
Winnowing Subspaces

Manfred K. Warmuth

MANFRED@CSE.UCSB.EDU

Computer Science Department, University of California - Santa Cruz

Abstract

We generalize the Winnow algorithm for learning disjunctions to learning subspaces of low rank. Subspaces are represented by symmetric projection matrices. The online algorithm maintains its uncertainty about the hidden low rank projection matrix as a symmetric positive definite matrix. This matrix is updated using a version of the Matrix Exponentiated Gradient algorithm that is based on matrix exponentials and matrix logarithms. As in the case of the Winnow algorithm, the bounds are logarithmic in the dimension n of the problem, but linear in the rank r of the hidden subspace. We show that the algorithm can be adapted to handle arbitrary matrices of any dimension via a reduction.

1. Introduction

Assume we want to label unit vectors as to whether they are “close to” a r dimensional subspace that is hidden from the learner. We will give a polynomial online algorithm that does this with a small number of prediction mistakes without ever producing an r dimensional subspace that determines the labels. Our algorithm is a matrix variant of the Winnow algorithm for learning disjunctions. In the noise-free case the number of mistakes of our algorithm is $O(r \log \frac{n}{r})$, where n is the dimension of the instances and r the rank of the hidden subspace. Whereas the original algorithm maintains a non-negative weight vector and is designed to learn as well as the best disjunction of size r , the new matrix version maintains a symmetric positive definite matrix as its parameter. This matrix summarizes the uncertainty about the hidden r dimensional subspace.

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

The algorithm and its mistake bounds are presented in an online setting. Standard conversion algorithms lead to bounds for batch learning models when the examples are drawn independently at random according to a fixed distribution (See e.g. (Littlestone, 1989; Cesa-Bianchi & Gentile, 2006)).

The original Winnow algorithm uses a linear threshold function as its hypothesis. In the matrix version the dot product between vectors is replaced by the trace between matrices. In this conference paper we only describe the algorithm in its simplest deterministic form and only prove bounds for the noise-free case. The algorithm and bound retain the case of disjunctions as a special case.

We leave it to the reader to adapt the fancier algorithms and bounds developed for the disjunction case to the matrix case: versions that can handle shifted and scaled versions of the problem as well as noisy data, a randomized version that has an expected mistake bound that is half of the mistake bound of the deterministic algorithms and modifications of the algorithms that learn well when the subspace shifts over time (Auer & Warmuth, 1998). Finally, the algorithms can also be kernelized (as was done in (Kuzmin & Warmuth, 2007) for principal component analysis).

Our first algorithm, Symmetric Matrix Winnow, learns well when the labels indicate closeness to a low rank subspace in \mathbb{R}^n . A subspace is represented as a (symmetric) projection matrix $\mathbf{P} = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^T$, where \mathbf{u}_i are the r orthogonal directions defining the subspace. The algorithm maintains its uncertainty about the hidden projection matrix as an $n \times n$ dimensional symmetric positive definite matrix. It assumes that the instance matrices are symmetric $n \times n$ matrices with eigenvalues in $[0, 1]$. In Figure 3 we visualize the regions produced by thresholding a trace between a projection matrix and an instance matrix of the form $\mathbf{x} \mathbf{x}^T$ where \mathbf{x} is a unit vector (or direction). At the end of the paper we give a second algorithm that can handle the case of arbitrary $m \times n$ dimensional instance matrices (with bounded maximum singular value) that

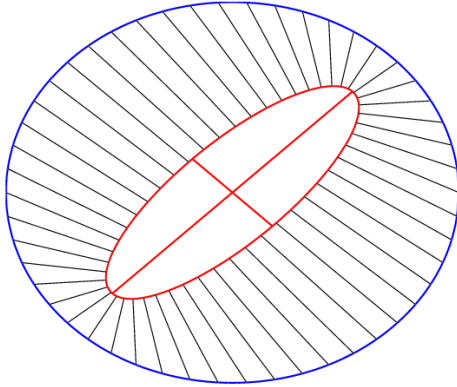


Figure 1: The symmetric matrix \mathbf{A} is depicted as the inner ellipse: The curve of the ellipse is plot of vector $\mathbf{A}\mathbf{u}$, where \mathbf{u} is unit vector. The lines connect the corresponding points on the outer unit circle and the ellipse. The eigenvectors correspond to the axes of the ellipse.

are labeled whether they are close to rank r matrix of dimension $m \times n$ that has r non-zero singular values of value one. The generalization is achieved via a reduction from the symmetric case.

1.1. Related Work

There has been a number of papers where the exponentiated gradient algorithms (Kivinen & Warmuth, 1997) for vector parameters have been generalized to matrix parameters. Instead of maintaining a vector of non-negative numbers, the algorithm now maintains a symmetric positive definite matrix.¹ The linear regression case and a generalization to Boosting appears in (Tsuda et al., 2005) and the expert setting of on-learning has been “lifted” to the matrix case in (Warmuth & Kuzmin, 2006a; Arora & Kale, 2007). Also in (Warmuth & Kuzmin, 2006b; Kuzmin & Warmuth, 2007) the methodology was used to derive on-line algorithms for principal component analysis.

Note that in all of the above cases where symmetric positive definite matrices are used as a parameter, there is always an original problem where the parameter of the algorithm and the instances are vectors. The original problem is retained as a special case when the parameter matrix and instance matrices are restricted to be diagonal matrices. Curiously enough the bounds are always the same for the general symmetric matrix case as for the vector/diagonal matrix case. This

¹Closely related is the case when the parameter vector is a *density matrix* (i.e. a symmetric positive definite matrix whose trace is normalized to one).

phenomenon has been dubbed the “free matrix lunch” (Warmuth, 2007). In particular, our bound for learning the Close to Subspace problem is the same as the bounds for the original disjunction problem.

In this paper we “lift” the Winnow algorithm for learning disjunctions to the symmetric matrix case. This is particularly interesting because it lets us learn well when the data is close to a low rank subspace. Many algorithms have been proposed for learning low-rank matrices. However, with the notable exception of principal component analysis, generalization bounds are often hard to obtain for these algorithms. In this paper we address a very simple problem where the target is a subspace of low rank (i.e. all eigenvalues defining the low-rank object are the same). In this simple case we are able to prove bounds and show how they carry over to general matrices of arbitrary shape. As for principal component analysis, the bounds are linear in the rank of the subspace but logarithmic in the feature dimension.

2. Notation and Some Linear Algebra

A *dyad* is an outer product $\mathbf{u}\mathbf{u}^\top$, where \mathbf{u} is a unit vector (i.e. a direction). Any linear combination of dyads (with real coefficients) is a symmetric matrix and any symmetric matrix \mathbf{A} of dimension n always has an eigendecomposition, i.e. it can be expressed as a linear combination of n orthogonal dyads: $\mathbf{A} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$, where the linear coefficients λ_i are the *real* eigenvalues and the \mathbf{u}_i are the orthonormal eigenvectors. We depict a symmetric matrix as an ellipse which is an affine transformation of the unit ball (See Figure 1). A dyad is a degenerate ellipse that has exactly one eigenvector with eigenvalue one and all other eigenvalues are zero. In other words a dyad can be viewed as an axis of length two that goes through the origin at the half point. The parameter of our algorithm will be a symmetric positive definite matrix, i.e. now the λ_i must be non-negative.

Alternatively, we can view the symmetric positive definite matrix \mathbf{A} as a covariance matrix of some random cost vector $\mathbf{c} \in \mathbb{R}^n$, i.e.

$$\mathbf{A} = \mathbb{E}((\mathbf{c} - \mathbb{E}(\mathbf{c}))(\mathbf{c} - \mathbb{E}(\mathbf{c}))^\top)$$

The variance along any vector \mathbf{u} is then

$$\begin{aligned} \mathbf{V}(\mathbf{c}^\top \mathbf{u}) &= \mathbb{E}((\mathbf{c}^\top \mathbf{u} - \mathbb{E}(\mathbf{c}^\top \mathbf{u}))^2) \\ &= \mathbf{u}^\top \underbrace{\mathbb{E}((\mathbf{c} - \mathbb{E}(\mathbf{c}))(\mathbf{c} - \mathbb{E}(\mathbf{c}))^\top)}_{\mathbf{A}} \mathbf{u}. \end{aligned}$$

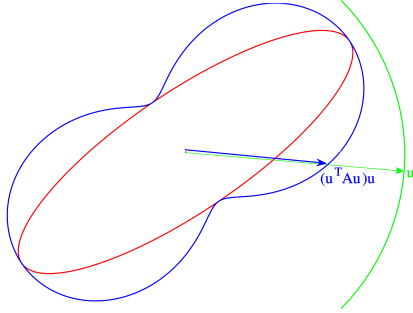


Figure 2: Plot of the variance: The outer figure eight is direction \mathbf{u} times the variance $\mathbf{u}^\top \mathbf{A} \mathbf{u}$. The inner red curve is matrix \mathbf{A} represented as an ellipse. The partial green arc is part of the unit circle.

We often express the variance as a trace²:

$$\mathbf{u}^\top \mathbf{A} \mathbf{u} = \text{tr}(\mathbf{u}^\top \mathbf{A} \mathbf{u}) = \text{tr}(\mathbf{A} \mathbf{u} \mathbf{u}^\top) \geq 0.$$

For an eigenvector, this variance equals the eigenvalue and touches the ellipse (See Figure 2).

3. Staying Close to a Subspace

Consider for example the 2 out of 5 literal monotone disjunction $v_1 \vee v_3$, which is represented by the bit vector $\mathbf{w} = (1, 0, 1, 0, 0)^\top$. The disjunction evaluates to true on the bit vector $\mathbf{x} = (0, 1, 1, 0, 0)^\top$ because $\mathbf{w} \cdot \mathbf{x} \geq \frac{1}{2}$.³ In general, an r literal disjunction is modeled as a bit vector with r ones and the disjunction is true if $\mathbf{d} \cdot \mathbf{x} \geq 1$. Note that if the vectors \mathbf{d} and \mathbf{x} are viewed as diagonal matrices, then the dot product becomes a trace: $\mathbf{d} \cdot \mathbf{x} = \text{tr}(\text{diag}(\mathbf{d}) \text{diag}(\mathbf{x}))$.

When we generalize our setup to matrices then disjunctions naturally correspond to subspaces. A subspace \mathbf{P} is defined by a rank r projection matrix, $\mathbf{P} = \sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^\top$, where the \mathbf{u}_i are r orthogonal directions defining the subspace. An instance is a positive definite symmetric matrix \mathbf{X} with eigenvalues in the range $[0,1]$ and \mathbf{P} labels \mathbf{X} as 1 if it is “close to” \mathbf{P} and 0 otherwise. Here closeness is measured with the trace $\text{tr}(\mathbf{P} \mathbf{X})$ which can be rewritten as a sum of variances over the r directions defining \mathbf{P} :

$$\text{tr}(\mathbf{P} \mathbf{X}) = \text{tr}\left(\sum_{i=1}^r \mathbf{u}_i \mathbf{u}_i^\top \mathbf{X}\right) = \sum_{i=1}^r \underbrace{\mathbf{u}_i^\top \mathbf{X} \mathbf{u}_i}_{\text{variance in dir. } \mathbf{u}_i}.$$

²The trace of a square matrix is the sum of the diagonal. For two such matrices \mathbf{A}, \mathbf{B} , $\text{tr}(\mathbf{A} \mathbf{B}) = \sum_{i,j} A_{i,j} B_{j,i} = \text{tr}(\mathbf{B} \mathbf{A})$. This inner product can also be denoted as $\mathbf{A}^\top \bullet \mathbf{B}$.

³The dot product with a positively labeled bit vector is ≥ 1 and with a negative bit vector is 0. The threshold of $\frac{1}{2}$ is right in the middle of the gap between the positive and negative bit vectors.

Since the eigenvalues of \mathbf{X} lie in $[0,1]$, the variance along each direction is at most one. The “generalized disjunction” evaluates to one if the total variance along all directions is above a threshold. In the simplest case the instance matrices \mathbf{X}_t are dyads $\mathbf{x}_t \mathbf{x}_t^\top$. See Figure 3 for a visualization of the regions produced by the threshold $\text{tr}(\mathbf{P} \mathbf{x} \mathbf{x}^\top) \geq \frac{1}{2}$ when \mathbf{P} is a projection matrix.

Close to Subspace Problem

Input: A sequence of examples (\mathbf{X}_t, y_t) , where the instances \mathbf{X}_t are symmetric matrices with eigenvalues in $[0,1]$, the labels y_t are ± 1 , and there is a projection matrix \mathbf{P} of rank r s.t. $\forall t$

$$\text{tr}(\mathbf{P} \mathbf{X}_t) = \begin{cases} \geq \frac{1}{2} & \text{if } y_t = 1 \\ 0 & \text{if } y_t = -1 \end{cases}$$

Output: A symmetric positive definite matrix \mathbf{W} of rank r s.t. $\forall t$

$$\text{tr}(\mathbf{W} \mathbf{X}_t) = \begin{cases} \geq \frac{1}{2} & \text{if } y_t = 1 \\ 0 & \text{if } y_t = -1 \end{cases}$$

We don’t know whether the Close to Subspace problem is NP-hard. However, this is easy to show if we augment the problem with a fixed set of additional linear constraints that must hold for both the consistent projection matrix \mathbf{P} and the output matrix \mathbf{W} .⁴

Theorem 1. *The Close to Subspace problem with linear side constraints is NP-hard.*

Proof. Given an instance of the Set Cover problem: a family of subsets S^i of a fixed set of n elements and a cover size r . The size r becomes the rank constraint of the corresponding CSP instance. Represent each set S^i as a positive example $(\mathbf{X}_i, +1)$ where \mathbf{X}_i is a binary diagonal instance matrix s.t. whose j diagonal element is 1 if the j -th element lies in S^i (and 0 otherwise). Add $O(n^2)$ additional constraints to the CSP instance that explicitly force all off-diagonal elements to be zero: $\text{tr}(\mathbf{W} \mathbf{X}^{i,j}) = 0$, where $\mathbf{X}^{i,j}$ is the $n \times n$ dimensional 0 matrix with positions (i, j) and (j, i) replaced by 1.

Clearly, any cover of size r can be encoded as a diagonal matrix \mathbf{W} with r ones and $n - r$ zeros along the diagonal. Also if a symmetric positive definite matrix \mathbf{W} has rank r and satisfies the side constraints then it must be a diagonal matrix with r non-zero entries in the diagonal. Wlog the non-zero elements are all 1’s and in this case \mathbf{W} corresponds to a cover. \square

⁴NP-hardness reductions from Maximum Cut are also possible (already for $r = 1$). In this case no side constraints are needed but the Close to Subspace problem needs to be scaled in various ways. Here we only give a simple reduction that uses side constraints.

Algorithm 1 Symmetric Matrix Winnow(η, θ, w_0, n)
 $\eta > 0, \theta, w_0 \in \mathbb{R}, n \in \mathbb{N}$

Initialize $\mathbf{W}_1 = w_0 \mathbf{I}_{n \times n}$

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

Receive instance \mathbf{X}_t with eigenvalues in $[0, 1]$

Predict with

$$\hat{y}_t = \begin{cases} +1 & \text{if } \text{tr}(\mathbf{W}_t \mathbf{X}_t) \geq \theta \\ -1 & \text{otherwise} \end{cases}$$

Receive label y_t

Update

$$\mathbf{W}_{t+1} = \begin{cases} \mathbf{W}_t & \text{if } \hat{y}_t = y_t \\ \exp(\log \mathbf{W}_t + \eta y_t \mathbf{X}_t) & \text{otherwise} \end{cases} \quad (1)$$

end for

4. The Algorithm

Consider the following online protocol: In each trial the algorithm is given an instance matrix \mathbf{X}_t . The algorithm then produces a prediction $\hat{y}_t \in \{0, 1\}$ based on its internal parameters and receives a label $y_t \in \{+1, -1\}$. Finally, the algorithm incurs a mistake if $\hat{y}_t \neq y_t$ and updates its parameters.

The Symmetric Matrix Winnow algorithm never outputs a projection matrix of low rank that predicts well. Instead it maintains a symmetric positive definite matrix \mathbf{W}_t (of full rank) as its parameter and produces a binary prediction \hat{y}_t based on thresholding $\text{tr}(\mathbf{W}_t \mathbf{X}_t)$. The parameter matrix is updated using the matrix log and exponential.⁵ As the original Winnow algorithm, the parameter is only updated when a prediction mistake occurs.⁶ Note that the matrix log and exponential of a symmetric matrix are computed by first computing the eigendecomposition of the argument matrix and then applying the function to all eigenvalues.

5. The Bound

We bound the number of mistakes made by our algorithm with similar methods as were used for the Winnow algorithm (see e.g. (Auer & Warmuth, 1998)).

⁵The normalized version of Symmetric Matrix Winnow scales the trace of the resulting matrix to the desired value.

⁶Such updates are called *conservative* in (Littlestone, 1988). The algorithm can be interpreted as a Matrix Exponentiated Gradient algorithm (Tsuda et al., 2005) wrt a hinge loss. In trials when no mistake occurs, the gradient of this hinge loss is zero, which results in no update (Gentile & Warmuth, 1998).

The original Winnow handles the special case when all matrices are diagonal. As customary for this type of analysis⁷ with symmetric positive definite matrix parameters, we use the quantum relative entropy (see e.g. (Nielsen & Chuang, 2000)) between such matrices \mathbf{V} and \mathbf{W} as our measure of progress:

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

As before, $\log \mathbf{A}$ of the symmetric positive definite matrix \mathbf{A} is defined as a spectral function which applies the logarithm to the eigenvalues of \mathbf{A} but leaves eigen-system unaltered.

The proof follows the methodology developed in (Tsuda et al., 2005; Warmuth & Kuzmin, 2006a) except that here we work with arbitrary symmetric positive definite matrices instead of density matrices which have the additional requirement that their trace is 1. We make use of the Golden Thompson inequality (see e.g. (Bhatia, 1997)), which holds for arbitrary symmetric matrices

$$\text{tr}(\exp(\mathbf{A} + \mathbf{B})) \leq \text{tr}(\exp(\mathbf{A}) \exp(\mathbf{B}))$$

We also need two lemmas from (Tsuda et al., 2005):

Lemma 1. *For any symmetric matrix \mathbf{X} , such that $\mathbf{0} \preceq \mathbf{X} \preceq \mathbf{I}$ and any constant $a \in \mathbb{R}$, the following holds:*

$$\exp(a\mathbf{X}) \preceq \mathbf{I} - (1 - \exp(a))\mathbf{X}.$$

Lemma 2. *For any positive semidefinite matrix \mathbf{A} and any symmetric matrices \mathbf{B} and \mathbf{C} , $\mathbf{B} \preceq \mathbf{C}$ implies $\text{tr}(\mathbf{A}\mathbf{B}) \leq \text{tr}(\mathbf{A}\mathbf{C})$.*

We are now ready to prove our main result. For the sake of simplicity this theorem only handles the noise-free case.

Theorem 2. *Given any sequence of examples (\mathbf{X}_t, y_t) such that the instances \mathbf{X}_t are symmetric positive definite matrices with eigenvalues in $[0, 1]$, the labels y_t are ± 1 , and there is a projection matrix $\mathbf{s.t.}$*

$$\text{tr}(\mathbf{P}\mathbf{X}_t) = \begin{cases} \geq \frac{1}{2} & \text{if } y_t = +1 \\ 0 & \text{if } y_t = -1 \end{cases} \quad (2)$$

Then the online algorithm Symmetric Matrix Winnow makes at most

$$7.18 r \ln \frac{n}{r}$$

mistakes on this sequence, when $\eta \approx 1.28$, $\theta \approx .19$ and $w_0 = r$.

Proof. We lower bound the per trial progress

$$\Delta(\mathbf{P}, \mathbf{W}_t) - \Delta(\mathbf{P}, \mathbf{W}_{t+1})$$

⁷Alternatively, $\text{tr}(\mathbf{W}_t)$ can be used as a potential.

towards a comparator matrix P , which is any projection matrix P satisfying the constraints (2). If no mistake occurs in trial t , then $\mathbf{W}_{t+1} = \mathbf{W}_t$ and the progress is zero. Otherwise, \mathbf{W}_{t+1} is updated to $\exp(\log \mathbf{W}_t + \eta y_t \mathbf{X}_t)$ and

$$\begin{aligned} & \Delta(\mathbf{P}, \mathbf{W}_t) - \Delta(\mathbf{P}, \mathbf{W}_{t+1}) \\ &= \text{tr}(\mathbf{P}(\log \mathbf{W}_{t+1} - \log \mathbf{W}_t) + \mathbf{W}_t - \mathbf{W}_{t+1}) \\ &\stackrel{(1)}{=} \eta y_t \text{tr}(\mathbf{P} \mathbf{X}_t) + \text{tr}(\mathbf{W}_t - \exp(\log \mathbf{W}_t + \eta y_t \mathbf{X}_t)) \\ &\stackrel{\text{G.Th.}}{\geq} \eta y_t \text{tr}(\mathbf{P} \mathbf{X}_t) + \text{tr}(\mathbf{W}_t - \exp(\log \mathbf{W}_t) \exp(\eta y_t \mathbf{X}_t)) \\ &= \eta y_t \text{tr}(\mathbf{P} \mathbf{X}_t) + \text{tr}(\mathbf{W}_t(\mathbf{I} - \exp(\eta y_t \mathbf{X}_t))) \\ &\stackrel{\text{Le.1\&2}}{\geq} \eta y_t \text{tr}(\mathbf{P} \mathbf{X}_t) + (1 - \exp(\eta y_t)) \text{tr}(\mathbf{W}_t \mathbf{X}_t) \end{aligned}$$

If $y_t = -1$ and $\hat{y}_t = +1$, then $\text{tr}(\mathbf{P} \mathbf{X}_t) = 0$, $1 - \exp(\eta y_t) > 0$ and $\text{tr}(\mathbf{W}_t \mathbf{X}_t) \geq \theta$. Thus the progress is lower bounded by

$$(1 - \exp(-\eta))\theta.$$

Similarly, if $y_t = +1$ and $\hat{y}_t = -1$, then $\text{tr}(\mathbf{P} \mathbf{X}_t) \geq \frac{1}{2}$, $1 - \exp(\eta y_t) < 0$ and $\text{tr}(\mathbf{W}_t \mathbf{X}_t) \leq \theta$. Now the progress is lower bounded by

$$\frac{\eta}{2} + (1 - \exp(\eta))\theta.$$

By setting both cases equal we obtain

$$\theta = \frac{\eta}{2(\exp(\eta) - \exp(-\eta))}.$$

With this choice of θ , the progress per mistake is $\frac{\eta}{2+2\exp(\eta)}$ and this is optimized at $\eta \approx 1.28$. This choice of η makes $\theta \approx .19$ and the progress $\approx .14$. By summing over trials we get

$$\underbrace{\Delta(\mathbf{P}, \mathbf{W}^1)}_{r \ln \frac{n}{r}} - \underbrace{\Delta(\mathbf{P}, \mathbf{W}^{T+1})}_{\geq 0} \geq .14M,$$

where M is the number of mistakes. This implies that $M \leq 7.18 r \log_2 \frac{n}{r}$. \square

TIGHTNESS OF THE BOUND: The Vapnik-Chervonenkis dimension of r out of n monotone disjunctions is at least $r \lceil \log_2 \frac{n}{r} \rceil$ (Littlestone, 1988). In other words, there are this many bit vectors \mathbf{x}_t of dimension n with the property that for each possible labeling there is a consistent monotone disjunction of r of the n variables. By feeding the instances $\text{diag}(\mathbf{x}_i)$ to any algorithm that solves the Close to Subspace problem and contradicting the prediction of the algorithm each time, we can force any such algorithm to make $r \lceil \log_2 \frac{n}{r} \rceil$ mistakes. The sequence

of examples will have a consistent disjunction of size r which is represented as a bit vector \mathbf{d} with r ones. This means that the r dimensional projection matrix $\text{diag}(\mathbf{d})$ with the identity eigensystem will label the matrix instances consistently. Thus the upper bound proven above for Symmetric Matrix Winnow cannot be improved by more than a constant.

TIME COMPLEXITY: Given that you have the eigendecomposition of \mathbf{W}_t , then the expensive part in computing \mathbf{W}_{t+1} is producing the eigendecomposition of the exponent $\log \mathbf{W}_t + \eta y_t \mathbf{X}_t$ of update (1). This costs $O(n^3)$ time per update. The algorithm of the next section is based on the singular value decomposition, which has similar complexity. In some applications one might get away with not updating at every trial. However, at this point the methods presented here are only applicable to small matrices. Nevertheless, they show what kind of bounds are possible if computation time is cheap.

NOISY CASE: In this paper we kept things simple and only focused on the noise-free case. The bounds for learning noisy linear threshold functions with Winnow when the loss function is the hinge loss (see e.g. (Gentile & Warmuth, 1998)) immediately carry over to the matrix case. For example this leads to bounds in terms of the hinge loss of the best subspace.

LINEAR SIDE CONSTRAINTS: Now the parameter matrix is projected onto the linear side constraints after the main update (1). Using the Bregman projection methods of (Herbster & Warmuth, 2001) one can show that the bounds stay the same provided that the target satisfies the constraints.⁸

GENERALIZATION BOUNDS: Any mistake bounded on-line algorithm can easily be converted to a generalization bound for the case when the instances are generated independently at random by a fixed but unknown distribution (Littlestone, 1989; Cesa-Bianchi & Gentile, 2006). One of the simplest conversion algorithm for a conservative on-line algorithm might be the following (Floyd & Warmuth, 1995). Draw a batch of m examples and in some default order on the instances do one pass over the examples with the mistake bounded on-line algorithm. If the on-line algorithm is guaranteed to make at most M mistakes then the subset of at most M examples where mistakes were made represents a consistent hypothesis. In other words a mistake bound of size M leads to a “compression scheme” of size M . This immediately implies that for any $\epsilon, \delta \in [0, 1]$, $O(\frac{1}{\epsilon}(M \log \frac{1}{\epsilon} + \log \frac{1}{\delta}))$ examples suffice

⁸Actually projections onto an arbitrary convex region can be used.

to guarantee the following about the produced consistent hypothesis: its probability of predicting wrongly is at most ϵ with probability at least $1 - \delta$.

6. Non-square Instance Matrices

Now the instances \mathbf{X} are non-square matrices of dimension $m \times n$, where $m \geq n$ for the purpose of discussion. We use the standard trick (see e.g. (Golub & Loan, 1996)) to embed⁹ such matrices into a symmetric matrix of dimension $(m + n) \times (m + n)$:

$$\begin{pmatrix} \mathbf{0} & \mathbf{X}^T \\ \mathbf{X} & \mathbf{0} \end{pmatrix}_{\substack{n \times n & n \times m \\ m \times n & m \times m}}$$

Note that $\mathbf{0}$ always denotes a square zero matrix of the appropriate dimension. From now on we also use bold greek letters to denote square diagonal matrices.

If \mathbf{X} has the singular value decomposition $\mathbf{U} \begin{pmatrix} \boldsymbol{\sigma} \\ \mathbf{0} \end{pmatrix} \mathbf{V}^T = \mathbf{U}_1 \boldsymbol{\sigma} \mathbf{V}^T$, where \mathbf{U}_1 consists of the first n columns of \mathbf{U} , then the corresponding symmetric matrix has the following eigendecomposition (Golub & Loan, 1996):

$$\begin{pmatrix} \mathbf{0} & \mathbf{X}^T \\ \mathbf{X} & \mathbf{0} \end{pmatrix}_{(m+n) \times (m+n)} = \mathbf{Q} \begin{pmatrix} \boldsymbol{\sigma} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{Q}^T = \widehat{\mathbf{Q}} \begin{pmatrix} \boldsymbol{\sigma} & \mathbf{0} \\ \mathbf{0} & -\boldsymbol{\sigma} \end{pmatrix} \widehat{\mathbf{Q}}^T_{(2n) \times (2n)}$$

where $\mathbf{Q} = \frac{1}{\sqrt{2}} \begin{pmatrix} \mathbf{v} & \mathbf{v} & \mathbf{0} \\ \mathbf{U}_1 & -\mathbf{U}_1 & \sqrt{2}\mathbf{U}_1 \end{pmatrix}_{(m+n) \times (m+n)}$
 and $\widehat{\mathbf{Q}} = \frac{1}{\sqrt{2}} \begin{pmatrix} \mathbf{v} & \mathbf{v} \\ \mathbf{U}_1 & -\mathbf{U}_1 \end{pmatrix}_{(m+n) \times (2n)}$.

We can now run Symmetric Matrix Winnow with the embedded instances $\begin{pmatrix} \mathbf{0} & \mathbf{X}_t^T \\ \mathbf{X}_t & \mathbf{0} \end{pmatrix}$ and the initial parameter matrix $\mathbf{W}_1 = w_0 \begin{pmatrix} \mathbf{0} & \mathbf{I}^T \\ \mathbf{I} & \mathbf{0} \end{pmatrix}$, where $w_0 = r$ and \mathbf{I} denotes the $m \times n$ dimensional identity matrix. Our update produces a symmetric positive definite matrix \mathbf{W}_t of the form

$$\frac{1}{2} \exp \begin{pmatrix} \mathbf{0} & \mathbf{R}_t^T \\ \mathbf{R}_t & \mathbf{0} \end{pmatrix} = \frac{1}{2} \begin{pmatrix} \tilde{\mathbf{v}} & \tilde{\mathbf{v}} \\ \tilde{\mathbf{U}}_1 & -\tilde{\mathbf{U}}_1 \end{pmatrix} \begin{pmatrix} \exp(\tilde{\boldsymbol{\sigma}}) & \mathbf{0} \\ \mathbf{0} & \exp(-\tilde{\boldsymbol{\sigma}}) \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{v}} & \tilde{\mathbf{v}} \\ \tilde{\mathbf{U}}_1 & -\tilde{\mathbf{U}}_1 \end{pmatrix}^T$$

where \mathbf{R}_t is a linear combination of the instances $\mathbf{X}_1 \dots \mathbf{X}_{t-1}$ that has the singular value decomposition $\tilde{\mathbf{U}}_1 \tilde{\boldsymbol{\sigma}} \tilde{\mathbf{V}}^T$. We can rewrite \mathbf{W}_t as

$$\frac{1}{2} \begin{pmatrix} \tilde{\mathbf{v}} \cosh(\tilde{\boldsymbol{\sigma}}) \tilde{\mathbf{V}}^T & \tilde{\mathbf{v}} \sinh(\tilde{\boldsymbol{\sigma}}) \tilde{\mathbf{U}}_1^T \\ \tilde{\mathbf{U}}_1 \sinh(\tilde{\boldsymbol{\sigma}}) \tilde{\mathbf{V}}^T & \tilde{\mathbf{U}}_1 \cosh(\tilde{\boldsymbol{\sigma}}) \tilde{\mathbf{U}}_1^T \end{pmatrix}^T$$

⁹The reduction is different from the EG^\pm reduction of (Kivinen & Warmuth, 1997) which would embed \mathbf{X} as $\begin{pmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & -\mathbf{X} \end{pmatrix}$, which is unfortunately not symmetric.

Algorithm 2 General Matrix Winnow(η, θ, w_0, m, n)
 $\eta > 0, \theta, w_0 \in \mathbb{R}, m, n \in \mathbb{N}$

Initialize $\mathbf{R}_1 = w_0 \mathbf{I}_{m \times n}$

for $t = 1$ to T do

 Receive instance \mathbf{X}_t with singular values in $[0, 1]_{m \times n}$

 Predict with

$$\hat{y}_t = \begin{cases} +1 & \text{if } \text{tr}(\mathbf{R}_t \mathbf{X}_t) \geq \theta \\ -1 & \text{otherwise} \end{cases}$$

 Receive label y_t

 Update

$$\mathbf{R}_{t+1} = \begin{cases} \mathbf{R}_t & \text{if } \hat{y}_t = y_t \\ \sinh(\text{arcsinh}(\mathbf{R}_t) + \eta \mathbf{X}_t) & \text{otherwise} \end{cases}$$

end for

where \sinh and \cosh are the diagonalized versions of the scalar sinh and cosh functions.¹⁰ This lets us write the trace which is used for forming the prediction as follows:

$$\text{tr} \left(\mathbf{W}_t \begin{pmatrix} \mathbf{0} & \mathbf{X}_t^T \\ \mathbf{X}_t & \mathbf{0} \end{pmatrix} \right) = \text{tr}(\tilde{\mathbf{U}}_1 \sinh(\tilde{\boldsymbol{\sigma}}) \tilde{\mathbf{V}}^T \mathbf{X}_t^T)$$

The General Matrix Winnow¹¹ algorithm 2 reformulates the entire algorithm sketched above using the original $m \times n$ dimensional instances \mathbf{X}_t and an $m \times n$ dimensional parameter matrix \mathbb{R}_t . For any $m \times n$ dimensional matrix \mathbf{A} with singular value decomposition $\mathbf{U} \boldsymbol{\sigma} \mathbf{V}^T$, the matrix function $\sinh(\mathbf{A})$ applies the scalar sinh function to the diagonal elements of $\tilde{\boldsymbol{\sigma}}$ and $\text{arcsinh}(\mathbf{A})$ is defined similarly based on the inverse of the scalar sinh function. The bounds for this algorithm follow from the symmetric case via the above embedding. The symmetric projection matrices of rank r are now *generalized projection matrices* of the form $\mathbf{P} = \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^T$, where $s_i \in \{+1, -1\}$ and the \mathbf{u}_i are a set of r orthogonal m -dimensional unit vector and the \mathbf{v}_i are a set of r orthogonal n -dimensional unit vectors.

7. Conclusion

It is unknown how to handle rank constraints. When linear side constraints are added, the Close to Subspace problem is hard even if we allow the output matrix \mathbf{W} to be any positive symmetric matrix of low rank. We follow the route taken by the original Winnow algorithm and bypass the hardness by never pro-

¹⁰ $\sinh(x) := \frac{\exp(x) - \exp(-x)}{2}$, $\cosh(x) := \frac{\exp(x) + \exp(-x)}{2}$.

¹¹The normalized version scales with the trace norm.

ducing a hypotheses of a certain form (here low rank). The algorithm maintains a positive definite parameter matrix of *full rank* instead. Nevertheless it is guaranteed to make few prediction mistakes if the instances are labeled based on a low rank projection matrix. The update of the algorithm is a special case of the Matrix Exponentiated Gradient algorithm and is motivated by regularizing the hinge loss on the linear trace with a quantum relative entropy. Note that the original Winnow is motivated along the same lines (see e.g. (Gentile & Warmuth, 1998)).

More generally the following picture is emerging. In the vector case the algorithms related to the original Winnow algorithm learn with a thresholded dot product between the weight vector and instance vector. In the noise-free case there has to be a gap between the positive and negative examples, and a threshold can be placed in this gap for the purpose of classifying the examples. For the original Winnow algorithm, the weight vectors are measured by the one-norm and this is also the norm by which the margin of the examples is normalized. For the instance vectors the range of the components (essentially the infinity norm) is crucial.

In the matrix case the dot product becomes a trace and the parameter matrices are measured by their trace norm (sum of singular values). This norm also naturally normalizes the margin of the examples. The instance matrices are naturally measured by the matrix 2-norm (maximum singular value).

Acknowledgment

Thanks to Dima Kuzmin for many brainstorming sessions and for preparing all the plots of this paper. Also thanks to Amin Coja-Oghlan and Elad Hazan for valuable discussions.

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.
- Auer, P., & Warmuth, M. K. (1998). Tracking the best disjunction. *Machine Learning*, *32*, 127–150. Earlier version in 36th FOCS, 1995.
- Bhatia, R. (1997). *Matrix analysis*. Berlin: Springer.
- Cesa-Bianchi, N., & Gentile, C. (2006). Improved risk tail bounds for on-line learning. *Advances in Neural Information Processing Systems 18 (NIPS 05)*. MIT Press.
- Floyd, S., & Warmuth, M. (1995). Sample compression, learnability, and the Vapnik-Chervonenkis dimension. *Machine Learning*, *21*, 269–304.
- Gentile, C., & Warmuth, M. K. (1998). Hinge loss and average margin. *In Advances in Neural Information Processing Systems 11 (NIPS*98)*.
- Golub, G. H., & Loan, C. F. V. (1996). *Matrix computations (third edition)*. The John Hopkins University Press.
- 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.
- Kuzmin, D., & Warmuth, M. K. (2007). Online Kernel PCA with entropic matrix updates. *ICML '07: Proceedings of the 24rd international conference on Machine learning*. Corvallis, Oregon: ACM Press.
- Littlestone, N. (1988). Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, *2*, 285–318.
- Littlestone, N. (1989). From on-line to batch learning. *Proceedings of the Second Annual Workshop on Computational Learning Theory* (pp. 269–284). San Mateo, CA: Morgan Kaufmann.
- Nielsen, M., & Chuang, I. (2000). *Quantum computation and quantum information*. Cambridge University Press.
- 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. (2007). When is there a free matrix lunch. *Proc. of the 20th Annual Conference on Learning Theory (COLT 07)*. San Diego, CA: Springer. Open problem, to appear.
- 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.

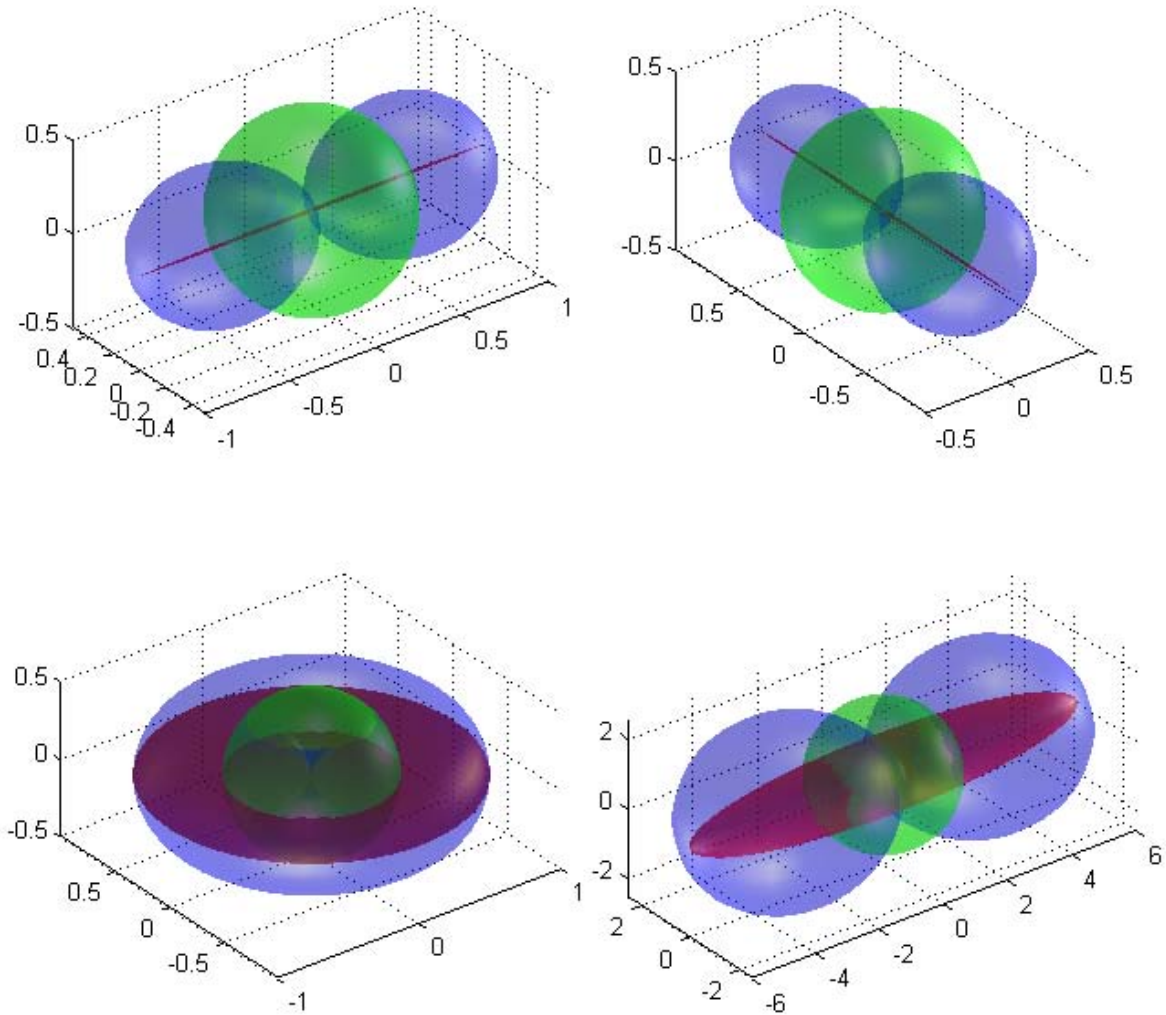


Figure 3: Depiction of a threshold of trace $\text{tr}(\mathbf{P}\mathbf{x}\mathbf{x}^\top)$: **Top left:** The projection matrix \mathbf{P} has dimension one. It is a single (red) dyad along the x -axis, i.e. $\mathbf{P} = \mathbf{e}_1\mathbf{e}_1^\top$, where $\mathbf{e}_1 = (1, 0, 0)^\top$. The blue double cloud is a plot of the variance $\text{tr}(\mathbf{P}\mathbf{x}\mathbf{x}^\top)$ for all unit directions \mathbf{x} . Recall that in dimension 2 the variance has figure eight shape (see blue curve in Figure 2). The green ball has radius $\frac{1}{2}$ and the threshold $\text{tr}(\mathbf{P}\mathbf{x}\mathbf{x}^\top) \geq \frac{1}{2}$ is satisfied if direction $\mathbf{x}\mathbf{x}^\top$ is close to the red axis, i.e. in the blue minus the green region. It is instructive to see what happens if the threshold (radius of the green ball) is raised to 1. Then the region labeled one (blue - green) just consists of the two tips of the red dyad $\mathbf{e}_1\mathbf{e}_1^\top$. That is, $\text{tr}(\mathbf{P}\mathbf{x}\mathbf{x}^\top) \geq 1$ iff $\mathbf{x} = \pm\mathbf{e}_1$. **Top right:** The same picture but now \mathbf{P} is the dyad $\mathbf{e}_2\mathbf{e}_2^\top$ corresponding to the y -axis. **Bottom left:** Now \mathbf{P} is the sum of the previous two dyads, i.e. $\mathbf{P} = \mathbf{e}_1\mathbf{e}_1^\top + \mathbf{e}_2\mathbf{e}_2^\top$ depicted as the red degenerate ellipse (disk) in the xy plane. The variance is a blue donut around this ellipse. The green ball again has radius $\frac{1}{2}$ and the threshold $\text{tr}(\mathbf{P}\mathbf{x}\mathbf{x}^\top) \geq \frac{1}{2}$ labels all directions $\mathbf{x}\mathbf{x}^\top$ as one that are close to the red ellipse, i.e. lie in the blue donut minus the green ball. If the radius of the green ball is raised to 1, then the blue minus the red region shrinks to the circle where the donut touches the disk. **Bottom right:** Now the inner red 3-dimensional ellipse is sum of three directions, i.e. $\mathbf{P} = 6\mathbf{e}_1\mathbf{e}_1^\top + \mathbf{e}_2\mathbf{e}_2^\top + \mathbf{e}_3\mathbf{e}_3^\top$. The variance is the blue handle and the threshold (green ball) is 2.