

DEPARTMENT OF STATISTICS  
University of Wisconsin  
1300 University Ave.  
Madison, WI 53706

TECHNICAL REPORT NO. 1107

May 6, 2005

# A Framework for Kernel Regularization with Application to Protein Clustering

Fan Lu<sup>1</sup>

Department of Statistics, University of Wisconsin, Madison, WI

Sündüz Keleş<sup>2</sup>

Department of Statistics and Department of Biostatistics and Medical Informatics,  
University of Wisconsin, Madison, WI

Stephen J. Wright<sup>3</sup>

Department of Computer Sciences, University of Wisconsin, Madison, WI

Grace Wahba<sup>1</sup>

Department of Statistics, Department of Computer Sciences and Department of  
Biostatistics and Medical Informatics, University of Wisconsin, Madison, WI

**Key Words and Phrases:** Regularized Kernel Estimation, positive definite matrices, noisy dissimilarity data, convex cone programming, protein clustering, globin family, support vector machines, classification.

---

<sup>1</sup>Research supported in part by NSF grant DMS0072292 and NIH grant EY09946.

<sup>2</sup>Research supported in part by Wisconsin Alumni Research Foundation.

<sup>3</sup>Research supported in part by NSF Grants ATM-0296033, CNS-0127857, CCF-0113051, ANI-0330538, DMS-0427689, CCF-0430504, and DOE grant DE-FG02-04ER25627.

# A Framework for Kernel Regularization with Application to Protein Clustering

Fan Lu<sup>1</sup>, Sündüz Keleş<sup>1,2</sup>, Stephen J. Wright<sup>3</sup> and Grace Wahba<sup>1,2,3</sup>

1. Department of Statistics

2. Department of Biostatistics and Medical Informatics

3. Department of Computer Sciences

University of Wisconsin, Madison, 53706, USA

May 6, 2005

## Abstract

We develop and apply a novel framework which is designed to extract information in the form of a positive definite kernel matrix from possibly crude, noisy, incomplete, inconsistent dissimilarity information between pairs of objects, obtainable in a variety of contexts. Any positive definite kernel defines a consistent set of distances, and the fitted kernel provides a set of coordinates in Euclidean space which attempt to respect the information available, while controlling for complexity of the kernel. The resulting set of coordinates are highly appropriate for visualization and as input to classification and clustering algorithms. The framework is formulated in terms of a class of optimization problems which can be solved efficiently using modern convex cone programming software. The power of the method is illustrated in the context of protein clustering based on primary sequence data. An application to the globin family of proteins resulted in a readily visualizable 3D sequence space of globins, where several sub-families and sub-groupings consistent with the literature were easily identifiable.

## 1 Introduction

It has long been recognized that symmetric positive definite kernels (hereinafter “kernels”) play a key role in function estimation [1][2], clustering and classification, dimension reduction, and other applications. Such kernels can be defined on essentially any conceivable domain of interest [3], originally function spaces and more recently, finite (but possibly large) collections of trees, graphs, images, DNA and protein sequences, microarray gene expression chips, and other objects. A kernel defines a distance metric between pairs of objects in the domain that admits an inner product. Thus they play a key role in the implementation of classification algorithms (via support vector machines (SVMs)) and clustering (via  $k$ -means algorithms, for example), along with their more classical role in function approximation and estimation, and the solution of ill-posed inverse problems [4]. Since the mid 90s, when the key role of these kernels became evident in SVMs [5][6][7], a massive literature has grown related to the use and choice of kernels in many domains of application, including, notably, computational biology [8]. A google search as of the date of this writing gave over *three million* hits on “Kernel Methods”, along with an ad from Google soliciting job applications from computer scientists!

Mathematically defined kernels, for example, spline kernels, radial basis functions and related positive definite functions defined on Euclidean space, have long been the workhorses in the field, generally with one or a few free parameters estimated from the data; see, for example [9]. A recent work [10] proposes estimating a kernel by optimizing a linear combination of prespecified kernels via a semidefinite programming approach. The reader may connect with the recent literature on kernel construction and use in a variety of contexts by going to the NIPS2004 website (<http://books.nips.cc/nips17.html>) or the book [11].

It is frequently possible to use expert knowledge or other information to obtain dissimilarity scores for pairs of objects, which serve as pseudo-distances between the objects. There are two problem types of interest. The first is to estimate full relative position information for a (training) set of objects in a space of preferably low dimension in order to visualize the data or to conduct further processing - typically, classification or clustering. One traditional approach for this purpose is multidimensional scaling (MDS) [12], which continues to be an active research area. The second problem is to place new objects in the space, given some dissimilarity information between them and some members of the training set, in the coordinate space of the training set.

This paper proposes regularized kernel estimation (RKE), a unified framework for solving both problems by fitting a positive definite kernel from possibly crude, noisy, incomplete, inconsistent, weighted, repetitious dissimilarity information, in a fully nonparametric approach, by solving a convex optimization problem with modern convex cone programming tools. The basic idea is to solve an optimization problem which trades off goodness of fit to the data and a complexity (shrinkage) penalty on the kernel which is used to fit the data - analogous to the well known bias-variance trade-off in the spline and ill-posed inverse literature, but not exactly the same. Within this framework, we provide an algorithm for placing new objects in the coordinate space of the training set.

The method can be used instead of MDS to provide a coherent set of coordinates for the given objects in few or many dimensions, without problems with local minima or (some) missing data. It can also be used to solve problems discussed in [10], but in a fully nonparametric way.

The feasibility of the RKE approach is demonstrated in the context of protein sequence clustering, by applying the method to global pairwise alignment scores of the heme-binding protein family of globins. In this example, we are already able to visualize the known globin subfamilies from a 3D plot of the training sequence coordinates that are obtained by the regularized kernel estimate. Furthermore, apparent sub-clusterings and outliers of the known globin subfamilies from the 3D plot reveal interesting observations consistent with the literature. Clustering of protein sequences from a family to identify subfamilies or clustering and classification of protein domains to determine protein function present one major application area for the novel framework presented here. However, we envision many more applications involving clustering and classification tasks in biological and non-biological data analysis, some of these are discussed in Section 5.

In Section 2, we present the general formulation of the problem and define the family of Regularized Kernel Estimates. Section 3 describes the formulation of RKE problems and the problem of placing test data in the coordinate space of training data as general convex cone problems. Also included is a brief discussion on tuning the parameters of the estimation procedure. Section 4 presents an application to the globin protein family to identify subfamilies and discusses the biological implication of the results. Examples of placing test data points in the coordinate system of training protein sequences are illustrated here. We conclude with a summary and discussion of future work in Section 5.

## 2 Dissimilarity Information and RKE

Given a set of  $N$  objects, suppose we have obtained a measure of dissimilarity,  $d_{ij}$ , for certain object pairs  $(i, j)$ . We introduce the class of Regularized Kernel Estimates (RKEs), which we define as solutions to optimization problems of the following form:

$$\min_{K \in S_N} \sum_{(i,j) \in \Omega} L(w_{ij}, d_{ij}, \hat{d}_{ij}(K)) + \lambda J(K), \quad (1)$$

where  $S_N$  is the convex cone of all real nonnegative definite matrices of dimension  $N$ ,  $\Omega$  is the set of pairs for which we utilize dissimilarity information, and  $L$  is some reasonable loss function, convex in  $\hat{d}_{ij}$ , where  $\hat{d}_{ij}$  is the dissimilarity induced by  $K$ .  $J$  is a convex kernel penalty (regularizing) functional, and  $\lambda$  is a tuning parameter balancing fit to the data and the penalty on  $K$ . The  $w_{ij}$  are weights that may, if desired, be associated with particular  $(i, j)$  pairs. The natural induced dissimilarity, which is a real squared distance admitting of an inner product, is  $\hat{d}_{ij} = K(i, i) + K(j, j) - 2K(i, j) = B_{ij} \cdot K$ , where  $K(i, j)$  is the  $(i, j)$  entry of  $K$  and  $B_{ij}$  is a symmetric matrix of dimension  $N$  with all elements 0 except  $B_{ij}(i, i) = B_{ij}(j, j) = 1$ ,  $B_{ij}(i, j) = B_{ij}(j, i) = -1$ . The inner (dot) product of two matrices of the same dimensions is defined as:  $A \cdot B = \sum_{i,j} A(i, j) \cdot B(i, j) \equiv \text{trace}(A^T B)$ . There are essentially no restrictions on the set of pairs other than requiring that the graph of the objects with pairs connected by edges be connected. A pair may have repeated observations, which just yield an additional term in (1) for each separate observation. If the pair set induces a connected graph, then the minimizer of (1) will have no local minima.

Although it is usually natural to require the observed dissimilarity information  $\{d_{ij}\}$  to satisfy  $d_{ij} \geq 0$  and  $d_{ij} = d_{ji}$ , the general formulation above does not require these properties to hold. The observed dissimilarity information may be incomplete (with the restriction noted), it may not satisfy the triangle inequality, or it may be noisy. It also may be crude, as for example when it encodes a small number of coded levels such as “very close”, “close”, “distant”, and “very distant”.

In this work we consider two special cases of the formulation (1), the first for its use in the application to be discussed.

## 3 Numerical Methods for RKE

In this section, we describe a specific formulation of the approach in Section 2, based on a linearly weighted  $l_1$  loss, and use the trace function in the regularization term to promote dimension reduction. The resulting problem is as follows:

$$\min_{K \succeq 0} \sum_{(i,j) \in \Omega} w_{ij} |d_{ij} - B_{ij} \cdot K| + \lambda \text{trace}(K). \quad (2)$$

We show how this formulation can be posed as a general conic optimization problem and also describe a “newbie” formulation in which the known solution to (2) for a set of  $N$  objects is augmented by the addition of one more object together with its dissimilarity data. A variant of (2), in which a quadratic loss function is used in place of the  $l_1$  loss function, is described in the appendix.

### 3.1 General Convex Cone Problem

We specify here the general convex cone programming problem. This problem, which is central to modern optimization research, involves some unknowns that are vectors in Euclidean space and others that are symmetric matrices. These unknowns are required to satisfy certain equality constraints and are also required to belong to cones of a certain type. The cones have the common feature that they all admit a self-concordant barrier function, which allows them to be solved by interior-point methods that are efficient in both theory and practice.

To describe the cone programming problem, we define some notation. Let  $\mathcal{R}^p$  be Euclidean  $p$ -space, and let  $P_p$  be the nonnegative orthant in  $\mathcal{R}^p$ , that is, the set of vectors in  $\mathcal{R}^p$  whose components are all nonnegative. We let  $Q_q$  be the second-order cone of dimension  $q$ , which is the set of vectors  $x = (x(1), \dots, x(q)) \in \mathcal{R}^q$  that satisfy the condition  $x(1) \geq [\sum_{i=2}^q x(i)^2]^{1/2}$ . We define  $S_s$  to be the cone of symmetric positive definite  $s \times s$  matrices of real numbers. Inner products between two vectors are defined in the usual way and we use the dot notation for consistency with the matrix inner product notation.

The general convex cone problem is then:

$$\begin{aligned} \min_{X_j, x_i, z} \quad & \sum_{j=1}^{n_s} C_j \cdot X_j + \sum_{i=1}^{n_q} c_i \cdot x_i + g \cdot z \\ \text{s.t.} \quad & \sum_{j=1}^{n_s} A_{rj} \cdot X_j + \sum_{i=1}^{n_q} a_{ri} \cdot x_i + g_r \cdot z = b_r, \quad \forall_r \\ & X_j \in S_{s_j} \quad \forall_j; \quad x_i \in Q_{q_i} \quad \forall_i; \quad z \in P_p. \end{aligned} \tag{3}$$

Here,  $C_j, A_{rj}$  are real symmetric matrices (not necessarily positive semidefinite) of dimension  $s_j$ ,  $c_i, a_{ri} \in \mathcal{R}^{q_i}$ ,  $g, g_r \in \mathcal{R}^p, b_r \in \mathcal{R}^1$ .

The solution of a general convex cone problem can be obtained numerically using publicly available software such as SDPT3 [13] and DSDP5 [14].

### 3.2 RKE with $l_1$ Loss

To convert the problem of equation (2) into a convex cone programming problem, without loss of generality, we let  $\Omega$  contain  $m$  distinct  $(i, j)$  pairs, which we index with  $r = 1, 2, \dots, m$ . Define  $I_N$  to be the  $N$ -dimensional identity matrix and  $e_{m,r}$  to be vector of length  $2m$  consisting of all zeros except for the  $r$ th element being 1 and  $(m+r)$ th element being  $-1$ . If we denote the  $r$ th element of  $\Omega$  as  $(i(r), j(r))$ , and with some abuse of the notation let  $i = i(r), j = j(r)$  and  $w \in P_{2m}$  with  $w(r) = w(r+m) = w_{i(r), j(r)}$ ,  $r = 1, \dots, m$ , we can formulate the problem of equation (2) as follows:

$$\begin{aligned} \min_{K \succeq 0, u \geq 0} \quad & w \cdot u + \lambda I_N \cdot K \\ \text{s.t.} \quad & d_{ij} - B_{ij} \cdot K + e_{m,r} \cdot u = 0, \quad \forall_r, \\ & K \in S_N, \quad u \in P_{2m}. \end{aligned} \tag{4}$$

### 3.3 ‘Newbie’ Formulation

We now consider the situation in which a solution  $K_N$  of (2) is known for some set of  $N$  objects. We wish to augment the optimal kernel (by one row and column), without changing any of its existing elements, to account for a new object. That is, we wish to find a new “pseudo-optimal” kernel  $\tilde{K}_{N+1}$  of the form

$$\tilde{K}_{N+1} = \begin{bmatrix} K_N & b^T \\ b & c \end{bmatrix} \succeq 0, \quad (5)$$

(where  $b \in \mathcal{R}^N$  and  $c$  is a scalar) that solves the following optimization problem:

$$\begin{aligned} \min_{c \geq 0, b} \sum_{i \in \Psi} w_i |d_{i,N+1} - B_{i,N+1} \cdot K_{N+1}| \\ \text{s.t. } b \in \text{Range}(K_N), \quad c - b^T K_N^+ b \geq 0, \end{aligned} \quad (6)$$

where  $K_N^+$  is the pseudo-inverse of  $K_N$  and  $\Psi$  is a subset of  $\{1, 2, \dots, N\}$  of size  $t$ . The quantities  $w_i$ ,  $i \in \Psi$  are the weights assigned to the dissimilarity data for the new point. The constraints in this problem are the necessary and sufficient conditions for  $\tilde{K}_{N+1}$  to be positive semidefinite.

Suppose that  $K_N$  has rank  $p < N$  and let  $K_N = \Gamma \Lambda \Gamma^T$ , where  $\Gamma_{N \times p}$  is the orthogonal matrix of non-zero eigenvectors and  $\Lambda$  is the  $p \times p$  matrix of positive eigenvalues of  $K_N$ . By introducing the variable  $\tilde{b}$  and setting  $b = \Gamma \Lambda^{1/2} \tilde{b}$ , we can ensure that the requirement  $b \in \text{Range}(K_N)$  is satisfied. We also introduce the scalar variable  $\tilde{c}$ , and enforce  $c \geq \tilde{c}^2$  by requiring that

$$Z \stackrel{\text{def}}{=} \begin{bmatrix} 1 & \tilde{c} \\ \tilde{c} & c \end{bmatrix} \in S_2. \quad (7)$$

Using these changes of variable, the condition  $c - b^T K_N^+ b \geq 0$  is implied by the condition

$$x \stackrel{\text{def}}{=} [\tilde{c} \quad \tilde{b}^T]^T \in Q_{p+1}.$$

Further we define the  $N \times (p+1)$  matrix  $\Sigma \stackrel{\text{def}}{=} [0_N \quad 2\Gamma \Lambda^{1/2}]$ , where  $0_N$  is the zero vector of length  $N$ , and let  $\Sigma_i$  be the row vector consisting of the  $p+1$  elements of row  $i$  of  $\Sigma$ . We use  $K_N(i, i)$  to denote the  $i$ th entry of  $K_N$  and define the weight vector  $w \in P_{2t}$  with components  $w(r) = w(t+r) = w_{i(r)}$ ,  $r = 1, \dots, t$ . We then replace problem (6) by the following equivalent convex cone program:

$$\begin{aligned} \min_{Z \succeq 0, u \geq 0, x} w \cdot u \\ \text{s.t. } \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \cdot Z = 1, \\ \begin{bmatrix} 0 & 0.5 \\ 0.5 & 0 \end{bmatrix} \cdot Z - \begin{bmatrix} 1 \\ o_p \end{bmatrix} \cdot x = 0, \\ d_{i,N+1} - K_N(i, i) - \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \cdot Z \\ + \Sigma_i \cdot x + e_{t,r} \cdot u = 0, \quad \forall r=1,2,\dots,t, \\ Z \in S_2, \quad x \in Q_{p+1}, \quad u \in P_{2t}, \end{aligned}$$

where  $i = i(r)$  as before. Note that the constraints on  $Z$  ensure that it has the form (7).

### 3.4 Choosing Elements of $\Omega$

If the dissimilarity information is symmetric (i.e.,  $d_{ij} = d_{ji}$ ), we can choose  $\Omega$  to be the subset of  $\{(i, j) : i < j\}$  for which information is available. However, the codes we use for solving formulation (4) [13, 14] require  $O(m^2)$  storage (where  $m$  is the size of  $\Omega$ ), which is prohibitive for the application we describe in Section 4. Hence, we define  $\Omega$  by randomly selecting a subset of the available dissimilarity information in a way that ensures that each object  $i$  appears with roughly the same frequency among the  $(i, j)$  pairs of  $\Omega$ . Specifically, for each  $i$ , we choose a fixed number  $k$  of pairs  $(i, j)$  with  $j \neq i$  (we call the objects  $j$  “buddies” of  $i$ ) and add either  $(i, j)$  or  $(j, i)$  to the set  $\Omega$ , reordering if necessary to ensure that the first index of the pair is smaller than the second. (It is possible that  $(j, i)$  has been placed in  $\Omega$  at an earlier stage.) We choose the parameter  $k$  sufficiently large that the solution of (4) does not vary noticeably with different random subsets.

The “newbie” formulation (6) is comparatively inexpensive to solve, so we take  $\Psi$  to be the complete set of objects for which dissimilarity information  $d_{i,N+1}$  is available.

### 3.5 Eigenanalysis, Tuning, Truncation

The left five panels of Figure 1 illustrate the effect of varying  $\lambda$  on the eigenvalues of the regularized estimate of  $K$  obtained by solving (4). The data is from the example to be discussed in Section 4 below, with  $N = 280$  objects and  $k = 55$  “buddies” for each of the  $N$  objects. Note that the vertical scale is in units of  $\log_{10} \lambda$ . As  $\lambda$  increases the smaller eigenvalues begin to shrink, although in this example there is a very broad range of values of  $\lambda$ , spanning several orders of magnitude, where the sensitivity to  $\lambda$  is barely visible. At  $\lambda = 10^{-8}$  the condition number of  $K$  is about  $10^3$ . As  $\lambda$  goes much past 200 in this example, the penalty on  $K$  dominates and the dissimilarity information in the data is suppressed.

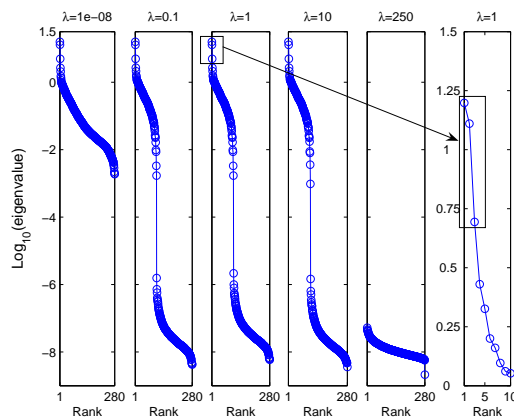


Figure 1: Left five panels: log scale eigensequence plots for five values of  $\lambda$ . As  $\lambda$  increases, smaller eigenvalues begin to shrink. Right panel: first ten eigenvalues of the  $\lambda = 1$  case displayed on a larger scale.

It is desirable to have a kernel with rank as low as possible while still respecting the data to an appropriate degree. Even if the rank of the regularized kernel estimate is not low, a low rank approximation obtained by setting all but a relatively small number of the largest eigenvalues to zero



might retain enough information to provide an efficient way of doing classification or clustering.

In the work described here, as well as in various simulation studies, we started with a very small positive  $\lambda$ , increased  $\lambda$  in a coarse log scale, and then experimented with retaining various numbers of eigenvalues to get a low rank approximation to  $K$ . The rightmost panel in Figure 1 shows the first 10 eigenvalues for the  $\lambda = 1$  case, in an expanded log scale. Natural breaks appear after both the second and the third eigenvalues. Setting all the eigenvalues of  $K$  after the largest  $p$  to 0 results in the  $\nu$ th coordinates of the  $j$ th object as  $x_j(\nu) = \sqrt{\lambda_\nu} \phi_\nu(j)$ ,  $\nu = 1, 2, \dots, p$ , where the  $\lambda_\nu, \phi_\nu$  are the first  $p$  eigenvalues and eigenvectors of  $K$  and  $\phi_\nu(j)$  is the  $j$  component of  $\phi_\nu$ . We remark that the coordinates of the  $N$  objects are always centered at the origin since the RKE estimate of  $K$  always has the constant vector as a 0 eigenvector. In the example discussed in Section 4 below with four classes of labeled objects, the choice of  $\lambda = 1$  and  $p = 3$  provided plots with a clear, informative clustering on the labels, that was verified from the science of the subject matter. We note that using the estimated  $K$  or a low rank version of it as the kernel in an SVM will result in linear classification boundaries in the object coordinates, (piecewise linear in the case of the multiclass SVM (MSVM) of [15]). It will be seen in the plots for labeled objects in Section 4 that piecewise linear classification boundaries in  $p = 3$  coordinates would apparently do quite well. However, that will not necessarily always be the case, and a more flexible workhorse kernel in the  $p$  object coordinates can be used. The MSVM [15] comes with a cross validation based method for choosing the MSVM tuning parameter(s) in a labeled training set. In principle, the parameters  $\lambda$  and  $p$  here can be incorporated in that method or other related methods, as additional MSVM parameters. Further examination of principled methods of choosing  $\lambda$  and  $p$  along with MSVM tuning parameter(s) will be deferred to later work.

## 4 Protein Clustering and Visualization with RKE

### 4.1 Background

One of the challenging problems of contemporary biology is inferring molecular functions of unannotated proteins. A widely used successful method of protein function prediction is based on sequence similarity. Statistically significant sequence similarity, which is typically based on a pairwise alignment score between two proteins, forms the basis for inferring the same function. Two major related problems exist for predicting function from sequence. The first problem is the clustering of large number of unlabeled protein sequences into subfamilies for the purpose of easing database searches and grouping similar proteins together. The second problem is concerned with assigning new unannotated proteins to the closest class, given the labeled or clustered training data. There is a substantial amount of literature addressing these two problems. [16] employ profile hidden Markov models (HMMs) for both problems. Clustering of proteins is obtained by a mixture of profile HMMs whereas assignment of new protein sequences to the clusters/classes is based on the likelihood of the new sequence under each of the cluster specific HMMs. Later, [17] addresses the second problem by first obtaining an explicit vector of features (Fisher scores) for each protein sequence and then utilizing a variant of SVMs, using a kernel called the Fisher kernel for classification purposes. The feature vector for each protein sequence is based on the likelihood scores of the input sequence evaluated at the corresponding maximum likelihood estimates of the HMM parameters fitted on the training data. More recently, [18] similarly uses SVMs for protein classification. However, in contrast to obtaining a feature vector by likelihood scores, they define a feature vector for each protein



sequence as a vector of its pairwise sequence similarity scores to all other proteins. Alternatively, [19] represent protein sequences as vectors in a high-dimensional feature space using a string-based feature map and train an SVM based on these vectors using a mismatch kernel. These latter works clearly illustrate the advantage of representing each protein sequence by a high-dimensional feature vector in some coordinate system and the power of kernel methods for protein classification. The RKE methodology presented here provides an efficient way to represent each protein sequence by a feature vector in a chosen coordinate system by utilizing the pairwise dissimilarity between protein sequences.

## 4.2 Data

We illustrate the utility of RKE methodology using a dataset of globins that was first analyzed in [16] by a profile HMM approach. The dataset, distributed with the HMMER2 software package [20], has a total of 630 globin sequences. The globin family is a large family of heme-containing proteins with many sub-families. It is mainly involved in binding and/or transportation of oxygen. For illustration purposes, we randomly choose 280 sequences from these data so that three large sub-classes of the globin family (alpha chains, beta chains, myoglobins) are included along with a heterogeneous class containing various types of chains. This selection resulted in a total of 112 “alpha-globins”, 101 “beta-globins”, 40 “myoglobins”, and 27 “globins” (the heterogeneous class). The proportion of sequences in each class were taken to be proportional to the class sizes in the original dataset.

## 4.3 Implementation of RKE

We used the RKE formulation of Section 3.2 for this application. The `Biconductor` package `pairseqsim` [21] was used to obtain global pairwise alignment scores for all pairs of  $N = 280$  sequences. This procedure gave a total of  $N(N - 1)/2 = 39060$  similarity scores, which we then normalized to map into the interval  $[0, 1]$ . We used one minus each of these numbers as the dissimilarity measure for each pair of sequences. During this process, alignment parameters were taken to be equal to the BLAST server [22] defaults. To construct the active index pair set  $\Omega$ , we used the procedure described in Section 3.4 with  $k = 55$  randomly chosen buddies for each protein sequence. The set  $\Omega$  thus contained approximately 14000 sequence pairs, corresponding to about 36% of the size of the complete index set. Replicated runs with  $k = 55$  proved to be nearly indistinguishable, as judged by examination of eigenvalue and 3D plots and the measure:  $\sum_{i < j} |\hat{d}_{ij1} - \hat{d}_{ij2}| / \sum_{i < j} \frac{1}{2}(\hat{d}_{ij1} + \hat{d}_{ij2})$ , where the third subscript in  $\hat{d}_{ijk}$  indexes different replicates (the above measure is typically about 5% for each pairwise comparison). The tuning parameter  $\lambda$  is set to 1 in the plots that follow later in this section.

## 4.4 Visualization of the globin sequence space and results

Figure 2 displays the 3D representation of the sequence space of 280 globins. This figure shows that the first three coordinates of the protein sequence space, corresponding to three largest eigenvalues, is already quite informative. The four main classes of the globin family are visually identifiable: The four colors red, blue, purple, and green represent alpha-globins, beta-globins, myoglobins, and globins, respectively.

Further investigation of this 3D plot reveals several interesting results. First, we observe that the five hemoglobin zeta chains, namely HBAZ\_HORSE, HBAZ\_HUMAN, HBAZ\_MOUSE, HBAZ\_PANTR, HBAZ\_PIG, shown by +, are located close to each other and are embedded within the alpha-globin cluster. Zeta-globin chains are alpha-like polypeptides and are synthesized in the yolk sac of the early embryo. It is well known that human zeta-globin polypeptide is more closely related to other mammalian embryonic alpha-like globins (i.e., zeta-globins) than to human alpha globins [23]. Furthermore, the zeta-globin gene in humans is a member of the alpha-globin gene cluster. Second we note that HBAM\_RANCA, which is represented by \* and is a hemoglobin alpha-type chain, seems to be isolated from the rest of the alpha-globin sequences. A possible explanation might be found in the structure of this protein. [24] note that the gene encoding this protein appeared through a gene duplication of hemoglobin and this took place near the time of the duplication that generated the alpha and beta chains. Our third observation is that the myoglobins MYG\_MUSAN, MYG\_THUAL, and MYG\_GALJA, denoted by □, - which are all fish myoglobins (*Mustelus antarcticus* (Gummy shark), *Thunnus albacares* (Yellowfin tuna), *Galeorhinus japonicus* (Shark)) - appear to be slightly separated from the rest of the myoglobin cluster. This is quite a remarkable observation because fish myoglobins are known to be structurally distinct from the mammalian myoglobins [26] and the RKE method nicely highlights this distinction on the basis of primary sequence data only. The 3D plot also reveals sub-clusters in the alpha-globin cluster. For example, all the 10 hemoglobin alpha-D chains (shown by △ in Figure 2) are clustered together within the alpha-globin cluster.

In a recent work, [27] provided a 3D plot of the protein structure space of 1898 chains. These authors utilized multi-dimensional scaling to project protein structures to a lower dimensional space based on the pairwise structural dissimilarity scores derived from 3D structures of proteins. Our application of RKE to the globin family, which is a few levels down from the top level of the protein structure hierarchy considered by [27], provide an analogous 3D plot for the sequence space of the globin family. It is quite encouraging that sub-protein domains of this family are readily distinguishable from the 3D embedding of the protein sequences. It is also worth mentioning that our current application is concerned only with pairwise sequence similarity, which can be obtained efficiently. However, clustering at the higher levels of the protein structure hierarchy is known to benefit enormously from 3D structural similarities and we plan to perform clustering at that level in future work (see Section 5 for details).

#### 4.5 Classification of new protein sequences.

We next illustrate how the newbie algorithm can be used to visualize unannotated protein sequences in the coordinate space of training data obtained by RKE. We used the following sequences as our test data: (1) HBAZ\_CAPHI: hemoglobin zeta chain from goat *Capra hircus*; (2) HBT\_PIG: Hemoglobin theta chain from pig *Sus scrofa*. Figure 3 displays the positions of these two test sequences with respect to 280 training sequences. We observe that HBAZ\_CAPHI (black circle) clusters nicely with the rest of the hemoglobin zeta chains, whereas HBT\_PIG (black star), which is an embryonic beta-type chain, is located closer to beta-globins. Additionally, we also used 17 Leghemoglobins (black triangles) as test data and found that these cluster tightly within the heterogeneous globin group. This observation is consistent with the results of [16], who also found a heterogeneous globin cluster with a tight sub-class of Leghemoglobins among their seven clusters obtained by a mixture of HMMs. These results indicate that RKE together with newbie algorithm provide a powerful means for clustering and classifying proteins.

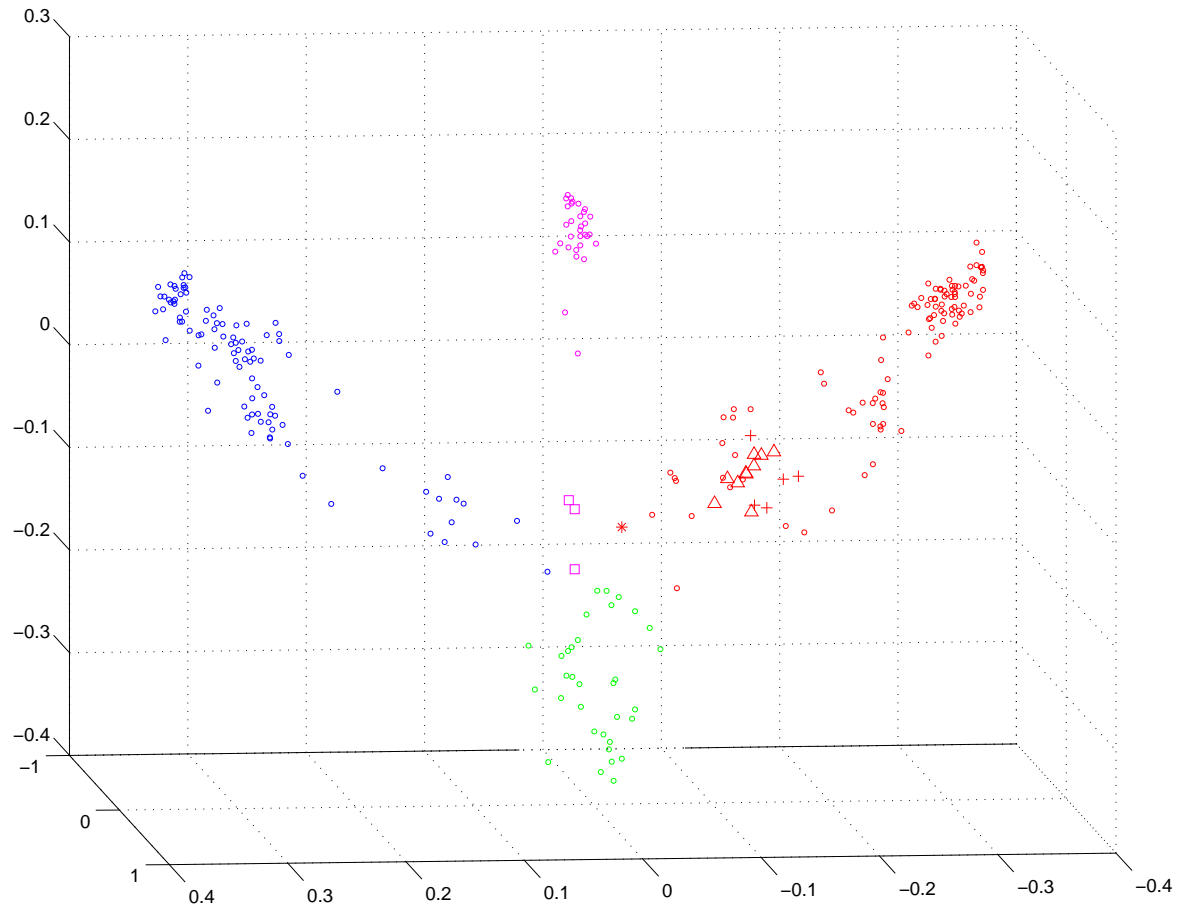


Figure 2: *3D representation of the sequence space for 280 proteins from the globin family.* Different subfamilies are encoded with different colors: Red symbols are alpha-globin subfamily, blue symbols are beta-globins, purple symbols represent myoglobin subfamily, and green symbols, scattered in the middle, are a heterogeneous group encompassing proteins from other small subfamilies within the globin family. Here, hemoglobin zeta chains are represented by the symbol +, fish myoglobins are marked by the symbol  $\square$ , and the diverged alpha-globin HBAM\_RANCA is shown by the symbol \*. Hemoglobin alpha-D chains, embedded within the alpha-globin cluster, are highlighted using the the symbol  $\triangle$ .

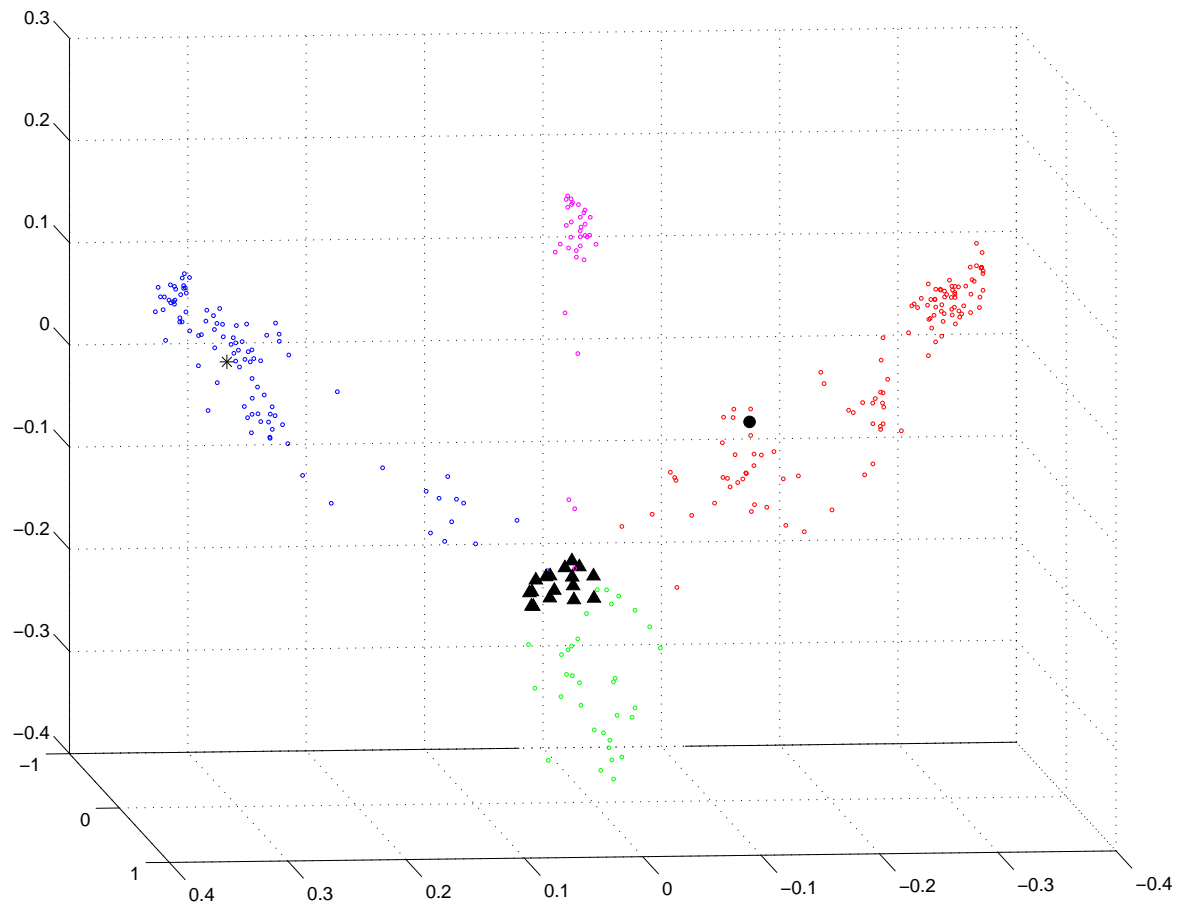


Figure 3: *Positioning test globin sequences in the coordinate system of 280 training sequences from the globin family. The newbie algorithm is used to locate one Hemoglobin zeta chain (black circle), one Hemoglobin theta chain (black star), and seventeen Leghemoglobins (black triangles) into the coordinate system of the training globin sequence data.*

## 5 Discussion

We have described a framework for estimating a regularized kernel (RKE methodology) from general dissimilarity information via the solution of a convex cone optimization problem. We have presented an application of the RKE methodology (including the “newbie” algorithm) to homology detection in the globin family of proteins. The most striking result here is perhaps the fact that a simple 3D plot is sufficient for visual identification of the subfamily information. However, in other applications, the plot coordinates (or higher dimensional coordinate vectors obtained by retaining more eigenvalues) may be used to build an automatic classification algorithm via the (principled) multicategory support vector machine (MSVM) [15]. That algorithm comes with a tuning method, it partitions the attribute space into regions for each training category, and it also comes with a method for signaling “none of the above”. Multicategory penalized likelihood estimates may also be used if there is substantial overlap of the data from different classes [9][28] [29][30].

A much harder problem in the context of protein classification and clustering is remote homology detection, that is, detecting homology in the presence of low sequence similarity. Since our framework accommodates an arbitrary notion of dissimilarities, we can easily take advantage of various types of dissimilarities such as presence or absence of discrete sequence motifs [31] and dissimilarities based on the primary, secondary, and tertiary structure [32], and obtain optimal kernels from each piece of information or data set. Without using labeled training sets, relations between a pair of kernels from different sources of information (or their lower rank approximations) can be quantified in various ways. A simple example is a measure of correlation:  $\sum_{ij} \hat{d}_{ij\alpha}^{s/2} \hat{d}_{ij\beta}^{s/2} / ((\sum_{ij} \hat{d}_{ij\alpha}^s)^{1/2} (\sum_{ij} \hat{d}_{ij\beta}^s)^{1/2})$  where  $\alpha$  and  $\beta$  index the different sources of information and  $s$  is a real number to be chosen. With labeled data, these kernels can further be examined and combined in an optimal way, as for example in [10], in the context of classification. As emphasized above, a striking feature of the presented methodology is the fact that it can exploit any type of dissimilarity measure, and data sets with missing information. These properties are clearly beneficial in biological data analysis, since many biologically relevant dissimilarities may not naturally result in positive semidefinite kernels (pairwise alignment scores, for example) which are essential for powerful classification methods such as SVMs.

Homology detection is one type of computational biology problem for which our framework offers rigorous, flexible tools. However, there are many other computational biology applications which can naturally be handled within this framework. Clustering of transcription factor position weight matrices (binding profiles) is one such application. With the increasing growth of transcription factor binding site databases, such as [33], a need for characterizing the space of DNA binding profiles and for developing tools to identify the class of newly estimated/studied profiles is emerging. A characterization of all available experimentally verified binding profiles such as in [33] might provide invaluable information regarding the possible class of binding profiles. Such information can then be utilized in supervised motif finding methods such as [34]. A natural dissimilarity measure for binding profiles is the Kullback-Leibler divergence. Clustering of the experimentally verified binding profiles based on a regularized kernel estimate of such dissimilarity measure might group binding profiles in a way that is consistent with the DNA binding domains of the transcription factors. We envision that this might generate a “protein binding profile space”, as the work of [27] generated a “protein structure space”.

Current work in progress includes both extension of the methodology and extension of the applications; in biology, the clustering of proteins at the top level of the protein hierarchy; and

in other contexts, medical images in particular. We are also investigating other choices of loss and penalty functionals in the noisy manifold unfolding problem (with Yi Lin). Future work of interest includes exploring the properties of alternatives provided here and their application in other contexts.

## Acknowledgments

The fourth author would like to acknowledge helpful conversations with Richard Gill, S. V. N. Vishwanathan, and David Callan. This work has been partly supported by NSF Grant DMS0072292 and NIH Grant EY09946 (Fan Lu and Grace Wahba), a WARF start-up grant from UW Madison (Sündüz Keleş), and NSF Grants ATM-0296033, CNS-0127857, CCF-0113051, ANI-0330538, DMS-0427689, CCF-0430504, and DOE grant DE-FG02-04ER25627 (Stephen Wright).

## References

- [1] Aronszajn, N. (1950) *Trans. Am. Math. Soc.* **68**, 337–404.
- [2] Kimeldorf, G. & Wahba, G. (1971) *J. Math. Anal. Applic.* **33**, 82–95.
- [3] Wahba, G. (1990) *Spline Models for Observational Data*. (SIAM). CBMS-NSF Regional Conference Series in Applied Mathematics, v. 59.
- [4] Wahba, G. (1977) *SIAM J. Numer. Anal.* **14**, 651–667.
- [5] Wahba, G. (1999) in *Advances in Kernel Methods-Support Vector Learning*, eds. Scholkopf, B., Burges, C., & Smola, A. (MIT Press), pp. 69–88.
- [6] Evgeniou, T., Pontil, M. & Poggio, T. (2000) *Advances in Computational Mathematics* **13**, 1–50.
- [7] Cristianini, N. & Shawe-Taylor, J. (2000) *An Introduction to Support Vector Machines*. (Cambridge University Press).
- [8] Scholkopf, B., Tsuda, K., & Vert, J-P. (2004) *Kernel Methods in Computational Biology*. (MIT Press).
- [9] Wahba, G. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 16524–16530.
- [10] Lanckriet, G., Cristianini, N., Bartlett, P., Ghaoui, L. E. & Jordan, M. (2004) *J. Mach. Learn. Res.* **5**, 27–72.
- [11] Shawe-Taylor, J. & Cristianini, N. (2004) *Kernel Methods for Pattern Analysis*. (Cambridge University Press).
- [12] Buja, A. & Swayne, D. (2002) *J. Classification* **19**, 7–43.
- [13] Tütüncü, R. H., Toh, K. C. & Todd, M. J. (2003) *Mathematical Programming* **95**, 189–217.

- [14] Benson, S. J. & Ye, Y. (2004) DSDP5: A software package implementing the dual-scaling algorithm for semidefinite programming, (Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL), Technical Report ANL/MCS-TM-255.
- [15] Lee, Y., Lin, Y. & Wahba, G. (2004) *J. Amer. Statist. Assoc.* **99**, 67–81.
- [16] Krogh, A., Brown, M., Mian, I. S., Sjolander, K. & Haussler, D. (1994) *J. Molecular Biology* **235**, 1501–1531.
- [17] Jaakkola, T., Diekhans, M. & Haussler, D. (2000) *J. Computational Biology* **7**, 95–114.
- [18] Liao, L. & Noble, W. S. (2003) *J. Computational Biology* **10**, 857–868.
- [19] Leslie, C., Eskin, E., Cohen, A., Weston, J. & Noble, W. S. (2004) *Bioinformatics* **20**, 467–476.
- [20] Eddy, S. R. (1998) *Bioinformatics* **14**, 755–763.
- [21] Gentleman, R. C., Carey, V. J., Bates, D. J., Bolstad, B. M., Dettling, M., Dudoit, S., Ellis, B., Gautier, L., Ge, Y., Gentry, J. *et al.* (2004) *Genome Biology* **5**, doi:10.1186/gb-2004-5-10-r80.
- [22] Atschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. (1990) *J. Molecular biology* **215**, 403–410.
- [23] Clegg, J. B. & Gagnon, J. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 6076–6080.
- [24] Maeda, N. & Fitch, W. M. (1982) *J. Biological Chemistry* **257**, 2806–2815.
- [25] Gasteiger, E., Gattiker, A., Hoogland, C., Ivanyi, I., Appel, R. D. & Bairoch, A. (2003) *Nucleic Acids Research* **31**, 3784–3788.
- [26] Cashon, R., Vayda, M. & Sidell, B. D. (1997) *Comparative Biochemistry and Physiology* **117B**, 613–620.
- [27] Hou, J., Jun, S., Zhang, C. & Kim, S. (2005) *Proc. Natl. Acad. Sci. USA* **102**, 3651–3656.
- [28] Lin, X. (1998) Ph.D. thesis (Department of Statistics, University of Wisconsin, Madison WI). Also Statistics Dept TR 1003 available via Grace Wahba’s website.
- [29] Wahba, G., Gu, C., Wang, Y. & Chappell, R. (1995) In *The Mathematics of Generalization, Santa Fe Institute Studies in the Sciences of Complexity, Proc. Vol. XX*, ed. Wolpert, D. (Addison-Wesley, Reading, MA), pp. 329–360.
- [30] Zhu, J. & Hastie, T. (2004) *Biostatistics* **5**, 329–340.
- [31] Ben-Hur, A. & Brutlag, D. (2003) *Bioinformatics* **19**, i26–i33.
- [32] Tang, C. L., Xie, L., Koh, I. Y. Y., Posy, S., Alexov, E. & Honig, B. (2003) *J. Molecular Biology* **334**, 1043–1062.



- [33] Sandelin, A., Alkema, W., Engström, P., Wasserman, W. & Lenhard, B. (2004) *Nucleic Acids Research* **32**, D91–D94.
- [34] Keleş, S., van der Laan, M. J., Dudoit, S., Xing, B. & Eisen, M. B. (2003) *Statistical Applications in Genetics and Molecular Biology* **2**, Article 5.

## Appendix

### Formulations with Square Loss Functions

We describe here the formulations of the RKE and Newbie Problems when a square loss function is used in place of an  $l_1$  loss function.

#### RKE Formulation

The sum-of-squares variant of (2) is

$$\min_{K \succeq 0} \sum_{(i,j) \in \Omega} w_{ij} (d_{ij} - B_{ij} \cdot K)^2 + \lambda \text{trace}(K). \quad (8)$$

We use  $r = 1, 2, \dots, m$  to index the elements of  $\Omega$ , as in Section 3.2, and obtain the following analog of (4):

$$\begin{aligned} \min \sum_{r=1}^m \begin{bmatrix} 0 & 0 \\ 0 & w(r) \end{bmatrix} \cdot X_r + \lambda I_N \cdot K \\ \text{s.t. } d_{ij} - B_{ij} \cdot K - \begin{bmatrix} 0 & 0.5 \\ 0.5 & 0 \end{bmatrix} \cdot X_r &= 0, \forall_r \\ \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \cdot X_r &= 1, \forall_r \\ K \succeq 0, X_r \succeq 0, &\quad \forall_r, \end{aligned} \quad (9)$$

where  $K \in S_N$  and  $X_r \in S_2$  for  $r = 1, 2, \dots, m$  and  $(i, j) = (i(r), j(r))$  as in Section 3.2. This is a pure positive semi-definite programming (SDP) problem. However, because of the number of constraints doubles, several times more space was needed for formulation (9) than for formulation (4). And in simulation experiments, we obtained similar results using the  $l_1$  and square loss formulations as  $\lambda$  varies. Hence, we decided to use the  $l_1$  formulation since it allowed us to solve larger problems.

#### Newbie Formulation

The Newbie Problem with least-squares objective is

$$\begin{aligned} \min \sum_{i \in \Psi} w_i (d_{i,N+1} - B_{i,N+1} \cdot K_{N+1})^2 \\ \text{s.t. } b \in \text{Range}(K_N), \quad c - b^T K_N^+ b \geq 0, \end{aligned}$$

where  $c \in \mathcal{R}$  and  $b \in \mathcal{R}^p$ . After performing the same transformations as in Subsection 3.3, we obtain the following formulation of this problem as a convex conic program:

$$\begin{aligned}
& \min \sum_{r=1}^t \begin{bmatrix} 0 & 0 \\ 0 & w(r) \end{bmatrix} \cdot X_r \\
& \text{s.t.} \quad \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \cdot Z = 1, \\
& \quad \quad \quad \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \cdot X_r = 1, \quad \forall r, \\
& \quad \quad \quad \begin{bmatrix} 0 & 0.5 \\ 0.5 & 0 \end{bmatrix} \cdot Z - \begin{bmatrix} 1 \\ o_p \end{bmatrix} \cdot x = 0, \\
& \quad \quad \quad \begin{bmatrix} 0 & 0.5 \\ 0.5 & 0 \end{bmatrix} \cdot X_r - \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \cdot Z \\
& \quad \quad \quad + \Sigma_i \cdot x + d_{i,N+1} - K_N(i, i) = 0, \quad \forall r, \\
& \quad \quad \quad Z \succeq 0, \quad X_r \succeq 0, \quad \forall r,
\end{aligned}$$

where

$$Z \in S_2, \quad X_r \in S_2, \quad \forall r \quad x \in Q_{p+1}.$$