# A New Algorithm for Tessellated Kernel Learning

**Brendon K. Colbert**
Arizona State University
Tempe, AZ 85287
brendon.colbert@asu.edu

**Matthew M. Peet**
Arizona State University
Tempe, AZ 85287
mpeet@asu.edu

## Abstract

The accuracy and complexity of machine learning algorithms based on kernel optimization are limited by the set of kernels over which they are able to optimize. An ideal set of kernels should: admit a linear parameterization (for tractability); be dense in the set of all kernels (for robustness); be universal (for accuracy). The recently proposed Tesselated Kernels (TKs) is currently the only known class which meets all three criteria. However, previous algorithms for optimizing TKs were limited to classification and relied on Semidefinite Programming (SDP) - limiting them to relatively small datasets. By contrast, the 2-step algorithm proposed here scales to 10,000 data points and extends to the regression problem. Furthermore, when applied to benchmark data, the algorithm demonstrates significant improvement in performance over Neural Nets and SimpleMKL with similar computation time.

## 1  Introduction

Kernel methods for classification and regression (and Support Vector Machines (SVMs) in particular) require selection of a kernel. Kernel Learning (KL) algorithms such as those found in [23, 21, 24] automate this task by finding the kernel, $k \in \mathcal{K}$ which optimizes an achievable metric such as the soft margin (for classification). The set of kernels, $k \in \mathcal{K}$, over which the algorithm can optimize, however, strongly influences the performance and robustness of the resulting classifier or predictor.

To understand how the choice of $\mathcal{K}$ influences performance and robustness, three properties were proposed in [4] to characterize the set $\mathcal{K}$ - tractability, density, and universality. Specifically, $\mathcal{K}$ is tractable if $\mathcal{K}$ is convex (or, preferably, a linear variety) - implying the KL problem is solvable using, e.g. [17, 10, 12, 16, 9]. The set $\mathcal{K}$ has the density property if, for any $\epsilon > 0$ and any positive kernel, $k^*$ there exists a $k \in \mathcal{K}$ where $\|k - k^*\| \leq \epsilon$. The density property implies the kernel will perform well on untrained data (robustness or generalizability). The set $\mathcal{K}$ has the universal property if any $k \in \mathcal{K}$ is universal - ensuring the classifier/predictor will perform arbitrarily well on large sets of training data.

In [4], the Tessellated Kernels (TKs) were shown to have all 3 properties, the first known such class of kernels. This work was based on a general framework for using positive matrices to parameterize positive kernels (as opposed to positive kernel matrices as in [12, 16, 15]). Unfortunately, however, the algorithms proposed in [4] were either based on SemiDefinite Programming (SDP) (thereby limiting the amount of training data) or used a randomized linear basis for the kernels (implying loss of density). Thus, while the algorithms in [4] outperformed all other methods (including deep learning) as measured by Test Set Accuracy (TSA), the computation times were not competitive. Furthermore, the results in [4] did not encompass the problem of regression.

In this paper, we extend the TK framework proposed in [4] to the problem of regression. The KL problem in regression has been studied using SDP in [16, 15] and Quadratic Programming (QP) in e.g. [17, 10]. However, neither of these previous works considered a set of kernels with both the tractability and the density property. By generalizing the Tessellated KL framework proposed in [4]

to the regression problem, we demonstrate significant increases in performance, as measured by Mean Square Error (MSE), and when compared to the results in [17, 10, 16].

In addition, we show that the SDP-based algorithm [4] for classification, and extended here to regression, can be decomposed into primal and dual sub-problems, $OPT\_A$ and $OPT\_P$ - similar to the approach taken in [17, 10]. Furthermore, we show that $OPT\_P$ (an SDP) admits an analytic solution using the Singular Value Decomposition (SVD) - an approach which allows us to consider higher dimensional feature spaces and more complex TKs. In addition, $OPT\_A$ is a convex QP and may be solved efficiently with achieved complexity which scales as $O(m^{2.16})$ where $m$ is the number of data points. We use a two-step algorithm on $OPT\_A$ and $OPT\_P$ and show that termination at $OPT\_A = OPT\_P$ is equivalent to global optimality. The resulting algorithm, then, does not require the use of SDP and, when applied to several standard test cases, is shown to retain the favorable TSA of [4] for classification, while offering improved MSE for regression, and competitive computation times as compared to other KL and deep learning algorithms.

## 2 An Ideal Set of Kernels for KL in Classification and Regression

Consider a generalized representation of the KL problem, which encompasses both classification and regression where (using the representor theorem [19]) the learned function is of the form $f_{\alpha,k}(z) = \sum_{i=1}^{m} \alpha_i k(x_i, z)$.

$$\min_{k \in \mathcal{K}} \min_{\alpha \in \mathbb{R}^m, b} \|f_{\alpha,k}\|^2 + C \sum_{i=1}^{m} l(f_{\alpha,k}, b)_{y_i, x_i} \tag{1}$$

Here $\|f_{\alpha,k}\|_X = \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j k(x_i, x_j)$ and $l(f_{\alpha,k}, b)_{y_i, x_i}$ is the loss function and is defined for SVM binary classification and SVM regression as $l^c(f_{\alpha,k}, b)_{y_i, x_i}$ and $l^r(f_{\alpha,k}, b)_{y_i, x_i}$, respectively, where

$$l^c(f_{\alpha,k}, b)_{y_i, x_i} = \max\{0, 1 - y_i(f_{\alpha,k}(x_i) - b)\} \text{ and } l(f_{\alpha,k}, b)^r_{y_i, x_i} = \max\{0, |y_i - (f_{\alpha,k}(x_i) - b)| - \epsilon\}.$$

The properties of the classifier/predictor, $f_{\alpha,k}$, resulting from Optimization Problem 1 will depend on the properties of the set $\mathcal{K}$, which is presumed to be a subset of the convex cone of all positive kernels. To understand how $\mathcal{K}$ influences the tractability of the optimization problem and the resulting fit, we consider three properties of the set, $\mathcal{K}$.

### 2.1 Tractability

We say a set of kernel functions, $\mathcal{K}$, is tractable if it can be represented using a countable basis.

**Definition 1.** *The set of kernels $\mathcal{K}$ is **tractable** if there exist a countable set $\{G_i(x, y)\}_i$ such that, for any $k \in \mathcal{K}$, there exists $N_G \in \mathbb{N}$ where $k(x, y) = \sum_{i=1}^{N_G} v_i G_i(x, y)$ for some $v \in \mathbb{R}^{N_G}$.*

Note the $G_i(x, y)$ need not be positive kernel functions. The tractable property is required for the KL problem to be tractable using algorithms for convex optimization.

### 2.2 Universality

Universal kernel functions always have positive definite (full rank) kernel matrices, implying that for arbitrary data $\{y_i, x_i\}_{i=1}^m$, there exists a function $f(z) = \sum_{i=1}^{m} \alpha_i k(x_i, z)$, such that $f(x_j) = y_j$ for all $j = 1, .., m$. Conversely, if a kernel is not universal, then exists a data set $\{x_i, y_i\}_{i=1}^m$ such that for any $\alpha \in \mathbb{R}^m$, there exists some $j \in \{1, \cdots, m\}$ such that $f(y_j) \neq \sum_{i=1}^{m} \alpha_i k(x_i, x_j)$. This ensures that SVMs using universal kernels can always benefit from additional training data, whereas non-universal kernels may saturate.

**Definition 2.** *A kernel $k : X \times X \to \mathbb{R}$ is said to be universal on the compact metric space $X$ if it is continuous and there exists an inner-product space $\mathcal{W}$ and feature map, $\Phi : X \to \mathcal{W}$ such that $k(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{W}}$ and where the unique Reproducing Kernel Hilbert Space (RKHS), $\mathcal{H} := \{f : f(x) = \langle v, \Phi(x) \rangle, v \in \mathcal{W}\}$ with associated norm $\|f\|_{\mathcal{H}} := \inf_v \{\|v\|_{\mathcal{W}} : f(x) = \langle v, \Phi(x) \rangle\}$ is dense in $\mathcal{C}(X) := \{f : X \to \mathbb{R} : f \text{ is continuous}\}$ where $\|f\|_{\mathcal{C}} := \sup_{x \in X} |f(x)|$.*

The following definition extends the universal property to a set of kernels.

**Definition 3.** *A set of kernel functions $\mathcal{K}$ has the universal property if every kernel function $k \in \mathcal{K}$ is universal.*

## 2.3 Density

The third property is density which distinguishes the TK class from other sets of kernel functions with the universal property. For instance consider a set containing a single Gaussian kernel function - which is clearly not ideal for kernel learning. The set containing a single Gaussian is tractable (it has only one element) and every member of the set is universal. However, it is not dense.

Considering SVM for classification, the KL problem determines the kernel $k \in \mathcal{K}$ for which we may obtain the maximum separation in the kernel-associated feature space. Increasing this separation distance makes the resulting classifier more robust (generalizable) [2]. The density property, then, ensures that the resulting KL algorithm will be maximally robust (generalizable) in the sense of separation distance.

Likewise, considering SVMs for regression, the KL problem finds the kernel $k \in \mathcal{K}$ which permits the "flattest" [20] function in feature space. In this case, the density property ensures that the resulting KL algorithm will be maximally robust (generalizable) in the sense of flatness.

These arguments motivate the following definition of the pointwise density property.

**Definition 4.** *The set of kernels $\mathcal{K}$ is said to be **pointwise dense** if for any positive kernel, $k^*$, any set of data $\{x_i\}_{i=1}^m$, and any $\epsilon > 0$, there exists $k \in \mathcal{K}$ such that $\|k(x_i, x_j) - k^*(x_i, x_j)\| \leq \epsilon$.*

## 3 A General Framework for Representation of Tractable Kernel Sets

Here we define a framework for constructing classes of tractable positive kernel functions and illustrate this approach on the class of General Polynomial Kernels.

**Proposition 5.** *Let $N$ be any bounded measurable function $N : X \times Y \to \mathbb{R}^q$ and $P \in \mathbb{R}^{q \times q}$ be a positive semidefinite matrix $P \geq 0$. Then*

$$k(x, y) = \int_X N(z, x)^T P N(z, y) dz \tag{2}$$

*is a positive kernel function.*

The proof for Proposition (5) may be found in [4].

**Lemma 6.** *Let $N$ be any bounded measurable function $N : X \times Y \to \mathbb{R}^q$ on compact $X$ and $Y$. Then the set of kernel functions*

$$\mathcal{K} := \left\{ k \mid k(x, y) = \int_X N(z, x)^T P N(z, y) dz, \ P \geq 0 \right\} \qquad \text{is tractable.} \tag{3}$$

For a given $N$, the map $P \mapsto k$ is linear. Specifically,

$$k(x, y) = \sum_{i=1}^q \sum_{j=1}^q P_{i,j} G_{i,j}(x, y) \ \text{where} \ G_{i,j}(x, y) = \int_X N_i(z, x) N_j(z, y) dz.$$

and thus by Definition 1 $\mathcal{K}$ is tractable.

In Subsection 3.1 we apply this framework to obtain Generalized Polynomial Kernels. In Subsection 4.1, we use the framework to obtain the TK class.

### 3.1 The Class of General Polynomial Kernels is Tractable

The class of General Polynomial Kernels (GPKs) is defined as the set of all polynomials, each of which is a positive kernel.

$$\mathcal{K}_P := \{k \in \mathbb{R}[x, y] \ : \ k \text{ is a positive kernel}\} \tag{4}$$

The GPK class is not universal, but is tractable, as per the following lemma.

**Lemma 7.** *$\mathcal{K}_P$ is tractable.*

*Proof.* Let $Z_d : \mathbb{R}^n \to \mathbb{R}^q$ be the vector of monomials of degree $d$ or less. From [4], we have that a polynomial $k$ of degree $2d$ is a positive polynomial kernel if and only if there exists some $P \geq 0$ such that $k(x, y) = Z_d(x)^T P Z_d(y)$. Now for any finite-dimensional subset of $\mathcal{K}_P$, let $d$ be the maximum degree over this subset and define $N(z, y) = Z_d(y)$. Then Lemma 6 implies that $\mathcal{K}_P$ is tractable. $\square$

This lemma implies that a representation of the form of Equation (2) is necessary and sufficient for a GPK to be positive. For convenience, we denote the set of GPK kernels of degree $d$ or less as follows [18].

$$\mathcal{K}_P^d := \{k \ : \ k(x, y) = Z_d(x)^T P Z_d(y) \ : \ P \geq 0\} \tag{5}$$

3

# 4 Tessellated Kernels: Tractable, Dense and Universal

In this section, we define the class of TK kernels and show it is tractable, dense, and universal.

## 4.1 Tessellated Kernels

Again, let $Z_d : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}^q$ be the vector of monomials of degree $d$. Define $\mathbf{I}$, the indicator function for the positive orthant, and the following choice of $N : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}^{2q}$ as

$$\mathbf{I}(z) = \begin{cases} 1 & z \geq 0 \\ 0 & \text{otherwise,} \end{cases} \quad \text{and} \quad N_T^d(z,x) = \begin{bmatrix} Z_d(z,x)\mathbf{I}(z-x) \\ Z_d(z,x)\mathbf{I}(x-z) \end{bmatrix} \tag{6}$$

where $z \geq 0$ means $z_i \geq 0$ for all $i$. We now define the set of TK kernels for $a < b \in \mathbb{R}^n$ as

$$\mathcal{K}_T^d := \left\{ k \ : \ k(x,y) = \int_a^b N_T^d(z,x)^T \ P \ N_T^d(z,y)dz, \ P \geq 0 \right\}, \ \mathcal{K}_T := \{k \ : \ k \in \mathcal{K}_T^d, \ d \in \mathbb{N}\}.$$

Kernels in the TK class are "Tessellated" in the sense that each datapoint defines a vertex which bisects each dimension of the domain of the resulting classifier/predictor - resulting in a tessellated partition of the feature space.

## 4.2 The Set of TK Kernels is Tractable

The class of TK kernels is prima facie in the form of Eqn. (3) in Lemma 6 and hence is tractable.

However, we will expand on this result by specifying the basis for the set of TK kernels, which will then be used in Section 5.

**Corollary 8.** *Suppose that for $a < b \in \mathbb{R}^n$, and $d \in \mathbb{N}$. We define the finite set $D_d := \{(\delta, \lambda) \in \mathbb{N}^{2n} : \|(\delta, \lambda)\|_1 \leq d\}$. Let $\{[\delta_i, \gamma_i]\}_{i=1}^q \subseteq D_d$ be some ordering of $D_d$ and define $Z_d(x,z)_j = x^{\delta_j} z^{\gamma_j}$ where $z^{\delta_j} x^{\gamma_j} := \prod_{i=1}^n z_i^{\delta_{j,i}} x_i^{\gamma_{j,i}}$. Now let $k$ be as defined in Eqn. (2) for some $P > 0$ and where $N$ is as defined in Eqn. (6). If we partition $P = \begin{bmatrix} Q & R \\ R^T & S \end{bmatrix}$ then we have,*

$$k(x,y) = \sum_{i,j=1}^q Q_{i,j} g_{i,j}(x,y) + R_{i,j} t_{i,j}(x,y) + R_{i,j}^T t_{i,j}(y,x) + S_{i,j} h_{i,j}(x,y)$$

*where $g_{i,j}, t_{i,j}, h_{i,j} : \mathbb{R}^{2n} \to \mathbb{R}$ are defined as*

$$g_{i,j}(x,y) := x^{\delta_i} y^{\delta_j} T(p^*(x,y), b, \gamma_{i,j} + \mathbf{1}), \ t_{i,j}(x,y) := x^{\delta_i} y^{\delta_j} T(x, b, \gamma_{i,j} + \mathbf{1}) - g_{i,j}(x,y), \text{ and}$$

$$h_{i,j}(x,y) := x^{\delta_i} y^{\delta_j} T(a, b, \gamma_i + \gamma_j + \mathbf{1}) - g_{i,j}(x,y) - t_{i,j}(x,y) - t_{i,j}(y,x),$$

*where $\mathbf{1} \in \mathbb{N}^n$ is the vector of ones, $p^* : \mathbb{R}^{2n} \to \mathbb{R}^n$ is defined elementwise as $p^*(x,y)_i = \max\{x_i, y_i\}$, and $T : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{N}^n \to \mathbb{R}$ is defined as*

$$T(x, y, \zeta) = \prod_{j=1}^n \left( \frac{y_j^{\zeta_j}}{\zeta_j} - \frac{x_j^{\zeta_j}}{\zeta_j} \right).$$

The proof of Corollary 8 can be found in [4].

## 4.3 The TK Class is Dense

The density property differentiates the set of TK kernels from other sets of kernel functions (e.g. a linear combination of Gaussian kernels of fixed bandwidths).

From [4] we have that the set of TK kernels satisfies the pointwise density property.

**Theorem 9.** *For any kernel matrix $K^*$ and any finite set $\{x_i\}_{i=1}^m$, there exists a $d \in \mathbb{N}$ and $k \in \mathcal{K}_T^d$ such that if $K_{i,j} = k(x_i, x_j)$, then $K = K^*$.*

In [4] an analytical solution, $K^*$, was found for the optimal trace-constrained kernel matrix that maximized the separation distance between two classes of points in the feature space. It was shown in this work that when $\{y_i\}$ has an equal number of positive and negative labels, $K^*$ contains an equal number of positive and negative elements - illustrating the importance of using kernels which are not pointwise positive (Gaussians are pointwise positive).

To illustrate the density property, then, we show how optimal GPK and TK kernels yield kernel matrices which approximate the analytic solution, $K^*$, of the optimal kernel matrix problem for a given set of data $\{x_i\}$ and labels $\{y_i\}$, while Gaussian kernels do not. Specifically, we consider the following optimization problem.

$$\min_{k \in \mathcal{K}} \|K - K^*\|_\infty \quad s.t. \ K_{i,j} = k(x_i, x_j) \qquad (7)$$

In these problems, the sets $\mathcal{K}$ will be: $\mathcal{K}_G^\gamma$ - the sum of $N$ Gaussians with bandwidths $\gamma_i$; $\mathcal{K}_P^d$ - the GPKs of degree $d$; and $\mathcal{K}_T^d$ - the TK kernels of degree $d$. More precisely, for bandwidths $\gamma \in \mathbb{R}^N$, we define $\mathcal{K}_G^\gamma := \left\{ k \ : \ k(x,y) = \sum_{i=1}^N \mu_i e^{\frac{||x-y||_2^2}{\gamma_i}} \ : \ \mu_i > 0 \right\}$.



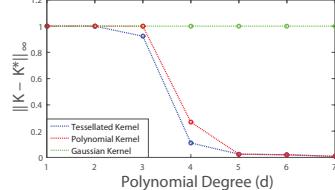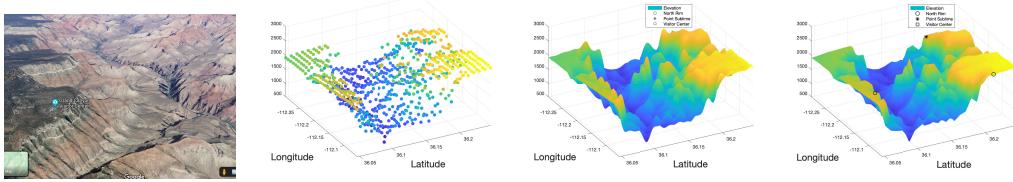Figure 1: The achieved objective ($\|K - K^*\|_\infty$) of Optimization Problem 7 for TK and GPK of degree $d$; and $m$ Gaussians with bandwidths in $[.01, 10]$. The number of bandwidths is selected so that the number of decision variables match in the Gaussian and TK cases.

Consider a spiral data set with 20 samples, using equal numbers of positive and negative labels. Fig. 1 shows the achieved objective value of Problem (7) for $\mathcal{K}_G^\gamma$, $\mathcal{K}_P^d$, and $\mathcal{K}_T^d$ as a function of the number of bandwidths (top $x$ axis - $N$ in $\mathcal{K}_G^\gamma$), polynomial degree (bottom $x$ axis - $d$ in $\mathcal{K}_P^d$, and $\mathcal{K}_T^d$). The $x$-axes of the plots are scaled to show equal numbers of decision variables. As expected, the case $\mathcal{K} = \mathcal{K}_\gamma^G$ saturates with an objective value significantly larger than the lower bound. The cases $\mathcal{K} = \mathcal{K}_P^d$ and $\mathcal{K} = \mathcal{K}_T^d$, meanwhile have almost no error at degree $d = 7$.

### 4.4  TK Kernels are Universal

Finally we discuss the universality property of the class of TK kernels which ensures that every TK function can fit the training data well.

The following theorem from [4] shows that any TK kernel with $P > 0$ is necessarily universal.

**Theorem 10.** *Suppose $k$ is as defined in Eqn. (2) for some $P > 0$, $d \in \mathbb{N}$ and $N$ as defined in Eqn. (6). Then $k$ is universal for $a < b \in \mathbb{R}^n$.*

This theorem implies that even if we use the subset of TK kernels defined by $d = 0$, this subset is still universal.

## 5  A New Algorithm for KL in Classification and Regression using TKs

In this section, we express the KL optimization problem for both classification and regression and break this optimization problem into two sub-problems which allow us to express the problem in primal and dual form. For convenience, we define the feasible sets for the sub-problems as

$$\mathcal{X} := \{P \in \mathbb{R}^{q \times q} \ : \ \text{trace}(P) = q, \ P > 0\}$$

$$\mathcal{Y}_c := \{\alpha \in \mathbb{R}^m \ : \ \sum_{i=1}^m \alpha_i y_i = 0, \ 0 \le \alpha_i \le C\}, \quad \mathcal{Y}_r := \{\alpha \in \mathbb{R}^m \ : \ \sum_{i=1}^m \alpha_i = 0, \ \alpha_i \in [-C, C]\}.$$

The common part of the objective is

$$O(\alpha, P) := -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \int_a^b N_T^d(z, x_i)^T P N_T^d(z, y_j) dz, \qquad (8)$$

while the unique parts of the objective are

$$\kappa_c(\alpha) := \sum_{i=1}^m \alpha_i \quad \text{and} \quad \kappa_r(\alpha) := -\epsilon \sum_{i=1}^m |\alpha_i| + \sum_{i=1}^m y_i \alpha_i.$$

Then the KL optimization problem ($OPT$) for TK kernels ($\odot$ being elementwise multiplication) is as follows for classification and regression, respectively.

$$OPT := \min_{P \in \mathcal{X}} \max_{\alpha \in \mathcal{Y}_c} \ O(\alpha \odot y, P) + \kappa_c(\alpha), \quad \text{and} \quad OPT := \min_{P \in \mathcal{X}} \max_{\alpha \in \mathcal{Y}_r} \ O(\alpha, P) + \kappa_r(\alpha).$$

**Primal Formulation:** We can formulate the primal problem ($OPT_P$) as

$$OPT_P = \min_{P \in \mathcal{X}} \max_{\alpha \in \mathcal{Y}} O(\alpha, P) + \kappa_r(\alpha) \text{ or } O(\alpha \odot y, P) + \kappa_c(\alpha) = \min_{P \in \mathcal{X}} \ OPT\_A(P) \qquad (9)$$

where for classification and regression, respectively,

$$OPT\_A(P) := \max_{\alpha \in \mathcal{Y}_c} \ O(\alpha \odot y, P) + \kappa_c(\alpha), \quad \text{and} \quad OPT\_A(P) := \max_{\alpha \in \mathcal{Y}_r} \ O(\alpha, P) + \kappa_r(\alpha).$$

5

(a) An image from Google Maps of a section of the Grand Canyon corresponding to (36.04, -112.05) latitude and (36.25, -112.3) longitude.

(b) Elevation data ($m = 750$) from [1] for a section of the Grand Canyon between (36.04, -112.05) latitude and (36.25, -112.3) longitude.

(c) Predictor using a hand-tuned Gaussian kernel trained on the elevation data in (b). The Gaussian predictor poorly represents the sharp edge at the north and south rim.

(d) Predictor from Algorithm 1 trained on the elevation data in (b). The TK predictor accurately represent the north and south rims of the canyon.

Figure 2: Subfigure (a) shows an 3D representation of the section of the Grand Canyon to be fitted. In (b) we plot elevation data of this section of the Grand Canyon. In (c) we plot the predictor for a hand-tuned Gaussian kernel. In (d) we plot the predictor from Algorithm 1 where $d = 2$.

**Dual Formulation:** Alternatively, we have the dual formulation ($OPT_D$).

$$OPT_D = \max_{\alpha \in \mathcal{Y}} OPT\_P(\alpha) \tag{10}$$

where $\mathcal{Y} = \mathcal{Y}_c$ for classification and $\mathcal{Y} = \mathcal{Y}_r$ regression. Likewise, for classification and regression, respectively,

$$OPT\_P(\alpha) := \min_{P \in \mathcal{X}} O(\alpha \odot y, P) + \kappa_c(\alpha) \quad \text{and} \quad OPT\_P(\alpha) := \min_{P \in \mathcal{X}} O(\alpha, P) + \kappa_r(\alpha).$$

**Lemma 11.** *For $\alpha \in \mathcal{Y}$, $P \in \mathcal{X}$, $OPT\_A(P) = OPT\_P(\alpha)$ if and only if: $\{\alpha, P\}$ solve $OPT$; $P$ solves $OPT_P$; and $\alpha$ solves $OPT_D$.*

*Proof.* For any minmax optimization problem with objective function $\phi$, we have

$$d^* = \max_{\alpha \in \mathcal{Y}} \min_{P \in \mathcal{X}} \phi(P, \alpha) \leq \min_{P \in \mathcal{X}} \max_{\alpha \in \mathcal{Y}} \phi(P, \alpha) = p^*,$$

and strong duality holds ($p^* - d^* = 0$) if $\mathcal{X}$ and $\mathcal{Y}$ are both convex and one is compact, $\phi(\cdot, \alpha)$ is convex for every $\alpha \in \mathcal{Y}$ and $\phi(P, \cdot)$ is concave for every $P \in \mathcal{X}$, and the function $\phi$ is continuous [7]. In our case, these conditions hold for both classification and regression where $\phi(P, \alpha) = O(\alpha, P) + \kappa_r(\alpha)$ or $O(\alpha \odot y, P) + \kappa_c(\alpha)$. Hence if $\alpha^*$ solves $OPT\_P$ and $P^*$ solves $OPT\_A$, then $\{\alpha^*, P^*\}$ solves $OPT$ and

$$OPT\_P(\alpha^*) = \max_{\alpha \in \mathcal{Y}} OPT\_P(\alpha) = \min_{P \in \mathcal{X}} OPT\_A(P) = OPT\_A(P^*).$$

Conversely, suppose $\alpha \in \mathcal{Y}$, $P \in \mathcal{X}$, then

$$OPT\_P(\alpha) \leq \max_{\alpha \in \mathcal{Y}} OPT\_P(\alpha) = OPT\_P(\alpha^*)$$
$$= OPT\_A(P^*) = \min_{P \in \mathcal{X}} OPT\_A(P) \leq OPT\_A(P).$$

Hence if $OPT\_A(P) = OPT\_P(\alpha)$, then $OPT\_A(P) = OPT\_A(P^*) = OPT\_P(\alpha^*) = OPT\_P(\alpha)$ and hence $P$ and $\alpha$ solve $OPT\_A$ and $OPT\_P$, respectively. $\square$

We propose Algorithm 1 as a two-step iterative algorithm for solving Optimization Problem (10).

### 5.1 Solving $OPT\_A(P)$

For a given $P > 0$, $OPT\_A(P)$ is a Quadratic Program (QP). General purpose QP solvers as applied to this problem have a worst-case complexity which scales as $O(m^3)$ [25] where $m$ is the number of data points. This computational complexity may be improved, however, by noting that the problem formulation is compatible with the representation defined in [3] for QPs

---

**Algorithm 1** Two Step TKL

Initialize $P = I$;
**while** $OPT\_P(\alpha_k) - OPT\_A(P_k) > \epsilon$ **do**
    $\alpha_{k+1} = OPT\_A(P_k)$
    $P_{k+1} = \frac{P_K + tOPT\_P(\alpha_{k+1})}{1+t}$ (select $t$ using line search)
    $k = k + 1$
**end while**

---

derived from SVM. In this case, the algorithm in LibSVM [3] can reduce the computational burden somewhat. This improved performance is illustrated in Figure 3 where we observe the achieved complexity scales as $O(m^{2.1})$. Note that for the 2-step algorithm proposed in this manuscript, solving the QP in $OPT\_A(P)$ is significantly slower that solving the Singular Value Decomposition (SVD) required for $OPT\_P(\alpha)$, which is defined in the following subsection. However, the achieved complexity of $O(m^{2.1})$ is also significantly faster than solving the large SDP, as described in [12], [16], and [4]. This complexity comparison will be further discussed in Section 6.
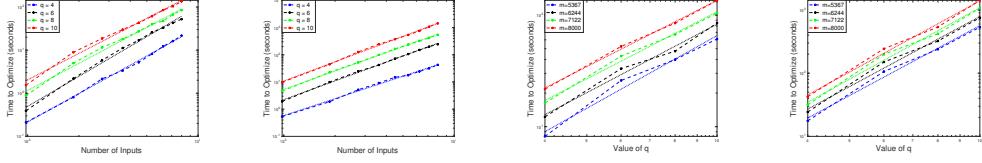
6

(a) Numerical complexity analysis of TKL for classification versus $m$.

(b) Numerical complexity analysis of TKL for regression versus $m$.

(c) Numerical complexity analysis of TKL for classification versus $q$.

(d) Numerical complexity analysis of TKL for regression versus $q$.

Figure 3: In (a) and (b) we plot log scale plots of the time taken to optimize TKL as the number of inputs change for $P \in \mathbb{R}^{q \times q}$. The line of best linear fit is plotted as a dotted line. In (c) and (d) we plot log scale plots of the time taken to optimize TKL as the value of $q$ changes for four different values of $m$.

## 5.2 Solving $OPT\_P(\alpha)$

For a given $\alpha$, $OPT\_P(\alpha)$ is an SDP. Fortunately, however, this SDP is structured so as to admit an analytic solution using the SVD. To solve $OPT\_P(\alpha)$ we minimize $O(\alpha, P)$ from Eq. (8) which, as per Corollary 8, is linear in $P$ and can be formulated as

$$OPT\_P(\alpha) := \min_{\substack{P \in \mathbb{R}^{q \times q} \\ \text{trace}(P)=q \\ P>0}} O(\alpha, P) := \min_{\substack{P \in \mathbb{R}^{q \times q} \\ \text{trace}(P)=q \\ P>0}} \text{trace}(C(\alpha)^T P) \tag{11}$$

where

$$C_{i,j}(\alpha) = \sum_{k,l=1}^m (\alpha_k y_k) G_{i,j}(x_k, x_l)(\alpha_l y_l), \qquad G_{i,j}(x,y) := \begin{cases} g_{i,j}(x,y) & \text{if } i \le \frac{q}{2}, j \le \frac{q}{2} \\ t_{i,j}(x,y) & \text{if } i \le \frac{q}{2}, j > \frac{q}{2} \\ t_{i,j}(y,x) & \text{if } i > \frac{q}{2}, j \le \frac{q}{2} \\ h_{i,j}(x,y) & \text{if } i > \frac{q}{2}, j > \frac{q}{2} \end{cases}$$

and $g$, $t$ and $h$ can be found in Corollary 8.

The following theorem gives an analytic solution for $OPT\_P$ using the SVD.

**Theorem 12.** *Let $C = V\Sigma V^T$ be the SVD of symmetric $C \in \mathbb{R}^{q \times q}$ and $v$ be the right singular vector corresponding to the minimum singular value of $C$. Then $P^* = qvv^T$ solves $OPT\_P$.*

*Proof.* Recall $OPT\_P$ has the form $\min_{P \in \mathbb{R}^{q \times q}} \text{trace}(C^T P)$ s.t. $P \ge 0$, $\text{trace}(P) = q$.

Denote the minimum singular value of $C$ as $\sigma_{\min}(C)$. Then for any feasible $P \in \mathcal{X}$, by [8] we have

$$\text{trace}(C^T P) \ge \sigma_{\min}(C) trace(P) = \sigma_{\min}(C)q.$$

Now consider $P = qvv^T \in \mathbb{R}^{q \times q}$. $P$ is feasible since $P \ge 0$, and $\text{trace}(P) = q$. Furthermore,

$$\text{trace}(CP) = q\,\text{trace}(V\Sigma V^T vv^T) = q\,\text{trace}(v^T V\Sigma V^T v) = q\,\sigma_{\min}(C)$$

as desired. $\square$

Note that the size of the SVD problem in $OPT\_P(\alpha)$ is $q^2$, which increases with the number of features, which is typically relatively small. As a result, we observe that the $OPT\_P$ step of Algorithm 1 is typically less computationally intense than the $OPT\_A$ step.

## 6 Complexity and Scalability of the New TK Kernel Learning Algorithm

We consider the computational complexity of Algorithm 1. If we define the number of data points used to learn the TK kernel function as $m$ and the size of $P$ as $q \times q$, then we find experimentally that the complexity of Algorithm 1 scales as approximately $O(m^{2.16}q^{2.23})$ for classification and $O(m^{2.24}q^{3.59})$ for regression as can be seen in Fig. 3. These results are lower with respect to $m$ than the value of $O(m^{2.6}q^{1.9})$ reported in [4] for binary classification. The values for classification and regression are both estimated using the data set: Combined Cycle Power Plant (CCPP) in [22, 11], containing 4 features and $m = 9568$ samples. In the case of classification, labels with value greater than or equal to the median of the were relabeled as 1, and those less than the median were relabeled as $-1$. Note that to study scalability in $q$, we varied the number of features in the dataset - thereby incrementing the size of the matrix $P \in \mathbb{R}^{q \times q}$.

Aside from improved scalability, the overall time required for Algorithm 1 is significantly reduced when compared with the algorithm in [4], improving by two orders of magnitude in some cases. This is illustrated for classification using four data sets in Table 1. This improved complexity is likely due to the lower overhead associated with QP and the SVD.

Table 1: We report the mean computation time (in seconds), along with standard deviation, for 30 trials comparing the SDP algorithm in [4] and the new TKL algorithm on several data sets. All tests are run on a computer with an Intel i7-5960X CPU at 3.00 GHz with 128 Gb of RAM.

| Method | Liver [6] | Cancer [13] | Heart [6] | Pima [6] |
|--------|-----------|-------------|-----------|----------|
| SDP | $95.75 \pm 2.68$ | $636.17 \pm 25.43$ | $221.67 \pm 29.63$ | $1211.66 \pm 27.01$ |
| TKL | $1.10 \pm 0.24$ | $8.20 \pm 0.36$ | $3.35 \pm 0.26$ | $12.66 \pm 0.44$ |

Table 2: Mean Square Error comparison for algorithms [TKL], [SimpleMKL] and [Neural Net]. In the data set column $m$ is the number of points in the training data set and $n$ is the number of features. All tests are run on a computer with an Intel i7-5960X CPU at 3.00 GHz with 128 Gb of RAM.

| Data Set | Method | Error | Time | Data Set | Method | Error | Time |
|----------|--------|-------|------|----------|--------|-------|------|
| CCPP [22, 11] | TKL | 9.70 | 1463.8 | Abalone [6] | TKL | 3.43 | 522.5 |
| $n = 4$, $m = 8000$ | SimpleMKL | 13.77 | 26097.1 | $n = 8$, $m = 4000$ | SimpleMKL | 4.28 | 1185.3 |
| $m_t = 1568$ | Neural Net | 15.00 | 850.4 | $m_t = 177$ | Neural Net | 8.72 | 483.4 |
| Airfoil [6] | TKL | 1.46 | 92.1 | Forest [5] | TKL | 2.05 | 7.6 |
| $n = 5$, $m = 1300$ | SimpleMKL | 3.63 | 1025.0 | $n = 10$, $m = 457$ | SimpleMKL | 2.07 | 0.8 |
| $m_t = 203$ | Neural Net | 4.28 | 61.3 | $m_t = 50$ | Neural Net | 6.40 | 117.7 |

## 7 Accuracy of the New TK Kernel Learning Algorithm for Regression

As expected, for classification, the accuracy of the new TK kernel learning algorithm (TKL) is identical to the analysis in [4].

For regression, we evaluate the accuracy of TKL when compared to other state of the art machine learning algorithms. Because the set of TK kernels is dense, for classification (as shown in [4]), TKL outperforms all existing algorithm with respect to TSA. For regression, the appropriate metric is Mean Square Error (MSE). The algorithms used in our comparison are as follows.

**[TKL]** Algorithm 1 with $d = 1$, $\epsilon = .1$ and we scale the data so that $x_i \in [0, 1]^n$, and then select $[a, b] = [0 - \delta, 1 + \delta]^n$, where $\delta > 0$ and $C$ are chosen by 5-fold cross-validation;

**[SimpleMKL]** We use SimpleMKL [17] with a standard selection of Gaussian and polynomial kernels with bandwidths arbitrarily chosen between .5 and 10 and polynomial degrees one through three - yielding approximately $13(n + 1)$ kernels. We set $\epsilon = .1$ as in TKL and $C$ is chosen by 5-fold cross-validation;

**[Neural Net]** We use a 3 layer neural network with 50 hidden layers using MATLABs (`feedforwardnet`) implementation and stopped learning after the error in a validation set decreased sequentially 50 times.

In Table 2, we see the average MSE on the test set for these three approaches as applied to randomly selected regression benchmark data sets where $n$ is the dimension of the data, $m$ is the number of training data and $m_t$ is the number of testing data points. In all cases except Forest, [TKL] had both a lower (or comparable) computation time and MSE than both SimpleMKL and Neural Net. In all cases, the MSE for TKL was significantly lower - illustrating the importance of the density property.

To further illustrate the importance of density property and the TKL framework for practical regression problems, we used elevation data from [1] to learn a TK kernel and associated SVM predictor representing the surface of the Grand Canyon in Arizona. This data set is particularly challenging due to the variety of geographical features. The result of the TKL algorithm can be seen in Figure 2(d).

## 8 Conclusion

We have extended the TK kernel learning framework to regression problems and proposed a faster algorithm for TK kernel learning which can be used for both classification and regression. The set of TK kernels is tractable, dense, and universal - implying that KL algorithms based on TK kernels are more robust - resulting in higher TSA for classification and lower MSE for regression. These three properties, combined with the improved computational complexity of the new algorithm, has resulted in a kernel learning framework which achieves both lower MSE and computation time when compared to both SimpleMKL and neural networks.

**Broader Impacts**

While machine learning algorithm have become very accurate in recent years, they perform poorly when faced with changes in the underlying process. As evidenced by Covid19, predictive models based on ML algorithms can be brittle [14]. The density property of the TK class ensures that the models generated using the algorithms described in this manuscript will be more robust to such changes in environment. Naturally, however, over-reliance on predictive models, without understanding of the process, can lead to negative outcomes, even if the models are robust.

# References

[1] J.J. Becker, D.T. Sandwell, W.H.F. Smith, J. Braud, B. Binder, J.L. Depner, D. Fabre, J. Factor, S. Ingalls, S.H. Kim, et al. Global bathymetry and elevation data at 30 arc seconds resolution: Srtm30_plus. *Marine Geodesy*, 32(4):355–371, 2009.

[2] B. Boehmke and B.M. Greenwell. *Hands-On Machine Learning with R*. CRC Press, 2019.

[3] C-C. Chang and C-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at `http://www.csie.ntu.edu.tw/~cjlin/libsvm`.

[4] B.K. Colbert and M.M. Peet. A convex parametrization of a new class of universal kernel functions. *Journal of Machine Learning Research*, 21(45):1–29, 2020.

[5] P. Cortez and A. Morais. A data mining approach to predict forest fires using meteorological data. 2007.

[6] D. Dua and C. Graff. UCI machine learning repository, 2017.

[7] K. Fan. Minimax theorems. *Proceedings of the National Academy of Sciences of the United States of America*, 39(1):42, 1953.

[8] Y. Fang, K.A. Loparo, and X. Feng. Inequalities for the trace of matrix product. *IEEE Transactions on Automatic Control*, 39(12):2489–2490, 1994.

[9] M. Gönen and E. Alpaydın. Multiple kernel learning algorithms. *Journal of Machine Learning Research*, 2011.

[10] A. Jain, S. Vishwanathan, and M. Varma. SPF-GMKL: generalized multiple kernel learning with a million kernels. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining*, 2012.

[11] H. Kaya, P. Tüfekci, and F.S. Gürgen. Local and global learning methods for predicting power of a combined gas & steam turbine. In *Proceedings of the international conference on emerging trends in computer and electronics engineering icetcee*, pages 13–18, 2012.

[12] G. Lanckriet, N. Cristianini, P. Bartlett, L. El Ghaoui, and M. Jordan. Learning the kernel matrix with semidefinite programming. *Journal of Machine Learning Research*, 2004.

[13] O.L. Mangasarian, R. Setiono, and W.H. Wolberg. Pattern recognition via linear programming: Theory and application to medical diagnosis. 1990.

[14] C. Mims. AI isn't magical and won't help you reopen your business. *Wall Street Journal*, May 2020.

[15] K. Ni, S. Kumar, and T. Nguyen. Learning the kernel matrix for superresolution. In *2006 IEEE Workshop on Multimedia Signal Processing*, pages 441–446, 2006.

[16] S. Qiu and T. Lane. Multiple kernel learning for support vector regression. *Computer Science Department, The University of New Mexico, Albuquerque, NM, USA, Tech. Rep*, page 1, 2005.

[17] A. Rakotomamonjy, F. R. Bach, S. Canu, and Y. Grandvalet. SimpleMKL. *Journal of Machine Learning Research*, 2008.

[18] B. Recht. *Convex Modeling with Priors*. PhD thesis, Massachusetts Institute of Technology, 2006.

[19] B. Schölkopf, R. Herbrich, and A.J. Smola. A generalized representer theorem. In *International conference on computational learning theory*, pages 416–426, 2001.

[20] A.J. Smola and B. Schölkopf. A tutorial on support vector regression. *Statistics and computing*, 14(3):199–222, 2004.

[21] S. Sonnenburg, G. Rätsch, S. Henschel, C. Widmer, J. Behr, A. Zien, F. De Bona, A. Binder, C. Gehl, and V. Franc. The shogun machine learning toolbox. *Journal of Machine Learning Research*, 11(60):1799–1802, 2010.

[22] P. Tüfekci. Prediction of full load electrical power output of a base load operated combined cycle power plant using machine learning methods. *International Journal of Electrical Power & Energy Systems*, 60:126–140, 2014.

[23] Z. Xu, R. Jin, H. Yang, I. King, and M.R. Lyu. Simple and efficient multiple kernel learning by group lasso. In *Proceedings of the 27th international conference on machine learning*, pages 1175–1182, 2010.

[24] H. Yang, Z. Xu, J. Ye, I. King, and M.R. Lyu. Efficient sparse generalized multiple kernel learning. *IEEE Transactions on neural networks*, 22(3):433–446, 2011.

[25] Y. Ye and E. Tse. An extension of karmarkar's projective algorithm for convex quadratic programming. *Mathematical programming*, 44(1-3):157–179, 1989.