

# Kernel Subspace and Feature Extraction

Xiangxiang Xu and Lizhong Zheng

Dept. EECS, MIT

Cambridge, MA 02139, USA

Email: {xuxx, lizhong}@mit.edu

**Abstract**—We study kernel methods in machine learning from the perspective of feature subspace. We establish a one-to-one correspondence between feature subspaces and kernels and propose an information-theoretic measure for kernels. In particular, we construct a kernel from Hirschfeld–Gebelein–Rényi maximal correlation functions, coined the maximal correlation kernel, and demonstrate its information-theoretic optimality. We use the support vector machine (SVM) as an example to illustrate a connection between kernel methods and feature extraction approaches. We show that the kernel SVM on maximal correlation kernel achieves minimum prediction error. Finally, we interpret the Fisher kernel as a special maximal correlation kernel and establish its optimality.

## I. INTRODUCTION

One main objective of machine learning is to obtain useful information from often high-dimensional data. To this end, it is a common practice to extract meaningful feature representations from original data and then process features [1]. Neural networks [2] and kernel methods [3]–[6] are two of the most representative approaches to map data into feature space. In neural networks, the features are represented as the outputs of hidden neurons in the network. In contrast, the feature mapping in kernel methods is defined by the used kernel, which is used implicitly and is often infinite dimensional. While kernel approaches require much fewer parameters and can obtain good empirical performance on certain tasks [7], the performance significantly relies on the choice of kernels. With many attempts to investigate kernel methods [6], [8], [9], there still lacks a theoretical understanding of the mechanism behind kernel methods, which restricts their applications on complicated data.

On the other hand, the feature extraction in deep neural networks has been studied recently by information-theoretic and statistical analyses [10], [11]. For example, it was shown in [10] that, the feature extracted by deep neural networks coincides with the most informative feature, which is essentially related to the classical Hirschfeld–Gebelein–Rényi (HGR) maximal correlation problem [12]–[14]. Such theoretical characterizations provide a better understanding of existing algorithms and have been shown useful in designing algorithms for multimodal learning tasks [15].

In this paper, our goal is to characterize kernel methods from the perspective of feature subspace and reveal its connection with other learning approaches. We first introduce the associated kernel with each given feature subspace, which we coin the *projection kernel*, to establish a correspondence between

kernel operations and geometric operations in feature subspaces. This connection allows us to study kernels methods via analyzing the corresponding feature subspaces. Specifically, we propose an information-theoretic measure for projection kernels, and demonstrate that the information-theoretically optimal kernel can be constructed from the HGR maximal correlation functions, coined the *maximal correlation kernel*. We further demonstrate that the support vector machine (SVM) with maximal correlation kernel can obtain the minimum prediction error, which justifies its optimality in learning tasks. Our analysis also reveals connections between SVM and other classification approaches including neural networks. Finally, we interpret the Fisher kernel, a classical kernel induced from parameterized distribution families [16], as a special case of maximal correlation kernels, thus demonstrating its optimality.

## II. PRELIMINARIES AND NOTATIONS

Throughout this paper, we use  $X, Y$  to denote two random variables with alphabets  $\mathcal{X}, \mathcal{Y}$ , and denote their joint distribution and marginals as  $P_{X,Y}$  and  $P_X, P_Y$ , respectively. We also use  $\mathbb{E}[\cdot]$  to denote the expectation with respect to  $P_{X,Y}$ .

### A. Feature Space

We adopt the notation convention introduced in [15], and let  $\mathcal{F}_{\mathcal{X}} \triangleq \{\mathcal{X} \rightarrow \mathbb{R}\}$  denote the feature space formed by the (one-dimensional) features of  $X$ , with the geometry defined as follows. The inner product  $\langle \cdot, \cdot \rangle_{\mathcal{F}_{\mathcal{X}}}$  on  $\mathcal{F}_{\mathcal{X}}$  is defined as  $\langle f_1, f_2 \rangle_{\mathcal{F}_{\mathcal{X}}} \triangleq \mathbb{E}_{P_X}[f_1(X)f_2(X)]$  for  $f_1, f_2 \in \mathcal{F}_{\mathcal{X}}$ . This induces a norm  $\|\cdot\|_{\mathcal{F}_{\mathcal{X}}}$  with  $\|f\|_{\mathcal{F}_{\mathcal{X}}} \triangleq \sqrt{\langle f, f \rangle_{\mathcal{F}_{\mathcal{X}}}}$  for  $f \in \mathcal{F}_{\mathcal{X}}$ . Then, for given  $f \in \mathcal{F}_{\mathcal{X}}$  and subspace  $\mathcal{G}$  of  $\mathcal{F}_{\mathcal{X}}$ , we denote the projection of  $f$  onto  $\mathcal{G}$  as

$$\Pi(f; \mathcal{G}) \triangleq \arg \min_{h \in \mathcal{G}} \|h - f\|_{\mathcal{F}_{\mathcal{X}}}. \quad (1)$$

In addition, for a  $d$ -dimensional feature  $f = (f_1, \dots, f_d)^T: \mathcal{X} \rightarrow \mathbb{R}^d$ , we use  $\text{span}\{f\} \triangleq \text{span}\{f_1, \dots, f_d\}$  to denote the subspace spanned by all dimensions. We also use  $\tilde{f}$  to denote the centered  $f$ , i.e.,  $\tilde{f}(x) \triangleq f(x) - \mathbb{E}_{P_X}[f(X)]$ , and denote  $\Lambda_f \triangleq \mathbb{E}_{P_X}[f(X)f^T(X)]$ .

### B. Kernel

Given  $\mathcal{X}$ ,  $\kappa: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a kernel on  $\mathcal{X}$ , if for all finite subset  $\mathcal{J} \subset \mathcal{X}$ , the  $|\mathcal{J}|$  by  $|\mathcal{J}|$  matrix  $[\kappa(x, x')]_{x \in \mathcal{J}, x' \in \mathcal{J}}$  is positive semidefinite. For each kernel  $\kappa$ , we define the associated functional operator  $\tau: \mathcal{F}_{\mathcal{X}} \rightarrow \mathcal{F}_{\mathcal{X}}$  as

$$[\tau(f)](x) \triangleq \mathbb{E}_{P_X}[\kappa(X, x)f(X)], \quad (2)$$

and we use  $\kappa \leftrightarrow \tau$  to denote the correspondence between  $\kappa$  and  $\tau$ . Furthermore, we define the centered kernel  $\tilde{\kappa}: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  as

$$\tilde{\kappa}(x, x') \triangleq \kappa(x, x') - \bar{\kappa}(x) - \bar{\kappa}(x') + \mathbb{E}_{P_X}[\bar{\kappa}(X)], \quad (3)$$

where we have defined  $\bar{\kappa}: (x \mapsto \mathbb{E}_{P_X}[\kappa(X, x)]) \in \mathcal{F}_X$ .

The following fact is the basis of the kernel trick in learning algorithms.

*Fact 1:* For each given kernel  $\kappa$ , there exist an inner product space  $\mathcal{V}$  with the inner product  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$ , and a mapping  $\nu: \mathcal{X} \rightarrow \mathcal{V}$ , such that  $\kappa(x, x') = \langle \nu(x), \nu(x') \rangle_{\mathcal{V}}$ .

*Remark 1:* Suppose  $\nu$  is one mapping for  $\kappa$  satisfying Fact 1. Then for the centered kernel  $\tilde{\kappa}$  [cf. (3)], we have  $\tilde{\kappa}(x, x') = \langle \tilde{\nu}(x), \tilde{\nu}(x') \rangle_{\mathcal{V}}$ , where  $\tilde{\nu}(x) \triangleq \nu(x) - \mathbb{E}_{P_X}[\nu(X)]$ .

In addition, we introduce the kernelized discriminative model (KDM) as follows.

*Definition 1 (Kernelized Discriminative Model):* For each kernel  $\kappa$ , we define its associated kernelized discriminative model  $P_{Y|X}^{(\kappa)}$  as

$$P_{Y|X}^{(\kappa)}(y|x) \triangleq P_Y(y) \left( 1 + \mathbb{E} \left[ \tilde{\kappa}(X, x) \middle| Y = y \right] \right). \quad (4)$$

Then, we use  $\hat{y}^{(\kappa)}$  to denote the maximum a posteriori (MAP) estimation induced from KDM  $P_{Y|X}^{(\kappa)}$ , i.e.,

$$\hat{y}^{(\kappa)}(x) \triangleq \arg \max_{y \in \mathcal{Y}} P_{Y|X}^{(\kappa)}(y|x). \quad (5)$$

The KDM can be regarded as a generalized probability distribution, since we have  $\sum_{y \in \mathcal{Y}} P_{Y|X}^{(\kappa)}(y|x) = 1$  for all  $x \in \mathcal{X}$  while  $P_{Y|X}^{(\kappa)}(y|x)$  can sometimes take negative values.

### C. Modal Decomposition, Maximal Correlation, and H-score

We first introduce the modal decomposition of joint distribution  $P_{X,Y}$  [11], [15].

*Proposition 1 (Modal Decomposition [11]):* For given  $P_{X,Y}$ , there exists  $K \leq \min\{|\mathcal{X}|, |\mathcal{Y}|\} - 1$ , such that

$$P_{X,Y}(x, y) = P_X(x)P_Y(y) \left( 1 + \sum_{i=1}^K \sigma_i f_i^*(x) g_i^*(y) \right), \quad (6)$$

where  $\sigma_1 \geq \sigma_2 \geq \sigma_K > 0$ , and  $\mathbb{E}[f_i^*(X)f_j^*(X)] = \mathbb{E}[g_i^*(Y)g_j^*(Y)] = \mathbb{1}_{\{i=j\}}$  for all  $1 \leq i, j \leq K$ , where  $\mathbb{1}_{\{\cdot\}}$  denotes the indicator function.

It can be shown that  $(f_i^*, g_i^*)$  pairs are the most correlated function pairs of  $X$  and  $Y$ , referred to as maximal correlation functions. We also denote  $\varrho \triangleq \sigma_1$ , known as the HGR maximal correlation [12]–[14] of  $X$  and  $Y$ , and define the  $K$ -dimensional feature  $f^*(x) \triangleq [f_1^*(x), \dots, f_K^*(x)]^T$ . In particular, when  $Y$  is binary, we have  $f^* = f_1^* \in \mathcal{F}_X$ .

It has been shown in [11] that the maximal correlation functions  $f_i^*, i = 1, \dots, K$  are the optimal features of  $X$  in inferring or estimating  $Y$ . In general, given a  $d$ -dimensional feature  $f$  of  $X$ , the effectiveness of  $f$  in inferring or estimating  $Y$  can be measured by its *H-score* [10], [11], defined as

$$\mathcal{H}(f) \triangleq \frac{1}{2} \cdot \mathbb{E} \left[ \left\| \mathbb{E} \left[ \Lambda_f^{-\frac{1}{2}} \tilde{f}(X) \middle| Y \right] \right\|^2 \right], \quad (7)$$

where  $\tilde{f}(x) \triangleq f(x) - \mathbb{E}[f(X)]$ . It can be verified that for all  $d$  and  $f: \mathcal{X} \rightarrow \mathbb{R}^d$ , we have

$$\mathcal{H}(f) \leq \mathcal{H}(f^*) = \frac{1}{2} \sum_{i=1}^K \sigma_i^2, \quad (8)$$

where  $\sigma_1, \dots, \sigma_K$  are as defined in (6).

### D. Binary Classification

We consider the binary classification problem which predicts binary label  $Y$  from the data variable  $X$ . For convenience, we assume  $Y$  takes values from  $\mathcal{Y} \triangleq \{-1, 1\}$ .

Suppose the training dataset contains  $n$  sample pairs  $\{(x_i, y_i)\}_{i=1}^n$  of  $(X, Y)$ , and let  $P_{X,Y}$  denote the corresponding empirical distribution, i.e.,

$$P_{X,Y}(x, y) \triangleq \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i=x, y_i=y\}}. \quad (9)$$

*1) Support Vector Machine:* The support vector machine (SVM) solves binary classification tasks by finding the optimal hyperplane that separates two classes with maximum margin [7]. Given  $d$ -dimensional feature mapping  $f: \mathcal{X} \rightarrow \mathbb{R}^d$ , the loss for SVM based on  $f$  can be written as

$$L_{\text{SVM}}(f, w, b; \lambda) \triangleq \mathbb{E}_{P_{X,Y}}[\ell_{\text{hinge}}(Y, \langle w, f(X) \rangle + b)] + \frac{\lambda}{2} \cdot \|w\|^2, \quad (10)$$

where  $w, b \in \mathbb{R}^d$  are the parameters of the hyperplane, where  $\lambda > 0$  is a hyperparameter of SVM, and where  $\ell_{\text{hinge}}: \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}$  denotes the hinge loss, defined as  $\ell_{\text{hinge}}(y, z) \triangleq (1 - yz)^+$  with  $x^+ \triangleq \max\{0, x\}$ .

Moreover, let  $(w_{\text{SVM}}, b_{\text{SVM}}) \triangleq \arg \min_{w, b} L_{\text{SVM}}(f, w, b; \lambda)$

and  $L_{\text{SVM}}^*(f; \lambda) \triangleq L_{\text{SVM}}(f, w_{\text{SVM}}, b_{\text{SVM}}; \lambda)$  denote the optimal parameters and the value of loss function, respectively. Then, the prediction of SVM is

$$\hat{y}_{\text{SVM}}(x; f, \lambda) \triangleq \text{sgn}(\langle w_{\text{SVM}}, f(x) \rangle + b_{\text{SVM}}), \quad (11)$$

where  $\text{sgn}(\cdot)$  denotes the sign function.

Specifically, for a given kernel  $\kappa$ , the prediction of the corresponding kernel SVM is  $\hat{y}_{\text{SVM}}^{(\kappa)}(x; \lambda) \triangleq \hat{y}_{\text{SVM}}(x; \nu, \lambda)$ , where  $\nu$  is any mapping given by Fact 1.

*2) Logistic Regression and Neural Networks:* Given  $d$ -dimensional feature  $f$  of  $X$ , the discriminative model of logistic regression is  $\tilde{P}_{Y|X}(y|x; f, w, b) \triangleq \text{sigmoid}(y \cdot (\langle w, f(x) \rangle + b))$ , where  $w \in \mathbb{R}^d, b \in \mathbb{R}$  are the weight and bias, respectively, and where  $\text{sigmoid}(\cdot)$  is defined as  $\text{sigmoid}(x) \triangleq \frac{1}{1 + \exp(-x)}$ .

Then, the loss of logistic regression is  $L_{\text{LR}}(f, w, b) \triangleq -\mathbb{E}[\log \tilde{P}_{Y|X}(Y|X; f, w, b)]$ , and the optimal parameters  $w_{\text{LR}}, b_{\text{LR}}$  are learned by minimizing the loss, i.e.,  $(w_{\text{LR}}, b_{\text{LR}}) \triangleq \arg \min_{w, b} L_{\text{LR}}(f, w, b)$ . The resulting decision rule is

$$\hat{y}_{\text{LR}}(x; f) \triangleq \arg \max_{y \in \mathcal{Y}} \tilde{P}_{Y|X}(y|x; f, w_{\text{LR}}, b_{\text{LR}})$$

<sup>1</sup>It is worth mentioning that the practical implementation of kernel SVM is typically done by solving a dual optimization problem without explicitly using  $\nu$ . See [17, Section 12] for detailed discussions.

$$= \text{sgn}(\langle w_{\text{LR}}, f(x) \rangle + b_{\text{LR}}). \quad (12)$$

The logistic regression is often used as the classification layer for multi-layer neural networks, where  $w$  and  $b$  correspond to weights and the bias term, respectively. In this case, the feature mapping  $f(\cdot)$  also takes a parameterized form, and the parameters of  $f(\cdot)$  are jointly learned with  $w$  and  $b$ .

Due to the space limitations, we omit most proofs in the rest of this paper, but refer the readers to the extended version of this paper [18] for the details.

### III. PROJECTION KERNEL AND INFORMATIVE FEATURES

In this section, we introduce a one-to-one correspondence between kernels and feature subspaces, and then characterize the informativeness of kernels by investigating the features in the associated subspaces.

#### A. Projection Kernel and Feature Subspace

We first introduce a family of kernels with one-to-one correspondence to feature subspace.

*Definition 2 (Projection Kernel):* Let  $\mathcal{G}$  denote a  $d$ -dimensional subspace of  $\mathcal{F}_X$  with a basis  $\{f_1, \dots, f_d\}$ . We use  $\kappa_{\mathcal{G}}: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  to denote the projection kernel associated with  $\mathcal{G}$ , defined as  $\kappa_{\mathcal{G}}(x, x') \triangleq f^{\text{T}}(x)\Lambda_f^{-1}f(x')$ , where we have defined  $f \triangleq (f_1, \dots, f_d)^{\text{T}}$  and  $\Lambda_f \triangleq \mathbb{E}[f(X)f^{\text{T}}(X)]$ .

With slight abuse of notation, we also denote  $\kappa_f \triangleq \kappa_{\text{span}\{f\}}$ , the projection kernel associated with  $\text{span}\{f\}$ .

Note that  $\kappa_{\mathcal{G}}$  is a valid kernel function, and the corresponding  $\nu$  mapping in Fact 1 can be chosen as  $\nu(x) = [f_1(x), \dots, f_d(x)]^{\text{T}}$  for any orthonormal basis  $\{f_1, \dots, f_d\}$  of  $\mathcal{G}$ . It turns out that the functional operators associated with projection kernels are projection operators in the feature space, which we formalize as follows.

*Property 1:* Let  $\tau \leftrightarrow \kappa_{\mathcal{G}}$  denote the operator corresponding to subspace  $\mathcal{G}$  [cf. (2)], then we have  $\tau(f) = \Pi(f; \mathcal{G})$  for all  $f \in \mathcal{F}_X$ .

Therefore, given a projection kernel  $\kappa$ , the associated subspace can be represented as  $\{f \in \mathcal{F}_X: \tau(f) = f\}$ , where  $\tau \leftrightarrow \kappa$  is the associated operator. This establishes a one-to-one correspondence between projection kernels and feature subspaces.

#### B. H-score and Informative Features

The projection kernel provides a connection between feature subspace and kernel, from which we can characterize subspace  $\mathcal{G}$  in terms of the corresponding kernel  $\kappa_{\mathcal{G}}$ . Specifically, we can represent the H-score [cf. (7)] of a feature  $f$  in terms of the projection kernel  $\kappa_f$ , formalized as follows.

*Proposition 2:* For all  $f$  with  $\text{span}\{f\} = \mathcal{G}$ , we have  $\mathcal{H}(f) = \frac{1}{2} \cdot (\mathbb{E}_{P_{X, X'}}[\kappa_{\mathcal{G}}(X, X')] - \mathbb{E}_{P_X P_{X'}}[\kappa_{\mathcal{G}}(X, X')])$ , where we have defined  $X'$  such that the joint distribution of  $X$  and  $X'$  is

$$P_{X, X'}(x, x') \triangleq \sum_{y \in \mathcal{Y}} P_Y(y) P_{X|Y=y}(x) P_{X|Y=y}(x'). \quad (13)$$

With slight abuse of notation, we can use  $\mathcal{H}(\mathcal{G})$  to denote the H-score corresponding to feature subspace  $\mathcal{G}$ . In particular,

we have the following characterization of  $\mathcal{H}(\mathcal{G})$  when  $Y$  is binary.

*Proposition 3:* Suppose  $Y$  is binary, and  $f^*$  is the maximal correlation function of  $P_{X, Y}$ . Then, for each subspace  $\mathcal{G}$  of  $\mathcal{F}_X$ , we have

$$\mathcal{H}(\mathcal{G}) = \frac{\rho^2}{2} \cdot \|\Pi(f^*; \mathcal{G})\|_{\mathcal{F}_X}^2 = \max_{f \in \mathcal{G}} \mathcal{H}(f) = \mathcal{H}(\Pi(f^*; \mathcal{G})). \quad (14)$$

From Proposition 3,  $\mathcal{H}(\mathcal{G})$  depends only on the projection of  $f^*$  onto  $\mathcal{G}$ , which is also the most informative feature in  $\mathcal{G}$ . In addition, note that since  $\|f^*\|_{\mathcal{F}_X} = 1$ ,  $\|\Pi(f^*; \mathcal{G})\|_{\mathcal{F}_X}^2$  is also the cosine value of the principal angle between  $f^*$  and  $\mathcal{G}$ . Therefore, we can interpret the H-score as a measure of the principal angle between the optimal feature  $f^*$  and the given subspace.

#### C. Maximal Correlation Kernel

Note that from (8),  $\mathcal{H}(f)$  is maximized when  $f$  takes the maximal correlation function  $f^*$ . Therefore, the subspace  $\text{span}\{f^*\}$  (and thus projection kernel  $\kappa_{f^*}$ ) is optimal in terms of the H-score measure. We will denote  $\kappa^* \triangleq \kappa_{f^*}$ , referred to as the *maximal correlation kernel*.

Specifically, the KDM (cf. Definition 1) of maximal correlation kernel  $\kappa^*$  coincides with the underlying conditional distribution  $P_{Y|X}$ , demonstrated as follows.

*Property 2:* For all  $x$  and  $y$ , we have  $P_{Y|X}(y|x) = P_{Y|X}^{(\kappa^*)}(y|x)$  and  $\hat{y}^{(\kappa^*)}(x) = \hat{y}_{\text{MAP}}(x)$ , where  $\hat{y}_{\text{MAP}}$  denotes the MAP estimation, i.e.,

$$\hat{y}_{\text{MAP}}(x) \triangleq \arg \max_{y \in \mathcal{Y}} P_{Y|X}(y|x). \quad (15)$$

As we will develop in the next section, the maximal correlation kernel also achieves the optimal performance in support vector machine.

### IV. SUPPORT VECTOR MACHINE ANALYSIS

In this section, we investigate support vector machine, a representative kernel approach for binary classification. Let  $(X, Y)$  denote the training data and corresponding label taken from  $\mathcal{Y} = \{-1, 1\}$ , with  $P_{X, Y}$  denoting the empirical distribution as defined in (9). Throughout this section, we will focus on the balanced dataset with

$$P_Y(-1) = P_Y(1) = \frac{1}{2}. \quad (16)$$

It can be verified that in this case, the MAP estimation [cf. (15)] can be expressed in terms of maximal correlation function.

*Property 3:* Under assumption (16), we can express the MAP estimation as  $\hat{y}_{\text{MAP}}(x) = \text{sgn}(f^*(x))$  for all  $x \in \mathcal{X}$ , where  $f^* \in \mathcal{F}_X$  is the maximal correlation function of  $P_{X, Y}$ .

#### A. SVM on Given Features

We first consider the SVM algorithm applied on a given feature representation  $f(X) \in \mathbb{R}^d$ , which can also be regarded as the kernel SVM on kernel  $\kappa(x, x') = \langle f(x), f(x') \rangle$ .

To begin, for each given feature  $f$  and  $\lambda > 0$ , let us define

$$\hat{L}(f; \lambda) \triangleq 1 - \frac{1}{2\lambda} \cdot \|\mathbb{E}[f(X)Y]\|^2.$$

Then we have the following characterization.

*Theorem 1:* For all given feature  $f$  and  $\lambda \geq 0$ , we have

$$\hat{L}(f; \lambda) \leq L_{\text{SVM}}^*(f; \lambda) \leq \hat{L}(f; \lambda) + \left(\frac{\lambda_{\text{T}}}{\lambda} - 1\right)^+, \quad (17)$$

where we have defined  $\lambda_{\text{T}} \triangleq M \cdot \|\mathbb{E}[f(X)Y]\|$  and  $M \triangleq \max_{x \in \mathcal{X}} \|\tilde{f}(x)\|$ , with  $\tilde{f}(x) \triangleq f(x) - \mathbb{E}[f(X)]$ , and where  $x^+ \triangleq \max\{0, x\}$ .

Specifically, when  $\lambda \geq \lambda_{\text{T}}$ , we have  $L_{\text{SVM}}^*(f; \lambda) = \hat{L}(f; \lambda)$ , which can be achieved by

$$w_{\text{SVM}} = \frac{1}{\lambda} \cdot \mathbb{E}[f(X)Y], \quad b_{\text{SVM}} = -\langle w_{\text{SVM}}, \mathbb{E}[f(X)] \rangle, \quad (18)$$

and the resulting SVM prediction is

$$\hat{y}_{\text{SVM}}(x; f, \lambda) = \text{sgn}\left(\left\langle \mathbb{E}[f(X)Y], \tilde{f}(x) \right\rangle\right) \quad (19)$$

$$= \arg \min_{y \in \mathcal{Y}} \|f(x) - \mathbb{E}[f(X)|Y=y]\|. \quad (20)$$

From Theorem 1, when  $\lambda \geq \lambda_{\text{T}}$ , the SVM decision  $\hat{y}_{\text{SVM}}(x; f, \lambda)$  does not depend on the value of  $\lambda$ . In the remaining, we will focus on the regime where  $\lambda \geq \lambda_{\text{T}}$ , and drop the  $\lambda$  in expressions whenever possible, e.g., we simply denote  $\hat{y}_{\text{SVM}}(x; f, \lambda)$  by  $\hat{y}_{\text{SVM}}(x; f)$ . As we will see soon, SVM can still obtain minimum prediction error in this regime, by using a good feature mapping  $f$  (or equivalently, a good kernel).

From (20), the SVM prediction can be interpreted as a nearest centroid classifier, where decision is based on comparing the distance between  $f(x)$  and the class centroids  $\mathbb{E}[f(X)|Y=y]$ ,  $y \in \mathcal{Y}$ . In addition, from

$$\begin{aligned} \mathbb{E}[f(X)Y] &= \mathbb{E}[Y \cdot \mathbb{E}[f(X)|Y]] \\ &= \frac{1}{2}(\mathbb{E}[f(X)|Y=1] - \mathbb{E}[f(X)|Y=-1]), \end{aligned}$$

we can interpret the SVM loss  $L_{\text{SVM}}^* = \hat{L}$  as measuring the distance between two class centroids.

Furthermore, when  $f$  is one-dimensional feature, we can rewrite (19) as

$$\hat{y}_{\text{SVM}}(x; f) = \text{sgn}\left(\left\langle \mathbb{E}[f(X)Y], \tilde{f}(x) \right\rangle\right) = \text{sgn}\left(\hat{f}(x)\right),$$

where  $\hat{f} \triangleq \Pi(f^*; \text{span}\{\tilde{f}\})$ . Therefore, the decision rule depends only the projection of  $f^*$  onto the subspace  $\text{span}\{\tilde{f}\}$ , which is also the most informative features on the subspace (cf. Proposition 3). Later on we will see a similar geometric illustration of kernel SVM.

Moreover, we can establish a connection between SVM loss and the H-score measure, formalized as the following corollary.

*Corollary 1:* Suppose  $\lambda \geq \lambda_{\text{T}}$ , then we have

$$1 - \frac{r_{\text{max}}}{\lambda} \cdot \mathcal{H}(\tilde{f}) \leq L_{\text{SVM}}^*(f; \lambda) \leq 1 - \frac{r_{\text{min}}}{\lambda} \cdot \mathcal{H}(\tilde{f}),$$

where  $r_{\text{max}}$  and  $r_{\text{min}}$  denote the maximum and minimum positive eigenvalues of the covariance matrix  $\Lambda_{\tilde{f}}$ , respectively. Specifically, if  $\Lambda_{\tilde{f}} = I$ , then we have  $L_{\text{SVM}}^*(f; \lambda) = 1 - \lambda^{-1} \cdot \mathcal{H}(\tilde{f})$ .

As a result, for each normalized feature  $f$  with covariance matrix  $\Lambda_{\tilde{f}} = I$ , the SVM loss  $L_{\text{SVM}}^*$  measures the informativeness of  $f$  in inferring the label  $Y$ .

## B. Kernel SVM

In practice, instead of applying SVM on a given or manually designed feature  $f$ , it is more often to directly implement SVM on a kernel  $\kappa$ . Similar to Theorem 1, we have the following characterization, from which we can interpret KDM as a probabilistic output for kernel SVM.

*Theorem 2:* For each given kernel  $\kappa$ , there exists a constant  $\lambda_{\text{T}} > 0$ , such that when  $\lambda \geq \lambda_{\text{T}}$ , the SVM prediction is  $\hat{y}_{\text{SVM}}^{(\kappa)}(x) = \text{sgn}([\tau(f^*)](x))$ , where  $\tau \leftrightarrow \tilde{\kappa}$  is the operator associated with centered kernel  $\tilde{\kappa}$  [cf. (2) and (3)]. In addition, the SVM prediction coincides with the KDM prediction (cf. Definition 1) obtained from  $\kappa$ , i.e., we have  $\hat{y}_{\text{SVM}}^{(\kappa)}(x) = \hat{y}^{(\kappa)}(x)$  for all  $x \in \mathcal{X}$ .

*Proof:* Let  $\mathcal{V}$  and  $\nu: \mathcal{X} \rightarrow \mathcal{V}$  denote the inner product space and mapping associated with kernel  $\kappa$  (cf. Fact 1), and let  $\tilde{\nu}(x) \triangleq \nu(x) - \mathbb{E}_{P_X}[\nu(X)]$ . Then, we have

$$\begin{aligned} \langle \mathbb{E}[\tilde{\nu}(X)Y], \tilde{\nu}(x) \rangle_{\mathcal{V}} &= \mathbb{E}[\langle \tilde{\nu}(X), \tilde{\nu}(x) \rangle_{\mathcal{V}} \cdot Y] \\ &= \mathbb{E}[\tilde{\kappa}(X, x) \cdot Y], \end{aligned} \quad (21)$$

which can be rewritten as

$$\begin{aligned} &\mathbb{E}[\tilde{\kappa}(X, x) \cdot Y] \\ &= \mathbb{E}_{P_{X,Y}}[\tilde{\kappa}(X, x) \cdot Y] \\ &= \mathbb{E}_{P_X P_Y}[\tilde{\kappa}(X, x) \cdot Y \cdot (1 + \varrho \cdot f^*(X) \cdot Y)] \\ &= \mathbb{E}[\tilde{\kappa}(X, x)] \cdot \mathbb{E}[Y] + \varrho \cdot \mathbb{E}[\tilde{\kappa}(X, x) f^*(X)] \cdot \mathbb{E}[Y^2] \\ &= \varrho \cdot \mathbb{E}[\tilde{\kappa}(X, x) f^*(X)] \\ &= \varrho \cdot [\tau(f^*)](x), \end{aligned}$$

where to obtain the second equality we have used the modal decomposition of  $P_{X,Y}$ .

Hence, from Theorem 1 we obtain

$$\begin{aligned} \hat{y}_{\text{SVM}}^{(\kappa)}(x) &= \hat{y}_{\text{SVM}}(x; \nu) = \text{sgn}(\langle \mathbb{E}[\tilde{\nu}(X)Y], \tilde{\nu}(x) \rangle) \\ &= \text{sgn}\left(\mathbb{E}[\tilde{\kappa}(X, x) \cdot Y]\right) \\ &= \text{sgn}([\tau(f^*)](x)). \end{aligned}$$

It remains only to establish the equivalence between  $\hat{y}_{\text{SVM}}^{(\kappa)}$  and KDM decision  $\hat{y}^{(\kappa)}$ . To this end, note that from (4) and the balanced dataset assumption (16), we have

$$\begin{aligned} P_{Y|X}^{(\kappa)}(y|x) &= P_Y(y) \left(1 + \mathbb{E}[\tilde{\kappa}(X, x) | Y=y]\right) \\ &= \frac{1}{2} \left(1 + y \cdot \mathbb{E}[\tilde{\kappa}(X, x) Y]\right) \end{aligned}$$

for all  $x \in \mathcal{X}, y \in \mathcal{Y}$ .

Hence, for all  $x \in \mathcal{X}$ ,

$$\begin{aligned} \hat{y}^{(\kappa)}(x) &= \arg \max_{y \in \mathcal{Y}} P_{Y|X}^{(\kappa)}(y|x) = \operatorname{sgn}\left(\mathbb{E}\left[\tilde{\kappa}(X, x)Y\right]\right) \\ &= \hat{y}_{\text{SVM}}^{(\kappa)}(x), \end{aligned}$$

which completes the proof.  $\blacksquare$

From Theorem 2, the final decision  $\hat{y}_{\text{SVM}}^{(\kappa)}$  depends on  $\kappa$  only through the centered kernel  $\tilde{\kappa}$ . Moreover, compare Theorem 2 with Property 3, kernel SVM prediction differs from MAP only in applying the operator  $\tau$  on  $f^*$ . In particular, when the maximal correlation function  $f^*$  is an eigenfunction of the corresponding operator  $\tau \leftrightarrow \tilde{\kappa}$ , i.e.,  $\tau(f^*) = c \cdot f^*$  for some  $c > 0$ , the SVM prediction coincides with the MAP prediction, i.e.,  $\hat{y}_{\text{SVM}}^{(\kappa)}(x) = \hat{y}_{\text{MAP}}(x)$  for all  $x \in \mathcal{X}$ .

If we restrict our attention to projection kernels, the kernel SVM decision can be further interpreted as a projection operation on the associated subspace. To see this, let  $\mathcal{G}$  denote a feature subspace of  $\mathcal{F}_X$  spanned by zero-mean features, then from Theorem 1 and Proposition 3, the kernel SVM loss for  $\kappa_{\mathcal{G}}$  is

$$1 - \frac{1}{\lambda} \cdot \mathcal{H}(\mathcal{G}) = 1 - \frac{\varrho^2}{2\lambda} \cdot \|\Pi(f^*; \mathcal{G})\|_{\mathcal{F}_X}^2,$$

which measures the principal angle between  $f^*$  and  $\mathcal{G}$ . In addition, the decision rule can be expressed as

$$\hat{y}_{\text{SVM}}^{(\kappa_{\mathcal{G}})}(x) = \operatorname{sgn}(\|\Pi(f^*; \mathcal{G})\|(x)), \quad (22)$$

From Proposition 3,  $\Pi(f^*; \mathcal{G})$  is also the most informative feature in  $\mathcal{G}$ . Therefore, kernel SVM on  $\kappa_{\mathcal{G}}$  is equivalent to first extracting the most informative feature in  $\mathcal{G}$ , and then using the extracted feature to make decision.

### C. Relationship to Other Classification Approaches

1) *Maximum a Posteriori (MAP) Estimation:* From (22), when the maximal correlation kernel  $\kappa^*$  is applied, the kernel SVM decision is  $\operatorname{sgn}(f^*(x))$ , which coincides with the MAP prediction (cf. Property 3). Since MAP achieves the minimum prediction error, kernel SVM on the maximal correlation kernel also obtains the minimum prediction error.

2) *Logistic Regression and Neural Networks:* We have interpreted SVM as extracting the most informative feature, where the informativeness is measured by H-score. The analysis in [10] has shown that logistic regression is also equivalent to maximizing the H-score, when  $X$  and  $Y$  are weakly independent. Indeed, we can show that SVM and logistic regression lead to the same prediction in a weak dependence regime, which we formalize as follows.

*Proposition 4:* Suppose  $\varrho = O(\epsilon)$  for some  $\epsilon > 0$ . For SVM and logistic regression applied on feature  $f: \mathcal{X} \rightarrow \mathbb{R}^d$  with covariance  $\Lambda_f = I_d$ , the optimal parameters satisfy

$$\begin{aligned} w_{\text{LR}} &= 2\lambda \cdot w_{\text{SVM}} + o(\epsilon), \\ b_{\text{LR}} &= 2\lambda \cdot b_{\text{SVM}} + o(\epsilon), \end{aligned}$$

where  $\lambda$  is the hyperparameter in SVM. In addition, we have  $\hat{y}_{\text{SVM}}(x; f) = \hat{y}_{\text{LR}}(x; f)$  for  $\epsilon$  sufficiently small.

*Remark 2:* Since H-score can also be directly maximized by implementing the maximal correlation regression [19], a similar connection holds for SVM and maximal correlation regression.

### V. FISHER KERNEL

We demonstrate that *Fisher kernel* [16], [20] can also be interpreted as a maximal correlation kernel.

Given a family of distributions  $\pi(\cdot; \theta)$  supported on  $\mathcal{X}$  and parameterized by  $\theta \in \mathbb{R}^m$ , suppose the score function  $s_{\theta}(x) \triangleq \frac{\partial}{\partial \theta} \log \pi(x; \theta)$  exists. Then, the Fisher kernel is defined as the projection kernel associated with the score function  $s_{\theta}$ , i.e.,  $\kappa_{s_{\theta}}$ .

Specifically, we consider classification tasks where the joint distribution between data variable  $X$  and label  $Y$  are a mixture of the parameterized forms. Suppose for each class  $Y = y \in \mathcal{Y}$ , the data variable  $X$  is generated from

$$P_{X|Y}(x|y) = \pi(x; \theta_y) \quad (23)$$

for some  $\theta_y \in \mathbb{R}^m$ . Then we have the following result.

*Theorem 3:* Suppose  $\|\theta_y\| < \epsilon$  for all  $y \in \mathcal{Y}$ , and let  $s(x) \triangleq s_0(x)$ . Then for the joint distribution  $P_{X,Y} = P_{X|Y}P_Y$  generated according to (23), we have

$$P_{X,Y}(x, y) = P_X(x)P_Y(y) \left(1 + \langle s(x), \tilde{\theta}_y \rangle\right) + o(\epsilon), \quad (24)$$

$$\kappa_s(x, x') = \kappa^*(x, x') + o(\epsilon), \quad (25)$$

where  $\tilde{\theta}_y \triangleq \theta_y - \mathbb{E}[\theta_Y]$  denotes the centered  $\theta_y$ , and where  $\kappa^*$  is the maximal correlation kernel defined on  $P_{X,Y}$ . In addition, the H-score of  $s$  satisfies

$$\mathcal{H}(s) = I(X; Y) + o(\epsilon^2), \quad (26)$$

where  $I(X; Y)$  denotes the mutual information between  $X$  and  $Y$ .

From (24), the score function  $s$  is equal to the maximal correlation function  $f^*$  of  $P_{X,Y}$  up to a linear transformation [cf. (6)], and we have  $P_{Y|X}(y|x) = P_{Y|X}^{(\kappa^*)}(y|x) = P_{Y|X}^{(\kappa_s)}(y|x) + o(\epsilon)$ . Therefore, Fisher kernel is the optimal kernel for tasks generated from (23).

### VI. CONCLUSION

In this paper, we study kernel methods from the perspective of feature subspace, where we demonstrate a connection between kernel methods and informative feature extraction problems. With SVM as an example, we illustrate the relationship between kernel methods and neural networks. The theoretical results can help guide practical kernel designs and incorporate kernel methods with feature-based learning approaches.

### ACKNOWLEDGMENTS

This work was supported in part by the National Science Foundation (NSF) under Award CNS-2002908 and the Office of Naval Research (ONR) under grant N00014-19-1-2621.

## REFERENCES

- [1] D. Storcheus, A. Rostamizadeh, and S. Kumar, "A survey of modern questions and challenges in feature extraction," in *Feature Extraction: Modern Questions and Challenges*. PMLR, 2015, pp. 1–18.
- [2] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [3] S. Mika, G. Ratsch, J. Weston, B. Scholkopf, and K.-R. Mullers, "Fisher discriminant analysis with kernels," in *Neural networks for signal processing IX: Proceedings of the 1999 IEEE signal processing society workshop (cat. no. 98th8468)*. Ieee, 1999, pp. 41–48.
- [4] F. R. Bach and M. I. Jordan, "Kernel independent component analysis," *Journal of machine learning research*, vol. 3, no. Jul, pp. 1–48, 2002.
- [5] T. Hofmann, B. Schölkopf, and A. J. Smola, "Kernel methods in machine learning," *The annals of statistics*, vol. 36, no. 3, pp. 1171–1220, 2008.
- [6] B. Schölkopf, "The kernel trick for distances," *Advances in neural information processing systems*, vol. 13, 2000.
- [7] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995.
- [8] J. Platt *et al.*, "Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods," *Advances in large margin classifiers*, vol. 10, no. 3, pp. 61–74, 1999.
- [9] J. Xu, R. Jenssen, A. Paiva, and I. Park, "A reproducing kernel hilbert space framework for itl," in *Information Theoretic Learning*. Springer, 2010, pp. 351–384.
- [10] X. Xu, S.-L. Huang, L. Zheng, and G. W. Wornell, "An information theoretic interpretation to deep neural networks," *Entropy*, vol. 24, no. 1, p. 135, 2022.
- [11] S.-L. Huang, A. Makur, G. W. Wornell, and L. Zheng, "On universal features for high-dimensional learning and inference," *arXiv preprint arXiv:1911.09105*, 2019.
- [12] H. O. Hirschfeld, "A connection between correlation and contingency," in *Proceedings of the Cambridge Philosophical Society*, vol. 31, no. 4, 1935, pp. 520–524.
- [13] H. Gebelein, "Das statistische problem der korrelation als variations-und eigenwertproblem und sein zusammenhang mit der ausgleichsrechnung," *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik*, vol. 21, no. 6, pp. 364–379, 1941.
- [14] A. Rényi, "On measures of dependence," *Acta Mathematica Academiae Scientiarum Hungarica*, vol. 10, no. 3-4, pp. 441–451, 1959.
- [15] X. Xu and L. Zheng, "Multivariate feature extraction," in *2022 58th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 2022, pp. 1–8.
- [16] T. S. Jaakkola, D. Haussler *et al.*, "Exploiting generative models in discriminative classifiers," *Advances in neural information processing systems*, pp. 487–493, 1999.
- [17] T. Hastie, R. Tibshirani, J. H. Friedman, and J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*. Springer, 2009, vol. 2.
- [18] X. Xu and L. Zheng, "Kernel subspace and feature extraction," *arXiv preprint arXiv:2301.01410*, 2023.
- [19] X. Xu and S.-L. Huang, "Maximal correlation regression," *IEEE Access*, vol. 8, pp. 26 591–26 601, 2020.
- [20] J. Shawe-Taylor, N. Cristianini *et al.*, *Kernel methods for pattern analysis*. Cambridge university press, 2004.