
Online Kernel PCA with Entropic Matrix Updates

Dima Kuzmin
Manfred K. Warmuth

DIMA@CSE.UCSC.EDU
MANFRED@CSE.UCSC.EDU

Computer Science Department, University of California - Santa Cruz

Abstract

A number of updates for density matrices have been developed recently that are motivated by relative entropy minimization problems. The updates involve a softmin calculation based on matrix logs and matrix exponentials. We show that these updates can be kernelized. This is important because the bounds provable for these algorithms are logarithmic in the feature dimension (provided that the 2-norm of feature vectors is bounded by a constant). The main problem we focus on is the kernelization of an online PCA algorithm which belongs to this family of updates.

1. Introduction

There are two main families of updates in machine learning when the parameter of the algorithm and the instances are vectors. The first family maintains a parameter vector that is a linear combination of the instances and for the second family the component-wise logarithm of the parameter vector is a linear combination of the instances. Family one is based on regularizing with the squared Euclidean distance and family two uses relative entropy as the regularization (Kivinen & Warmuth, 1997). Both families have their advantages. Family one can be kernelized, i.e. the instances \mathbf{x} can be expanded into a feature vector $\phi(\mathbf{x})$ and the algorithm can be computed efficiently provided the kernel function $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$ can be computed efficiently. The advantage of family two is that it allows for bounds that are logarithmic in the dimension of the instances (see e.g. (Warmuth & Vishwanathan, 2005) for an extended discussion).

More recently the entropic family of updates has been

Appearing in *Proceedings of the 24th International Conference on Machine Learning*, Corvallis, OR, 2007. Copyright 2007 by the author(s)/owner(s).

generalized to the case when the instances are symmetric matrices \mathbf{X} and the parameter is a density matrix (Tsuda et al., 2005; Warmuth & Kuzmin, 2006a; Warmuth & Kuzmin, 2006b; Arora & Kale, 2007). The regularization is now the quantum relative entropy for density matrices instead of the regular relative entropy for probability vectors, and the matrix logarithm of the density matrix parameter is essentially a linear combination of the instance matrices. We show in this paper that if the instances \mathbf{X} are outer products of feature vectors,¹ i.e. $\mathbf{X} = \begin{matrix} & \phi(\mathbf{x}) & \\ \phi(\mathbf{x})^\top & & \end{matrix}$, then the matrix generalization of the entropic updates can be kernelized. If there is an efficient kernel function available, then the most expensive operation in computing the generalized updates is the computation of the eigendecomposition of the kernel matrix.

2. Offline Kernel PCA

Principal Component Analysis (PCA) is a well-known algorithm that chooses a low-dimensional linear subspace that best approximates the data (See Offline PCA Algorithm 1).

The Offline PCA Algorithm can be kernelized (see Offline Kernel PCA Algorithm 2): That is, the instances \mathbf{x} can be transformed into feature vectors $\phi(\mathbf{x})$, the individual features of the instances never need to be accessed, and the entire computation can be done via calls to a kernel function (Schölkopf et al., 1998). For the sake of simplicity, we assume in this paper that the data is centered (in feature space) and deal with uncentered case in the full paper.

In the kernel setting, the m columns of the data matrix \mathbf{X} are the expanded instances $\phi(\mathbf{x}_i) \in \mathbb{R}^N$, where the number of features N is typically much larger than the dimension n of the original instances \mathbf{x}_i . If we ran PCA in feature space, then we would need the eigendecomposition of $\mathbf{C} = \mathbf{X}\mathbf{X}^\top$ and this is too expensive when N is large. However there is a related matrix

¹Such outer products are called *dyads*

Algorithm 1 Offline PCA Algorithm

Input: Data matrix \mathbf{X} $n \times m$ having the data points $\mathbf{x}_i \in \mathbb{R}^n$ as columns, dimension r
 Compute covariance matrix $\mathbf{C} = \mathbf{X}\mathbf{X}^\top$ $n \times n$
 Compute eigendecomposition of the covariance matrix

$$\mathbf{C} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$$

Let \mathbf{U}_r be the column matrix of eigenvectors for the r largest eigenvalues. Also, $\mathbf{P} = \mathbf{U}_r \mathbf{U}_r^\top$ $n \times n$

Projection: $\mathbf{U}_r^\top \mathbf{x}$ is the r -dimensional compression vector for \mathbf{x} and $\mathbf{P}\mathbf{x}$ is the projection into the original space

Approximation error:

$$\|\mathbf{x} - \mathbf{P}\mathbf{x}\|_2^2 = \text{tr}((\mathbf{I} - \mathbf{P})\mathbf{x}\mathbf{x}^\top)$$

$\mathbf{K} = \mathbf{X}^\top \mathbf{X}$ $m \times m$ of the same rank called the *kernel matrix* because $\mathbf{K}_{ij} = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$. Kernel PCA gets away with computing only the eigendecomposition of this smaller matrix and the projection onto the principal components can also be obtained via kernel computations. The following lemma relates the eigensystems of the two matrices.

Lemma 1. *If the kernel matrix \mathbf{K} has eigenvalues λ_i and eigenvectors \mathbf{u}_i , then the covariance matrix \mathbf{C} has the same eigenvalues λ_i and orthonormal eigenvectors $\frac{\mathbf{X}\mathbf{u}_i}{\sqrt{\lambda_i}}$.*

Proof. Since

$$\mathbf{C}\mathbf{X}\mathbf{u}_i = \mathbf{X}\mathbf{X}^\top \mathbf{X}\mathbf{u}_i = \mathbf{X}\mathbf{K}\mathbf{u}_i = \lambda_i \mathbf{X}\mathbf{u}_i,$$

the eigenvalues λ_i for \mathbf{K} are also eigenvalues for \mathbf{C} . The orthogonality of the eigenvectors \mathbf{u}_i and \mathbf{u}_j of \mathbf{K} implies that $\mathbf{X}\mathbf{u}_i$ and $\mathbf{X}\mathbf{u}_j$ are also orthogonal:

$$(\mathbf{X}\mathbf{u}_i)^\top (\mathbf{X}\mathbf{u}_j) = \mathbf{u}_i^\top \mathbf{X}^\top \mathbf{X}\mathbf{u}_j = \mathbf{u}_i^\top \mathbf{K}\mathbf{u}_j = \lambda_j \mathbf{u}_i^\top \mathbf{u}_j = 0.$$

Finally, the $\frac{1}{\sqrt{\lambda_i}}$ factor makes our eigenvectors have norm 1:

$$\left\| \frac{\mathbf{X}\mathbf{u}_i}{\sqrt{\lambda_i}} \right\|_2^2 = \frac{1}{\lambda_i} (\mathbf{X}\mathbf{u}_i)^\top (\mathbf{X}\mathbf{u}_i) = \frac{1}{\lambda_i} \mathbf{u}_i^\top \mathbf{K}\mathbf{u}_i = \frac{1}{\lambda_i} \lambda_i \mathbf{u}_i^\top \mathbf{u}_i.$$

We conclude that the vectors $\frac{\mathbf{X}\mathbf{u}_i}{\sqrt{\lambda_i}}$ form an orthonormal eigensystem for \mathbf{C} . \square

Algorithm 2 Offline Kernel PCA Algorithm

Input: Data points $\mathbf{x}_1, \dots, \mathbf{x}_m$, dimension r , kernel function k
 Compute kernel matrix $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ $m \times m$
 Compute eigendecomposition of the kernel matrix

$$\mathbf{K} = \sum_{i=1}^m \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$$

Let \mathbf{U}_r be the matrix with the top r eigenvectors as columns and let $\mathbf{\Lambda}_r$ be the diagonal matrix of the top r eigenvalues

$$\text{Compute } \tilde{\mathbf{P}} = \mathbf{U}_r \mathbf{\Lambda}_r^{-1} \mathbf{U}_r^\top$$

For a new data point \mathbf{x} compute the kernel vector

$$\mathbf{y}_i = k(\mathbf{x}_i, \mathbf{x})$$

Projection: $\mathbf{\Lambda}_r^{-1} \mathbf{U}_r^\top \mathbf{y}$ is the r -dimensional compression vector for \mathbf{x}

Approximation error:

$$k(\mathbf{x}, \mathbf{x}) - \text{tr}(\tilde{\mathbf{P}}\mathbf{y}\mathbf{y}^\top)$$

3. Online Kernel PCA Algorithm

Online algorithms often maintain a parameter which expresses the uncertainty of the algorithm over which choice is good. For example, in the expert setting, Weighted Majority Algorithm maintains a probability vector over the experts. The algorithm hedges its bets by using a softmin function on the experts' losses so far, since committing to the best expert can be easily exploited by an adversary.

In the case of PCA the uncertainty over the r -dimensional subspace can be represented as a density matrix. The Online PCA Algorithm in (Warmuth & Kuzmin, 2006b) uses this idea and has good worst-case loss bounds. It uses the density matrix update based on matrix logs and exponentials like the one used in (Tsuda et al., 2005; Warmuth & Kuzmin, 2006a; Arora & Kale, 2007) but with the additional constraint that the eigenvalues are upper bounded (capped). Density matrices are symmetric, positive definite matrices of trace one. Thus their eigenvalues form a probability vector and the whole matrix can be seen as a mixture of the dyads formed by its eigenvectors.

A subspace of dimension r is characterized by a projection matrix \mathbf{P} of rank r . We need to work with the complementary subspace which is characterized by

the projection matrix $\mathbf{I} - \mathbf{P}$ of rank $n - r$. The matrix $\mathbf{I} - \mathbf{P}$ is the sum of $n - r$ orthogonal dyads and the scaled projection matrix $\frac{1}{n-r}\mathbf{P}$ is a density matrix with $n - r$ non-zero eigenvalues of value $\frac{1}{n-r}$. The Online PCA Algorithm maintains a mixture of scaled rank $n - r$ projection matrices and such mixtures are density matrices with the constraint that each eigenvalue is at most $\frac{1}{n-r}$. Crucially, every convex combination of rank $n - r$ projection matrices is a constrained density matrix and every such density matrix can be written as convex combination of at most n rank $n - r$ projection matrices, thus making capped density matrices the natural objects for representing uncertainty over subspaces.

The parameter matrix \mathbf{W}^t of the Online PCA Algorithm at trial t is a capped density matrix. The parameter update involves matrix logs and exponentials as well as a step that caps the eigenvalues of the updated matrix. At each trial the online algorithm produces a rank r projection matrix \mathbf{P}^t for compressing the current data point \mathbf{x}^t by decomposing the capped density matrix \mathbf{W}^t into a mixture of scaled projection matrices of rank $n - r$. From this mixture it is easy to probabilistically choose a single projection matrix of rank $n - r$ and then choose its complement as the projection matrix \mathbf{P}^t used for compressing at trial t . The decomposition and the capping can be done efficiently via algorithms 4 and 5, respectively. Note that capping in this context doesn't mean that we just set any components above the threshold to the threshold value, since this would result in a vector that does not sum to one anymore. More details on the Decomposition Algorithm 4 are provided in (Warmuth & Kuzmin, 2006b) and a linear time implementation of the Capping Algorithm 5 is given in (Herbster & Warmuth, 2001). The Online PCA Algorithm 3 gives a summary of all the steps. Here a *corner* is a diagonal matrix with $n - r$ of the diagonal elements set to $\frac{1}{n-r}$ and the rest set to zero.

An important property of the standard PCA problem is that the 2-norm approximation error is linear in the matrix parameter space:

$$\|\mathbf{x} - \mathbf{P}\mathbf{x}\|_2^2 = \text{tr}((\mathbf{I} - \mathbf{P})\mathbf{x}\mathbf{x}^\top)$$

This means that the expected loss of the online algorithm that chooses a subspace probabilistically based on the decomposed density matrix $\mathbf{W} = \frac{1}{n-r} \sum_i p_i (\mathbf{I} - \mathbf{P}_i)$ will be equal to the linear loss of the matrix parameter \mathbf{W} :

$$\sum_i p_i \|\mathbf{x} - \mathbf{P}_i \mathbf{x}\|_2^2 = (n - r) \text{tr}(\mathbf{W}\mathbf{x}\mathbf{x}^\top).$$

(We need to multiply by $n - r$ because the density

Algorithm 3 Online PCA Algorithm

Input: Dimension r and initial density matrix \mathbf{W}^1 with eigenvalues at most $\frac{1}{n-r}$, learning rate $\eta > 0$
for $t = 1$ to T **do**
 Perform eigendecomposition $\mathbf{W}^t = \mathbf{U}\boldsymbol{\lambda}\mathbf{U}^\top$
 Decompose $\boldsymbol{\lambda}$ into a mixture of n corners using the Decomposition Algorithm 4: $\boldsymbol{\lambda} = \sum_{j=1}^n p_j \mathbf{c}_j$, where the p_j are the mixture coefficients and \mathbf{c}_j the corners.
 Draw a corner $\mathbf{c} = \mathbf{c}_j$ with probability p_j
 Form a matrix corner $\mathbf{R} = \mathbf{U} \text{diag}(\mathbf{c}) \mathbf{U}^\top$
 Form a rank r projection matrix

$$\mathbf{P}^t = \mathbf{I} - (n - r)\mathbf{R}$$

Let \mathbf{U}_r^t be a column matrix of those r eigenvectors that correspond to zero entries in \mathbf{c}

Receive data instance vector \mathbf{x}^t

Projection: $(\mathbf{U}_r^t)^\top \mathbf{x}^t$ is the r -dimensional compression vector for \mathbf{x}^t

Approximation error:

$$\|\mathbf{x}^t - \mathbf{P}^t \mathbf{x}^t\|_2^2 = \text{tr}((\mathbf{I} - \mathbf{P}^t) \mathbf{x}^t (\mathbf{x}^t)^\top)$$

Expected approximation error:

$$(n - r) \text{tr}(\mathbf{W}^t \mathbf{x}^t (\mathbf{x}^t)^\top)$$

Update:

$$\widehat{\mathbf{W}}^t = \frac{\exp(\log \mathbf{W}^t - \eta \mathbf{x}^t (\mathbf{x}^t)^\top)}{Z^t},$$

where Z^t normalizes the trace to 1

\mathbf{W}^{t+1} is obtained from $\widehat{\mathbf{W}}^t$ by capping the eigenvalues to at most $\frac{1}{n-r}$ using the Capping Algorithm 5.

end for

Algorithm 4 Decomposition Algorithm

Input probability vector \mathbf{w} , $\max(\mathbf{w}) \leq 1/k$

repeat

Pick a k -corner \mathbf{r} whose components correspond to non-zero components of \mathbf{w} and contain all the components of \mathbf{w} that are equal to $\frac{|\mathbf{w}|}{k}$

Let s be the smallest of the k components in \mathbf{w}

Let l be the largest of the remaining $n - k$ components

$$\mathbf{w} := \mathbf{w} - \overbrace{\min(m s, |\mathbf{w}| - m l)}^p \mathbf{r}$$

output $p\mathbf{r}$

until $\mathbf{w} = \mathbf{0}$

Algorithm 5 Capping Algorithm

Input ω probability vector, $1/k$ capping constant
 Let ω^\downarrow index the vector in decreasing order, i.e.
 $\omega_1^\downarrow = \max(\omega)$
if $\max(\omega) \leq 1/k$ **then**
 Return ω
end if
 $i := 1$
repeat
 Set first i largest components to $1/k$ and normal-
 ize the rest:
 $\bar{\omega} := \omega$
 $\bar{\omega}_j^\downarrow := 1/k, j = 1 \dots i$
 $\bar{\omega}_j^\downarrow := \frac{k-i}{k} \frac{\bar{\omega}_j^\downarrow}{\sum_{l=i+1}^n \bar{\omega}_l^\downarrow}, j = i + 1 \dots n$
 $i := i + 1$
until $\max(\bar{\omega}) \leq 1/k$
 Return $\bar{\omega}$

matrix is a mixture of scaled projection matrices.) One might want to solve a PCA-like problem, where the approximation error is measured by something other than a 2-norm. Online learning techniques exist for dealing with very many different loss functions (see e.g. (Kivinen & Warmuth, 1999)), but only for the case of linear loss will the expectation of the losses equal the loss of the expected parameter. Thus, for other loss functions the algorithm might have good bounds, but it won't produce an actual low-dimensional subspace as a prediction.

Our new Online Kernel PCA Algorithm uses techniques from the Offline Kernel PCA Algorithm and makes use of the relationship between the offline and online PCA algorithms. Note that offline PCA picks the eigenvectors corresponding to the top k eigenvalues of the covariance matrix. Online PCA can be seen as a different way of choosing k eigenvectors from the eigensystem of the data covariance matrix: Probabilistically choose a rank $n - k$ subspace based on a softmin function of the eigenvalues and then use the complementary subspace of rank k as the subspace for compression. Online kernel PCA uses a softmin function on the kernel matrix instead of the covariance matrix.

First note that if we drop the capping constraint and initialize with $\mathbf{W}^1 = \frac{1}{n} \mathbf{I}$, then the density matrix at step t of the Online PCA Algorithm 3 has the form

$$\mathbf{W}^t = \frac{\exp(-\eta \sum_{i=1}^{t-1} \mathbf{x}^i (\mathbf{x}^i)^\top)}{Z^t}.$$

Thus in this case, the parameter matrix \mathbf{W}^t has the same eigensystem as the covariance matrix and its eigenvalues are obtained by applying a softmin func-

tion to the eigenvalues of the covariance matrix. By Lemma 1 which relates the eigensystems of the covariance and kernel matrices, all the information needed to compute \mathbf{W}^t is given by the eigendecomposition of the smaller kernel matrix. The capped version² of the above \mathbf{W}^t will be the crucial parameter matrix implicitly maintained by our Online Kernel PCA Algorithm 6. Our algorithm differs from the Offline Kernel PCA Algorithm only in the way it chooses k eigenvectors from the kernel matrix for forming the projection matrix: instead of taking the eigenvectors corresponding to the maximum k eigenvalues, it picks $n - k$ eigenvectors probabilistically based on a softmin function of the eigenvalues and then chooses the complementary k eigenvectors as the projection matrix. This more complicated procedure assures the expected compression error for all sequences of examples is not too much larger than the compression error of the offline algorithm.

The next section discusses how the proof of the original loss bound for online PCA needs to be altered to work with this kernelized version.

4. Relative Loss Bounds for Online Kernel PCA

We begin by recalling the motivation for the online PCA update as a solution of some optimization problem. The regularization used in this optimization is the quantum relative entropy for density matrices, which is defined as:

$$\Delta(\mathbf{V}, \mathbf{W}) = \text{tr}(\mathbf{V}(\log \mathbf{V} - \log \mathbf{W})) - \text{tr}(\mathbf{V}) + \text{tr}(\mathbf{W}).$$

The update of the Online PCA Algorithm 3 is determined as

$$\mathbf{W}^{t+1} = \underset{\substack{\text{tr}(\mathbf{W})=1 \\ \mathbf{W} \preceq \frac{1}{N-r} \mathbf{I}}}{\text{argmin}} \left(\Delta(\mathbf{W}, \mathbf{W}^t) + \eta \text{tr}(\mathbf{W} \mathbf{x}^t (\mathbf{x}^t)^\top) \right). \tag{1}$$

We can't kernelize the above update because the iterative capping complicates things.³ Instead we always go back to the beginning, which allows us to define

²Note that in the Online Kernel PCA Algorithm 6 the kernel matrix is capped with $\frac{1}{\epsilon-r}$ whereas in the original Online PCA Algorithm the covariance matrix is capped with $\frac{1}{N-r}$. It can be shown that the two types of capping have the same effect.

³Specifically, the accumulated capping Lagrangians can make the eigenvectors of the resulting density matrix different from the eigenvectors of the data covariance matrix. If you look at equation (4) in the proof of Theorem 1 later in this section, you can deduce the form of \mathbf{W}^t for the above update: $\mathbf{W}^t \sim \exp(-\eta \sum \mathbf{x}^i (\mathbf{x}^i)^\top - \sum \Gamma^i)$. In general, this is no longer a spectral function of the data covariance

Algorithm 6 Online Kernel PCA Algorithm

Input: dimension r , kernel function k , initial r data points $\mathbf{x}^1, \dots, \mathbf{x}^r \in \mathbb{R}^n$, learning rate η

for $t = r + 1$ to T **do**

 Compute kernel matrix on data so far:

$$\mathbf{K}_{ij}^t = k(\mathbf{x}^i, \mathbf{x}^j)$$

$(t-1) \times (t-1)$

 Eigendecompose the kernel matrix

$$\mathbf{K}^t = \sum_{i=1}^t \lambda_i \mathbf{u}_i \mathbf{u}_i^\top = \mathbf{U}^t \mathbf{\Lambda}^t (\mathbf{U}^t)^\top$$

 Transform the eigenvalues:

$$\bar{\lambda}_i \leftarrow \frac{e^{-\eta \lambda_i}}{Z}, \quad Z = \sum_{j=1}^{t-1} e^{-\eta \lambda_j}$$

 Cap values in vector $\bar{\lambda}$ to at most $\frac{1}{t-1-r}$ using the Capping Algorithm 5

 Decompose $\bar{\lambda}$ as $\sum_j p_j \mathbf{c}_j$, where \mathbf{c}_j are corners of size $t-1-r$ using the Decomposition Algorithm 4

 Draw the corner $\mathbf{c} = \mathbf{c}_j$ with probability p_j

 Form matrix \mathbf{U}_r^t with the columns being $(t-1) \times (r)$

 those eigenvectors of \mathbf{K}^t that correspond to zero values in the selected corner \mathbf{c} .

 Let $\mathbf{\Lambda}_r^t$ be the diagonal matrix of the eigenvalues of \mathbf{K}^t associated with eigenvectors in \mathbf{U}_r^t .

$$\tilde{\mathbf{P}}^t = \mathbf{U}_r^t (\mathbf{\Lambda}_r^t)^{-1} (\mathbf{U}_r^t)^\top$$

$$\tilde{\mathbf{W}}^t = \mathbf{U}^t (\mathbf{I} - (t-1-r) \text{diag}(\bar{\lambda})) (\mathbf{\Lambda}^t)^{-1} (\mathbf{U}^t)^\top$$

 Receive data instance vector \mathbf{x}^t

 Compute kernel vector $\mathbf{y}_i^t = k(\mathbf{x}_i, \mathbf{x}^t)$

Projection: $(\mathbf{\Lambda}_r^t)^{-1} (\mathbf{U}_r^t)^\top \mathbf{y}^t$ is the r -dimensional compression vector for \mathbf{x}^t

Approximation error:

$$k(\mathbf{x}^t, \mathbf{x}^t) - \text{tr}(\tilde{\mathbf{P}}^t \mathbf{y}^t (\mathbf{y}^t)^\top)$$

Expected approximation error:

$$k(\mathbf{x}^t, \mathbf{x}^t) - \text{tr}(\tilde{\mathbf{W}}^t \mathbf{y}^t (\mathbf{y}^t)^\top)$$

end for

\mathbf{W}^{t+1} by a single capping:

$$\mathbf{W}^{t+1} = \underset{\substack{\text{tr}(\mathbf{W})=1 \\ \mathbf{W} \preceq \frac{1}{N-r} \mathbf{I}}}{\text{argmin}} \left(\Delta(\mathbf{W}, \frac{1}{N} \mathbf{I}) + \eta \sum_{i=1}^t \text{tr}(\mathbf{W} \mathbf{x}^i (\mathbf{x}^i)^\top) \right). \tag{2}$$

We will use \mathbf{C}^{t+1} to denote the covariance matrix used in obtaining \mathbf{W}^{t+1} , i.e $\mathbf{C}^{t+1} = \sum_{i=1}^t \mathbf{x}^i (\mathbf{x}^i)^\top$. The resulting density matrix used by the algorithm looks like this:

$$\mathbf{W}^{t+1} = \text{cap} \left(\frac{\exp(-\eta \mathbf{C}^{t+1})}{Z} \right)$$

This is a spectral function of \mathbf{C}^{t+1} and thus allows for easy kernelization. With some additional effort, we were able to prove the same loss bound for the update (2) as was proven in (Warmuth & Kuzmin, 2006b) for the update (1). This is done in the following theorem, which says that the loss incurred by the algorithm over a sequence of trials is not much larger than the loss incurred by the offline algorithm that is given all of the data in advance. In our case this offline algorithm is the standard PCA for the entire data. Note that with our representation of mixtures of subspaces as density matrices, the expected loss of algorithm (i.e. the expected approximation error) at trial t has the form $(N-r) \text{tr}(\mathbf{W}^t \mathbf{x}^t (\mathbf{x}^t)^\top)$. The expectation is wrt the internal randomization of the algorithm, the bound itself makes no probabilistic assumptions on the data.

Theorem 1. *For any sequence of data $\mathbf{x}^1, \dots, \mathbf{x}^T \in \mathbb{R}^n$, such that the squared 2-norm of each instance is bounded by 1^4 , positive learning rate η and arbitrary capped comparator density matrix \mathbf{V} , the following relative loss bound holds:*

$$\sum_{t=1}^T (N-r) \text{tr}(\mathbf{W}^t \mathbf{x}^t (\mathbf{x}^t)^\top) \leq (N-r) \frac{\eta \sum_{t=1}^T \text{tr}(\mathbf{V} \mathbf{x}^t (\mathbf{x}^t)^\top) + \Delta(\mathbf{V}, \frac{1}{N} \mathbf{I})}{1 - \exp(-\eta)}$$

Proof. As was mentioned before, the update we use is different from online PCA described previously, and thus the proof of (Warmuth & Kuzmin, 2006b) does not apply. The old proof relied on bounding the progress of the algorithm in terms of the quantum relative entropy and then using generalized Pythagorean matrix. On the other hand, capping once only affects the eigenvalues of the capped matrix as shown in (Warmuth & Kuzmin, 2006b).

⁴The bound can be generalized based on the assumption that the squared norm of the expanded instances $\phi(\mathbf{x}_t)$ is upper bounded by some constant Q .

theorem for Bregman projections to deal with the capping constraints. Here we adopt a somewhat different approach which makes use of the dual optimization problem of the update (2). Let $F^t(\mathbf{W})$ be the objective function of this update, i.e.

$$F^t(\mathbf{W}) = \Delta(\mathbf{W}, \frac{1}{N}\mathbf{I}) + \eta \operatorname{tr}(\mathbf{W}\mathbf{C}^t)$$

Following (Kivinen & Warmuth, 1999), we use the optimum value of $F^t(\mathbf{W})$ as a potential, i.e.

$$P^t := \min_{\substack{\operatorname{tr}(\mathbf{W})=1 \\ \mathbf{W} \preceq \frac{1}{N}\mathbf{I}}} F^t(\mathbf{W})$$

We would like to lower bound the per trial drop of the potential in terms of the loss of the algorithm.

$$P^{t+1} - P^t \quad (3)$$

We use strong duality (Boyd & Vandenberghe, 2004) to compute the potential as a maximum of the dual problem. The purpose of going to the dual problem is only to help us analyze the algorithm, the algorithm itself is fully determined by the optimization problem in (2) and kernelization of the resulting density matrix computations. We begin by formulating a Lagrangian with δ as the dual variable for the trace constraint and the symmetric positive definite matrix $\mathbf{\Gamma}$ as the dual variable for the capping constraint. For the relevant derivatives see (Tsuda et al., 2005).

$$L^t(\mathbf{W}, \mathbf{\Gamma}, \delta) = \operatorname{tr}(\mathbf{W}(\log \mathbf{W} + \ln N \mathbf{I})) - \operatorname{tr}(\mathbf{W}) - \ln N + \eta \operatorname{tr}(\mathbf{W}\mathbf{C}^t) + \delta(\operatorname{tr}(\mathbf{W}) - 1) + \operatorname{tr}((\mathbf{W} - \frac{1}{r}\mathbf{I})\mathbf{\Gamma})$$

$$\frac{\partial L^t}{\partial \mathbf{W}} = \log \mathbf{W} + \ln N \mathbf{I} + \eta \mathbf{C}^t + \delta \mathbf{I} + \mathbf{\Gamma}.$$

The extra constant $\ln N$ can be rolled into δ , i.e. $\delta' = \delta + \ln N$. Setting the above to zero we get:

$$\mathbf{W} = \exp(-\eta \mathbf{C}^t - \delta' \mathbf{I} - \mathbf{\Gamma}). \quad (4)$$

We substitute this form of \mathbf{W} into L^t and enforce the $\operatorname{tr}(\mathbf{W}) = 1$ constraint by choosing $\delta' = \ln \operatorname{tr}(\exp(-\eta \mathbf{C}^t - \mathbf{\Gamma}))$. After numerous simplifications we arrive at the dual optimization problem

$$\max_{\mathbf{\Gamma} \succeq \mathbf{0}} L^t(\mathbf{\Gamma}) = -\ln \operatorname{tr} \left(\exp(-\eta \mathbf{C}^t - \mathbf{\Gamma}) - 1 - \frac{1}{r} \operatorname{tr}(\mathbf{\Gamma}) \right)$$

Let $\mathbf{\Gamma}^t$ be the maximizer of the dual optimization problem for L^t and define $\mathbf{\Gamma}^{t+1}$ similarly. Strong duality holds for our problem. We now want to lower bound the rewritten drop of the potential (3) by the same expression that was used in (Warmuth & Kuzmin, 2006a) for their proof:

$$L^{t+1}(\mathbf{\Gamma}^{t+1}) - L^t(\mathbf{\Gamma}^t) \geq -\ln \operatorname{tr}(\exp(\log \mathbf{W}^t - \eta \mathbf{x}^t (\mathbf{x}^t)^\top)). \quad (5)$$

We begin by substituting the solution

$$\mathbf{W}^t = \exp(-\eta \mathbf{C}^t - \ln \operatorname{tr}(\exp(-\eta \mathbf{C}^t - \mathbf{\Gamma}^t))\mathbf{I} - \mathbf{\Gamma}^t)$$

into the rhs and rewrite it as follows:

$$-\ln \operatorname{tr}(\exp(-\eta \mathbf{C}^{t+1} - \mathbf{\Gamma}^t)) + \ln \operatorname{tr}(\exp(-\eta \mathbf{C}^t - \mathbf{\Gamma}^t)) = L^{t+1}(\mathbf{\Gamma}^t) - L^t(\mathbf{\Gamma}^t).$$

By substituting this rhs into (5) and simplifying, we get that inequality (5) holds iff

$$L^{t+1}(\mathbf{\Gamma}^{t+1}) - L^{t+1}(\mathbf{\Gamma}^t) \geq 0.$$

However this inequality trivially hold, since $\mathbf{\Gamma}^{t+1}$ maximizes $L^{t+1}(\mathbf{\Gamma})$.

We now lower bound the rhs of (5) in term of the loss of the algorithm in the usual way (Warmuth & Kuzmin, 2006a):

$$\text{r.h.s.} \geq \operatorname{tr}(\mathbf{W}^t \mathbf{x}^t (\mathbf{x}^t)^\top) (1 - e^{-\eta})$$

Finally we sum the drop of the potential (5) over t . Telescoping occurs and we get the following bound:

$$\begin{aligned} P^{T+1} - \overbrace{P^1}^0 &= L^{T+1}(\mathbf{\Gamma}^{T+1}) - L^1(\mathbf{\Gamma}^1) \\ &\geq \sum_{t=1}^T \operatorname{tr}(\mathbf{W}^t \mathbf{x}^t (\mathbf{x}^t)^\top) (1 - e^{-\eta}) \end{aligned}$$

Since

$$P^{T+1} \leq \Delta(\mathbf{V}, \frac{1}{n}) + \eta \operatorname{tr}(\mathbf{V}\mathbf{C}^{T+1}),$$

the bound follows. \square

We now set η as a function of N , r , the maximum 2-norm Q of the instances $\phi(\mathbf{x}_t)$ and an upper bound on the loss of the best r -projection \hat{L} (Freund & Schapire, 1997):

$$\eta = \frac{\ln(1 + \sqrt{\frac{2r \ln \frac{N}{r}}{\hat{L}}})}{Q^2}. \quad (6)$$

This tuning of η results in the following bound for the Online Kernel PCA Algorithm, that holds for all sequences with 2-norm of expanded instances $\phi(\mathbf{x}^t)$ bounded by Q and that have the loss of the best r -projection bounded by \hat{L} :

$$\begin{aligned} &(\text{expected loss of alg.}) - (\text{loss best } r\text{-projection}) \\ &\leq Q \sqrt{2\hat{L} r \ln \frac{N}{r}} + rQ^2 \ln \frac{N}{r}. \quad (7) \end{aligned}$$

Note that the dependence on the feature dimension N is logarithmic. Therefore we can accommodate a large number of features without degrading the bound

substantially. For example for the polynomial kernel of degree d , $\ln N = d \ln n$, and if the feature vectors have bounded 2-norm then the bound is still reasonable. The current analysis is not directly applicable to the infinite-dimensional kernels. To handle that case we need a different way to bound the entropy of the comparator $\Delta(\mathbf{V}, \frac{1}{n} \mathbf{I})$. Intuitively, the comparator \mathbf{V} will only use the finite dimensional subspace spanned by the past expanded instances and thus the entropy remains bounded (probably by $O(\log T)$), but we haven't worked out the details yet.

5. Experiments

We implemented the Online Kernel PCA Algorithm and tested it on a simple synthetic dataset generated as follows: A point was picked at random on a cone defined by the equation $x^2 + y^2 - z^2 = 0$. These three dimensions were embedded into \mathbb{R}^{20} by padding with zeros. Finally Gaussian noise was added to all dimensions. The value of z was constrained to the range $[-0.6 : 0.6]$, which had the effect of bounding the 2-norm in feature space to at most 1. The resulting dataset is visualized in Figure 1.

Since the cone is a quadratic surface, PCA without kernelization is not able to pick up the 2-dimensional structure of this dataset. A polynomial kernel of degree 2 (i.e., $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^2$) is more suitable (the dimension N of the feature space in this case is 400). Figure 2 shows the total loss of our Online Kernel PCA Algorithm 6 and compares it to the total loss of the Offline Kernel PCA Algorithm 2 whose projection matrix is formed based on the entire dataset. As we can see the online algorithm is not too much worse. We also plotted the upper bound of Theorem 1 which not very tight on this random data set. (Note that our bounds hold for worst-case sequences.) Figure 3 shows the average "regret" per time step of our Online Kernel PCA Algorithm, i.e. the difference in total losses between the online and offline kernel PCA algorithms divided by the number of iterations so far.

While implementing the algorithm we found some numerical stability issues that need to be addressed. First, the theoretical guarantee for the algorithm allows us to assume an arbitrary bound Q on the 2-norm of expanded instance vectors $\phi(\mathbf{x}_t)$. In practice, however, large norms of the data instances lead to overly large eigenvalues of the kernel matrix, which then cause problems during exponentiation for the softmax computation. This difficulty is specific to our Online Kernel PCA Algorithm 6, which at trail t uses the kernel matrix formed by the past $t - 1$ instances and the eigenvalues of this matrix may diverge. Appropriate

scaling of the instance avoids this problem.

Second, the Capping Algorithm 5 may encounter the case where the remaining portion of the probability vector is all zeros due to the exponentiation of large numbers. In that case, a division by zero can occur. The proper approach in this case is to set enough of the remaining zeros to the capping constant, so that the whole vector sums to one again.

We also did some preliminary experiments using human face datasets and the Gaussian kernel. While we don't yet have a theoretical guarantee for infinite-dimensional feature spaces, the algorithm still seems to work in practice.

6. Conclusion

The question of whether kernels can be applied to updates that are motivated with relative entropy regularizations has generated great interest (See e.g. discussion in (Warmuth & Vishwanathan, 2005)). Only very few cases have been found for applying kernels with entropic updates when the instances are vectors (Takimoto & Warmuth, 2003). In this paper we show how the matrix generalizations of these updates can be kernelized when the instance matrices are of the form $\phi(\mathbf{x}_t)\phi(\mathbf{x}_t)^\top$. We chose the online PCA problem as our example problem for working out the details. The crucial property needed is that the exponent must be a linear combination of the outer products $\phi(\mathbf{x}_t)\phi(\mathbf{x}_t)^\top$. The other crucial property is that the exponential of this linear combination only depends on the eigenvalues of the linear combination and keeps the eigensystem unchanged. All applications of the Matrix Exponentiated Gradient Algorithm (Tsuda et al., 2005; Warmuth & Kuzmin, 2006a; Warmuth & Kuzmin, 2006b) have these properties and generalizations are straightforward to families of algorithms defined by other Bregman divergences.

In this paper we only considered bounds compared to the best fixed off-line comparator. However these online algorithms can be adapted to the case when the comparator shifts with time (See e.g. (Bousquet & Warmuth, 2002)). A thorough experimental analysis of these extensions would be useful.



Figure 1. Our synthetic dataset is a sample of a 3-dimensional cone embedded in \mathbb{R}^{20} with added Gaussian noise. We depict the key 3 of the 20 dimensions.

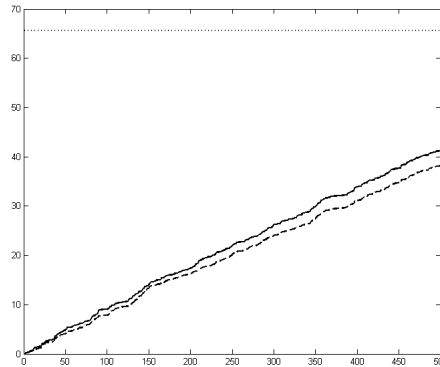


Figure 2. Plots of the total losses the algorithms and our bound versus the iteration number. Lowest (dashed) line: total loss of the Offline Kernel PCA Algorithm 2 (based on the entire dataset). Solid line right above: total loss of our Online Kernel PCA Algorithm with tuned learning rate (6) (for this dataset $\eta = 1.5$). Flat dotted line: tuned theoretical bound (7) on the total loss of the Online Kernel PCA Algorithm 6.

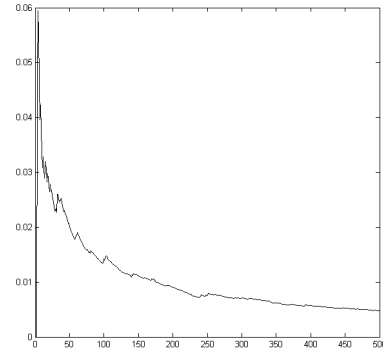


Figure 3. The average per-trial regret of the online algorithm converges to zero (Here η was set to 1).

References

- Arora, S., & Kale, S. (2007). A combinatorial primal-dual approach to semidefinite programs. *Proc. 39th Annual ACM Symposium on Theory of Computing*. ACM. To appear.
- Bousquet, O., & Warmuth, M. K. (2002). Tracking a small set of experts by mixing past posteriors. *Journal of Machine Learning Research*, 3, 363–396.
- Boyd, S., & Vandenberghe, L. (2004). *Convex optimization*. Cambridge University Press.
- Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to Boosting. *Journal of Computer and System Sciences*, 55, 119–139.
- Herbster, M., & Warmuth, M. K. (2001). Tracking the best linear predictor. *Journal of Machine Learning Research*, 1, 281–309.
- Kivinen, J., & Warmuth, M. K. (1997). Additive versus exponentiated gradient updates for linear prediction. *Information and Computation*, 132, 1–64.
- Kivinen, J., & Warmuth, M. K. (1999). Averaging expert predictions. *Computational Learning Theory, 4th European Conference, EuroCOLT '99, Nordkirchen, Germany, March 29-31, 1999, Proceedings* (pp. 153–167). Springer.
- Schölkopf, B., Smola, A. J., & Müller, K.-R. (1998). Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10, 1299–1319.
- Takimoto, E., & Warmuth, M. K. (2003). Path kernels and multiplicative updates. *Journal of Machine Learning Research*, 4, 773–818.
- Tsuda, K., Rätsch, G., & Warmuth, M. K. (2005). Matrix exponentiated gradient updates for on-line learning and Bregman projections. *Journal of Machine Learning Research*, 6, 995–1018.
- Warmuth, M. K., & Kuzmin, D. (2006a). Online variance minimization. *Proceedings of the 19th Annual Conference on Learning Theory (COLT 06)*. Pittsburgh: Springer.
- Warmuth, M. K., & Kuzmin, D. (2006b). Randomized PCA algorithms with regret bounds that are logarithmic in the dimension. *Advances in Neural Information Processing Systems 19 (NIPS 06)*. MIT Press.
- Warmuth, M. K., & Vishwanathan, S. (2005). Leaving the span. *Proceedings of the 18th Annual Conference on Learning Theory (COLT 05)*. Bertinoro, Italy: Springer. Journal version: <http://www.cse.ucsc.edu/~manfred/pubs/span.pdf>.