

1 **R1:** “*non-conservative flag*” : In **practice** the algorithm should **always** be run non-conservatively to obtain a regret
2 bound as that is the more flexible/noise-tolerant bound. We **only** include the conservative case for presentation issues as
3 the bound is *simpler*, analogous to Novikoff’s (perceptron) bound, and in order to give “easy” upper and lower bound
4 insights in that case (e.g., lines 244-249, 350-353). We will update the manuscript to reflect this advice.

5 “*For Theorem 1, ... non-conservative = 0? ... don’t depend on the value of γ at all, so why does it show up in the bounds?*”
6 : See line 172 - 173, 1) the algorithm’s predictions depends on γ through η ($\eta = \gamma$). 2) the bound fails hold if for the
7 “true” comparator matrix U the margin complexity is large i.e., $\text{mc}(U) > 1/\gamma$.

8 “*modification upon MEG ... (e.g. were any new proof techniques required).*” : See 109-112 which then refers to Appendix
9 B.2 (esp. lines 520-536) which discusses our new proof techniques wrt MEG.

10 **R2:** “*Why does $\text{mc}(U)$... resembles the margin in SVM*”: In lines 2-8 of the abstract we discuss the term $1/\gamma^2$ which
11 serves as the algorithmic proxy for $\text{mc}^2(U)$ i.e., the tightest bound is obtained if $1/\gamma^2 = \text{mc}^2(U)$ and it is in those
12 lines that we make the explicit connection to the “SVM margin.” This initial discussion is then built on in lines 54-61
13 where we further refer to the three references [5,6,7].

14 “*The latent block structure seems not well-motivated from a practical part. In Netflix, it should be a five rated problem ... ,*
15 *the latent block structure may be a strong assumption and not easy to be tested in reality.*” : See lines 40-53 (esp. 51-53).
16 Observe that the latent block structure (LBS) assumption is actually a *weaker* assumption than the common low-rank
17 assumption. Note that Netflix does not meet the low-rank assumption as it is experimentally known that human ratings
18 are only *ordinally* qualitative. Finally the definition of LBS trivially generalizes to categorical data as do our algorithms.

19 **R3:** “*Lower bounds*” : We agree that there is value in formalizing the setting. Our casual mistake lower bounds
20 generalize to lower bounds for regret using [Agnostic Online Learning, 2009, Ben-David et al; Lem. 14].

21 **R4:** “*(c) The exposition ... Section 4.1, can be simplified ... matrices M and N directly in Section 4.1 rather than*
22 *arriving at them from a graph perspective.*” Yes, an alternate presentation is to assume feature vectors associated with
23 each row and column and let M and N be the inverse of the gram (kernel) matrices (also see 251-256). However
24 we focus on the graph perspective because of the smaller bound on \mathcal{D} wrt graphs given in Thm 3. Eq (7). We note in
25 the batch setting different methods of using graph side information were considered in [15,16] (lines 87-94). In the
26 inductive setting (Sec. 5.1) we give an example of using non-graph side-information. Finally we note that 4.1 serves as
27 the necessary background for 4.2 which has two important observations: 1) a bound from [22] can be recovered and
28 extended, 2) An example of a class of $k \times k$ biclustered matrices (whose rank are k) for which max norm is $O(1)$.

29 “*(d), (e)*” : We will move related work and use forward references for undefined notation.

30 “*(f) What are $\mathcal{R}_M, \mathcal{R}_N, \mathcal{R}_L$...*” : See line 129. The “*Iff*” (147-148) .. follows from definition of $\|U\|_{\max}$.

31 “*(2) The transductive setting seems unnatural. Why is it reasonable to assume that M and N will be available*
32 *beforehand?*” : We argue side-information is the norm rather than the exception. I.e., in the Netflix example we
33 may have demographic info on the *users* as well as categorization, actor lists, etc on the *movies*. M and N may
34 then be constructed from feature vectors by selecting kernels and inverting the kernel matrices. Or as in graph-based
35 semi-supervised learning, a graph may be constructed using the feature vectors and using the corresponding Laplacian.

36 “*(3) ... experiments (4) ...time complexity ...*” : We agree that experiments would be useful and that the time complexity
37 is large. However, we note that natural heuristics include maintaining only a low-rank approximation to \bar{W}^t and
38 maintaining a fixed number of indices in \mathbb{U} (e.g., decaying old indices) for Algs. 1 and 2, respectively.

39 “*... counterintuitive ... inversely proportional to gamma. ... Shouldn’t more mistakes be made when the margin*
40 *requirements are higher?*” : See lines 3-4 of the Abstract. This is the usual intuition behind perceptron, SVM, and
41 other “margin margin” classifiers. I.e., the further that data in classes are apart the easier it is separate and thus (online)
42 fewer mistakes (batch) better generalization as opposed to the case where the classes are arbitrarily near one another.

43 “*How is the max norm block invariant? I believe that as the sizes of the matrices are different, they will be on different*
44 *scales.*” : **Theorem: max-norm is block invariant.** *Pf. Sketch.* We first show $\|X\|_{\max} \geq \|RXC^T\|_{\max}$ for all
45 $m, k, n, \ell \in \mathbb{N}^+$ with $m \geq k, n \geq \ell, R \in \mathcal{B}^{m,k}$ and $C \in \mathcal{B}^{n,\ell}$ (where $\mathcal{B}^{m,k}$ is the set of all $m \times k$ block expansion
46 matrices (cf lines 140-143)). WLOG. assume X is $k \times \ell$. If $PQ^T = X$ then $(RP)(CQ)^T = RXC^T$. Observe
47 that $\max_{i \in [k]} \|P_i\| = \max_{s \in [m]} \|(RP)_s\|$ since every row in P is duplicated by (1+) rows in (RP) and there are no
48 distinct rows in (RP) that are not in P . Recall $\|X\|_{\max} := \min_{PQ^T = X} \{\max_{1 \leq i \leq m} \|P_i\| \times \max_{1 \leq j \leq n} \|Q_j\|\}$ then
49 since for every decomposition $PQ^T = X$ there exists a decomposition $(RP)(CQ)^T = RXC^T$ thus $\|X\|_{\max} \geq$
50 $\|RXC^T\|_{\max}$. We have $\|X\|_{\max} \leq \|RXC^T\|_{\max}$ since trivially the max-norm of a sub-matrix cannot be larger than
51 that of the matrix ■.

52 “*Where is the minimum in (7)?*” : See line 205. The minimum is before the large ‘{’ in (7) and is over the matrices
53 R, C, U^* s.t. $U = RU^*C^T$.