

Online Variance Minimization^{*}

Manfred K. Warmuth and Dima Kuzmin

Computer Science Department
University of California, Santa Cruz
{manfred, dima}@cse.ucsc.edu

Abstract. We design algorithms for two online variance minimization problems. Specifically, in every trial t our algorithms get a covariance matrix \mathcal{C}_t and try to select a parameter vector \mathbf{w}_t such that the total variance over a sequence of trials $\sum_t \mathbf{w}_t^\top \mathcal{C}_t \mathbf{w}_t$ is not much larger than the total variance of the best parameter vector \mathbf{u} chosen in hindsight. Two parameter spaces are considered - the probability simplex and the unit sphere. The first space is associated with the problem of minimizing risk in stock portfolios and the second space leads to an online calculation of the eigenvector with minimum eigenvalue. For the first parameter space we apply the Exponentiated Gradient algorithm which is motivated with a relative entropy. In the second case the algorithm maintains a mixture of unit vectors which is represented as a density matrix. The motivating divergence for density matrices is the quantum version of the relative entropy and the resulting algorithm is a special case of the Matrix Exponentiated Gradient algorithm. In each case we prove bounds on the additional total variance incurred by the online algorithm over the best offline parameter.

1 Introduction

In one of the simplest settings of learning with expert advice [FS97], the learner has to commit to a probability vector \mathbf{w} over the experts at the beginning of each trial. It then receives a loss vector \mathbf{l} and incurs loss $\mathbf{w} \cdot \mathbf{l} = \sum_i w_i l_i$. The goal is to design online algorithms whose total loss over a sequence of trials is close to loss of the best expert in all trials, i.e. the total loss of the online algorithm $\sum_t \mathbf{w}_t \cdot \mathbf{l}_t$ should be close to the total loss of the best expert chosen in hindsight, which is $\inf_i \sum_t l_{t,i}$, where t is the trial index.

In this paper we investigate online algorithms for minimizing the total variance over a sequence of trials. Instead of receiving a loss vector \mathbf{l} in each trial, we now receive a covariance matrix \mathcal{C} of a random loss vector \mathbf{l} , where $\mathcal{C}(i, j)$ is the covariance between l_i and l_j at the current trial. Intuitively the loss vector provides first-order information (means), whereas covariance matrices give second order information. The variance/risk of the loss for probability vector \mathbf{w} when the covariance matrix is \mathcal{C} can be expressed as $\mathbf{w}^\top \mathcal{C} \mathbf{w} = \mathbf{Var}(\mathbf{w} \cdot \mathbf{l})$. Our goal

^{*} Supported by NSF grant CCR 9821087. Some of this work was done while visiting National ICT Australia in Canberra.

is to minimize the total variance over a sequence of trials: $\sum_t \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t$. More precisely, we want algorithms whose total variance is close to the total variance of the best probability vector \mathbf{u} chosen in hindsight, i.e. the total variance of the algorithm should be close to $\inf_{\mathbf{u}} \mathbf{u}^\top (\sum_t \mathbf{C}_t) \mathbf{u}$ (where the minimization is over the probability simplex).

In a more general setting one actually might want to optimize trade-offs between first-order and second order terms: $\mathbf{w} \cdot \mathbf{l} + \gamma \mathbf{w}^\top \mathbf{C} \mathbf{w}$, where $\gamma \geq 0$ is a risk-aversion parameter. Such problems arise in Markowitz portfolio optimization (See e.g. discussion in [BV04], Section 4.4). For the sake of simplicity, in this paper we focus on minimizing the variance by itself.

We develop an algorithm for the above online variance minimization problem. The parameter space is the probability simplex. We use the Exponentiated Gradient algorithm for solving this problem since it maintains a probability vector. The latter algorithm is motivated and analyzed using the relative entropy between probability vectors [KW97]. The bounds we obtain are similar to the bounds of the Exponentiated Gradient algorithm when applied to linear regression.

In the second part of the paper we focus on the same online variance minimization problem, but now the parameter space that we compare against is the unit sphere of direction vectors instead of the probability simplex and the total loss of the algorithm is to be close to $\inf_{\mathbf{u}} \mathbf{u}^\top (\sum_t \mathbf{C}_t) \mathbf{u}$, where the minimization is over unit vectors. The solution of the offline problem is an eigenvector that corresponds to a minimum eigenvalue of the total covariance $\sum_t \mathbf{C}_t$.

Note that the variance $\mathbf{u}^\top \mathbf{C} \mathbf{u}$ can be rewritten using the trace operator: $\mathbf{u}^\top \mathbf{C} \mathbf{u} = \text{tr}(\mathbf{u}^\top \mathbf{C} \mathbf{u}) = \text{tr}(\mathbf{u} \mathbf{u}^\top \mathbf{C})$. The outer product $\mathbf{u} \mathbf{u}^\top$ for unit \mathbf{u} is called a *dyad* and the offline problem can be reformulated as minimizing trace of a product of a dyad with the total covariance matrix: $\inf_{\mathbf{u}} \text{tr}(\mathbf{u} \mathbf{u}^\top (\sum_t \mathbf{C}_t))$ (where \mathbf{u} is unit length).¹

In the original experts setting, the offline problem involved a minimum over experts. Now its a minimum over dyads and the best dyad corresponds to an eigenvector with minimum eigenvalue. The algorithm for the original expert setting maintains its uncertainty over which expert is best as a probability vector \mathbf{w} , i.e. w_i is the current belief that expert i is best. This algorithm is the Continuous Weighted Majority (WMC) [LW94] (which was reformulated as the Hedge algorithm in [FS97]). It uses exponentially decaying weights $w_{t,i} = \frac{e^{-\eta \sum_{q=1}^{t-1} l_{q,i}}}{Z_t}$, where Z_t is a normalization factor.

In the generalized setting we need to maintain uncertainty over dyads. The natural parameter space is therefore mixtures of dyads which are called density matrices in statistical physics (symmetric positive definite matrices of trace one). Note that the vector of eigenvalues of such matrices is a probability vector. Using the methodology of [TRW05, War05] we develop a matrix version of the Weighted Majority algorithm for solving our second online variance minimization problem.

¹ In this paper we upper bound the total variance of our algorithm, whereas the generalized Bayes rule of [War05, WK06] is an algorithm for which the sum of the negative logs of the variances is upper bounded.

Now the density matrix parameter has the form $\mathbf{W}_t = \frac{\exp(-\eta \sum_{q=1}^{t-1} \mathbf{C}_q)}{Z_t}$, where \exp is the matrix exponential and Z_t normalizes the trace of the parameter matrix to one. When the covariance matrices \mathbf{C}_q are the diagonal matrices $\text{diag}(\mathbf{l}_q)$ then the matrix update becomes the original expert update. In other words the original update may be seen as a special case of the new matrix update when the eigenvectors are fixed to the standard basis vectors and are not updated.

The original weighted majority type update may be seen as a softmin calculation, because as $\eta \rightarrow \infty$, the parameter vector \mathbf{w}_t puts all of its weight on $\arg \min_i \sum_{q=1}^{t-1} l_{q,i}$. Similarly, the generalized update is a soft eigenvector calculation for the eigenvectors with the minimum eigenvalue.

What replaces the loss $\mathbf{w} \cdot \mathbf{l}$ of the algorithm in the more general context? The dot product for matrices is a trace and we use the generalized loss $\text{tr}(\mathbf{W}\mathbf{C})$. If the eigendecomposition of the parameter matrix \mathbf{W} consists of the eigenvectors \mathbf{w}_i and associated eigenvalues ω_i then this loss can be rewritten as

$$\text{tr}(\mathbf{W}\mathbf{C}) = \text{tr}\left(\left(\sum \omega_i \mathbf{w}_i \mathbf{w}_i^\top\right) \mathbf{C}\right) = \sum_i \omega_i \mathbf{w}_i^\top \mathbf{C} \mathbf{w}_i$$

In other words it may be seen as an expected variance along the eigenvectors \mathbf{w}_i that is weighted by the eigenvalues ω_i . Curiously enough, this trace is also a quantum measurement, where \mathbf{W} represents a mixture state of a particle and \mathbf{C} the instrument (See [War05, WK06] for additional discussion). Again the dot product $\mathbf{w} \cdot \mathbf{l}$ is the special case when the eigenvectors are the standard basis vectors, i.e.

$$\text{tr}(\text{diag}(\mathbf{w}) \text{diag}(\mathbf{l})) = \text{tr}\left(\left(\sum w_i \mathbf{e}_i \mathbf{e}_i^\top\right) \text{diag}(\mathbf{l})\right) = \sum_i w_i \mathbf{e}_i^\top \text{diag}(\mathbf{l}) \mathbf{e}_i = \sum_i w_i l_i.$$

The new update is motivated and analyzed using the quantum relative entropy (due to Umegaki, see e.g. [NC00]) instead of the standard relative entropy (also called Kullback-Leibler divergence). The analysis is a fancier version of the original online loss bound for WMC that uses the Golden-Thompson inequality and some lemmas developed in [TRW05].

2 Variance Minimization over the Probability Simplex

2.1 Definitions

In this paper we only consider symmetric matrices. Such matrices always have an eigendecomposition of the form $\mathbf{W} = \mathbf{W}\boldsymbol{\omega}\mathbf{W}^\top$, where \mathbf{W} is an orthogonal matrix of eigenvectors and $\boldsymbol{\omega}$ is a diagonal matrix of the corresponding eigenvalues. Alternatively, the decomposition can be written as $\mathbf{W} = \sum_i \omega_i \mathbf{w}_i \mathbf{w}_i^\top$, with the ω_i being the eigenvalues and the \mathbf{w}_i the eigenvectors. Note that the dyads $\mathbf{w}_i \mathbf{w}_i^\top$ are square matrices of rank one.

Matrix \mathcal{M} is called *positive semidefinite* if for all vectors \mathbf{w} we have $\mathbf{w}^\top \mathcal{M} \mathbf{w} \geq 0$. This is also written as a generalized inequality $\mathcal{M} \succeq \mathbf{0}$. In eigenvalue terms this

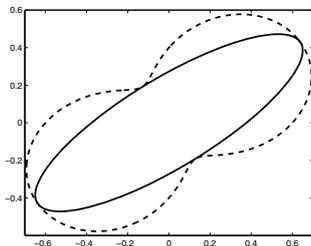


Fig. 1. An ellipse \mathbf{C} in \mathbb{R}^2 : The eigenvectors are the directions of the axes and the eigenvalues their lengths from the origin. Ellipses are weighted combinations of the one-dimensional degenerate ellipses (dyads) corresponding to the axes. (For unit \mathbf{w} , the dyad $\mathbf{w}\mathbf{w}^\top$ is a degenerate one-dimensional ellipse which is a line between $-\mathbf{w}$ and \mathbf{w}). The solid curve of the ellipse is a plot of direction vector $\mathbf{C}\mathbf{w}$ and the outer dashed figure eight is direction \mathbf{w} times the variance $\mathbf{w}^\top \mathbf{C}\mathbf{w}$. At the eigenvectors, this variance equals the eigenvalues and the figure eight touches the ellipse.

means that all eigenvalues of matrix are ≥ 0 . A matrix is *strictly positive definite* if all eigenvalues are > 0 . In what follows we will drop the semi- prefix and call any matrix $\mathbf{M} \succeq \mathbf{0}$ simply positive definite.

Let \mathbf{l} be a random vector, then $\mathbf{C} = \mathbf{E}((\mathbf{l} - \mathbf{E}(\mathbf{l}))(\mathbf{l} - \mathbf{E}(\mathbf{l}))^\top)$ is its *covariance matrix*. It is symmetric and positive definite. For any other vector \mathbf{w} we can compute the variance of the dot product $\mathbf{l}^\top \mathbf{w}$ as follows:

$$\begin{aligned} \text{Var}(\mathbf{l}^\top \mathbf{w}) &= \mathbf{E}((\mathbf{l}^\top \mathbf{w} - \mathbf{E}(\mathbf{l}^\top \mathbf{w}))^2) \\ &= \mathbf{E}(((\mathbf{l}^\top - \mathbf{E}(\mathbf{l}^\top))\mathbf{w})^\top ((\mathbf{l}^\top - \mathbf{E}(\mathbf{l}^\top))\mathbf{w})) \\ &= \mathbf{E}(\mathbf{w}^\top (\mathbf{l} - \mathbf{E}(\mathbf{l}))(\mathbf{l} - \mathbf{E}(\mathbf{l}))^\top \mathbf{w}) \\ &= \mathbf{w}^\top \mathbf{C}\mathbf{w}. \end{aligned}$$

A covariance matrix can be depicted as an ellipse $\{\mathbf{C}\mathbf{w} : \|\mathbf{w}\|_2 = 1\}$ centered at the origin. The eigenvectors of \mathbf{C} form the axes of the ellipse and eigenvalues are the lengths of the axes from the origin (See Figure 1 taken from [War05]).

For two probability vectors \mathbf{u} and \mathbf{w} (e.g. vectors whose entries are nonnegative and sum to one) their relative entropy (or Kullback-Leibler divergence) is given by:

$$d(\mathbf{u}, \mathbf{w}) = \sum_{i=1}^n u_i \log \frac{u_i}{w_i}.$$

We call this a divergence (and not a distance) since its not symmetric and does not satisfy the triangle inequality. It is however nonnegative and convex in both arguments.

2.2 Risk Minimization

The problem of minimizing the variance when the direction \mathbf{w} lies in the probability simplex is connected to risk minimization in stock portfolios. In Markowitz

portfolio theory, vector \mathbf{p} denotes the relative price change of all assets in a given trading period. Let \mathbf{w} be a probability vector that specifies the proportion of our capital invested into each asset (assuming short positions are not allowed). Then the relative capital change after a trading period is the dot product $\mathbf{w} \cdot \mathbf{p}$. If \mathbf{p} is a random vector with known or estimated covariance matrix \mathbf{C} , then the variance of the capital change for our portfolio is $\mathbf{w}^\top \mathbf{C} \mathbf{w}$. This variance is clearly associated with the risk of our investment. Our problem is then to “track” the performance of minimum risk portfolio over a sequence of trading periods.

2.3 Algorithm and Motivation

Let us reiterate the setup and the goal for our algorithm. On every trial t it must produce a probability vector \mathbf{w}_t . It then gets a covariance matrix \mathbf{C}_t and incurs a loss equal to the variance $\mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t$. Thus for a sequence of T trials the total loss of the algorithm will be $L_{\text{alg}} = \sum_{t=1}^T \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t$. We want this loss to be comparable to the total variance of the best probability vector \mathbf{u} chosen in hindsight, i.e. $L_{\mathbf{u}} = \min_{\mathbf{u}} \mathbf{u}^\top \left(\sum_{t=1}^T \mathbf{C}_t \right) \mathbf{u}$, where \mathbf{u} lies in the probability simplex. This offline problem is a quadratic optimization problem with non-negativity constraints which does not have a closed form solution. However we can still prove bounds for the online algorithm.

The natural choice for an online algorithm for this problem is the Exponentiated Gradient algorithm of [KW97] since it maintains a probability vector as its parameter. Recall that for a general loss function $L_t(\mathbf{w}_t)$, the probability vector of Exponentiated Gradient algorithm is updated as

$$w_{t+1,i} = \frac{w_{t,i} e^{-\eta(\nabla L_t(\mathbf{w}_t))_i}}{\sum_i w_{t,i} e^{-\eta(\nabla L_t(\mathbf{w}_t))_i}}.$$

This update is motivated by considering the tradeoff between the relative entropy divergence to the old probability vector and the current loss, where $\eta > 0$ is the tradeoff parameter:

$$\mathbf{w}_{t+1} \approx \arg \min_{\mathbf{w} \text{ prob. vec.}} d(\mathbf{w}, \mathbf{w}_t) + \eta L_t(\mathbf{w}),$$

where \approx comes from the fact that the gradient at \mathbf{w}_{t+1} that should appear in the exponent is approximated by the gradient at \mathbf{w}_t (See [KW97] for more discussion). In our application, $L_t(\mathbf{w}_t) = \frac{1}{2} \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t$ and $\nabla L_t(\mathbf{w}_t) = \mathbf{C}_t \mathbf{w}_t$, leading to the following update:

$$w_{t+1,i} = \frac{w_{t,i} e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i}}{\sum_{i=1}^n w_{t,i} e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i}}.$$

2.4 Proof of Relative Loss Bounds

We now use the divergence $d(\mathbf{u}, \mathbf{w})$ that motivated the update as a measure of progress in the analysis.

Lemma 1. Let \mathbf{w}_t be the weight vector of the algorithm before trial t and let \mathbf{u} be an arbitrary comparison probability vector. Also, let r be the bound on the range of elements in covariance matrix \mathbf{C}_t , specifically let $\max_{i,j} |\mathbf{C}_t(i,j)| \leq \frac{r}{2}$. For any constants a and b such that $0 < a \leq \frac{b}{1+rb}$ and a learning rate $\eta = \frac{2b}{1+rb}$ we have:

$$a \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - b \mathbf{u}^\top \mathbf{C}_t \mathbf{u} \leq d(\mathbf{u}, \mathbf{w}_t) - d(\mathbf{u}, \mathbf{w}_{t+1}).$$

Proof. The proof given in Appendix A follows the same outline as Lemma 5.8 of [KW97] which gives an inequality for the Exponentiated Gradient algorithm when applied to linear regression. \square

Lemma 2. Let $\max_{i,j} |\mathbf{C}_t(i,j)| \leq \frac{r}{2}$ as before. Then for arbitrary positive c and learning rate $\eta = \frac{2c}{r(c+1)}$, the following bound holds:

$$L_{\text{alg}} \leq (1+c)L_{\mathbf{u}} + \left(1 + \frac{1}{c}\right) r d(\mathbf{u}, \mathbf{w}_1).$$

Proof. Let $b = \frac{c}{r}$, then for $a = \frac{b}{rb+1} = \frac{c}{r(c+1)}$ and $\eta = 2a = \frac{2c}{r(c+1)}$, we can use the inequality of Lemma 1 and obtain:

$$\frac{c}{c+1} \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - c \mathbf{u}^\top \mathbf{C}_t \mathbf{u} \leq r(d(\mathbf{u}, \mathbf{w}_t) - d(\mathbf{u}, \mathbf{w}_{t+1})).$$

Summing over the trials t results in:

$$\frac{c}{c+1} L_{\text{alg}} - c L_{\mathbf{u}} \leq r(d(\mathbf{u}, \mathbf{w}_1) - d(\mathbf{u}, \mathbf{w}_{t+1})) \leq r d(\mathbf{u}, \mathbf{w}_1).$$

Now the statement of the lemma immediately follows. \square

The following theorem describes how to choose the learning rate for the purpose of minimizing the upper bound:

Theorem 1. Let $\mathbf{C}_1, \dots, \mathbf{C}_T$ be an arbitrary sequence of covariance matrices such that $\max_{i,j} |\mathbf{C}_t(i,j)| \leq \frac{r}{2}$ and assume that $\mathbf{u}^\top \sum_{t=1}^T \mathbf{C}_t \mathbf{u} \leq L$. Then running our algorithm with uniform start vector $\mathbf{w}_1 = (\frac{1}{n}, \dots, \frac{1}{n})$ and learning rate $\eta = \frac{2\sqrt{L \log n}}{r\sqrt{\log n} + \sqrt{rL}}$ leads to the following bound:

$$L_{\text{alg}} \leq L_{\mathbf{u}} + 2\sqrt{rL \log n} + r \log n.$$

Proof. By Lemma 2 and since $d(\mathbf{u}, \mathbf{w}_1) \leq \log n$:

$$L_{\text{alg}} \leq L_{\mathbf{u}} + cL + \frac{r \log n}{c} + r \log n.$$

By differentiating we see that $c = \sqrt{\frac{r \log n}{L}}$ minimizes the r.h.s. and substituting this choice of c gives the bound of the theorem. \square

3 Variance Minimization over the Unit Sphere

3.1 Definitions

The *trace* $\text{tr}(\mathcal{A})$ of a square matrix \mathcal{A} is the sum of its diagonal elements. It is invariant under a change of basis transformation and thus it is also equal to the sum of eigenvalues of the matrix. The trace generalizes the normal dot product between vectors to the space of matrices, i.e. $\text{tr}(\mathcal{A}\mathcal{B}) = \text{tr}(\mathcal{B}\mathcal{A}) = \sum_{i,j} \mathcal{A}(i,j)\mathcal{B}(i,j)$. The trace is also a linear operator, that is $\text{tr}(a\mathcal{A} + b\mathcal{B}) = a\text{tr}(\mathcal{A}) + b\text{tr}(\mathcal{B})$. Another useful property of the trace is its cycling invariance, i.e. $\text{tr}(\mathcal{A}\mathcal{B}\mathcal{C}) = \text{tr}(\mathcal{B}\mathcal{C}\mathcal{A}) = \text{tr}(\mathcal{C}\mathcal{A}\mathcal{B})$. A particular instance of this is the following manipulation: $\mathbf{u}^\top \mathcal{A} \mathbf{u} = \text{tr}(\mathbf{u}^\top \mathcal{A} \mathbf{u}) = \text{tr}(\mathcal{A} \mathbf{u} \mathbf{u}^\top)$.

Dyads have trace one because $\text{tr}(\mathbf{u}\mathbf{u}^\top) = \mathbf{u}^\top \mathbf{u} = 1$. We generalize mixtures or probability vectors to *density matrices*. Such matrices are mixtures of any number of dyads, i.e. $\mathcal{W} = \sum_i \alpha_i \mathbf{u}_i \mathbf{u}_i^\top$ where $\alpha_j \geq 0$ and $\sum_i \alpha_i = 1$. Equivalently, density matrices are arbitrary symmetric positive definite matrices of trace one. Any density matrix \mathcal{W} can be decomposed into a sum of exactly n dyads corresponding to the orthogonal set of its eigenvectors \mathbf{w}_i , i.e. $\mathcal{W} = \sum_{i=1}^n \omega_i \mathbf{w}_i \mathbf{w}_i^\top$ where the vector $\boldsymbol{\omega}$ of the n eigenvalues must be a probability vector. In quantum physics density matrices over the field of complex numbers represent the mixed state of a physical system.

We also need the matrix generalizations of the exponential and logarithm operations. Given the decomposition of a symmetric matrix $\mathcal{A} = \sum_i \alpha_i \mathbf{a}_i \mathbf{a}_i^\top$, the matrix exponential and logarithm denoted as **exp** and **log** are computed as follows:

$$\mathbf{exp}(\mathcal{A}) = \sum_i e^{\alpha_i} \mathbf{a}_i \mathbf{a}_i^\top, \quad \mathbf{log}(\mathcal{A}) = \sum_i \log \alpha_i \mathbf{a}_i \mathbf{a}_i^\top$$

In other words, the exponential and the logarithm are applied to the eigenvalues and the eigenvectors remain unchanged. Obviously, the matrix logarithm is only defined when the matrix is strictly positive definite. In analogy with the exponential for numbers, one would expect the following equality to hold: $\mathbf{exp}(\mathcal{A} + \mathcal{B}) = \mathbf{exp}(\mathcal{A}) \mathbf{exp}(\mathcal{B})$. However this is only true when the symmetric matrices \mathcal{A} and \mathcal{B} commute, i.e. $\mathcal{A}\mathcal{B} = \mathcal{B}\mathcal{A}$, which occurs iff both matrices share the same eigensystem. On the other hand, the following trace inequality, called the Golden-Thompson inequality, holds for arbitrary symmetric matrices:

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

The following quantum relative entropy is a generalization of the classical relative entropy to density matrices due to Umegaki (see e.g. [NC00]):

$$\Delta(\mathcal{U}, \mathcal{W}) = \text{tr}(\mathcal{U}(\log \mathcal{U} - \log \mathcal{W})).$$

We will also use generalized inequalities for the cone of positive definite matrices: $\mathcal{A} \preceq \mathcal{B}$ if $\mathcal{B} - \mathcal{A}$ positive definite.

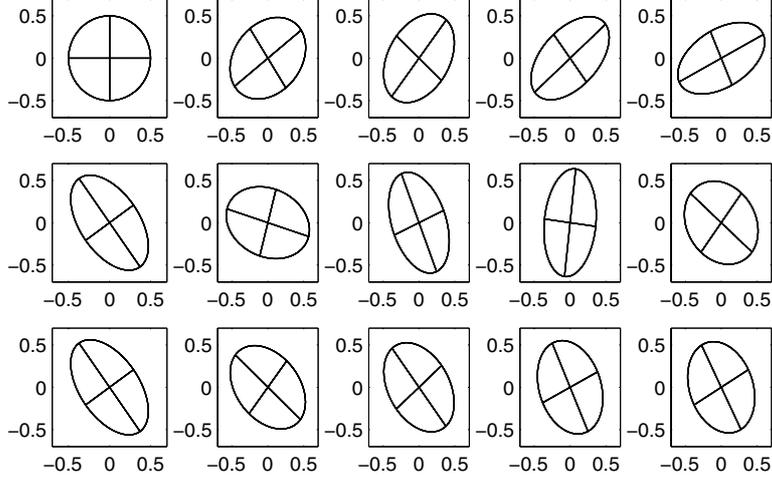


Fig. 2. The figure depicts a sequence of updates for the density matrix algorithm when the dimension is 2. All 2-by-2 matrices are represented as ellipses. The top row shows the density matrices \mathcal{W}_t chosen by the algorithm. The middle row shows the covariance matrix \mathbf{C}_t received in that trial. Finally, the bottom row is the average $\mathbf{C}_{\leq t} = \frac{\sum_{q=1}^t \mathbf{C}_q}{t}$ of all covariance matrices so far. By the update (1), $\mathcal{W}_{t+1} = \frac{\exp(-\eta t \mathbf{C}_{\leq t})}{Z_t}$, where Z_t is a normalization. Therefore, $\mathbf{C}_{\leq t}$ in the third row has the same eigensystem as the density matrix \mathcal{W}_{t+1} in the next column of the first row. Note the tendency of the algorithm to try to place more weight on the minimal eigenvalue of the covariance average. Since the algorithm is not sure about the future, it does not place the full weight onto that eigenvalue but hedges its bets instead and places some weight onto the other eigenvalues as well.

3.2 Applications

We develop online algorithms that perform as well as the eigenvector associated with a minimum (or maximum) eigenvalue. It seems that online versions of principal component analysis and other spectral methods can also be developed using the methodology of this paper. For instance, spectral clustering methods of [CSTK01] use a similar form of loss.

3.3 Algorithm and Motivation

As before we briefly review our setup. On each trial t our algorithm chooses a density matrix \mathcal{W}_t described as a mixture $\sum_i \omega_{t,i} \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top$. It then receives a covariance matrix \mathbf{C}_t and incurs a loss equal to the expected variance of its mixture:

$$\text{tr}(\mathcal{W}_t \mathbf{C}_t) = \text{tr}\left(\left(\sum_i \omega_{t,i} \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top\right) \mathbf{C}_t\right) = \sum_i \omega_{t,i} \mathbf{w}_{t,i}^\top \mathbf{C}_t \mathbf{w}_{t,i}.$$

On a sequence of T trials the total loss of the algorithm will be $L_{\text{alg}} = \sum_{t=1}^T \text{tr}(\mathcal{W}_t \mathbf{C}_t)$. We want this loss to be not too much larger than the

total variance of best unit vector \mathbf{u} chosen in hindsight, i.e. $L_{\mathbf{u}} = \text{tr}(\mathbf{u}\mathbf{u}^\top \sum_t \mathbf{C}_t) = \mathbf{u}^\top (\sum_t \mathbf{C}_t) \mathbf{u}$. The set of dyads is not a convex set. We therefore close it by using convex combinations of dyads (i.e. density matrices) as our parameter space. The best offline parameter is still a single dyad:

$$\min_{\mathbf{u} \text{ dens.mat.}} \text{tr}(\mathbf{u}\mathbf{C}) = \min_{\mathbf{u} : \|\mathbf{u}\|_2=1} \mathbf{u}^\top \mathbf{C} \mathbf{u}$$

Curiously enough our, loss $\text{tr}(\mathbf{W}\mathbf{C})$ has interpretation in quantum mechanics as the expected outcome of measuring a physical system in mixture state \mathbf{W} with instrument \mathbf{C} . Let \mathbf{C} be decomposed as $\sum_i \gamma_i \mathbf{c}_i \mathbf{c}_i^\top$. The eigenvalues γ_i are the possible numerical outcomes of measurement. When measuring a pure state specified by unit vector \mathbf{u} , the probability of obtaining outcome γ_i is given as $(\mathbf{u} \cdot \mathbf{c}_i)^2$ and the expected outcome is $\text{tr}(\mathbf{u}\mathbf{u}^\top \mathbf{C}) = \sum_i (\mathbf{u} \cdot \mathbf{c}_i)^2 \gamma_i$. For a mixed state \mathbf{W} we have the following double expectation:

$$\text{tr}(\mathbf{W}\mathbf{C}) = \text{tr} \left(\left(\sum_i \omega_i \mathbf{w}_i \mathbf{w}_i^\top \right) \left(\sum_j \gamma_j \mathbf{c}_j \mathbf{c}_j^\top \right) \right) = \sum_{i,j} (\mathbf{w}_i \cdot \mathbf{c}_j)^2 \gamma_i \omega_j,$$

where the matrix of measurement probabilities $(\mathbf{w}_i \cdot \mathbf{c}_j)^2$ is a doubly stochastic matrix. Note also, that for the measurement interpretation the matrix \mathbf{C} does not have to be positive definite, but only symmetric. The algorithm and the proof of bounds in fact work fine for this case, but the meaning of the algorithm when \mathbf{C} is not a covariance matrix is less clear, since despite all these connections our algorithm does not seem to have the obvious quantum-mechanical interpretation. Our update clearly is not a unitary evolution of the mixture state and a measurement does not cause a collapse of the state as is the case in quantum physics. The question of whether this type of algorithm is still doing something quantum-mechanically meaningful remains intriguing. See also [War05, WK06] for additional discussion.

To derive our algorithm we use the trace expression for expected variance as our loss and replace the relative entropy with its matrix generalization. The following optimization problem produces the update:

$$\mathbf{W}_{t+1} = \underset{\mathbf{W} \text{ dens.mat.}}{\text{arg min}} \Delta(\mathbf{W}, \mathbf{W}_t) + \eta \text{tr}(\mathbf{W}\mathbf{C}_t)$$

Using a Lagrangian that enforces the trace constraint [TRW05], it is easy to solve this constrained minimization problem:

$$\mathbf{W}_{t+1} = \frac{\exp(\log \mathbf{W}_t - \eta \mathbf{C}_t)}{\text{tr}(\exp(\log \mathbf{W}_t - \eta \mathbf{C}_t))} = \frac{\exp(-\eta \sum_{q=1}^t \mathbf{C}_q)}{\text{tr}(\exp(-\eta \sum_{q=1}^t \mathbf{C}_q))}. \quad (1)$$

Note that for the second equation we assumed that $\mathbf{W}_1 = \frac{1}{n} \mathbf{I}$. The update is a special case of the Matrix Exponentiated Gradient update with the linear loss $\text{tr}(\mathbf{W}\mathbf{C}_t)$.

3.4 Proof Methodology

For the sake of clarity, we begin by recalling the proof of the worst-case loss bound for the Continuous Weighted Majority (WMC)/Hedge algorithm in the expert advice setting [LW94]. In doing so we clarify the dependence of the algorithm on the range of the losses. The update of that algorithm is given by:

$$w_{t+1,i} = \frac{w_{t,i}e^{-\eta l_{t,i}}}{\sum_i w_{t,i}e^{-\eta l_{t,i}}} \quad (2)$$

The proof always starts by considering the progress made during the update towards any comparison vector/parameter \mathbf{u} in terms of the motivating divergence for the algorithm, which in this case is the relative entropy:

$$d(\mathbf{u}, \mathbf{w}_t) - d(\mathbf{u}, \mathbf{w}_{t+1}) = \sum_i u_i \log \frac{w_{t+1,i}}{w_{t,i}} = -\eta \mathbf{u} \cdot \mathbf{l}_t - \log \sum_i w_{t,i} e^{-\eta l_{t,i}}.$$

We assume that $l_{t,i} \in [0, r]$, for $r > 0$, and use the inequality $\beta^x \leq 1 - (1 - \beta)^{\frac{x}{r}}$, for $x \in [0, r]$, with $\beta = e^{-\eta}$:

$$d(\mathbf{u}, \mathbf{w}_t) - d(\mathbf{u}, \mathbf{w}_{t+1}) \geq -\eta \mathbf{u} \cdot \mathbf{l}_t - \log\left(1 - \frac{\mathbf{w}_t \cdot \mathbf{l}_t}{r}(1 - e^{-\eta r})\right),$$

We now apply $\log(1 - x) \leq -x$:

$$d(\mathbf{u}, \mathbf{w}_t) - d(\mathbf{u}, \mathbf{w}_{t+1}) \geq -\eta \mathbf{u} \cdot \mathbf{l}_t + \frac{\mathbf{w}_t \cdot \mathbf{l}_t}{r}(1 - e^{-\eta r}),$$

and rewrite the above to

$$\mathbf{w}_t \cdot \mathbf{l}_t \leq \frac{r(d(\mathbf{u}, \mathbf{w}_t) - d(\mathbf{u}, \mathbf{w}_{t+1})) + \eta r \mathbf{u} \cdot \mathbf{l}_t}{1 - e^{-\eta r}}$$

Here $\mathbf{w}_t \cdot \mathbf{l}_t$ is the loss of the algorithm at trial t and $\mathbf{u} \cdot \mathbf{l}_t$ is the loss of the probability vector \mathbf{u} which serves as a comparator.

So far we assumed that $l_{t,i} \in [0, r]$. However, it suffices to assume that $\max_i l_{t,i} - \min_i l_{t,i} \leq r$. In other words, the individual losses can be positive or negative, as long as their range is bounded by r . For further discussion pertaining to the issues with losses having different signs see [CBMS05]. As we shall observe below, the requirement on the range of losses will become a requirement on the range of eigenvalues of the covariance matrices.

Define $\tilde{l}_{t,i} := l_{t,i} - \min_i l_{t,i}$. The update remains unchanged when the shifted losses $\tilde{l}_{t,i}$ are used in place of the original losses $l_{t,i}$ and we immediately get the inequality

$$\mathbf{w}_t \cdot \tilde{\mathbf{l}}_t \leq \frac{r(d(\mathbf{u}, \mathbf{w}_t) - d(\mathbf{u}, \mathbf{w}_{t+1})) + \eta r \mathbf{u} \cdot \tilde{\mathbf{l}}_t}{1 - e^{-\eta r}}.$$

Summing over t and dropping the $d(\mathbf{u}, \mathbf{w}_{t+1}) \geq 0$ term results in a bound that holds for any \mathbf{u} and thus for the best \mathbf{u} as well:

$$\sum_t \mathbf{w}_t \cdot \tilde{\mathbf{l}}_t \leq \frac{rd(\mathbf{u}, \mathbf{w}_t) + \eta r \sum_t \mathbf{u} \cdot \tilde{\mathbf{l}}_t}{1 - e^{-\eta r}}.$$

We can now tune the learning rate following [FS97]: if $\sum_t \mathbf{u} \cdot \tilde{\mathbf{l}}_t \leq \tilde{L}$ and $d(\mathbf{u}, \mathbf{w}_1) \leq D \leq \ln n$, then with $\eta = \frac{\log(1+\sqrt{2D/\tilde{L}})}{r}$ we get the bound

$$\sum_t \mathbf{w}_t \cdot \tilde{\mathbf{l}}_t \leq \sum_t \mathbf{u} \cdot \tilde{\mathbf{l}}_t + \sqrt{2r\tilde{L}D} + rd(\mathbf{u}, \mathbf{w}_1),$$

which is equivalent to

$$\underbrace{\sum_t \mathbf{w}_t \cdot \mathbf{l}_t}_{L_{\text{alg}}} \leq \underbrace{\sum_t \mathbf{u} \cdot \mathbf{l}_t}_{L_u} + \sqrt{2r\tilde{L}D} + rd(\mathbf{u}, \mathbf{w}_1).$$

Note that \tilde{L} is defined wrt the tilde versions of the losses and the update as well as the above bound is invariant under shifting the loss vectors \mathbf{l}_t by arbitrary constants. If the loss vectors \mathbf{l}_t are replaced by gain vectors, then the minus sign in the exponent of the update becomes a plus sign. In this case the inequality above is reversed and the last two terms are subtracted instead of added.

3.5 Proof of Relative Loss Bounds

In addition to the Golden-Thompson inequality we will need lemmas 2.1 and 2.2 from [TRW05]:

Lemma 3. *For any symmetric \mathcal{A} , such that $\mathbf{0} \preceq \mathcal{A} \preceq \mathbf{I}$ and any $\rho_1, \rho_2 \in \mathbb{R}$ the following holds:*

$$\exp(\mathcal{A}\rho_1 + (\mathbf{I} - \mathcal{A})\rho_2) \preceq \mathcal{A}e^{\rho_1} + (\mathbf{I} - \mathcal{A})e^{\rho_2}.$$

Lemma 4. *For any positive semidefinite \mathcal{A} and any symmetric \mathcal{B}, \mathcal{C} , $\mathcal{B} \preceq \mathcal{C}$ implies $\text{tr}(\mathcal{A}\mathcal{B}) \leq \text{tr}(\mathcal{A}\mathcal{C})$.*

We are now ready to generalize the WMC bound to matrices:

Theorem 2. *For any sequence of covariance matrices $\mathcal{C}_1, \dots, \mathcal{C}_T$ such that $\mathbf{0} \preceq \mathcal{C}_t \preceq r\mathbf{I}$ and for any learning rate η , the following bound holds for arbitrary density matrix \mathcal{U} :*

$$\text{tr}(\mathcal{W}_t \mathcal{C}_t) \leq \frac{r(\Delta(\mathcal{U}, \mathcal{W}_t) - \Delta(\mathcal{U}, \mathcal{W}_{t+1})) + \eta r \text{tr}(\mathcal{U}\mathcal{C}_t)}{1 - e^{-r\eta}}.$$

Proof. We start by analyzing the progress made towards the comparison matrix \mathcal{U} in terms of quantum relative entropy:

$$\begin{aligned} \Delta(\mathcal{U}, \mathcal{W}_t) - \Delta(\mathcal{U}, \mathcal{W}_{t+1}) &= \text{tr}(\mathcal{U}(\log \mathcal{U} - \log \mathcal{W}_t)) - \text{tr}(\mathcal{U}(\log \mathcal{U} - \log \mathcal{W}_{t+1})) \\ &= -\text{tr} \left(\mathcal{U} \left(\log \mathcal{W}_t + \log \frac{\exp(\log \mathcal{W}_t - \eta \mathcal{C}_t)}{\text{tr}(\exp(\log \mathcal{W}_t - \eta \mathcal{C}_t))} \right) \right) \\ &= -\eta \text{tr}(\mathcal{U}\mathcal{C}_t) - \log(\text{tr}(\exp(\log \mathcal{W}_t - \eta \mathcal{C}_t))). \end{aligned} \tag{3}$$

We will now bound the log of trace term. First, the following holds via the Golden-Thompson inequality:

$$\mathrm{tr}(\exp(\log \mathbf{W}_t - \eta \mathbf{C}_t)) \leq \mathrm{tr}(\mathbf{W}_t \exp(-\eta \mathbf{C}_t)). \quad (4)$$

Since $\mathbf{0} \preceq \frac{\mathbf{C}_t}{r} \preceq \mathbf{I}$, we can use Lemma 3 with $\rho_1 = -\eta r$, $\rho_2 = 0$:

$$\exp(-\eta \mathbf{C}_t) \preceq \mathbf{I} - \frac{\mathbf{C}_t}{r}(1 - e^{-\eta r}).$$

Now multiply both sides on the left with \mathbf{W}_t and take a trace. The inequality is preserved according to Lemma 4:

$$\mathrm{tr}(\mathbf{W}_t \exp(-\eta \mathbf{C}_t)) \leq \left(1 - \frac{\mathrm{tr}(\mathbf{W}_t \mathbf{C}_t)}{r}(1 - e^{-r\eta})\right).$$

Taking logs of both sides we have:

$$\log(\mathrm{tr}(\mathbf{W}_t \exp(-\eta \mathbf{C}_t))) \leq \log\left(1 - \frac{\mathrm{tr}(\mathbf{W}_t \mathbf{C}_t)}{r}(1 - e^{-r\eta})\right). \quad (5)$$

To bound the log expression on the right we use inequality $\log(1 - x) \leq -x$:

$$\log\left(1 - \frac{\mathrm{tr}(\mathbf{W}_t \mathbf{C}_t)}{r}(1 - e^{-r\eta})\right) \leq -\frac{\mathrm{tr}(\mathbf{W}_t \mathbf{C}_t)}{r}(1 - e^{-r\eta}). \quad (6)$$

By combining inequalities (4-6), we obtain the following bound on the log trace term:

$$-\log(\mathrm{tr}(\exp(\log \mathbf{W}_t - \eta \mathbf{C}_t))) \geq \frac{\mathrm{tr}(\mathbf{W}_t \mathbf{C}_t)}{r}(1 - e^{-r\eta}).$$

Plugging this into equation (3) we obtain

$$r(\Delta(\mathbf{u}, \mathbf{W}_t) - \Delta(\mathbf{u}, \mathbf{W}_{t+1})) + \eta r \mathrm{tr}(\mathbf{u} \mathbf{C}_t) \geq \mathrm{tr}(\mathbf{W}_t \mathbf{C}_t)(1 - e^{-r\eta}),$$

which is the inequality of the theorem. \square

Note that our density matrix update (1) is invariant wrt the variable change $\widetilde{\mathbf{C}}_t = \mathbf{C}_t - \lambda_{\min}(\mathbf{C}_t)\mathbf{I}$. Therefore by the above theorem, the following inequality holds whenever $\lambda_{\max}(\mathbf{C}_t) - \lambda_{\min}(\mathbf{C}_t) \leq r$:

$$\mathrm{tr}(\mathbf{W}_t \widetilde{\mathbf{C}}_t) \leq \frac{r(\Delta(\mathbf{u}, \mathbf{W}_t) - \Delta(\mathbf{u}, \mathbf{W}_{t+1})) + \eta r \mathrm{tr}(\mathbf{u} \widetilde{\mathbf{C}}_t)}{1 - e^{-r\eta}}.$$

We can now sum over trials and tune the learning rate as done at the end of Section 3.4. If $\sum_t \mathrm{tr}(\mathbf{u} \widetilde{\mathbf{C}}_t) \leq \tilde{L}$ and $\Delta(\mathbf{u}, \mathbf{W}_1) \leq D$, with $\eta = \frac{\log(1 + \sqrt{\frac{2D}{\tilde{L}}})}{r}$ we get the bound:

$$\underbrace{\sum_t \mathrm{tr}(\mathbf{W}_t \mathbf{C}_t)}_{L_{\mathrm{alg}}} \leq \underbrace{\sum_t \mathrm{tr}(\mathbf{u} \mathbf{C}_t)}_{L_{\mathbf{u}}} + \sqrt{2r\tilde{L}D} + r\Delta(\mathbf{u}, \mathbf{W}_1).$$

4 Conclusions

We presented two algorithms for online variance minimization problems. For the first problem, the variance was measured along a probability vector. It would be interesting to combine this work with the online algorithms considered in [HSSW98, Cov91] that maximize the return of the portfolio. It should be possible to design online algorithms that minimize a trade off between the return of the portfolio (first order information) and the variance/risk. Note that it is easy to extend the portfolio vector to maintain short positions: Simply keep two weights w_i^+ and w_i^- per component as is done in the EG[±] algorithm of [KW97].

In our second problem the variance was measured along an arbitrary direction. We gave a natural generalization of the WMC/Hedge algorithm to the case when the parameters are density matrices. Note that in this paper we upper bounded the sum of the expected variances over trials, whereas in [War05, WK06] a Bayes rule for density matrices was given for which a lower bound was provided on the product of the expected variances over trials.²

Much work has been done on exponential weight updates for the experts. In particular, algorithms have been developed for shifting experts by combining the exponential updates with an additive “sharing update” [HW98]. In preliminary work we showed that these techniques easily carry over to the density matrix setting. This includes the more recent work on the “sharing to the past average” update, which introduces a long-term memory [BW02].

Appendix A

Proof of Lemma 1

Begin by analyzing the progress towards the comparison vector \mathbf{u} :

$$\begin{aligned} d(\mathbf{u}, \mathbf{w}_t) - d(\mathbf{u}, \mathbf{w}_{t+1}) &= \sum u_i \log \frac{u_i}{w_{t,i}} - \sum u_i \log \frac{u_i}{w_{t+1,i}} \\ &= \sum u_i \log w_{t+1,i} - \sum u_i \log w_{t,i} \\ &= \sum u_i \log \frac{w_{t,i} e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i}}{\sum w_{t,i} e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i}} - \sum u_i \log w_{t,i} \\ &= \sum u_i \log w_{t,i} - \eta \sum u_i (\mathbf{C}_t \mathbf{w}_t)_i - \\ &\quad - \log \left(\sum w_{t,i} e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i} \right) - \sum u_i \log w_{t,i} \\ &= -\eta \sum u_i (\mathbf{C}_t \mathbf{w}_t)_i - \log \left(\sum w_{t,i} e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i} \right) \end{aligned}$$

Thus, our bound is equivalent to showing $F \leq 0$ with F given as:

$$F = a \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - b \mathbf{u}^\top \mathbf{C}_t \mathbf{u} + \eta \mathbf{u}^\top \mathbf{C}_t \mathbf{w}_t + \log \left(\sum w_{t,i} e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i} \right)$$

² This amounts to an upper bound on the sum of the negative logarithms of the expected variances.

We proceed by bounding the log term. The assumption on the range of elements of \mathbf{C}_t and the fact that \mathbf{w}_t is a probability vector allows us to conclude that $\max_i(\mathbf{C}_t \mathbf{w}_t)_i - \min_i(\mathbf{C}_t \mathbf{w}_t)_i \leq r$, since $(\mathbf{C}_t \mathbf{w}_t)_i = \sum_j \mathbf{C}_t(i, j) \mathbf{w}_t(j)$. Now, assume that l is a lower bound for $(\mathbf{C}_t \mathbf{w}_t)_i$, then we have that $l \leq (\mathbf{C}_t \mathbf{w}_t)_i \leq l + r$, or $0 \leq \frac{(\mathbf{C}_t \mathbf{w}_t)_i - l}{r} \leq 1$. This allows us to use the inequality $a^x \leq 1 - x(1 - a)$ for $a \geq 0$ and $0 \leq x \leq 1$. Let $a = e^{-\eta r}$:

$$e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i} = e^{-\eta l} (e^{-\eta r})^{\frac{(\mathbf{C}_t \mathbf{w}_t)_i - l}{r}} \leq e^{-\eta b} \left(1 - \frac{(\mathbf{C}_t \mathbf{w}_t)_i - l}{r} (1 - e^{-\eta r}) \right)$$

Using this inequality we obtain:

$$\log \left(\sum_i w_{t,i} e^{-\eta(\mathbf{C}_t \mathbf{w}_t)_i} \right) \leq -\eta l + \log \left(1 - \frac{\mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - l}{r} (1 - e^{-\eta r}) \right)$$

This gives us $F \leq G$, with G given as:

$$G = a \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - b \mathbf{u}^\top \mathbf{C}_t \mathbf{u} + \eta \mathbf{u}^\top \mathbf{C}_t \mathbf{w}_t - \eta l + \log \left(1 - \frac{\mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - l}{r} (1 - e^{-\eta r}) \right)$$

It is sufficient to show that $G \leq 0$. Let $\mathbf{z} = \sqrt{\mathbf{C}_t} \mathbf{u}$. Then $G(\mathbf{z})$ becomes:

$$G(\mathbf{z}) = -b \mathbf{z}^\top \mathbf{z} + \eta \mathbf{z}^\top \sqrt{\mathbf{C}_t} \mathbf{w}_t + \text{constant}.$$

The function $G(\mathbf{z})$ is concave quadratic and is maximized at:

$$\frac{\partial G}{\partial \mathbf{z}} = -2b \mathbf{z} + \eta \sqrt{\mathbf{C}_t} \mathbf{w}_t = 0, \quad \mathbf{z} = \frac{\eta}{2b} \sqrt{\mathbf{C}_t} \mathbf{w}_t$$

We substitute this value of \mathbf{z} into G and get $G \leq H$, where H is given by:

$$H = a \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t + \frac{\eta^2}{4b} \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - \eta l + \log \left(1 - \frac{\mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - l}{r} (1 - e^{-\eta r}) \right).$$

Since $l \leq (\mathbf{C}_t \mathbf{w}_t)_i \leq l + r$, then obviously so is $\mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t$, since weighted average stays within the bounds. Now we can use the inequality $\log(1 - p(1 - e^q)) \leq pq + \frac{q^2}{8}$, for $0 \leq p \leq 1$ and $q \in \mathbb{R}$:

$$\log \left(1 - \frac{\mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - l}{r} (1 - e^{-\eta r}) \right) \leq -\eta \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t + \eta l + \frac{\eta^2 r^2}{8}.$$

We get $H \leq S$, where S is given as:

$$\begin{aligned} S &= a \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t + \frac{\eta^2}{4b} \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t - \eta \mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t + \frac{\eta^2 r^2}{8} \\ &= \frac{\mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t}{4b} (4ab + \eta^2 - 4b\eta) + \frac{\eta^2 r^2}{8}. \end{aligned}$$

By our assumptions $\mathbf{w}_t^\top \mathbf{C}_t \mathbf{w}_t \leq \frac{r}{2}$, and therefore:

$$S \leq Q = \eta^2 \left(\frac{r^2}{8} + \frac{r}{8b} \right) - \frac{\eta r}{2} + \frac{ar}{2}$$

We want to make this expression as small as possible, so that it stays below zero. To do so we minimize it over η :

$$2\eta\left(\frac{r^2}{8} + \frac{r}{8b}\right) - \frac{r}{2} = 0, \quad \eta = \frac{2b}{rb+1}$$

Finally we substitute this value of η into Q and obtain conditions on a , so that $Q \leq 0$ holds:

$$a \leq \frac{b}{rb+1}$$

This concludes the proof. □

References

- [BV04] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [BW02] O. Bousquet and M. K. Warmuth. Tracking a small set of experts by mixing past posteriors. *J. of Machine Learning Research*, 3(Nov):363–396, 2002.
- [CBMS05] Nicolo Cesa-Bianchi, Yishay Mansour, and Gilles Stoltz. Improved second-order bounds for prediction with expert advice. In *Proceedings of the 18th Annual Conference on Learning Theory (COLT 05)*, pages 217–232. Springer, June 2005.
- [Cov91] T. M. Cover. Universal portfolios. *Mathematical Finance*, 1(1):1–29, 1991.
- [CSTK01] Nello Cristianini, John Shawe-Taylor, and Jaz Kandola. Spectral kernel methods for clustering. In *Advances in Neural Information Processing Systems 14*, pages 649–655. MIT Press, December 2001.
- [FS97] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, August 1997.
- [HSSW98] D. Helmbold, R. E. Schapire, Y. Singer, and M. K. Warmuth. On-line portfolio selection using multiplicative updates. *Mathematical Finance*, 8(4):325–347, 1998.
- [HW98] M. Herbster and M. K. Warmuth. Tracking the best expert. *Journal of Machine Learning*, 32(2):151–178, August 1998.
- [KW97] J. Kivinen and M. K. Warmuth. Additive versus exponentiated gradient updates for linear prediction. *Information and Computation*, 132(1):1–64, January 1997.
- [LW94] N. Littlestone and M. K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108(2):212–261, 1994.
- [NC00] M.A. Nielsen and I.L. Chuang. *Quantum Computation and Quantum Information*. Cambridge University Press, 2000.
- [TRW05] K. Tsuda, G. Rätsch, and M. K. Warmuth. Matrix exponentiated gradient updates for on-line learning and Bregman projections. *Journal of Machine Learning Research*, 6:995–1018, June 2005.
- [War05] Manfred K. Warmuth. Bayes rule for density matrices. In *Advances in Neural Information Processing Systems 18 (NIPS 05)*. MIT Press, December 2005.
- [WK06] Manfred K. Warmuth and Dima Kuzmin. A Bayesian probability calculus for density matrices. Unpublished manuscript, March 2006.