

Online Linear Regression using Burg Entropy

Prateek Jain, Brian Kulis, and Inderjit Dhillon

Technical Report TR-07-08
University of Texas at Austin
Austin, TX 78712

February 14, 2007

Abstract

We consider the problem of online prediction with a linear model. In contrast to existing work in online regression, which regularizes based on squared loss or KL-divergence, we regularize using divergences arising from the Burg entropy. We demonstrate regret bounds for our resulting online gradient-descent algorithm; to our knowledge, these are the first online bounds involving Burg entropy. We extend this analysis to the matrix case, where our algorithm employs LogDet-based regularization, and discuss an application to online metric learning. We demonstrate empirically that using Burg entropy for regularization is useful in the presence of noisy data.

1 Introduction

Recently, online learning has received significant attention, and several recent results have demonstrated strong provable online bounds. In the online linear regression problem, considered in this paper, the algorithm receives an instance \mathbf{x}_t ($\mathbf{x}_t \in \mathbb{R}_+^m$) at time t and predicts $\hat{y}_t = \mathbf{w}_t \cdot \mathbf{x}_t$; \mathbf{w}_t is the current model used by the algorithm. We denote w_t^i as the i -th component of the model at the t -th step. The error incurred at step t is $l_t(\mathbf{w}_t) = (y_t - \hat{y}_t)^2$, where y_t is the true/target distance, and the goal is to minimize the total loss over all time steps $\sum_t l_t(\mathbf{w}_t)$.

If there is no correlation between the input \mathbf{x}_t and the output y_t , then the cumulative loss of the algorithm may be unbounded. Hence, online algorithms are traditionally compared to the performance of the best possible offline solution, where all the input and output instances are provided beforehand. The goal is to prove a relative loss bound on how the online algorithm compares to the optimal offline algorithm. Given a T -trial sequence $S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_T, y_T)\}$, the optimal offline solution is given by

$$\mathbf{u} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{t=1}^T l_t(\mathbf{w})$$

s.t. $w^i > 0 \quad \forall i$

A typical approach to solving the online linear regression problem uses the standard gradient descent method with regularization for minimizing the loss at each step. Specifically, the algorithm tries to minimize the following function at each step:

$$U(\mathbf{w}) = \overbrace{D(\mathbf{w}, \mathbf{w}_t)}^{\text{Regularization Term}} + \eta_t \overbrace{L(y_t, \mathbf{w} \cdot \mathbf{x}_t)}^{\text{Loss Term}}, \quad (1)$$

where η_t is learning rate at t -th step and $d(\mathbf{w}, \mathbf{w}_t)$ is some distance measure between the new weight vector \mathbf{w} and the current weight vector \mathbf{w}_t .

In the above objective function, the regularization term favors weight vectors which are “close” to the current model \mathbf{w}_t , representing a tendency for conservativeness. On the other hand, the loss term is minimized

when \mathbf{w} is updated to exactly satisfy $y_t = \hat{y}_t$; hence the loss term has a tendency to satisfy target outputs for recent examples. The tradeoff between these two quantities is handled by the learning rate η_t , an important parameter in online learning problems.

The performance of any method following this approach will depend heavily on the choice of the regularization term and the loss function. Kivinen et. al[6, 8] show that the performance of any regularization term in turn depends heavily on the problem domain. For example, for sparse inputs, relative entropy as a regularization term often works better than squared Euclidean distance. For dense inputs, relative entropy’s performance is generally very poor compared to that of squared Euclidean distance.

In this paper, we study the performance of online linear regression when the Burg entropy is considered as the regularization term:

$$D_{Burg}(\mathbf{w}_{t+1}, \mathbf{w}_t) = \sum_{i=1}^N \left(\frac{w_{t+1}^i}{w_t^i} - \log \left(\frac{w_{t+1}^i}{w_t^i} \right) \right).$$

The relative Burg entropy (or Itakura-Saito distance) is a special case of a *Bregman divergence*, and such measures have been shown to be useful distance measures for various machine learning tasks [1] (relative entropy and squared Euclidean distance are other special cases of Bregman divergences). The Burg entropy in particular has recently been shown to be useful for matrix approximations and kernel learning [7]. Our study of Burg entropy is also motivated by the fact that online information-theoretic metric learning reduces to online linear regression problem with the LogDet divergence—the natural generalization of the Burg entropy as regularization term [3]. The results in this paper complement those of our work on information-theoretic metric learning.

We present an algorithm for using Burg relative entropy for regularization and prove an associated regret bound. We also generalize the above approach to the matrix case when the input matrices are rank-one; this is precisely the case needed for information-theoretic metric learning. In section 5, we empirically show how the Burg entropy performs compared to other standard regularization terms.

2 The Vector Case

In this study, we select the Bregman divergence corresponding to the Burg entropy as our regularization term for (1):

$$D_{Burg}(\mathbf{w}_{t+1}, \mathbf{w}_t) = \sum_{i=1}^N \left(\frac{w_{t+1}^i}{w_t^i} - \log \left(\frac{w_{t+1}^i}{w_t^i} \right) \right).$$

We derive Algorithm 1 from standard gradient descent; in particular, if we differentiate $U(\mathbf{w})$, we obtain the update given as (2) in Algorithm 1.

Algorithm 1 Burg Gradient (BG) Algorithm

- 1: **Initialize:** Set $\eta_t = \frac{1}{4mR^2}$ and $w_0^i = \frac{1}{N}, \forall i$.
- 2: **Prediction:** Upon receiving the t -th instance \mathbf{x}_t , give the prediction $\hat{y}_t = \mathbf{w}_t \cdot \mathbf{x}_t$.
- 3: **Update:** Upon receiving the t -th outcome y_t , update the weights according to the rule,

$$w_{t+1}^i = ((w_t^i)^{-1} + 2\eta_t(\hat{y}_t - y_t)x_t^i)^{-1}, \tag{2}$$

where

$$\eta_t = \begin{cases} n_{t-1} & \text{if } \hat{y}_t - y_t > 0 \\ \min \left(n_{t-1}, \frac{1}{-2(\hat{y}_t - y_t)x_t^i} \left(\frac{1}{w_t^i} - 1 \right) \right) & \text{if } \hat{y}_t - y_t < 0. \end{cases} \tag{3}$$

Note that the divergence D_{Burg} is only defined over positive-valued weight vectors. To maintain positivity, we use an adaptive learning rate η_t , which guarantees that the resulting vector \mathbf{w} is positive. Assuming $\eta_{t-1} > 0$, then $\eta_t > 0$ since η_t is the minimum of two positive quantities. Also, if $\hat{y}_t - y_t > 0$ then clearly

$\mathbf{w}_{t+1} > 0$ (assuming $\mathbf{w}_t > 0$). On the other hand, if $\hat{y}_t - y_t < 0$, then since $\eta_t \leq \frac{1}{-2(\hat{y}_t - y_t)x_t^i} \left(\frac{1}{w_t^i} - 1 \right)$, we can conclude that $\mathbf{w}_{t+1} > 0$. Hence, we maintain positivity.

2.1 Regret Bounds

Now we demonstrate regret bounds for the proposed algorithm. Let the total loss of Algorithm 1 be $Loss_{BG} = \sum_t l_t(\mathbf{w}_t)$, and let $Loss_{\mathbf{u}} = \sum_t l_t(\mathbf{u})$ be the loss of the best offline solution. First we bound the loss of the algorithm at a single trial in terms of the loss of a comparison vector \mathbf{u} at that trial and the progress of the algorithm towards \mathbf{u} .

Lemma 1.

$$a_t(\hat{y}_t - y_t)^2 - b_t(\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 \leq D_{Burg}(\mathbf{u}, \mathbf{w}_t) - D_{Burg}(\mathbf{u}, \mathbf{w}_{t+1}),$$

where a_t and b_t are positive constants and \mathbf{u} is the optimal offline solution.

Proof.

$$\begin{aligned} D_{Burg}(\mathbf{u}, \mathbf{w}_t) - D_{Burg}(\mathbf{u}, \mathbf{w}_{t+1}) &= \sum_{i=1}^N \left(u^i \left(\frac{1}{w_t^i} - \frac{1}{w_{t+1}^i} \right) + \ln \frac{w_t^i}{w_{t+1}^i} \right) \\ &= -2\eta_t(\hat{y}_t - y_t) \sum_{i=1}^N u^i x_t^i + \sum_{i=1}^N \ln(1 + 2\eta_t(\hat{y}_t - y_t)w_t^i x_t^i) \\ &= -2\eta_t(\hat{y}_t - y_t) \mathbf{u} \cdot \mathbf{x}_t + \sum_{i=1}^N \ln(1 + 2\eta_t(\hat{y}_t - y_t)w_t^i x_t^i). \end{aligned}$$

Consider the following two cases:

1. $\hat{y}_t - y_t \geq 0$: Since $w_t^i > 0$ and $x_t^i > 0$, it follows that $2\eta_t(\hat{y}_t - y_t)w_t^i x_t^i \geq 0$.
2. $\hat{y}_t - y_t < 0$: Since $\eta_t < \frac{1}{-2(\hat{y}_t - y_t)x_t^i} \left(\frac{1}{w_t^i} - 1 \right)$, it follows that $2\eta_t(\hat{y}_t - y_t)w_t^i x_t^i > \frac{1}{w_t^i} - 1$. Now, $w_t^i > 0$ and so $2\eta_t(\hat{y}_t - y_t)w_t^i x_t^i > -1$.

Thus, in either case, $2\eta_t(\hat{y}_t - y_t)w_t^i x_t^i > -1$. We can therefore apply the inequality

$$\ln(1 + z) \geq z - z^2/2,$$

which holds for all $z > -1$. Hence

$$D_{Burg}(\mathbf{u}, \mathbf{w}_t) - D_{Burg}(\mathbf{u}, \mathbf{w}_{t+1}) \geq -2\eta_t(\hat{y}_t - y_t) \mathbf{u} \cdot \mathbf{x}_t + \sum_{i=1}^N (2\eta_t(\hat{y}_t - y_t)w_t^i x_t^i - 2\eta_t^2(\hat{y}_t - y_t)^2 (w_t^i x_t^i)^2).$$

Assume $x_t^i \leq R$ and $\sum_{i=1}^N w_t^i \leq m$; the above inequality simplifies to

$$D_{Burg}(\mathbf{u}, \mathbf{w}_t) - D_{Burg}(\mathbf{u}, \mathbf{w}_{t+1}) \geq -2\eta_t(\hat{y}_t - y_t) \mathbf{u} \cdot \mathbf{x}_t + 2\eta_t(\hat{y}_t - y_t)\hat{y}_t - 2\eta_t^2(\hat{y}_t - y_t)^2 mR^2.$$

Let $r = \mathbf{u} \cdot \mathbf{x}_t$. Proving the lemma amounts to showing that

$$D_{Burg}(\mathbf{u}, \mathbf{w}_t) - D_{Burg}(\mathbf{u}, \mathbf{w}_{t+1}) \geq -2\eta_t(\hat{y}_t - y_t)r + 2\eta_t(\hat{y}_t - y_t)\hat{y}_t - 2\eta_t^2(\hat{y}_t - y_t)^2 mR^2 \quad (4)$$

$$\geq a_t(\hat{y}_t - y_t)^2 - b_t(r - y_t)^2, \quad (5)$$

for some positive constants a and b . Consider the function

$$F(r) = a_t(\hat{y}_t - y_t)^2 - b_t(r - y_t)^2 + 2\eta_t(\hat{y}_t - y_t)r - 2\eta_t(\hat{y}_t - y_t)\hat{y}_t + 2\eta_t^2(\hat{y}_t - y_t)^2 mR^2.$$

Equation 5 is equivalent to $F(r) \leq 0, \forall r$. It can be easily seen that $F(r)$ is maximized when $r = y_t + \frac{a_t}{b_t}(\hat{y}_t - y_t)$. Substituting for r in $F(r)$ and simplifying, we get:

$$a_t(\hat{y}_t - y_t)^2 + \frac{\eta_t^2}{b_t}(\hat{y}_t - y_t)^2 - 2\eta_t(\hat{y}_t - y_t)^2 + 2\eta_t^2(\hat{y}_t - y_t)^2 mR^2.$$

Hence, we need to prove that

$$0 \geq a_t(\hat{y}_t - y_t)^2 + \frac{\eta_t^2}{b_t}(\hat{y}_t - y_t)^2 - 2\eta_t(\hat{y}_t - y_t)^2 + 2\eta_t^2(\hat{y}_t - y_t)^2 mR^2.$$

Since $(\hat{y}_t - y_t)^2 \geq 0$, this amounts to showing

$$0 \geq a_t b_t + \eta_t^2 - 2\eta_t b_t + 2\eta_t^2 mR^2 b_t = Q(\eta_t).$$

Now, $Q(\eta_t)$ is minimized for $\eta_t = \frac{b_t}{1+2mR^2 b_t}$, which implies $b_t = \frac{\eta_t}{1-2\eta_t mR^2}$. Since, $\eta_t \leq \eta_0 = \frac{1}{4mR^2}$, we know that $b_t > 0$. $Q(\eta_t) < 0$ iff $a_t \leq \frac{b_t}{1+2mR^2 b_t} = \eta_t$. Hence, for $0 \leq a_t \leq \eta_t$ and $b_t = \frac{\eta_t}{1-2\eta_t mR^2}$, Lemma 1 holds. \square

Using Lemma 1, we can sum the loss over all time steps to get an overall bound on the loss of the Burg entropy-based online learning algorithm.

Theorem 1.

$$L_{BG} \leq \frac{1}{2\eta_T mR^2} L_{\mathbf{u}} + \frac{1}{\eta_T} D_{Burg}(\mathbf{u}, \mathbf{w}_0),$$

where L_{BG} is the loss incurred by the Burg Gradient descent algorithm and $L_{\mathbf{u}}$ is the loss incurred by the optimal offline algorithm.

Proof. By Lemma 1, for each trial t ,

$$\eta_t(\hat{y}_t - y_t)^2 - \frac{\eta_t}{1-2\eta_t mR^2}(\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 \leq D_{Burg}(\mathbf{u}, \mathbf{w}_t) - D_{Burg}(\mathbf{u}, \mathbf{w}_{t+1}).$$

Hence, adding over all trials we get

$$\sum_{t=1}^T \eta_t(\hat{y}_t - y_t)^2 \leq \sum_{t=1}^T \frac{\eta_t}{1-2\eta_t mR^2}(\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + D_{Burg}(\mathbf{u}, \mathbf{w}_0) - D_{Burg}(\mathbf{u}, \mathbf{w}_T)$$

Thus, since $\eta_t \leq \eta_T$ for all t , we have:

$$L_{BG} \leq \frac{1}{2\eta_T mR^2} L_{\mathbf{u}} + \frac{1}{\eta_T} D_{Burg}(\mathbf{u}, \mathbf{w}_0).$$

\square

Note that similar bounds can be obtained using the framework presented in Gordon[5], but the resulting bound is substantially weaker than our bound and the proof is much more involved.

3 The Matrix Case

In the matrix case, at the t -th step, the algorithm receives an instance X_t ($X_t \in \mathfrak{R}_+^{m \times m}$) and predicts $\hat{y}_t = \text{tr}(W_t X_t)$. W_t is the current model used by the algorithm. The error incurred at the t -th step is $l_t(W_t) = (y_t - \hat{y}_t)^2$.

Similarly to the vector case, for a given T -trial sequence $S = \{(X_1, y_1), (X_2, y_2), \dots, (X_T, y_T)\}$, the optimal offline solution is given by

$$U = \underset{W}{\operatorname{argmin}} \sum_{t=1}^T l_t(W) \\ \text{s.t. } W \succ 0.$$

As a result, the goal here is to compare the loss of the online algorithm with the optimal offline solution, just as in the vector case. For the matrix case, we consider the LogDet divergence as the regularization in the objective. The LogDet divergence is the generalization of D_{Burg} to matrices, given by

$$D_{ld}(W, W_t) = \text{tr}(W W_t^{-1}) - \log \det(W W_t^{-1}).$$

3.1 Rank-One Input

We consider the special case where each input X_t is a rank-1 matrix. This assumption holds in the case of information-theoretic metric learning (discussed later), where the inputs are rank-one constraints. We further assume that input is symmetric; thus, the input matrix X_t can be written as

$$X_t = \mathbf{z}_t \mathbf{z}_t^T.$$

Using LogDet as the regularization term, we derive Algorithm 2 analogously as in the vector case. We set η_t carefully at each step so that the algorithm ensures positive definiteness of the weight matrix W_t . Lemma 2 proves this invariant.

Algorithm 2 Matrix Burg Gradient (BG) Algorithm

- 1: **Initialize:** Set $\eta_t = \frac{1}{4R^2}$ and $W_0 = \frac{1}{N}I, \forall i$.
- 2: **Prediction:** Upon receiving the t -th instance X_t , give the prediction $\hat{y}_t = \text{tr}(W_t X_t)$.
- 3: **Update:** Upon receiving the t -th outcome y_t , update the weights according to the rule

$$W_{t+1} = (W_t^{-1} + 2\eta_t(\hat{y}_t - y_t)X_t)^{-1}, \quad (6)$$

where

$$\eta_t = \begin{cases} n_{t-1} & \text{if } \hat{y}_t - y_t > 0 \\ \min\left(n_{t-1}, \frac{1}{2(y_t - \hat{y}_t)} \left(\frac{1}{\text{tr}((I + (W_t^{-1} - I)^{-1})W_t X_t)} \right)\right) & \text{if } \hat{y}_t - y_t < 0. \end{cases} \quad (7)$$

Lemma 2. *Algorithm 2 preserves the invariant $0 \prec W_t \prec I$.*

Proof. We prove this by induction over t . The base case holds trivially for $N > 1$. By the induction hypothesis, $0 \prec W_t \prec I$. Assume $X_t = \mathbf{z}\mathbf{z}^T$ as X_t is a symmetric positive definite rank-one matrix. Thus, using the update rule and the Sherman Morrison-Woodbury formula[4]:

$$W_{t+1} = W_t - \frac{2\eta_t(\hat{y}_t - y_t)W_t \mathbf{z}\mathbf{z}^T W_t}{1 + 2\eta_t(\hat{y}_t - y_t)\mathbf{z}^T W_t \mathbf{z}}. \quad (8)$$

We break up the proof into the following two cases:

Case 1. $\hat{y}_t - y_t \geq 0$: Consider the update $W_{t+1} = (W_t^{-1} + 2\eta_t(\hat{y}_t - y_t)X_t)^{-1}$. Since $\eta_t > 0$ and $\hat{y}_t - y_t > 0$, then $2\eta_t(\hat{y}_t - y_t)X_t \succ 0$. Now, W_{t+1}^{-1} is sum of two symmetric positive definite matrices, and so W_{t+1} is positive definite.

Using Equation 8,

$$I - W_{t+1} = I - W_t + \frac{2\eta_t(\hat{y}_t - y_t)W_t \mathbf{z}\mathbf{z}^T W_t}{1 + 2\eta_t(\hat{y}_t - y_t)\mathbf{z}^T W_t \mathbf{z}}.$$

Since, W_t is symmetric,

$$W_t \mathbf{z}\mathbf{z}^T W_t = W_t \mathbf{z}\mathbf{z}^T W_t^T = \mathbf{v}\mathbf{v}^T \succ 0,$$

where $\mathbf{v} = W_t \mathbf{z}$. Also, $\eta_t > 0$ and $\hat{y}_t - y_t > 0$, and thus $\frac{2\eta_t(\hat{y}_t - y_t)W_t \mathbf{z}\mathbf{z}^T W_t}{1 + 2\eta_t(\hat{y}_t - y_t)\mathbf{z}^T W_t \mathbf{z}} \succ 0$. Furthermore, $W_t \prec I$, and so $I - W_t \succ 0$. Thus, $I - W_{t+1}$ is sum of two symmetric positive definite matrices, implying $I - W_{t+1} \succ 0$. We conclude that $W_{t+1} \prec I$.

Case 2. $\hat{y}_t - y_t < 0$: By Equation 8,

$$\begin{aligned} W_{t+1} &= W_t - \frac{2\eta_t(\hat{y}_t - y_t)W_t \mathbf{z}\mathbf{z}^T W_t}{1 + 2\eta_t(\hat{y}_t - y_t)\mathbf{z}^T W_t \mathbf{z}} \\ &= W_t + \frac{2\eta_t(y_t - \hat{y}_t)W_t \mathbf{z}\mathbf{z}^T W_t}{1 - 2\eta_t(y_t - \hat{y}_t)\mathbf{z}^T W_t \mathbf{z}}. \end{aligned}$$

The update rule guarantees that $\eta_t < \frac{1}{2(y_t - \hat{y}_t)\text{tr}(W_t X_t)}$. Since

$$\text{tr}(W_t X_t) = \text{tr}(W_t z z^T) = \text{tr}(z^T W_t z) = z^T W_t z, \quad (9)$$

we know that $1 - 2\eta_t(y_t - \hat{y}_t)z^T W_t z > 0$, implying that W_{t+1} is sum of two positive definite matrices hence $W_{t+1} \succ 0$. Write W_t as its eigendecomposition $W_t = Q\Sigma Q^T$. Since W_t is positive definite matrix, Σ is a diagonal matrix with positive diagonal. Using this decomposition, the update for W_{t+1} is expressed as

$$W_{t+1} = Q\Sigma Q^T - \frac{2\eta_t(\hat{y}_t - y)Q\Sigma Q^T z z^T Q\Sigma Q^T}{1 + 2\eta_t(\hat{y}_t - y)z^T Q\Sigma Q^T z}. \quad (10)$$

Since $W_t \prec I$, then $I - \Sigma \succ 0$. Using Equation 10, we have that

$$\begin{aligned} I - W_{t+1} &= I - Q\Sigma Q^T - \frac{2\eta_t(y_t - \hat{y}_t)Q\Sigma Q^T z z^T Q\Sigma Q^T}{1 + 2\eta_t(\hat{y}_t - y_t)z^T Q\Sigma Q^T z} \\ &= Q \left(I - \Sigma - \frac{2\eta_t(y_t - \hat{y}_t)\Sigma Q^T z z^T Q\Sigma}{1 + 2\eta_t(\hat{y}_t - y_t)z^T Q\Sigma Q^T z} \right) Q^T \\ &= Q(I - \Sigma)^{\frac{1}{2}} \left(I - \frac{2\eta_t(y_t - \hat{y}_t)(I - \Sigma)^{-\frac{1}{2}}\Sigma Q^T z z^T Q\Sigma(I - \Sigma)^{-\frac{1}{2}}}{1 + 2\eta_t(\hat{y}_t - y_t)z^T Q\Sigma Q^T z} \right) (I - \Sigma)^{\frac{1}{2}} Q^T. \end{aligned}$$

Let $G = Q(I - \Sigma)^{\frac{1}{2}}$. If the middle term is symmetric positive definite, then it can be written as EE^T . This would lead to $I - W_{t+1} = GEE^T G^T = GE(GE)^T \succ 0$, since for any A , $AA^T \succ 0$.

Thus to prove that $W_{t+1} \prec I$, we need to prove that

$$\left(I - \frac{2\eta_t(y_t - \hat{y}_t)(I - \Sigma)^{-\frac{1}{2}}\Sigma Q^T z z^T Q\Sigma(I - \Sigma)^{-\frac{1}{2}}}{1 + 2\eta_t(\hat{y}_t - y_t)z^T Q\Sigma Q^T z} \right) \succ 0. \quad (11)$$

To do this, we calculate the eigenvalues of this matrix, and show that they are greater than zero. Let $v = (I - \Sigma)^{-\frac{1}{2}}\Sigma Q^T z$. Then $(I - \Sigma)^{-\frac{1}{2}}\Sigma Q^T z z^T Q\Sigma(I - \Sigma)^{-\frac{1}{2}} = vv^T$. Since, vv^T is a rank-one matrix, it has exactly one positive eigenvalue, and that eigenvalue is equal to its trace $\text{tr}((I - \Sigma)^{-\frac{1}{2}}\Sigma Q^T z z^T Q\Sigma(I - \Sigma)^{-\frac{1}{2}})$. Expanding the trace yields

$$\begin{aligned} \text{tr}((I - \Sigma)^{-\frac{1}{2}}\Sigma Q^T z z^T Q\Sigma(I - \Sigma)^{-\frac{1}{2}}) &= \text{tr}(Q\Sigma(I - \Sigma)^{-1}\Sigma Q^T z z^T) \\ &= \text{tr}(Q(\Sigma^{-1} - I)^{-1}\Sigma Q^T z z^T) \\ &= \text{tr}((W_t^{-1} - I)^{-1}W_t X_t). \end{aligned}$$

To prove Equation 11, we need to show that the eigenvalues are positive; equivalently, we can show that

$$1 - \frac{2\eta_t(y_t - \hat{y}_t)\text{tr}((W_t^{-1} - I)^{-1}W_t X_t)}{1 - 2\eta_t(y - \hat{y}_t)\text{tr}(W_t X_t)} > 0. \quad (12)$$

By the update rule, $\eta_t < \frac{1}{2(y_t - \hat{y}_t)\text{tr}(W_t X_t)}$, and hence the denominator of the second term is positive. Also, $\eta_t < \frac{1}{2(y_t - \hat{y}_t)} \left(\frac{1}{\text{tr}((I + (W_t^{-1} - I)^{-1})W_t X_t)} \right)$, so Equation 12 holds. Thus $W_{t+1} \prec I$. \square

3.2 Regret Bound

Similar to the vector case, we bound loss at each step incurred by algorithm in terms of loss incurred by the optimal offline solution. Assume $\|z\|_2 \leq R$.

Lemma 3.

$$a_t(\hat{y}_t - y_t)^2 - b_t(\text{tr}(U X_t) - y_t)^2 \leq D_{id}(U, W_t) - D_{id}(U, W_{t+1}),$$

where a_t and b_t are positive constants and U is the optimal offline solution.

Proof.

$$D_{ld}(U, W_t) - D_{ld}(U, W_{t+1}) = \log\left(\frac{\det(W_t)}{\det(W_{t+1})}\right) + \text{tr}((W_t^{-1} - W_{t+1}^{-1})U).$$

Since $\log(\det(A)) = \text{tr}(\log(A))$ and $\text{tr}(AB) = \text{tr}(BA)$, we have that

$$\begin{aligned} D_{ld}(U, W_t) - D_{ld}(U, W_{t+1}) &= \text{tr}(\log(W_t) - \log(W_{t+1})) + \text{tr}((W_t^{-1} - W_{t+1}^{-1})U) \\ &= \text{tr}(\log(W_t W_{t+1}^{-1})) - 2\eta_t(\hat{y}_t - y_t)\text{tr}(X_t U) \\ &= \text{tr}(\log(I + 2\eta_t(\hat{y}_t - y_t)W_t X_t)) - 2\eta_t(\hat{y}_t - y_t)\text{tr}(U X_t). \end{aligned}$$

Using the Taylor series expansion and ignoring lower order terms,

$$D_{ld}(U, W_t) - D_{ld}(U, W_{t+1}) \geq 2\eta_t(\hat{y}_t - y_t)\text{tr}(W_t X_t) - 2\eta_t^2(\hat{y}_t - y_t)^2\text{tr}((W_t X_t)^2) - 2\eta_t(\hat{y}_t - y_t)\text{tr}(U X_t).$$

We write X as $\mathbf{z}\mathbf{z}^T$, so $\text{tr}((W_t X_t)^2) = \text{tr}(W_t \mathbf{z}\mathbf{z}^T W_t \mathbf{z}\mathbf{z}^T) = (\mathbf{z}^T W_t \mathbf{z})\text{tr}(W \mathbf{z}\mathbf{z}^T) = \text{tr}(W_t X_t)^2$. Hence,

$$\begin{aligned} D_{ld}(U, W_t) - D_{ld}(U, W_{t+1}) &\geq 2\eta_t(\hat{y}_t - y_t)\text{tr}(W_t X_t) - 2\eta_t^2(\hat{y}_t - y_t)^2\text{tr}((W_t X_t)^2) - 2\eta_t(\hat{y}_t - y_t)\text{tr}(U X_t) \\ &\geq 2\eta_t(\hat{y}_t - y_t)\text{tr}(W_t X_t) - 2\eta_t^2(\hat{y}_t - y_t)^2\text{tr}(W_t X_t)^2 - 2\eta_t(\hat{y}_t - y_t)\text{tr}(U X_t). \end{aligned}$$

Since $W_t \prec I$ by Lemma 2, $\text{tr}(W_t X_t) = \mathbf{z}^T W_t \mathbf{z} \leq \mathbf{z}^T \mathbf{z} \leq R^2$. As a result,

$$\begin{aligned} D_{ld}(U, W_t) - D_{ld}(U, W_{t+1}) &\geq 2\eta_t(\hat{y}_t - y_t)\text{tr}(W_t X_t) - 2\eta_t^2(\hat{y}_t - y_t)^2 R^2 - 2\eta_t(\hat{y}_t - y_t)\text{tr}(U X_t) \\ &= 2\eta_t(\hat{y}_t - y_t)\hat{y}_t - 2\eta_t^2(\hat{y}_t - y_t)^2 R^2 - 2\eta_t(\hat{y}_t - y_t)r, \end{aligned}$$

where $r = \text{tr}(U X_t)$. Thus, to prove Lemma 3, we need to prove the same inequality as in Equation 5, with mR^2 replaced by R^2 . Hence, for $0 \leq a_t \leq \eta_t$ and $b_t = \frac{\eta_t}{1 - 2\eta_t R^2}$, Lemma 3 holds. \square

As in the vector case, we can obtain the final bound by summing the loss at each step.

Theorem 2.

$$L_{BG} \leq \frac{1}{2\eta_T R^2} L_U + \frac{1}{\eta_T} H_{ld}(U),$$

where L_{MBG} is the loss incurred by the Matrix Burg Gradient descent algorithm, L_U is loss incurred by optimal offline algorithm and $H_{ld}(U) = -\log \det(U)$.

Proof. The proof is very similar to that of the vector case. \square

4 Application: Online Metric Learning

Selecting an appropriate distance measure (or metric) is fundamental to many learning algorithms such as k -means, nearest neighbor searches, and others. However, choosing such a measure is highly problem-specific and ultimately dictates the success—or failure—of the learning algorithm. To this end, there have been several recent approaches that attempt to learn distance functions. These methods work by exploiting distance information that is intrinsically available in many learning settings. For example, in the problem of semi-supervised clustering, points are constrained to be either similar (i.e. the distance between them should be relatively small) or dissimilar (the distance should be larger). In information retrieval settings, constraints between pairs of distances can be gathered from click-through feedback. In fully supervised settings, constraints can be inferred so that points in the same class have smaller distances to each other than to points in different classes.

Given a set of n points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ in \mathbb{R}^d , we seek a positive definite matrix A which parameterizes the Mahalanobis distance¹:

$$d_A(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T A (\mathbf{x}_i - \mathbf{x}_j). \quad (13)$$

¹The term ‘Mahalanobis distance’ frequently refers to the squared Mahalanobis distance, even though it is $\sqrt{d_A(\mathbf{x}, \mathbf{y})}$ that is a metric. This difference does not affect our method, hence we drop the qualification *squared* when referring to (13)

We can quantify the measure of “closeness” between two Mahalanobis matrices A and A_0 via a natural information-theoretic approach. There exists a simple bijection (up to a scaling function) between the set of Mahalanobis distances and the set of equal-mean multivariate Gaussian distributions (without loss of generality, we can assume the Gaussians have zero-mean $\boldsymbol{\mu}$). Given a Mahalanobis distance parameterized by A , we express its corresponding multivariate Gaussian as $p(\mathbf{x}; A) = \frac{1}{Z} \exp(-\frac{1}{2}d_A(\mathbf{x}, \boldsymbol{\mu}))$, where Z is a normalizing constant. Using this bijection, we measure the distance between two Mahalanobis distance functions parameterized by A and A_0 by the (differential) relative entropy between their corresponding multivariate Gaussians:

$$\text{KL}(p(\mathbf{x}; A) \| p(\mathbf{x}; A_0)) = \int p(\mathbf{x}; A) \log \frac{p(\mathbf{x}; A)}{p(\mathbf{x}; A_0)} d\mathbf{x}. \quad (14)$$

The distance (14) provides a well-founded measure of “closeness” between two Mahalanobis distance functions and forms the basis of information-theoretic metric learning. It is known that the differential relative entropy between two multivariate Gaussians can be expressed as the convex combination of a Mahalanobis distance between mean vectors and the LogDet divergence between the covariance matrices. Assuming the means to be the same, we have

$$\text{KL}(p(\mathbf{x}; A) \| p(\mathbf{x}; A_0)) = \frac{1}{2} D_{ld}(A, A_0), \quad (15)$$

We therefore select $d(A, A_t) = D_{ld}(A, A_t)$ as the regularizer for (1), using the above equivalence 15 to obtain the update

$$A_{t+1} = \underset{A}{\text{argmin}} D_{ld}(A, A_t) + \eta_t (d_t - \hat{d}_t)^2.$$

We can view the problem of online metric learning in the following manner: at each time step, we receive a matrix X_t and a target y_t . These encode the distance between two vectors \mathbf{x}_i and \mathbf{x}_j . If we let $\mathbf{z} = \mathbf{x}_i - \mathbf{x}_j$, then $\text{tr}(A_t X_t) = \mathbf{z}^T A_t \mathbf{z} = (\mathbf{x}_i - \mathbf{x}_j)^T A_t (\mathbf{x}_i - \mathbf{x}_j) = d_{A_t}(\mathbf{x}_i, \mathbf{x}_j)$. Thus, the prediction $\text{tr}(A_t X_t)$ captures the Mahalanobis distance between the vectors \mathbf{x}_i and \mathbf{x}_j using the current Mahalanobis matrix A_t . The resulting loss $(\text{tr}(A_t X_t) - y_t)^2$, where y_t is the target distance between the two vectors, tells us how close the current Mahalanobis matrix A_t is in terms of satisfying the given distance constraint.

As a result, the rank-one online regression algorithm using the LogDet divergence can be used as method for learning a Mahalanobis matrix A , where at every time step two points, along with their target distance, are presented to the algorithm. The regret bound proven in the previous section applies directly to this case.

See [2, 3] for further details on information-theoretic metric learning, including the analysis of an offline algorithm.

5 Experimental Results

In section 2.1 and 3.2, we showed worst case bounds for our algorithms. To simulate the online learning problem, we now generate simple artificial data. First, we generate a random vector \mathbf{u} drawn from S^{m-1} . Then, we generate some random inputs \mathbf{x}_t , also drawn from S^{m-1} . We output y as $\mathbf{u} \cdot \mathbf{x}_t$ mixed with some noise, i.e.,

$$y = (\mathbf{u} \cdot \mathbf{x}_t)(1 + \nu Z),$$

where Z is an uniformly distributed random number over $[-0.5, 0.5]$ and ν is the noise level.

We generate this artificial data for $N = 100$ and $T = 10000$ iterations. Figure 1 compares performance of the Burg entropy-based algorithm with those using relative entropy and squared Euclidean distance as the regularization term, when the noise level is set to 0. Clearly, the Burg algorithm converges slower compared to the other methods. But, Figures 2 and 3 show that as noise increases, the Burg performance improves in comparison to both relative entropy and squared Euclidean distance. Figure 4 shows that when noise is around 30% then the Burg algorithm performs comparatively to relative entropy, and squared Euclidean performs poorly.

6 Conclusion

In this paper, we introduced online regression based on using a Burg entropy-based regularization, leading to new online regret bounds. An extension to the rank-one matrix case was also analyzed. The main application

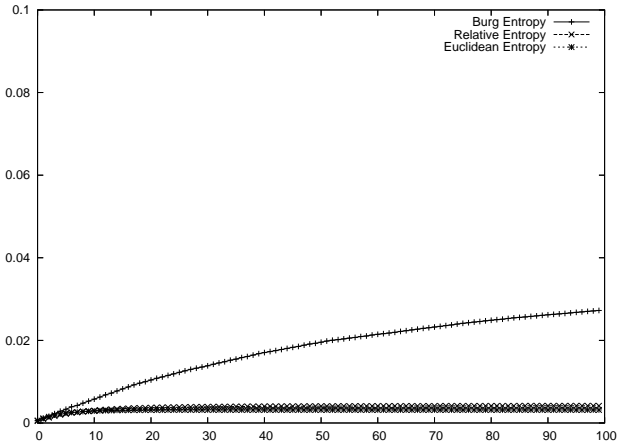


Figure 1: Comparison of performance of burg, relative and euclidean distance measure with noise=0%

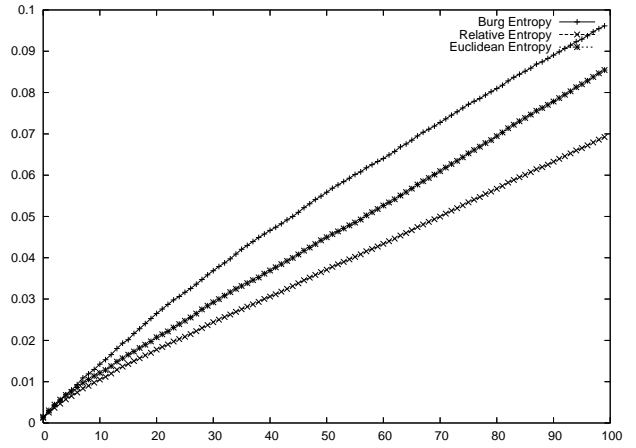


Figure 2: Comparison of performance of burg, relative and euclidean distance measure with noise=10%

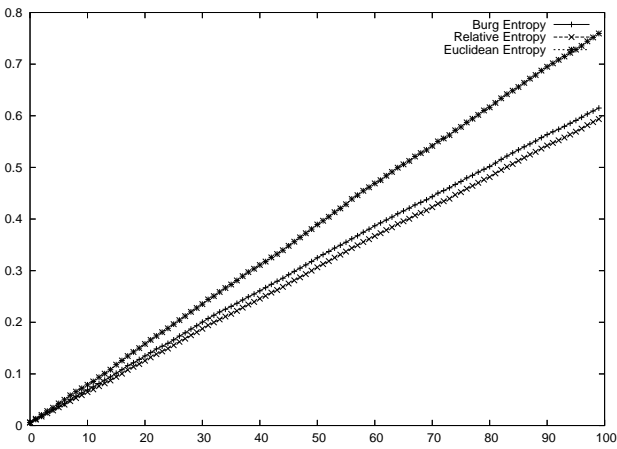


Figure 3: Comparison of performance of burg, relative and euclidean distance measure with noise=20%

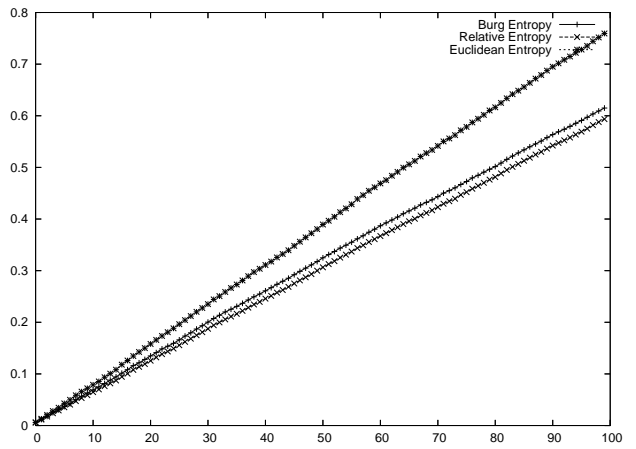


Figure 4: Comparison of performance of burg, relative and euclidean distance measure with noise=30%

is to online metric learning, which can be cast as an online regression problem for rank-one matrices; this metric learning formulation has an information-theoretic motivation, and results in [2] have shown that it outperforms existing methods. Future work includes improving the bounds so that they do not depend on the learning rate, and extending the matrix analysis to the general case.

References

- [1] A. BANERJEE, S. MERUGU, I. S. DHILLON, AND J. GHOSH, *Clustering with bregman divergences.*, in SDM, 2004, pp. 234–245.
- [2] J. DAVIS, B. KULIS, P. JAIN, S. SRA, AND I. S. DHILLON, *Information-theoretic metric learning*, in International Conference on Metric Learning, 2007, p. to appear.
- [3] J. DAVIS, B. KULIS, S. SRA, AND I. DHILLON, *Information-theoretic metric learning*, in NIPS 2006 Workshop on Learning to Compare Examples, 2006.
- [4] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Johns Hopkins Series in the Mathematical Sciences, Johns Hopkins Univ. Press, second ed., 1989.
- [5] G. J. GORDON, *Regret bounds for prediction problems*, in COLT, 1999, pp. 29–40.
- [6] J. KIVINEN AND M. K. WARMUTH, *Exponentiated gradient versus gradient descent for linear predictors*, Inf. Comput., 132 (1997), pp. 1–63.
- [7] B. KULIS, M. SUSTIK, AND I. S. DHILLON, *Learning low-rank kernel matrices.*, in ICML, 2006, pp. 505–512.
- [8] K. TSUDA, G. RAETSCH, AND M. K. WARMUTH, *Matrix exponentiated gradient updates of online learning and bregman projection*, Journal of Machine Learning Research, 6 (2005).