
Mistake Bounds for Binary Matrix Completion

Mark Herbster
University College London
Department of Computer Science
London WC1E 6BT, UK
m.herbster@cs.ucl.ac.uk

Stephen Pasteris
University College London
Department of Computer Science
London WC1E 6BT, UK
s.pasteris@cs.ucl.ac.uk

Massimiliano Pontil
Istituto Italiano di Tecnologia
16163 Genoa, Italy
and
University College London
Department of Computer Science
London WC1E 6BT, UK
m.pontil@cs.ucl.ac.uk

Abstract

We study the problem of completing a binary matrix in an online learning setting. On each trial we predict a matrix entry and then receive the true entry. We propose a Matrix Exponentiated Gradient algorithm [1] to solve this problem. We provide a mistake bound for the algorithm, which scales with the *margin complexity* [2, 3] of the underlying matrix. The bound suggests an interpretation where each row of the matrix is a prediction task over a finite set of objects, the columns. Using this we show that the algorithm makes a number of mistakes which is comparable up to a logarithmic factor to the number of mistakes made by the Kernel Perceptron with an optimal kernel in hindsight. We discuss applications of the algorithm to predicting as well as the best biclustering and to the problem of predicting the labeling of a graph without knowing the graph in advance.

1 Introduction

We consider the problem of predicting *online* the entries in an $m \times n$ binary matrix U . We formulate this as the following game: *nature* queries an entry (i_1, j_1) ; the *learner* predicts $\hat{y}_1 \in \{-1, 1\}$ as the matrix entry; *nature* presents a label $y_1 = U_{i_1, j_1}$; *nature* queries the entry (i_2, j_2) ; the *learner* predicts \hat{y}_2 ; and so forth. The learner's goal is to minimize the total number of mistakes $M = |\{t : \hat{y}_t \neq y_t\}|$. If nature is adversarial, the learner will always mispredict, but if nature is regular or simple, there is hope that a learner may make only a few mispredictions.

In our setting we are motivated by the following interpretation of matrix completion. Each of the m rows represents a task (or binary classifier) and each of the n columns is associated with an object (or input). A task is the problem of predicting the binary label of each of the objects. For a single task, if we were given a kernel matrix between the objects we could then use the Kernel Perceptron algorithm to sequentially label the objects and this algorithm would incur $\mathcal{O}(1/\lambda^2)$ mistakes, where λ is the margin of the best linear classifier in the inner product space induced by the kernel. Unfortunately, in our setup, we do not know a good kernel in advance. However, we will show that a remarkable property of our algorithm is that it enjoys, up to logarithmic factors, a mistake bound of $\mathcal{O}(1/\gamma^2)$ per task, where γ is the largest possible margin (over the choice of the kernel) which is achieved on all tasks.

The problem of predicting online the labels of a finite set of objects under the assumption that the similarity between objects can be described by a graph was introduced in [4], building upon earlier work in the batch setting [5, 6]. In this and later research the common assumption is that two objects are similar if there is an edge in the graph connecting them and the aim is to predict well when there are few edges between objects with disagreeing labels. Lower bounds and an optimal algorithm (up to logarithmic factors) for this problem were given in [7, 8]. The problem of predicting well when the graph is unknown was previously addressed in [9, 10]. That research took the approach that when receiving a vertex to predict, edges local to that vertex were then revealed. In this paper we take a different approach - the graph structure is never revealed to the learner. Instead, we have a number of tasks over the same unknown graph, and the hope is to perform comparably to the case in which the graph is known in advance.

The general problem of matrix completion has been studied extensively in the batch statistical i.i.d. setting, see for example [11, 12, 13] and references therein. These studies are concerned either with Rademacher bounds or statistical oracle inequalities, both of which are substantially different from the focus of the present paper. In the online mistake-bound setting a special form of matrix completion was previously considered as the problem of learning a binary relation [14, 15] (see Section 5). In a more general online setting, with minimal assumptions on the loss function [16, 17] bounded the regret of the learner in terms of the *trace-norm* of the underlying matrix. Instead our bounds are with respect to the *margin complexity* of the matrix. As a result, although our bounds have a more restricted applicability they have the advantage that they become non-trivial after only $\Theta(n)$ matrix entries¹ are observed as opposed to the required $\Theta(n^{3/2})$ in [16] and $\Theta(n^{7/4})$ in [17]. The notion of margin complexity in machine learning was introduced in [2] where it was used to study the learnability of concept classes via linear embeddings and further studied in [3], where it was linked to the γ_2 norm. Here we adopt the terminology in [11] and refer to the γ_2 norm as the max-norm. The margin complexity seems to be a more natural parameter as opposed to the trace-norm for the 0-1 loss as it only depends on the signs of the underlying comparator matrix. To the best of our knowledge the bounds contained herein are the first online matrix completion bounds in terms of the margin complexity.

To obtain our results, we use an online matrix multiplicative weights algorithm, e.g., see [1, 18, 17, 19] and references therein. These kinds of algorithms have been applied in a number of learning scenarios, including online PCA [20], online variance minimization [21], solving SDPs [18], and online prediction with switching sequences [22]. These algorithms update a new hypothesis matrix on each trial by trading off fidelity to the previous hypothesis and the incorporation of the new label information. The tradeoff is computed as an approximate spectral regularization via the quantum relative entropy (see [1, Section 3.1]). The particular matrix multiplicative weights algorithm we apply is Matrix Winnow [19]; we adapt this algorithm and its mistake bound analysis for our purposes via selection of comparator, threshold, and appropriate “progress inequalities.”

The paper is organized as follows. In Section 2 we introduce basic notions used in the paper. In Section 3 we present our algorithm and derive a mistake bound, also comparing it to related bounds in the literature. In Section 4 we observe that our algorithm is able to exploit matrix structure to perform comparably to the Kernel Perceptron with the best kernel known in advance. Finally, in Section 5 we discuss the example of biclustered matrices, and argue that our bound is optimal up to a polylogarithmic factor. The appendix contains proofs of the results only stated in the main body of the paper, and other auxiliary results.

2 Preliminaries

We denote the set of the first m positive integers as $\mathbb{N}_m = \{1, \dots, m\}$. We denote the inner product of vectors $\mathbf{x}, \mathbf{w} \in \mathbb{R}^n$ as $\langle \mathbf{x}, \mathbf{w} \rangle = \sum_{i=1}^n x_i w_i$ and the norm as $|\mathbf{w}| = \sqrt{\langle \mathbf{w}, \mathbf{w} \rangle}$. We let $\mathbb{R}^{m \times n}$ be the set of all $m \times n$ real-valued matrices. If $\mathbf{X} \in \mathbb{R}^{m \times n}$ then \mathbf{X}_i denotes the i -th n -dimensional row vector and the (i, j) entry in \mathbf{X} is X_{ij} . The trace of a square matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$ is $\text{Tr}(\mathbf{X}) = \sum_{i=1}^n X_{ii}$. The trace norm of a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$ is $\|\mathbf{X}\|_1 = \text{Tr}(\sqrt{\mathbf{X}^\top \mathbf{X}})$, where $\sqrt{\cdot}$ indicates the unique positive square root of a positive semi-definite matrix. For every matrix $\mathbf{U} \in \{-1, 1\}^{m \times n}$, we define $\text{SP}(\mathbf{U}) = \{\mathbf{V} \in \mathbb{R}^{m \times n} : \forall_{ij} V_{ij} U_{ij} > 0\}$, the set of matrices which

¹For simplicity we assume $m \in \Theta(n)$.

are sign consistent with U . We also define $\text{SP}^1(U) = \{V \in \mathbb{R}^{m \times n} : \forall_{ij} V_{ij} U_{ij} \geq 1\}$, that is the set of matrices which are sign consistent to U with a margin of at least one.

The max-norm (or γ_2 norm [3]) of a matrix $U \in \mathbb{R}^{m \times n}$ is defined by the formula

$$\|U\|_{\max} := \inf_{PQ^\top = U} \left\{ \max_{1 \leq i \leq m} |P_i| \max_{1 \leq j \leq n} |Q_j| \right\}, \quad (1)$$

where the infimum is over all matrices $P \in \mathbb{R}^{m \times k}$ and $Q \in \mathbb{R}^{n \times k}$ and every integer k . The *margin complexity* of a matrix $U \in \mathbb{R}^{m \times n}$ is

$$\text{mc}(U) := \inf_{PQ^\top \in \text{SP}(U)} \max_{ij} \frac{|P_i||Q_j|}{|\langle P_i, Q_j \rangle|}.$$

This quantity plays a central role in the analysis of our algorithm. If we interpret the rows of U as m different binary classification tasks, and the columns as a finite set of objects which we wish to label, the “min-max” margin with respect to an embedding is smallest of the m maximal margins over the tasks. The quantity $1/\text{mc}(U)$ is then the maximum “min-max” margin with respect to all possible embeddings. Specifically, the rows of matrix P represent the “weights” of the binary classifiers and the rows of matrix Q the “input vectors” associated with the objects. The quantity $\frac{|\langle P_i, Q_j \rangle|}{|P_i||Q_j|}$ is the margin of the i -th classifier on the j -th input. Observe that margin complexity depends only on the sign pattern of the matrix and not the magnitudes. The margin complexity is equivalently $\text{mc}(U) = \min_{V \in \text{SP}^1(U)} \|V\|_{\max}$, see e.g., [3, Lemma 3.1].

In our online setting we are concerned with predicting an (*example*) sequence $((i_1, j_1), y_1), \dots, ((i_T, j_T), y_T) \in (\mathbb{N}_m \times \mathbb{N}_n) \times \{-1, 1\}$. A sequence must be *consistent*, that is, given examples $((i, j), y)$ and $((i', j'), y')$ if $(i, j) = (i', j')$ then $y = y'$. We define the set of sign-consistent matrices with a sequence \mathcal{S} as $\text{cons}(\mathcal{S}) := \{M \in \mathbb{R}^{m \times n} : 0 < y M_{ij}, ((i, j), y) \in \mathcal{S}\}$. We extend the notion of margin complexity to sequences via $\text{mc}(\mathcal{S}) := \inf_{U \in \text{cons}(\mathcal{S})} \text{mc}(U)$.

The number of margin violations in a sequence \mathcal{S} at complexity γ is defined to be,

$$\text{merr}(\mathcal{S}, \gamma) := \inf_{PQ^\top \in \text{cons}(\mathcal{S})} \left| \left\{ ((i, j), y) \in \mathcal{S} : \frac{|P_i||Q_j|}{|\langle P_i, Q_j \rangle|} > \frac{1}{\gamma} \right\} \right|. \quad (2)$$

In particular, note that $\text{merr}(\mathcal{S}, \gamma) = 0$ if $\gamma \leq \frac{1}{\text{mc}(\mathcal{S})}$.

Finally, we introduce the following quantity, which plays a central role in the amortized analysis of our algorithm.

Definition 2.1. *The quantum relative entropy of symmetric positive semidefinite square matrices A and B is*

$$\Delta(A, B) := \text{Tr}(A \log(A) - A \log(B) + B - A).$$

3 Algorithm and Analysis

Algorithm 1 presents an adaptation of the Matrix Exponentiated Gradient algorithm [1, 17, 18, 19] to our setting. This algorithm is a matrix analog of the Winnow algorithm [19]; we refer to the above papers for more insights into this family of algorithms.

The following theorem provides a mistake bound for the algorithm.

Theorem 3.1. *The number of mistakes, M , on sequence \mathcal{S} made by the Algorithm 1 with parameter $0 < \gamma \leq 1$ is upper bounded by*

$$M \leq c \left[(m+n) \log(m+n) \frac{1}{\gamma^2} + \text{merr}(\mathcal{S}, \gamma) \right], \quad (3)$$

where $c = 1/(3 - e) \leq 3.55$ and the quantity $\text{merr}(\mathcal{S}, \gamma)$ is given in equation (2).

Proof. Given $U \in \mathbb{R}^{m \times n}$, let $P \in \mathbb{R}^{m \times k}$ and $Q \in \mathbb{R}^{n \times k}$ be such that $PQ^\top = U$. For every $i \in \mathbb{N}_m$, we denote by P_i the i -th row vector of P and for every $j \in \mathbb{N}_n$, we denote by Q_j the j -th row vector of Q . We construct the $(m+n) \times k$ matrix

$$R := \text{diag} \left(\frac{1}{|P_1|}, \dots, \frac{1}{|P_m|}, \frac{1}{|Q_1|}, \dots, \frac{1}{|Q_n|} \right) \begin{bmatrix} P \\ Q \end{bmatrix}$$

Algorithm 1 Predicting a binary matrix.

Parameters: Learning rate $0 < \gamma \leq 1$.

Initialization: $\mathbf{W}^{(0)} \leftarrow \frac{\mathbf{I}}{(m+n)}$, where \mathbf{I} is the $(m+n) \times (m+n)$ identity matrix.

For $t = 1, \dots, T$

- Get pair $(i_t, j_t) \in \mathbb{N}_m \times \mathbb{N}_n$.
- Define $\mathbf{X}^{(t)} := \frac{1}{2}(\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t})(\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t})^\top$, where \mathbf{e}_k is the k -th basis vector of \mathbb{R}^{m+n} .
- Predict

$$\hat{y}_t = \begin{cases} 1 & \text{if } \text{Tr}(\mathbf{W}^{(t-1)} \mathbf{X}^{(t)}) \geq \frac{1}{m+n}, \\ -1 & \text{otherwise.} \end{cases}$$

- Receive label $y_t \in \{-1, 1\}$ and if $\hat{y}_t \neq y_t$ update

$$\mathbf{W}^{(t)} \leftarrow \exp\left(\log\left(\mathbf{W}^{(t-1)}\right) + \frac{\gamma}{2}(y_t - \hat{y}_t)\mathbf{X}^{(t)}\right).$$

and construct $\tilde{\mathbf{U}} := \left(\frac{1}{m+n}\right)\mathbf{R}\mathbf{R}^\top$. Define matrix $\mathbf{X}^{(t)} := \frac{1}{2}(\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t})(\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t})^\top$, where \mathbf{e}_k is the k -th basis vector of \mathbb{R}^{m+n} .

Note that $\text{Tr}(\mathbf{X}^{(t)}) = 1$, $\text{Tr}(\tilde{\mathbf{U}}) = 1$ (since every row of \mathbf{R} is normalized) and

$$\begin{aligned} \text{Tr}(\tilde{\mathbf{U}}\mathbf{X}^{(t)}) &= \frac{1}{n+m} \text{Tr}((\mathbf{R}\mathbf{R}^\top) \frac{1}{2}(\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t})(\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t})^\top) \\ &= \frac{1}{2(n+m)} (\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t})^\top \mathbf{R}\mathbf{R}^\top (\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t}) \\ &= \frac{1}{2(n+m)} (\mathbf{R}^\top (\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t}))^\top (\mathbf{R}^\top (\mathbf{e}_{i_t} + \mathbf{e}_{m+j_t})) \\ &= \frac{1}{2(n+m)} \left(\frac{\mathbf{P}_{i_t}}{|\mathbf{P}_{i_t}|} + \frac{\mathbf{Q}_{j_t}}{|\mathbf{Q}_{j_t}|} \right) \left(\frac{\mathbf{P}_{i_t}}{|\mathbf{P}_{i_t}|} + \frac{\mathbf{Q}_{j_t}}{|\mathbf{Q}_{j_t}|} \right)^\top \\ &= \frac{1}{(n+m)} \left(1 + \frac{\langle \mathbf{P}_{i_t}, \mathbf{Q}_{j_t} \rangle}{|\mathbf{P}_{i_t}| |\mathbf{Q}_{j_t}|} \right). \end{aligned}$$

For a trial t we say there is a *margin violation* if $\frac{|\mathbf{P}_{i_t}| |\mathbf{Q}_{j_t}|}{|\langle \mathbf{P}_{i_t}, \mathbf{Q}_{j_t} \rangle|} > \frac{1}{\gamma}$. Let M^- denote the number of mistakes made in trials with margin violations and let M^+ denote the number of mistakes made in trials without margin violations.

From Lemma A.3 in the appendix we have

$$\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) \geq \frac{\gamma}{2}(y_t - \hat{y}_t) \text{Tr}(\tilde{\mathbf{U}}\mathbf{X}^{(t)}) + \left(1 - e^{\frac{\gamma}{2}(y_t - \hat{y}_t)}\right) \text{Tr}(\mathbf{W}^{(t-1)}\mathbf{X}^{(t)}),$$

then substituting in the above we have that

$$\begin{aligned} \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) &\geq \frac{\gamma}{2}(y_t - \hat{y}_t) \frac{1}{n+m} \left(1 + \frac{\langle \mathbf{P}_{i_t}, \mathbf{Q}_{j_t} \rangle}{|\mathbf{P}_{i_t}| |\mathbf{Q}_{j_t}|} \right) \\ &\quad + \left(1 - e^{\frac{\gamma}{2}(y_t - \hat{y}_t)}\right) \text{Tr}(\mathbf{W}^{(t-1)}\mathbf{X}^{(t)}). \end{aligned}$$

To further simplify the above we use Lemma A.4 presented in the appendix, which gives

$$\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) \geq \begin{cases} (c' - 1) \frac{1}{n+m} \gamma^2, & \text{if there is a margin violation,} \\ c' \frac{1}{n+m} \gamma^2, & \text{otherwise.} \end{cases}$$

where $c' = 3 - e$.

Using a telescoping sum, this gives

$$\begin{aligned}\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(0)}) &\geq \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(0)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(T)}) \geq M^+ c' \frac{1}{n+m} \gamma^2 + M^- (c' - 1) \frac{1}{n+m} \gamma^2 \\ &= (c' M^+ - (1 - c') M^-) \frac{1}{n+m} \gamma^2\end{aligned}$$

and hence

$$M^+ \leq \frac{1}{c' \frac{1}{n+m} \gamma^2} \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(0)}) + \frac{1 - c'}{c'} M^-.$$

We conclude that

$$M = M^+ + M^- \leq \frac{1}{c' \frac{1}{n+m} \gamma^2} \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(0)}) + \frac{1}{c'} M^-.$$

We also have that

$$\begin{aligned}\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(0)}) &= \text{Tr}(\tilde{\mathbf{U}} \log(\tilde{\mathbf{U}})) - \text{Tr}(\tilde{\mathbf{U}} \log(\mathbf{W}^{(0)})) + \text{Tr}(\mathbf{W}^{(0)}) - \text{Tr}(\tilde{\mathbf{U}}) \\ &= \text{Tr}(\tilde{\mathbf{U}} \log(\tilde{\mathbf{U}})) - \text{Tr}(\tilde{\mathbf{U}} \log(\mathbf{W}^{(0)})) + 1 - 1 \\ &= \text{Tr}(\tilde{\mathbf{U}} \log(\tilde{\mathbf{U}})) - \text{Tr}(\tilde{\mathbf{U}} \log(\mathbf{W}^{(0)})).\end{aligned}$$

Write the eigen-decomposition of $\tilde{\mathbf{U}}$ as $\sum_{i=1}^{m+n} \lambda_i \alpha_i \alpha_i^T$. Now we have $\sum_{i=1}^{m+n} \lambda_i = \text{Tr}(\tilde{\mathbf{U}}) = 1$ so all eigenvalues λ_i are in the range $[0, 1]$ meaning $\log(\lambda_i) \leq 0$ so $\lambda_i \log(\lambda_i) < 0$ which are the eigenvalues of $\tilde{\mathbf{U}} \log(\tilde{\mathbf{U}})$ meaning that $\text{Tr}(\tilde{\mathbf{U}} \log(\tilde{\mathbf{U}})) \leq 0$. Also, $\log(\mathbf{W}^{(0)}) = \log(\frac{1}{n+m}) \mathbf{I}$ so $\tilde{\mathbf{U}} \log(\mathbf{W}^{(0)}) = \log(\frac{1}{n+m}) \tilde{\mathbf{U}}$ and hence $-\text{Tr}(\tilde{\mathbf{U}} \log(\mathbf{W}^{(0)})) = -\log(\frac{1}{n+m}) \text{Tr}(\tilde{\mathbf{U}}) = \log(m+n)$. So by the above we have

$$\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(0)}) \leq \log(m+n)$$

and hence putting together we get

$$M \leq \frac{m+n}{c' \gamma^2} \log(m+n) + \frac{1}{c'} M^-.$$

□

Observe that in the simplifying case when we have no margin errors ($\text{merr}(\mathcal{S}, \gamma) = 0$) and the learning rate is $\gamma := \frac{1}{\text{mc}(\mathcal{S})}$ we have that the number of mistakes of Algorithm 1 is bounded by $\tilde{\mathcal{O}}((n+m) \text{mc}^2(\mathcal{S}))$. More generally although the learning rate is fixed in advance, we may use a “doubling trick” to avoid the need to tune the γ .

Corollary 3.2. *For any value of γ^* the number of mistakes M made by the following algorithm:*

DOUBLING ALGORITHM:

Set $\kappa \leftarrow \sqrt{2}$ and loop over

1. Run Algorithm 1 with $\gamma = \frac{1}{\kappa}$ until it has made $\lceil 2c(m+n) \log(m+n) \kappa^2 \rceil$ mistakes
 2. Set $\kappa \leftarrow \kappa \sqrt{2}$
-

is upper bounded by

$$M \leq 12c \left[(m+n) \log(m+n) \frac{1}{(\gamma^*)^2} + \text{merr}(\mathcal{S}, \gamma^*) \right],$$

with $c = 1/(3 - e) \approx 3.55$.

See the appendix for a proof. We now compare our bound to other online learning algorithms for matrix completion. The algorithms of [16, 17] address matrix completion in a significantly more general setting. Both algorithms operate with weak assumptions on the loss function, while our algorithm is restricted to the 0–1 loss (mistake counting). Those papers present regret bounds, whereas we apply the stronger assumption that there exists a consistent predictor. As a regret bound is not possible for a deterministic predictor with the 0–1 loss, we compare Theorem 3.1 to their

bound when their algorithm is allowed to predict $\hat{y} \in [-1, 1]$ and uses absolute loss. For clarity in our discussion we will assume that $m \in \Theta(n)$.

Under the above assumptions, the regret bound in [17, Corollary 7] becomes $2\sqrt{\|\mathbf{U}\|_1(m+n)^{1/2} \log(m+n)T}$. For simplicity we consider the simplified setting in which each entry is predicted, that is $T = mn$; then absorbing polylogarithmic factors, their bound is $\tilde{O}(n^{5/4}\|\mathbf{U}\|_1^{1/2})$. From Theorem 3.1 we have a bound of $\tilde{O}(n \text{mc}^2(\mathbf{U}))$. Using [11, Theorem 10], we may upper bound the margin complexity in terms of the trace norm,

$$\text{mc}(\mathbf{U}) \leq 3 \min_{\mathbf{V} \in \text{SP}^1(\mathbf{U})} \|\mathbf{V}\|_1^{1/3} \leq 3\|\mathbf{U}\|_1^{1/3}. \quad (4)$$

Substituting this into Theorem 3.1 our bound is $\tilde{O}(n\|\mathbf{U}\|_1^{2/3})$. Since the trace norm may be bounded as $n \leq \|\mathbf{U}\|_1 \leq n^{3/2}$, both bounds become vacuous when $\|\mathbf{U}\|_1 = n^{3/2}$, however if the trace norm is bounded away from $n^{3/2}$, the bound of Theorem 3.1 is smaller by a polynomial factor. An aspect of the bounds which this comparison fails to capture is the fact that since [17, Corollary 7] is a regret bound it will degrade more smoothly under adversarial noise than Theorem 3.1.

The algorithm in [16] is probabilistic and the regret bound is of $\tilde{O}(\|\mathbf{U}\|_1\sqrt{n})$. Unlike [17], the setting of [16] is transductive, that is each matrix entry is seen only once, and thus less general. If we use the upper bound from [11, Theorem 10] as in the discussion of [17] then [16] improves uniformly on our bound and the bound in [17]. However, using this upper bound oversimplifies the comparison as $1 \leq \text{mc}^2(\mathbf{U}) \leq n$ while $n \leq \|\mathbf{U}\|_1 \leq n^{3/2}$ for $\mathbf{U} \in \{-1, 1\}^{m \times n}$. In other words we have been very conservative in our comparison; the bound (4) may be loose and our algorithm may often have a much smaller bound. A specific example is provided by the class of (k, ℓ) -biclustered matrices (see also the discussion in Section 5 below) where $\text{mc}^2(\mathbf{U}) \leq \min(k, \ell)$, in which case bound becomes nontrivial after $\tilde{\Theta}(\min(k, \ell)n)$ examples while the bounds in [16] and [17] become nontrivial after at least $\tilde{\Theta}(n^{3/2})$ and $\tilde{\Theta}(n^{7/4})$ examples, respectively.

With respect to computation our algorithm on each trial requires a single eigenvalue decomposition of a PSD matrix, whereas the algorithm of [17] requires multiple eigenvalue decompositions per trial. Although [16] does not discuss the complexity of their algorithm beyond the fact that it is polynomial, in [17] it is conjectured that it requires at a minimum $\Theta(n^4)$ time per trial.

4 Comparison to the Best Kernel Perceptron

In this section, we observe that Algorithm 1 has a mistake bound that is comparable to Novikoff’s bound [23] for the Kernel Perceptron with an optimal kernel in hindsight. To explain our observation, we interpret the rows of matrix \mathbf{U} as m different binary classification tasks, and the columns as a finite set of objects which we wish to label; think for example of users/movies matrix in recommendation systems. If we solve the tasks independently using a Kernel Perceptron algorithm, we will make $O(1/\gamma^2)$ mistakes per task, where γ is the largest margin of a consistent hypothesis. If every task has a margin larger than γ we will make $O(m/\gamma^2)$ mistakes in total. This algorithm and the parameter γ crucially depend on the kernel used: if there exists a kernel which makes γ large for all (or most of) the tasks, then the Kernel Perceptron will incur a small number of mistakes on all (or most of) the tasks. We now argue that our bound mimics this “oracle”, without knowing in advance the kernel. Without loss of generality, we assume $m \geq n$ (otherwise apply the same reasoning below to matrix \mathbf{U}^\top). In this scenario, Theorem 3.1 upper bounds the number of mistakes as

$$O\left(\frac{m \log m}{\gamma^2}\right)$$

where γ is chosen so that $\text{merr}(\mathcal{S}, \gamma) = 0$. To further illustrate our idea, we define the *task complexity* of a matrix $\mathbf{U} \in \mathbb{R}^{m \times n}$ as

$$\tau(\mathbf{U}) = \min \{h(\mathbf{V}) : \mathbf{V} \in \text{SP}^1(\mathbf{U})\}$$

where

$$h(\mathbf{V}) = \inf_{\mathbf{K} \succ 0} \max_{1 \leq i \leq m} \mathbf{V}_i \mathbf{K}^{-1} \mathbf{V}_i^\top \max_{1 \leq j \leq n} K_{jj}. \quad (5)$$

Note that the quantity $\mathbf{V}_i \mathbf{K}^{-1} \mathbf{V}_i^\top \max_{1 \leq j \leq n} K_{jj}$ is exactly the bound in Novikoff’s Theorem on the number of mistakes of the Kernel Perceptron on the i -th task with kernel \mathbf{K} . Hence the quantity

$h(\mathbf{V})$ represents the best upper bound on the number of mistakes made by a Kernel Perceptron on the worst (since we take the maximum over i) task.

Proposition 4.1. *For every $\mathbf{U} \in \mathbb{R}^{m \times n}$, it holds that $\text{mc}^2(\mathbf{U}) = \tau(\mathbf{U})$.*

Proof. The result follows by Lemma A.6 presented in the appendix and by the formula $\text{mc}(\mathbf{U}) = \min_{\mathbf{V} \in \text{SP}^1(\mathbf{U})} \|\mathbf{V}\|_{\max}$, see, e.g., [3, Lemma 3.1]. \square

Returning to the interpretation of the bound in Theorem 3.1, we observe that if no more than r out of the m tasks have margin smaller than a threshold λ then in Algorithm 1 setting parameter $\gamma = \lambda$, Theorem 3.1 gives a bound of

$$O\left(\frac{(m-r)\log m}{\lambda^2} + rn\right).$$

Thus we essentially “pay” linearly for every object in a difficult task. Since we assume $n \leq m$, provided r is small the bound is “robust” to the presence of bad tasks.

We specialize the above discussion to the case that each of the m tasks is a binary labeling of an unknown underlying connected graph $\mathcal{G} := (\mathcal{V}, \mathcal{E})$ with n vertices and assume that $m \geq n$. We let $\mathbf{U} \in \{-1, 1\}^{m \times n}$ be the matrix, the rows of which are different binary labelings of the graph. For every $i \in \mathbb{N}_m$, we interpret \mathbf{U}_i , the i -th row of matrix \mathbf{U} , as the i -th labeling of the graph and let Φ_i be the corresponding cutsize, namely, $\Phi_i := |\{(j, j') \in \mathcal{E} : U_{ij} \neq U_{ij'}\}|$ and define $\Phi_{\max} := \max_{1 \leq i \leq m} \Phi_i$. In order to apply Theorem 3.1, we need to bound the margin complexity of \mathbf{U} . Using the above analysis (Proposition 4.1), this quantity is upper bounded by

$$\text{mc}^2(\mathbf{U}) \leq \max_{1 \leq i \leq m} \mathbf{U}_i \mathbf{K}^{-1} \mathbf{U}_i^\top \max_{1 \leq j \leq n} K_{jj}. \quad (6)$$

We choose the kernel $\mathbf{K} := \mathbf{L}^+ + (R\mathbf{1}\mathbf{1}^\top)$, where \mathbf{L} is the graph Laplacian of \mathcal{G} , the vector $\mathbf{1}$ has all components equal to one, and $R = \max_j L_{jj}^+$. Since the graph is connected then $\mathbf{1}$ is the only eigenvector of \mathbf{L} with zero eigenvalue. Hence \mathbf{K} is invertible and $\mathbf{K}^{-1} = \mathbf{L} + (R\mathbf{1}\mathbf{1}^\top)^+ = \mathbf{L} + (Rn\frac{1}{\sqrt{n}}\mathbf{1}\mathbf{1}^\top\frac{1}{\sqrt{n}})^+ = \mathbf{L} + \frac{1}{Rn^2}\mathbf{1}\mathbf{1}^\top$. Then using the formula $\Phi_i = \frac{1}{4}\mathbf{U}_i\mathbf{L}\mathbf{U}_i^\top$ we obtain from (6) that

$$\text{mc}^2(\mathbf{U}) \leq \max_{1 \leq i \leq m} \left(4\Phi_i + \frac{1}{R}\right) R.$$

Theorem 3.1 then gives a bound of $M \leq O((1 + \Phi_{\max}R)m \log m)$. The quantity R may be further upper bounded by the graph resistance diameter, see for example [24].

5 Biclustering and Near Optimality

The problem of learning a (k, ℓ) -binary-biclustered matrix, corresponds to the assumption that the row indices and column indices represent k and ℓ distinct object types and that there exists a binary relation on these objects which determines the matrix entry. Formally we have the following

Definition 5.1. *The class of (k, ℓ) -binary-biclustered matrices is defined as*

$$\mathbb{B}_{k,\ell}^{m,n} = \{\mathbf{U} \in \mathbb{R}^{m \times n} : \mathbf{r} \in \mathbb{N}_k^m, \mathbf{c} \in \mathbb{N}_\ell^n, \mathbf{F} \in \{-1, 1\}^{k \times \ell}, U_{ij} = F_{r_i c_j}, i \in \mathbb{N}_m, j \in \mathbb{N}_n\}.$$

The intuition is that a matrix is (k, ℓ) -biclustered if after a permutation of the rows and columns the resulting matrix is a $k \times \ell$ grid of rectangles and all entries in a given rectangle are either 1 or -1 . The problem of determining a (k, ℓ) -biclustered matrix with a minimum number of “violated” entries given a subset of entries was shown to be NP-hard in [25]. Thus although we do not give an algorithm that provides a biclustering, we provide a bound in terms of the best consistent biclustering.

Lemma 5.2. *If $\mathbf{U} \in \mathbb{B}_{k,\ell}^{m,n}$ then $\text{mc}^2(\mathbf{U}) \leq \min(k, \ell)$.*

Proof. We use Proposition 4.1 to upper bound $\text{mc}^2(\mathbf{U})$ by $h(\mathbf{U})$, where the function h is given in equation (5). We further upper bound $h(\mathbf{U})$ by choosing a kernel matrix in the underlying optimization problem. By Definition 5.1, there exists $\mathbf{r} \in \mathbb{N}_k^m, \mathbf{c} \in \mathbb{N}_\ell^n$ and $\mathbf{F} \in \{-1, 1\}^{k \times \ell}$

such that $U_{ij} = F_{r_i c_j}$, for every $i \in \mathbb{N}_m$ and every $j \in \mathbb{N}_n$. Then we choose the kernel matrix $\mathbf{K} = (K_{jj'})_{1 \leq j, j' \leq n}$ such that

$$K_{jj'} := \delta_{c_j c_{j'}} + \epsilon \delta_{jj'}$$

One verifies that $\mathbf{U}_i \mathbf{K}^{-1} \mathbf{U}_i^\top \leq \ell$ for every $i \in \{1, \dots, m\}$, hence by taking the limit for $\epsilon \rightarrow 0$ Proposition 4.1 gives that $\text{mc}^2(\mathbf{U}) \leq \ell$. By the symmetry of our construction we can swap ℓ with k , giving the bound. \square

Using this lemma with Theorem 3.1 gives us the following upper bound on the number of mistakes.

Corollary 5.3. *The number of mistakes of Algorithm 1 applied to sequences generated by a (k, ℓ) -binary-biclustered matrix is upper bounded by $\mathcal{O}(\min(k, \ell)(m+n) \log(m+n))$.*

A special case of the setting in this corollary was first studied in the mistake bound setting in [14]. In [15] the bound was improved and generalized to include robustness to noise (for simplicity we do not compare in the noisy setting). In both papers the underlying assumption is that there are k distinct row types and no restrictions on the number of columns thus $\ell = n$. In this case they obtained an upper bound of $kn + \min(\frac{m^2}{2e} \log_2 e, m\sqrt{3n \log_2 k})$. Comparing the two bounds we can see that when $k < n^{\frac{1}{2}-\epsilon}$ the bound in Corollary 5.3 improves over [15, Corollary 1] by a polynomial factor and on other hand when $k \geq n^{\frac{1}{2}}$ we are no worse than a polylogarithmic factor.

We now establish that the mistake bound (3) is tight up to a poly-logarithmic factor.

Theorem 5.4. *Given an online algorithm \mathcal{A} that predicts the entries of a matrix $\mathbf{U} \in \{-1, 1\}^{m \times n}$ and given an $\ell \in \mathbb{N}_n$ there exists a sequence \mathcal{S} constructed by an adversary with margin complexity $\text{mc}(\mathcal{S}) \leq \sqrt{\ell}$. On this sequence the algorithm \mathcal{A} will make at least $\ell \times m$ mistakes.*

See the appendix for a proof.

6 Conclusion

In this paper, we presented a Matrix Exponentiated Gradient algorithm for completing the entries of a binary matrix in an online learning setting. We established a mistake bound for this algorithm, which is controlled by the margin complexity of the underlying binary matrix. We discussed improvements of the bound over related bounds for matrix completion. Specifically, we noted that our bound requires fewer examples before it becomes non-trivial, as compared to the bounds in [16, 17]. Here we require only $\tilde{\Theta}(m+n)$ examples as opposed to the required $\tilde{\Theta}((m+n)^{3/2})$ in [16] and $\tilde{\Theta}((m+n)^{7/4})$, respectively. Thus although our bound is more sensitive to noise, it captures structure more quickly in the underlying matrix. When interpreting the rows of the matrix as binary tasks, we argued that our algorithm performs comparably (up to logarithmic factors) to the Kernel Perceptron with the optimal kernel in retrospect. Finally, we highlighted the example of completing a biclustered matrix and noted that this is instrumental in showing the optimality of the algorithm in Theorem 5.4.

We observed that Algorithm 1 has a per trial computational cost which is smaller than currently available algorithms for matrix completion with online guarantees. In the future it would be valuable to study if improvements in this computation are possible by exploiting the special structure in our algorithm. Furthermore, it would be very interesting to study a modification of our analysis to the case in which the tasks (rows of matrix \mathbf{U}) grow over time, a setting which resembles the lifelong learning frameworks in [26, 27].

Acknowledgements. We wish to thank the anonymous reviewers for their useful comments. This work was supported in part by EPSRC Grants EP/P009069/1, EP/M006093/1, and by the U.S. Army Research Laboratory and the U.K. Defence Science and Technology Laboratory and was accomplished under Agreement Number W911NF-16-3-0001. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Army Research Laboratory, the U.S. Government, the U.K. Defence Science and Technology Laboratory or the U.K. Government. The U.S. and U.K. Governments are authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

References

- [1] K. Tsuda, G. Rätsch, and M.K. Warmuth. Matrix exponentiated gradient updates for on-line learning and bregman projection. *Journal of Machine Learning Research*, 6:995–1018, 2005.
- [2] S. Ben-David, N. Eiron, and H. U. Simon. Limitations of learning via embeddings in euclidean half spaces. *Journal of Machine Learning Research*, 3:441–461, 2003.
- [3] N. Linial, S. Mendelson, G. Schechtman, and A. Shraibman. Complexity measures of sign matrices. *Combinatorica*, 27(4):439–463, 2007.
- [4] M. Herbster, M. Pontil, and L. Wainer. Online learning over graphs. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 305–312, 2005.
- [5] X. Zhu, Z. Ghahramani, and J. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *Proc. 20th International Conference on Machine Learning*, pages 912–919, 2003.
- [6] M. Belkin and P. Niyogi. Semi-supervised learning on riemannian manifolds. *Machine Learning*, 56:209–239, 2004.
- [7] N. Cesa-Bianchi, C. Gentile, and F. Vitale. Fast and optimal prediction of a labeled tree. In *Proceedings of the 22nd Annual Conference on Learning Theory*, 2009.
- [8] N. Cesa-Bianchi, C. Gentile, F. Vitale, and G. Zappella. Random spanning trees and the prediction of weighted graphs. *Journal of Machine Learning Research*, 14(1):1251–1284, 2013.
- [9] N. Cesa-Bianchi, C. Gentile, and F. Vitale. Predicting the labels of an unknown graph via adaptive exploration. *Theoretical Computer Science*, 412(19):1791–1804, 2011.
- [10] C. Gentile, M. Herbster, and S. Pasteris. Online similarity prediction of networked data from known and unknown graphs. In *Proceedings of the 26th Annual Conference on Learning Theory*, 2013.
- [11] N. Srebro and A. Shraibman. Rank, trace-norm and max-norm. In *Proceedings of the 18th Annual Conference on Learning Theory*, pages 545–560, 2005.
- [12] E. J. Candès and T. Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Trans. Inf. Theor.*, 56(5):2053–2080, May 2010.
- [13] A. Maurer and M. Pontil. Excess risk bounds for multitask learning with trace norm regularization. In *Proceedings of The 27th Conference on Learning Theory (COLT)*, pages pages 55–76, 2013.
- [14] S. A. Goldman, R. L. Rivest, and R. E. Schapire. Learning binary relations and total orders. *SIAM J. Comput.*, 22(5), 1993.
- [15] S. A. Goldman and M. K. Warmuth. Learning binary relations using weighted majority voting. In *Proceedings of the 6th Annual Conference on Computational Learning Theory*, pages 453–462, 1993.
- [16] N. Cesa-Bianchi and O. Shamir. Efficient online learning via randomized rounding. In *Advances in Neural Information Processing Systems 24*, pages 343–351, 2011.
- [17] E. Hazan, S. Kale, and S. Shalev-Shwartz. Near-optimal algorithms for online matrix prediction. In *Proc. 23rd Annual Conference on Learning Theory*, volume 23:38.1-38.13. JMLR W&CP, 2012.
- [18] S. Arora and S. Kale. A combinatorial, primal-dual approach to semidefinite programs. In *Proceedings of the 29th Annual ACM Symposium on Theory of Computing*, pages 227–236, 2007.
- [19] M.K. Warmuth. Winoing subspaces. In *Proceedings of the 24th International Conference on Machine Learning*, pages 999–1006, 2007.
- [20] J. Nie, W. Kotłowski, and M. K. Warmuth. Online PCA with optimal regrets. In *Proceedings of the 24th International Conference on Algorithmic Learning Theory*, pages 98–112, 2013.
- [21] M. K. Warmuth and D. Kuzmin. Online variance minimization. *Machine Learning*, 87(1):1–32, 2012.
- [22] M. Herbster, S. Pasteris, and S. Pontil. Predicting a switching sequence of graph labelings. *Journal of Machine Learning Research*, 16:2003–2022, 2015.
- [23] A.B. Novikoff. On convergence proofs on perceptrons. In *Proceedings of the Symposium on the Mathematical Theory of Automata*, pages 615–622, 1962.
- [24] M. Herbster and M. Pontil. Prediction on a graph with a perceptron. In *Advances in Neural Information Processing Systems 19*, pages 577–584, 2006.
- [25] S. Wulff, R. Uerner, and S. Ben-David. Monochromatic bi-clustering. In *Proc. 30th International Conference on Machine Learning*, volume 28, pages 145–153. JMLR W&CP, 2013.
- [26] P. Alquier, T.-T. Mai, and M. Pontil. Regret bounds for lifelong learning. *Preprint*, 2016.
- [27] M.-F. Balcan, A. Blum, and S. Vempala. Efficient representations for lifelong learning and autoencoding. In *Proc. 28th Conference on Learning Theory*, pages 191–210, 2015.
- [28] R. Bhatia. *Matrix Analysis*. Springer Verlag, New York, 1997.

A Appendix

In this appendix we give proof of some of the results only stated in the main body of the paper and collect some auxiliary results. The first result is the well known Golden-Thompson Inequality, whose proof can be found, for example, in [28].

Lemma A.1. *For any symmetric matrices \mathbf{A} and \mathbf{B} we have that*

$$\text{Tr}(\exp(\mathbf{A} + \mathbf{B})) \leq \text{Tr}(\exp(\mathbf{A}) \exp(\mathbf{B})). \quad (7)$$

The next result is taken from [1].

Lemma A.2. *If $A > 0$ with eigenvalues in $[0, 1]$ and $a \in \mathbb{R}$, then*

$$(1 - e^a)\mathbf{A} \preceq \mathbf{I} - \exp(a\mathbf{A}).$$

The next lemma is useful for the analysis of Matrix Winnow, see [1, 19, 22].

Lemma A.3.

$$\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) \geq \frac{\gamma}{2}(y_t - \hat{y}_t) \text{Tr}(\tilde{\mathbf{U}} \mathbf{X}^{(t)}) + \left(1 - e^{\frac{\gamma}{2}(y_t - \hat{y}_t)}\right) \text{Tr}(\mathbf{W}^{(t-1)} \mathbf{X}^{(t)})$$

Proof. We observe that

$$\begin{aligned} \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) &= \text{Tr}(\tilde{\mathbf{U}} \log \mathbf{W}^{(t)} - \tilde{\mathbf{U}} \log \mathbf{W}^{(t-1)}) + \text{tr} \mathbf{W}^{(t-1)} - \text{tr} \mathbf{W}^{(t)} \\ &= \frac{\gamma}{2}(y_t - \hat{y}_t) \text{Tr}(\tilde{\mathbf{U}} \mathbf{X}^{(t)}) + \text{tr} \mathbf{W}^{(t-1)} - \text{tr} \left(e^{\log \mathbf{W}^{(t-1)} + \frac{\gamma}{2}(y_t - \hat{y}_t) \mathbf{X}^{(t)}} \right) \\ &\geq \frac{\gamma}{2}(y_t - \hat{y}_t) \text{Tr}(\tilde{\mathbf{U}} \mathbf{X}^{(t)}) + \text{tr} \mathbf{W}^{(t-1)} - \text{tr} \left(e^{\log \mathbf{W}^{(t-1)}} e^{\frac{\gamma}{2}(y_t - \hat{y}_t) \mathbf{X}^{(t)}} \right) \\ &= \frac{\gamma}{2}(y_t - \hat{y}_t) \text{Tr}(\tilde{\mathbf{U}} \mathbf{X}^{(t)}) + \text{tr} \left(\mathbf{W}^{(t-1)} \left(\mathbf{I} - e^{\frac{\gamma}{2}(y_t - \hat{y}_t) \mathbf{X}^{(t)}} \right) \right) \\ &\geq \frac{\gamma}{2}(y_t - \hat{y}_t) \text{Tr}(\tilde{\mathbf{U}} \mathbf{X}^{(t)}) + (1 - e^{\frac{\gamma}{2}(y_t - \hat{y}_t)}) \text{Tr}(\mathbf{W}^{(t-1)} \mathbf{X}^{(t)}) \end{aligned}$$

where the second equality follows by the update formula in Algorithm 1, the first inequality follows by Golden-Thompson Inequality (Lemma A.1), and the last inequality follows by Lemma A.2. \square

The next lemma further analyzes the quantity $\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)})$ distinguishing between the case in which there is a margin violation or not.

Lemma A.4. *Let $c' = 3 - e$. If $y_t \neq \hat{y}_t$, it holds that*

$$\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) \geq \begin{cases} (c' - 1) \frac{1}{n+m} \gamma^2, & \text{if there is a margin violation,} \\ c' \frac{1}{n+m} \gamma^2, & \text{otherwise.} \end{cases}$$

Proof. If $y_t = 1, \hat{y}_t = -1$ and there is not a margin violation then

$$\begin{aligned} \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) &\geq \gamma \frac{1}{n+m} \left(1 + \frac{\langle \mathbf{P}_{i_t}, \mathbf{Q}_{j_t} \rangle}{|\mathbf{P}_{i_t}| |\mathbf{Q}_{j_t}|} \right) + \frac{1}{n+m} (1 - e^\gamma) \\ &\geq \frac{1}{n+m} (\gamma + \gamma^2 + 1 - e^\gamma) \\ &\geq c' \frac{1}{n+m} \gamma^2. \end{aligned}$$

And if $y_t = -1$, $\hat{y}_t = 1$ and there is not a margin violation then

$$\begin{aligned}\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) &\geq -\gamma \frac{1}{n+m} \left(1 + \frac{\langle \mathbf{P}_{i_t}, \mathbf{Q}_{j_t} \rangle}{|\mathbf{P}_{i_t}| |\mathbf{Q}_{j_t}|}\right) + \frac{1}{n+m} (1 - e^{-\gamma}) \\ &\geq \frac{1}{n+m} (-\gamma + \gamma^2 + 1 - e^{-\gamma}) \\ &\geq c' \frac{1}{n+m} \gamma^2.\end{aligned}$$

If $y_t = 1$, $\hat{y}_t = -1$ and there is a margin violation then

$$\begin{aligned}\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) &\geq \gamma \frac{1}{n+m} \left(1 + \frac{\langle \mathbf{P}_{i_t}, \mathbf{Q}_{j_t} \rangle}{|\mathbf{P}_{i_t}| |\mathbf{Q}_{j_t}|}\right) + \frac{1}{n+m} (1 - e^{-\gamma}) \\ &\geq \frac{1}{n+m} (\gamma + \gamma \cdot 0 + 1 - e^{-\gamma}) \\ &= \frac{1}{n+m} (\gamma + 1 - e^{-\gamma}) \\ &\geq (c' - 1) \frac{1}{n+m} \gamma^2.\end{aligned}$$

And if $y_t = -1$, $\hat{y}_t = 1$ and there is a margin violation then

$$\begin{aligned}\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) &\geq -\gamma \frac{1}{n+m} \left(1 + \frac{\langle \mathbf{P}_{i_t}, \mathbf{Q}_{j_t} \rangle}{|\mathbf{P}_{i_t}| |\mathbf{Q}_{j_t}|}\right) + \frac{1}{n+m} (1 - e^{-\gamma}) \\ &\geq -\frac{1}{n+m} (\gamma - \gamma \cdot 0 - 1 + e^{-\gamma}) \\ &= \frac{1}{n+m} (-\gamma + 1 - e^{-\gamma}) \\ &\geq (c' - 1) \frac{1}{n+m} \gamma^2.\end{aligned}$$

So on a mistaken trial t without a margin violation we have

$$\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) \geq c' \frac{1}{n+m} \gamma^2$$

and on a mistaken trial t with a margin violation we have

$$\Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t-1)}) - \Delta(\tilde{\mathbf{U}}, \mathbf{W}^{(t)}) \geq (c' - 1) \frac{1}{n+m} \gamma^2.$$

The constant c' is chosen via Lemma A.5 below. □

Lemma A.5. For every $x \in [0, 1]$ it holds that

$$\min(x^2 - x + 1 - e^{-x}, x^2 + x + 1 - e^x) \geq (3 - e)x^2.$$

Proof. Let $f(x) := x^2 - x + 1 - e^{-x} - (3 - e)x^2$ and $g(x) := x^2 + x + 1 - e^x - (3 - e)x^2$.

We first show that $f(x) \geq 0$. Taking a derivative we have that $f'(x) := e^{-x} + (2e - 4)x$. Since $f'(0) = 0$ and $f'(x)$ is strictly increasing for $x \geq 0$ this implies that $f(x) \geq 0$ for $x \geq 0$.

Second, we show that $g(x) \geq 0$. Taking derivatives we have that $g'(x) = 1 + (2e - 4)x - e^x$ and $g''(x) = (2e - 4) - e^x$. Since g'' is strictly decreasing, g' is concave and has at most two zero-crossings. Therefore g has at most two critical points. One critical point is at 0 and the other is at an $x^* \in [0.6, 0.8]$. Since $g''(0)$ is positive, 0 is a minimum of g and since $g''(x) < 0$ for $0.6 < x$, x^* is a maximum of g . Thus g is unimodal in $[0, 1]$ with minimum values at $g(0) = g(1) = 0$. □

The following provides an alternative formulation of the max-norm.

Lemma A.6. For every $\mathbf{V} \in \mathbb{R}^{m \times n}$, we have that $\|\mathbf{V}\|_{\max}^2 = h(\mathbf{V})$.

Proof. Let $f(\mathbf{P}, \mathbf{Q}) = \max_{1 \leq i \leq m} |\mathbf{P}_i|^2 \max_{1 \leq j \leq n} |\mathbf{Q}_j|^2$. Then we have that

$$\|\mathbf{V}\|_{\max}^2 = \inf_{\mathbf{P}\mathbf{Q}^\top = \mathbf{V}} f(\mathbf{P}, \mathbf{Q}) \quad (8)$$

where the infimum is over all real matrices $\mathbf{P} \in \mathbb{R}^{m \times k}$, $\mathbf{Q} \in \mathbb{R}^{n \times k}$ and every integer k . For any $\mathbf{K} \succ 0$, set $\mathbf{P} = \mathbf{V}\sqrt{\mathbf{K}^{-1}}$ and $\mathbf{Q} = \sqrt{\mathbf{K}}$. Note that $\mathbf{P}\mathbf{Q}^\top = \mathbf{V}$ and, for every $i \in \mathbb{N}_m$, $|\mathbf{P}_i|^2 = \mathbf{V}_i \mathbf{K}^{-1} \mathbf{V}_i^\top$. Moreover, for any $j \in \mathbb{N}_n$, it holds $|\mathbf{Q}_j|^2 = K_{jj}$. This shows that every feasible point in the optimization problem (5) maps to a feasible point for the optimization problem (8) and that the two objective functions take the same value at these points. Hence, we have shown, for every $\mathbf{V} \in \mathbb{R}^{m \times n}$, that

$$\|\mathbf{V}\|_{\max}^2 \leq h(\mathbf{V}). \quad (9)$$

To show the reverse inequality, let $k \in \mathbb{N}$, let $\mathbf{P} \in \mathbb{R}^{m \times k}$ and $\mathbf{Q} \in \mathbb{R}^{n \times k}$ be a solution of the optimization problem (8). Set $\mathbf{K} = \mathbf{Q}\mathbf{Q}^\top + \epsilon \mathbf{I}$ for $\epsilon > 0$. Then for every $j \in \mathbb{N}_n$, $K_{jj} = |\mathbf{Q}_j|^2 + \epsilon$ and, for every $i \in \mathbb{N}_m$, we have that

$$\mathbf{V}_i \mathbf{K}^{-1} \mathbf{V}_i^\top = \mathbf{e}_i^\top \mathbf{V} \mathbf{K}^{-1} \mathbf{V}^\top \mathbf{e}_i = \mathbf{e}_i^\top \mathbf{P} \mathbf{Q}^\top \mathbf{K}^{-1} \mathbf{Q} \mathbf{P}^\top \mathbf{e}_i \leq \mathbf{P}_i \mathbf{P}_i^\top = |\mathbf{P}_i|^2$$

where we used the fact that $\mathbf{Q}^\top (\mathbf{Q}\mathbf{Q}^\top + \epsilon \mathbf{I})^{-1} \mathbf{Q} \prec \mathbf{I}$. Thus

$$\max_{1 \leq i \leq m} \mathbf{V}_i \mathbf{K}^{-1} \mathbf{V}_i^\top \max_{1 \leq j \leq n} K_{jj} < \max_{1 \leq i \leq m} |\mathbf{P}_i|^2 \max_{1 \leq j \leq n} (|\mathbf{Q}_j|^2 + \epsilon)$$

and taking the limit of $\epsilon \rightarrow 0$, we conclude that $h(\mathbf{V}) \leq \|\mathbf{V}\|_{\max}^2$. This inequality, combined with inequality (9) proves the result. \square

Proof of Corollary 3.2. For any $\gamma > 0$ define $\mathcal{A}(\gamma) := c(m+n) \log(m+n) \frac{1}{\gamma^2}$, $\mathcal{B}(\gamma) := c \text{merr}(\mathcal{S}, \gamma)$, and $\mathcal{M}(\gamma) := \mathcal{A}(\gamma) + \mathcal{B}(\gamma)$ which is equal to the mistake bound for Algorithm 1 with learning rate γ . Let η be such that for all $\gamma \leq \eta$ we have $\mathcal{A}(\gamma) \geq \mathcal{B}(\gamma)$ and for all $\gamma > \eta$ we have $\mathcal{A}(\gamma) \leq \mathcal{B}(\gamma)$ which is defined and positive since $\mathcal{A}(\cdot)$ is continuous and monotonic decreasing (and limiting to 0) and $\mathcal{B}(\cdot)$ is monotonic non-decreasing (with $\mathcal{B}(0) = 0$). Since $\mathcal{A}(\cdot)$ is continuous let ϵ be such that $\epsilon > 0$ and $\mathcal{A}(\eta + \epsilon) \geq \mathcal{A}(\eta) - \frac{1}{2}$. For $\gamma \leq \eta + \epsilon$ we have $\mathcal{M}(\gamma) \geq \mathcal{A}(\gamma) \geq \mathcal{A}(\eta + \epsilon)$ and for $\gamma \geq \eta + \epsilon$ we have $\mathcal{M}(\gamma) \geq \mathcal{B}(\gamma) \geq \mathcal{B}(\eta + \epsilon) \geq \mathcal{A}(\eta + \epsilon)$ so in either case we have $\mathcal{M}(\gamma) \geq \mathcal{A}(\eta + \epsilon) \geq \mathcal{A}(\eta) - \frac{1}{2}$.

Let z be the minimum integer power of $\sqrt{2}$ that is greater than or equal to $1/\eta$. i.e. $z := \min\{\sqrt{2}^i : i \in \mathbb{N}, \sqrt{2}^i \geq 1/\eta\}$. Since $1/z \leq \eta$ we have $\mathcal{B}(1/z) \leq \mathcal{B}(\eta) \leq \mathcal{A}(\eta) \leq \mathcal{A}(1/z)$ so $\mathcal{M}(1/z) = \mathcal{A}(1/z) + \mathcal{B}(1/z) \leq 2\mathcal{A}(1/z) = 2c(m+n) \log(m+n) z^2$ and hence the doubling algorithm above terminates at some $\kappa \leq z$. The total number of mistakes, M , made by doubling algorithm above is then bounded above by

$$\begin{aligned} & (2\mathcal{A}(2^{-1/2}) + 1) + (2\mathcal{A}(2^{-2/2}) + 1) + (2\mathcal{A}(2^{-3/2}) + 1) + \dots + (2\mathcal{A}(1/z) + 1) \\ & \leq 3\mathcal{A}(2^{-1/2}) + 3\mathcal{A}(2^{-2/2}) + \dots + 3\mathcal{A}(1/z) \\ & \leq 3c(m+n) \log(m+n) (2 + 2^2 + \dots + z^2) \\ & \leq 6c(m+n) \log(m+n) z^2 - 6 \\ & = 6\mathcal{A}(1/z) - 6. \end{aligned}$$

We also have $z \leq \frac{\sqrt{2}}{\eta}$ so $\mathcal{A}(1/z) \leq 2\mathcal{A}(\eta)$ and putting together we get $M \leq 12\mathcal{A}(\eta) - 6$ so since, by above, $\mathcal{A}(\eta) \leq \mathcal{M}(\eta) + \frac{1}{2}$ for all $\eta > 0$ we're done. \square

Proof of Theorem 5.4. Observe that for the set of (k, ℓ) -biclustered $m \times n$ matrices we have that,

$$k \times \ell \leq \text{VC-dimension}(\mathbb{B}_{k,\ell}^{m,n})$$

as each of $k \times \ell$ ‘‘tiles’’ may be labeled independently of the others. Thus there exists an adversary that may force $m \times \ell$ mistakes via a constructed sequence \mathcal{S} that is consistent with some (m, ℓ) -biclustered matrix \mathbf{U} . By Lemma 5.2 we have that $\text{mc}^2(\mathbf{U}) \leq \min(m, \ell)$ and hence $\text{mc}(\mathcal{S}) \leq \sqrt{\ell}$. \square