s margin and dimension bounds applied to this family do not lead to non-trivial values.

Viewed from this perspective, our results can be interpreted as follows: one should *not* use a kernel that is invariant under a set T of transformations if T does *not* reflect symmetries in the data. The kernel becomes very poor especially when the family $\{f_t : t \in T\}$ contains much “orthogonality” (which is sort of the opposite of collapsing to a singleton or to a family of highly correlated functions) because Forster’s bounds, applied to pairwise (almost) orthogonal functions, are extremely strong.

This interpretation makes clear that our results are not particularly surprising but, on the other hand, quantify (in terms of small margin and large dimension bounds) in a meaningful and rigorous fashion an existing mismatch between a kernel and the (missing or existing) symmetries in the data.

7.2 Open Problems

Haasdonk and Burkhardt [8] consider two notions of invariance: “simultaneous invariance” and “total invariance”. Simultaneous invariance very much corresponds to the notion of invariance that we discussed in Section 3 so that our margin and dimension bounds apply. Total invariance is a stronger notion so that our bounds apply more than ever. But the obvious challenge is to find *stronger* margin and dimension bounds for *totally* invariant kernels.

The basic idea behind our paper is roughly as follows. For a family of kernels (e.g., polynomial kernels), we argue that the existence a “good representation” for a particular target function implies the existence of a “good representation” for a whole family of target functions (so that classical margin and dimension bounds can be brought into play). We think that invariance under a group operation (the notion considered in this paper) is just the first obvious thing one should consider. We would like to develop more versatile techniques that, while following the same basic idea, lead to strong margin and dimension bounds for a wider class of kernels.

References

- [1] Shai Ben-David, Nadav Eiron, and Hans Ulrich Simon. Limitations of learning via embeddings in euclidean half-spaces. *Journal of Machine Learning Research*, 3:441–461, 2002.

- [2] Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, pages 144–152, 1992.
- [3] Nello Cristianini and John Shawe-Taylor. *An Introduction to Support Vector Machines*. Cambridge University Press, 2000.
- [4] Jürgen Forster. A linear lower bound on the unbounded error communication complexity. *Journal of Computer and System Sciences*, 65(4):612–625, 2002.
- [5] Jürgen Forster, Matthias Krause, Satyanarayana V. Lokam, Rustam Mubarakzjanov, Niels Schmitt, and Hans Ulrich Simon. Relations between communication complexity, linear arrangements, and computational complexity. In *Proceedings of the 21st Annual Conference on the Foundations of Software Technology and Theoretical Computer Science*, pages 171–182, 2001.
- [6] Jürgen Forster, Niels Schmitt, Hans Ulrich Simon, and Thorsten Suttrop. Estimating the optimal margins of embeddings in euclidean half spaces. *Machine Learning*, 51(3):263–281, 2003.
- [7] Jürgen Forster and Hans Ulrich Simon. On the smallest possible dimension and the largest possible margin of linear arrangements representing given concept classes. *Theoretical Computer Science*, 350(1):40–48, 2006.
- [8] Bernard Haasdonk and Hans Burkhardt. Invariant kernel functions for pattern analysis and machine learning. *Machine Learning*, 68(1):35–61, 2007.
- [9] Jaz S. Kandola, John Shawe-Taylor, and Nello Cristianini. Learning semantic similarity. In *Advances in Neural Information Processing Systems 15*, pages 657–664. MIT Press, 2003.
- [10] Eike Kiltz and Hans Ulrich Simon. Threshold circuit lower bounds on cryptographic functions. *Journal of Computer and System Sciences*, 71(2):185–212, 2005.
- [11] Risi I. Kondor and John D. Lafferty. Diffusion kernels on graphs and other discrete input spaces. In *Proceedings of the 19th International Conference on Machine Learning*, pages 315–322, 2002.
- [12] Christina Leslie, Eleazar Eskin, and William S. Noble. The spectrum kernel: A string kernel for SVM protein classification. In *Pacific Symposium on Biocomputing*, pages 564–575, 2002.
- [13] Nathan Linial, Shahar Mendelson, Gideon Schechtman, and Adi Shraibman. Complexity measures of sign matrices. *Combinatorica*. To appear.
- [14] Nati Linial and Adi Shraibman. Lower bounds in communication complexity based on factorization norms. In *Proceedings of the 39th Annual Symposium on Theory of Computing*, pages 699–708, 2007.
- [15] Alexander A. Razborov and Alexander A. Sherstov. The sign-rank of AC^0 . Personal Communication.
- [16] Ken Sadohara. Learning of boolean functions using support vector machines. In *Proceedings of the 12th International Conference on Algorithmic Learning Theory*, pages 106–118, 2001.
- [17] John Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [18] Eiji Takimoto and Manfred K. Warmuth. Pathe kernels and multiplicative updates. *Journal of Machine Learning Research*, 4:773–818, 2003.
- [19] Vladimir Vapnik, Christopher J. C. Burges, Bernhard Schoelkopf, and R. Lyons. A new method for constructing artificial neural networks. Interim ARPA Technical Report, AT&T Bell Laboratories, 1995.