

## A. Unregularized LEML with Squared $L_2$ Loss Recovers CPLST

**Claim 1.** If  $\ell(\mathbf{y}, f(\mathbf{x}; Z)) = \|\mathbf{y} - Z^T \mathbf{x}\|_2^2$  and  $\lambda = 0$ , then

$$V_X \Sigma_X^{-1} M_k = \arg \min_{Z: \text{rank}(Z) \leq k} \|Y - XZ\|_F^2,$$

where  $X = U_X \Sigma_X V_X^T$  is the thin SVD decomposition of  $X$ , and  $M_k$  is the rank- $k$  truncated SVD of  $M \equiv U_X^T Y$ .

*Proof of Claim 1.* Let  $X = U_X \Sigma_X V_X^T$  be the thin SVD decomposition of  $X$ , and  $M_k$  be the rank- $k$  truncated SVD approximation of  $U_X^T Y$ . We have

$$\begin{aligned} \arg \min_{Z: \text{rank}(Z) \leq k} \|Y - XZ\|_F &= \arg \min_{Z: \text{rank}(Z) \leq k} \|(U_X U_X^T)(Y - XZ) + (I - U_X U_X^T)(Y - XZ)\|_F \\ &= \arg \min_{Z: \text{rank}(Z) \leq k} \|(U_X U_X^T)(Y - XZ) + (I - U_X U_X^T)(Y - XZ)\|_F^2 \\ &= \arg \min_{Z: \text{rank}(Z) \leq k} \|U_X U_X^T(Y - XZ)\|_F^2 + \|(I - U_X U_X^T)(Y - XZ)\|_F^2 \\ &= \arg \min_{Z: \text{rank}(Z) \leq k} \|U_X^T(Y - XZ)\|_F^2 \\ &= \arg \min_{Z: \text{rank}(Z) \leq k} \|U_X^T(Y - XZ)\|_F \\ &= \arg \min_{Z: \text{rank}(Z) \leq k} \|U_X^T Y - \Sigma_X V_X^T Z\|_F \\ &= V_X \Sigma_X^{-1} M_k. \end{aligned}$$

The second and the fifth inequalities follow from the fact that the  $(\cdot)^2$  is an increasing function. The third equality follows from the Pythagorean theorem since  $U_X U_X^T$  constitutes an orthonormal projection. Since  $U_X U_X^T X = X$  as  $U_X^T U_X = I_r$ , where  $r$  is the rank of  $X$ , we have  $(I - U_X U_X^T)(Y - XZ) = (I - U_X U_X^T)Y$ . Since the last term does not depend on the variable  $Z$ , it can be removed from consideration and the fourth equality follows. The sixth equality follows due to the same reason as  $U_X^T X = I_r \Sigma_X V_X^T = \Sigma_X V_X^T$ .

For the last equality, first of all note that  $Z = V_X \Sigma_X^{-1} M_k$  is a feasible solution to the problem since  $\text{rank}(V_X \Sigma_X^{-1} M_k) \leq \text{rank}(M_k) \leq k$  by definition of  $M_k$ . Next, notice that for any feasible  $Z'$ , since  $\text{rank}(\Sigma_X V_X^T Z') \leq \text{rank}(Z') \leq k$ , we have  $\|U_X^T Y - \Sigma_X V_X^T Z'\|_F \geq \|U_X^T Y - M_k\|_F$ , again by the definition of  $M_k$ . The result follows since by  $V_X^T V_X = I_r$ , we have  $\Sigma_X V_X^T (V_X \Sigma_X^{-1} M_k) = M_k$ .  $\square$

**Claim 2.** The solution to (3) is equivalent to  $Z^{CPLST} = W_{CPLST} H_{CPLST}^T$  which is the closed form solution for the CPLST scheme, i.e.,

$$\begin{aligned} (W_{CPLST}, H_{CPLST}) &= \arg \min_{\substack{W \in \mathbb{R}^{d \times k} \\ H \in \mathbb{R}^{L \times k}}} \|XW - YH\|_F^2 + \|Y - YHH^T\|_F^2, \\ \text{s.t. } &H^T H = I_k. \end{aligned} \quad (13)$$

*Proof of Claim 2.* Let  $U_k[A] \Sigma_k[A] V_k[A]$  be the rank- $k$  truncated SVD approximation of a matrix  $A$ . In (Chen & Lin, 2012), the authors show that the closed form solution to (13) is

$$\begin{aligned} H_C &= V_k[Y^T X X^\dagger Y], \\ W_C &= X^\dagger Y H_C, \end{aligned}$$

where  $X^\dagger$  is the pseudo inverse of  $X$ . It follows from  $X^\dagger = V_X \Sigma_X^{-1} U_X^T$  that  $Y^T X X^\dagger Y = Y^T U_X U_X^T Y = M^T M$  and  $V_k[Y^T X X^\dagger Y] = V_k[M]$ . Thus, we have

$$\begin{aligned} Z^{CPLST} &= W_C H_C^T \\ &= X^\dagger Y H_C H_C^T \\ &= V_X^T \Sigma_X^{-1} U_X^T Y V_k[M] V_k[M]^T \\ &= V_X^T \Sigma_X^{-1} M V_k[M] V_k[M]^T \\ &= V_X^T \Sigma_X^{-1} M_k \end{aligned} \quad \square$$

## B. Algorithm Details

### B.1. Derivative Computations for Various Losses

Note that for the logistic and  $L_2$ -hinge loss in Table 5,  $Y_{ij}$  is assumed to be  $-1, +1$  instead of  $\{0, 1\}$ . Note that although  $L_2$ -hinge loss is not twice-differentiable, the sub-differential of  $\frac{\partial}{\partial b} \ell(a, b)$  still can be used for TRON to solve (6).

Table 5. Computation of  $\ell'(a, b)$  and  $\ell''(a, b)$  for different loss functions.

	$\ell(a, b)$	$\frac{\partial}{\partial b} \ell(a, b)$	$\frac{\partial^2}{\partial b^2} \ell(a, b)$
Squared loss	$\frac{1}{2}(a - b)^2$	$b - a$	1
Logistic loss	$\log(1 + e^{-ab})$	$\frac{-a}{1 + e^{-ab}}$	$\frac{-a^2 e^{-ab}}{(1 + e^{-ab})^2}$
$L_2$ -hinge loss	$(\max(0, 1 - ab))^2$	$-2a \max(0, 1 - ab)$	$2 \cdot \mathcal{I}[ab < 1]$

### B.2. Conjugate Gradient for Squared Loss

In Algorithm 3, we show the detailed conjugate gradient procedure used to solve (6) when the squared loss is used. Note that  $\nabla^2 g(\mathbf{w})$  is invariant to  $\mathbf{w}$  as (6) is a quadratic problem due to the squared loss function.

---

**Algorithm 3** Conjugate gradient for solving (6) with the squared loss

---

- Set initial  $\mathbf{w}_0, \mathbf{r}_0 = -\nabla g(\mathbf{w}_0), \mathbf{d}_0 = \mathbf{r}_0$ .
  - For  $t = 0, 1, 2, \dots$ 
    - If  $\|\mathbf{r}_t\|$  is small enough, then stop the procedure and return  $\mathbf{w}_t$ .
    - $\alpha_t = \frac{\mathbf{r}_t^T \mathbf{r}_t}{\mathbf{d}_t^T \nabla^2 g(\mathbf{w}_0) \mathbf{d}_t}$
    - $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{d}_t$
    - $\mathbf{r}_{t+1} = \mathbf{r}_t - \alpha_t \nabla^2 g(\mathbf{w}_0) \mathbf{d}_t$
    - $\beta_t = \frac{\mathbf{r}_{t+1}^T \mathbf{r}_{t+1}}{\mathbf{r}_t^T \mathbf{r}_t}$
    - $\mathbf{d}_{t+1} = \mathbf{r}_{t+1} + \beta_t \mathbf{d}_t$
-

## C. Analyzing Trace Norm-bounded Predictors

In this section, we shall provide a proof of Theorems 3 and 4. Our proof shall proceed by demonstrating a uniform convergence style bound for the empirical losses. More precisely, we shall show, for both trace norm as well as Frobenius regularizations, that with high probability, we have

$$\mathcal{L}(\hat{Z}) \leq \hat{\mathcal{L}}(\hat{Z}) + \epsilon.$$

Suppose  $Z^* \in \arg \min_{r(Z) \leq \lambda} \mathcal{L}(Z)$ , then a similar analysis will allow us to show, again with high probability,

$$\hat{\mathcal{L}}(Z^*) \leq \mathcal{L}(Z^*) + \epsilon.$$

Combining the two along with the fact that  $\hat{Z}$  is the empirical risk minimizer i.e.  $\hat{\mathcal{L}}(\hat{Z}) \leq \hat{\mathcal{L}}(Z^*)$  will yield the announced claim in the following form:

$$\mathcal{L}(\hat{Z}) \leq \mathcal{L}(Z^*) + 2\epsilon.$$

Thus, in the sequel, we shall only concentrate on proving the aforementioned uniform convergence bound. We shall denote the regularized class of predictors as  $\mathcal{Z} = \{Z \in \mathbb{R}^{d \times L}, r(Z) \leq \lambda\}$ , where  $r(Z) = \|Z\|_{\text{tr}}$  or  $r(Z) = \|Z\|_F$ . We shall also use the following shorthand for the loss incurred by the predictor on a specific label  $l \in [L]$ :  $\ell(\mathbf{y}_i^l, Z_l, \mathbf{x}) := \ell(\mathbf{y}_i^l, f^l(\mathbf{x}; Z))$ , where  $Z_l$  denotes the  $l^{\text{th}}$  column of the matrix  $Z$ .

We shall perform our analysis in several steps outlined below:

1. Step 1: In this step we shall show, by an application of McDiarmid's inequality, that with high probability, the excess risk of the learned predictor can be bounded by bounding the expected suprēmus deviation of empirical risks from population risks over the set of predictors in the class  $\mathcal{Z}$ .
2. Step 2: In this step we shall show that the expected suprēmus deviation can be bounded by a Rademacher average term.
3. Step 3: In this step we shall reduce the estimation of the Rademacher average term to the estimation of the spectral norm of a random matrix that we shall describe.
4. Step 4: Finally, we shall use tools from random matrix theory to bound the spectral norm of the random matrix.

We now give details of each of the steps in the following subsections:

### C.1. Step 1: Bounding Excess Risk by Expected Suprēmus Deviation

We will first analyze the case  $s = 1$  and will later show how to extend the analysis to  $s > 1$ . In this case, we receive  $n$  training points  $(\mathbf{x}_i, \mathbf{y}_i)$  and for each training point  $\mathbf{x}_i$ , we get to see the value of a random label  $l_i \in [L]$  i.e. we get to see the true value of  $\mathbf{y}_i^{l_i}$ . Thus, for any predictor  $Z \in \mathcal{Z}$ , the observed training loss is given by

$$\hat{\mathcal{L}}(Z) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i).$$

The population risk functional is given by

$$\mathcal{L}(Z) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}, l)} [\ell_l(\mathbf{y}^l, f^l(\mathbf{x}; Z))] = \mathbb{E}_{(\mathbf{x}, \mathbf{y}, l)} [\ell_l(\mathbf{y}^l, Z_l, \mathbf{x})]$$

We note here that our subsequent analysis shall hold even for non uniform distributions for sampling the labels. The definition of the population risk functional incorporates this. In case we have a uniform distribution over the labels, the above definition reduces to

$$\mathcal{L}(Z) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}, l)} [\ell_l(\mathbf{y}^l, Z_l, \mathbf{x})] = \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i, \tilde{l}_i)} \left[ \frac{1}{n} \sum_{i=1}^n \ell(\tilde{\mathbf{y}}_i^{\tilde{l}_i}, Z_{\tilde{l}_i}, \tilde{\mathbf{x}}_i) \right]$$

Given the above, we now analyze the excess risk i.e. the difference between the observed training loss  $\hat{\mathcal{L}}(\hat{Z})$  and the population risk  $\mathcal{L}(\hat{Z})$ .

$$\begin{aligned} \mathcal{L}(\hat{Z}) - \hat{\mathcal{L}}(\hat{Z}) &\leq \sup_{Z \in \mathcal{Z}} \left\{ \mathcal{L}(Z) - \hat{\mathcal{L}}(Z) \right\} \\ &= \sup_{Z \in \mathcal{Z}} \underbrace{\left\{ \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i, \tilde{l}_i)} \left[ \frac{1}{n} \sum_{i=1}^n \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{\tilde{l}_i}, \tilde{\mathbf{x}}_i) \right] - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right\}}_{g((\mathbf{x}_1, \mathbf{y}_1, l_1), \dots, (\mathbf{x}_n, \mathbf{y}_n, l_n))} \end{aligned}$$

Since all the label-wise loss functions are bounded, an arbitrary change in any  $(\mathbf{x}_i, \mathbf{y}_i)$  or any  $l_i$  should not perturb the expression  $g((\mathbf{x}_1, \mathbf{y}_1, l_1), \dots, (\mathbf{x}_n, \mathbf{y}_n, l_n))$  by more than  $\mathcal{O}(\frac{1}{n})$ . Thus, by an application of McDiarmid's inequality, we have, with probability at least  $1 - \delta$ ,

$$\mathcal{L}(\hat{Z}) - \hat{\mathcal{L}}(\hat{Z}) \leq \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i} [g((\mathbf{x}_1, \mathbf{y}_1, l_1), \dots, (\mathbf{x}_n, \mathbf{y}_n, l_n))] + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right)$$

Thus, we conclude that the excess risk of the learned predictor can be bounded by calculating the expected suprēmus deviation of empirical risks from population risks.

## C.2. Step 2: Bounding Expected Suprēmus Deviation by a Rademacher Average

We now analyze the expected suprēmus deviation. We have

$$\begin{aligned} &\mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i} [g((\mathbf{x}_1, \mathbf{y}_1, l_1), \dots, (\mathbf{x}_n, \mathbf{y}_n, l_n))] \\ &= \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i, \tilde{l}_i)} \left[ \frac{1}{n} \sum_{i=1}^n \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{\tilde{l}_i}, \tilde{\mathbf{x}}_i) \right] - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right\} \right] \\ &\leq \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i, \tilde{l}_i)} \left[ \frac{1}{n} \sum_{i=1}^n \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{\tilde{l}_i}, \tilde{\mathbf{x}}_i) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) \right] \right\} \right] \\ &\quad + \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) \right] - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right\} \right] \\ &= \mathbb{E}_{l_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i, \tilde{l}_i)} \left[ \frac{1}{n} \sum_{i=1}^n \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{\tilde{l}_i}, \tilde{\mathbf{x}}_i) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) \right] \right\} \right] \\ &\quad + \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) \right] - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right\} \right] \\ &\leq \mathbb{E}_{(l_i, \tilde{l}_i)} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{\tilde{l}_i}, Z_{\tilde{l}_i}, \tilde{\mathbf{x}}_i) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) \right] \right\} \right] \\ &\quad + \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i, (\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right\} \right] \\ &= \mathbb{E}_{(l_i, \tilde{l}_i), \epsilon_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i \left( \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{\tilde{l}_i}, Z_{\tilde{l}_i}, \tilde{\mathbf{x}}_i) \right] - \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) \right] \right) \right\} \right] \\ &\quad + \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i, (\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i), \epsilon_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i \left( \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) - \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right) \right\} \right] \\ &\leq 2 \mathbb{E}_{l_i, \epsilon_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i)} \left[ \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) \right] \right\} \right] + 2 \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i, \epsilon_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right\} \right] \end{aligned}$$

$$\begin{aligned}
 &\leq 2 \mathbb{E}_{(\tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i), l_i, \epsilon_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i \ell(\tilde{\mathbf{y}}_i^{l_i}, Z_{l_i}, \tilde{\mathbf{x}}_i) \right\} \right] + 2 \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i, \epsilon_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right\} \right] \\
 &\leq 4 \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i, \epsilon_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i \ell(\mathbf{y}_i^{l_i}, Z_{l_i}, \mathbf{x}_i) \right\} \right] \leq \frac{4C}{n} \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i), l_i, \epsilon_i} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \sum_{i=1}^n \epsilon_i \langle Z_{l_i}, \mathbf{x}_i \rangle \right\} \right] \\
 &= \frac{4C}{n} \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sup_{Z \in \mathcal{Z}} \langle Z, X_{\epsilon}^{\mathbf{1}} \rangle \right],
 \end{aligned}$$

where for any  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X}$ ,  $\mathbf{1} \in [L]^n$  and  $\epsilon \in \{-1, +1\}^n$ , we define the matrix  $X_{\epsilon}^{\mathbf{1}}$  as follows:

$$X_{\epsilon}^{\mathbf{1}} := \left[ \sum_{i \in I_1} \epsilon_i \mathbf{x}_i \quad \sum_{i \in I_2} \epsilon_i \mathbf{x}_i \quad \dots \quad \sum_{i \in I_L} \epsilon_i \mathbf{x}_i \right]$$

where for any  $l \in [L]$ , we define  $I_l := \{i : l_i = l\}$ . Note that in the last second last inequality we have used the contraction inequality for Rademacher averages (see [Ledoux & Talagrand, 2002](#), proof of Theorem 4.12). We also note that the above analysis also allows for separate label-wise loss functions, so long as they are all bounded and  $C$ -Lipschitz. For any matrix predictor class  $\mathcal{Z}$ , we define its Rademacher complexity as follows:

$$\mathcal{R}_n(\mathcal{Z}) := \frac{1}{n} \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sup_{Z \in \mathcal{Z}} \langle Z, X_{\epsilon}^{\mathbf{1}} \rangle \right]$$

We have thus established that with high probability,

$$\mathcal{L}(\hat{Z}) - \hat{\mathcal{L}}(\hat{Z}) \leq 4C\mathcal{R}_n(\mathcal{Z}) + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right).$$

We now establish that the same analysis also extends to situations wherein, for each training point we observe values of  $s$  labels instead. Thus, for each  $\mathbf{x}_i$ , we observe values for labels  $l_i^1, \dots, l_i^s$ . In this case the empirical loss is given by

$$\hat{\mathcal{L}}(Z) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^s \ell(\mathbf{y}_i^{l_i^j}, Z_{l_i^j}, \mathbf{x}_i)$$

The change in any  $\mathbf{x}_i$  leads to a perturbation of at most  $\mathcal{O}\left(\frac{s}{n}\right)$  whereas the change in any  $l_i^j$  leads to a perturbation of  $\mathcal{O}\left(\frac{1}{n}\right)$ . Thus the sum of squared perturbations is bounded by  $\frac{2s^2}{n}$ . Thus on application of the McDiarmid's inequality, we will be able to bound the excess risk by the following expected sup̄remus deviation term

$$\mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i, l_i^j)} \left[ \sup_{Z \in \mathcal{Z}} \left\{ s \mathbb{E}_{(\mathbf{x}, \mathbf{y}, l)} \left[ \ell_l(\mathbf{y}^l, Z_l, \mathbf{x}) \right] - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^s \ell(\mathbf{y}_i^{l_i^j}, Z_{l_i^j}, \mathbf{x}_i) \right\} \right]$$

plus a quantity that behaves like  $\mathcal{O}\left(s\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right)$ . We analyze the expected sup̄remus deviation term below:

$$\begin{aligned}
 &\mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i, l_i^j)} \left[ \sup_{Z \in \mathcal{Z}} \left\{ s \mathbb{E}_{(\mathbf{x}, \mathbf{y}, l)} \left[ \ell_l(\mathbf{y}^l, Z_l, \mathbf{x}) \right] - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^s \ell(\mathbf{y}_i^{l_i^j}, Z_{l_i^j}, \mathbf{x}_i) \right\} \right] \\
 &= \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i, l_i^j)} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \sum_{j=1}^s \left( \mathbb{E}_{(\mathbf{x}, \mathbf{y}, l)} \left[ \ell_l(\mathbf{y}^l, Z_l, \mathbf{x}) \right] - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i^{l_i^j}, Z_{l_i^j}, \mathbf{x}_i) \right) \right\} \right] \\
 &\leq \sum_{j=1}^s \mathbb{E}_{(\mathbf{x}_i, \mathbf{y}_i, l_i^j)} \left[ \sup_{Z \in \mathcal{Z}} \left\{ \mathbb{E}_{(\mathbf{x}, \mathbf{y}, l)} \left[ \ell_l(\mathbf{y}^l, Z_l, \mathbf{x}) \right] - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i^{l_i^j}, Z_{l_i^j}, \mathbf{x}_i) \right\} \right] \\
 &\leq \sum_{j=1}^s \frac{4C}{n} \mathbb{E}_{X, \mathbf{1}^j, \epsilon} \left[ \sup_{Z \in \mathcal{Z}} \langle Z, X_{\epsilon}^{\mathbf{1}^j} \rangle \right] = \frac{4Cs}{n} \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sup_{Z \in \mathcal{Z}} \langle Z, X_{\epsilon}^{\mathbf{1}} \rangle \right] = 4Cs\mathcal{R}_n(\mathcal{Z})
 \end{aligned}$$

and thus, it just suffices to prove bounds for the case where a single label is observed per point. As an aside, we note that the case  $s = 1$  resembles that of multi-task learning. However, multi-task learning is typically studied in a different learning model and mostly uses group regularization that is distinct from ours.

### C.3. Step 3: Estimating the Rademacher Average

We will now bound the following quantity:

$$\mathcal{R}_n(\mathcal{Z}) = \frac{1}{n} \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sup_{Z \in \mathcal{Z}} \langle Z, X_\epsilon^1 \rangle \right]$$

where  $X_\epsilon^1$  is as defined above. Approaches to bounding such Rademacher average terms usually resort to Martingale techniques (Kakade et al., 2008) or use of tools from convex analysis (Kakade et al., 2012) and decompose the Rademacher average term. However, such decompositions shall yield suboptimal results in our case. Our proposed approach will, instead involve an application of Hölder’s inequality followed by an application from results from random matrix theory to bound the spectral norm of a random matrix.

For simplicity of notation, for any  $l \in [L]$ , we denote  $V_l = \sum_{i \in I_l} \epsilon_i x_i$  and  $V := X_\epsilon^1 = [V_1 \ V_2 \ \dots \ V_L]$ . Also, for any  $l \in [L]$ , let  $n_l = |I_l|$  denote the number of training points for which values of the  $l^{\text{th}}$  label was observed i.e.  $n_l = \sum_{i=1}^n \mathbf{1}_{l_i=l}$ .

#### C.3.1. DISTRIBUTION INDEPENDENT BOUND

We apply Hölder’s inequality to get the following result:

$$\frac{1}{n} \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sup_{Z \in \mathcal{Z}} \langle Z, X_\epsilon^1 \rangle \right] \leq \frac{1}{n} \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sup_{Z \in \mathcal{Z}} \|Z\|_{\text{tr}} \|X_\epsilon^1\|_F \right] \leq \frac{1}{n} \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \lambda \|X_\epsilon^1\|_2 \right] \leq \frac{\lambda}{n} \sqrt{\mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \|X_\epsilon^1\|_2^2 \right]}$$

Then the following bound can be derived in a straightforward manner:

$$\begin{aligned} \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \|X_\epsilon^1\|_2^2 \right] &\leq \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \|X_\epsilon^1\|_F^2 \right] = \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sum_{l=1}^L \|V_l\|_2^2 \right] = \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sum_{l=1}^L \left\| \sum_{i \in I_l} \epsilon_i x_i \right\|_2^2 \right] \\ &= \mathbb{E}_{X, \mathbf{1}, \epsilon} \left[ \sum_{l=1}^L \sum_{i \in I_l} \|x_i\|_2^2 + \sum_{i \neq j \in I_l} \epsilon_i \epsilon_j \langle x_i, x_j \rangle \right] \\ &\leq \mathbb{E}_{\mathbf{1}} \left[ \sum_{l=1}^L n_l \mathbb{E} \left[ \|x\|_2^2 \right] \right] \leq \mathbb{E}_{\mathbf{1}} \left[ \sum_{l=1}^L n_l \right] = n \end{aligned}$$

where we have assumed, without loss of generality that  $\mathbb{E}_{x \sim \mathcal{D}} \left[ \|x\|_2^2 \right] \leq 1$ . This proves

$$\mathcal{R}_n(\mathcal{Z}) \leq \frac{\lambda}{\sqrt{n}},$$

which establishes Theorem 3. Note that the same analysis holds if  $Z$  is Frobenius norm regularized since we can apply the Hölder’s inequality for Frobenius norm instead and still get the same Rademacher average bound.

#### C.3.2. TIGHTER BOUNDS FOR TRACE NORM REGULARIZATION

Notice that in the above analysis, we did not exploit the fact that the top singular value of the matrix  $X_\epsilon^1$  could be much smaller than its Frobenius norm. However, there exist distributions where trace norm regularization enjoys better performance guarantees over Frobenius norm regularization. In order to better present our bounds, we model the data distribution  $\mathcal{D}$  on  $\mathcal{X}$  (or rather its marginal) more carefully. Let  $X := \mathbb{E} \left[ x x^\top \right]$  and suppose the distribution  $\mathcal{D}$  satisfies the following conditions:

1. The top singular value of  $X$  is  $\|X\|_2 = \sigma_1$
2. The matrix  $X$  has trace  $\text{tr}(X) = \Sigma$

3. The distribution on  $\mathcal{X}$  is sub-Gaussian i.e. for some  $\eta > 0$ , we have, for all  $\mathbf{v} \in \mathbb{R}^d$ ,

$$\mathbb{E} \left[ \exp(\mathbf{x}^\top \mathbf{v}) \right] \leq \exp \left( \|\mathbf{v}\|_2^2 \eta^2 / 2 \right)$$

In order to be consistent with previous results, we shall normalize the vectors  $x$  so that they are unit-norm *on expectation*. Since  $\mathbb{E} \left[ \|x\|_2^2 \right] = \text{tr}(X) = \Sigma$ , we wish to bound the Rademacher average as

$$\mathcal{R}_n(\mathcal{Z}) \leq \frac{1}{n\sqrt{\Sigma}} \mathbb{E}_{X,1,\epsilon} \left[ \sup_{Z \in \mathcal{Z}} \langle Z, X_\epsilon^1 \rangle \right]$$

In this case, it is possible to apply the Hölder's inequality as

$$\frac{1}{n\sqrt{\Sigma}} \mathbb{E}_{X,1,\epsilon} \left[ \sup_{Z \in \mathcal{Z}} \langle Z, X_\epsilon^1 \rangle \right] \leq \frac{1}{n\sqrt{\Sigma}} \mathbb{E}_{X,1,\epsilon} \left[ \sup_{Z \in \mathcal{Z}} \|Z\|_{\text{tr}} \|X_\epsilon^1\|_2 \right] \leq \frac{1}{n\sqrt{\Sigma}} \mathbb{E}_{X,1,\epsilon} \left[ \lambda \|X_\epsilon^1\|_2 \right] \leq \frac{\lambda}{n\sqrt{\Sigma}} \sqrt{\mathbb{E}_{X,1,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \right]}$$

Thus, in order to bound  $\mathcal{R}_n(\mathcal{Z})$ , it suffices to bound  $\mathbb{E}_{X,1,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \right]$ . In this case, since our object of interest is the spectral norm of the matrix  $X_\epsilon^1$ , we expect to get much better guarantees, for instance, in case the training points  $x \in \mathcal{X}$  are being sampled from some (near) isotropic distribution. We note that Frobenius norm regularization will not be able to gain any advantage in these situations since it would involve the Frobenius norm of the matrix  $X_\epsilon^1$  (as shown in the previous subsection) and thus, cannot exploit the fact that the spectral norm of this matrix is much smaller than its Frobenius norm.

#### C.4. Step 4: Calculating the Spectral norm of a Random Matrix

To bound  $\mathbb{E}_{X,1,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \right]$ , we first make some simplifications (we will take care of the normalizations later). For any  $l \in [L]$ , let the probability of the value for label  $l$  being observed be  $p_l \in (0, 1]$  such that  $\sum_l p_l = 1$ . Also let  $P = \max_{l \in [L]} p_l$  and  $p = \min_{l \in [L]} p_l$ . Call the event  $\mathcal{E}_{\max}$  as the event when  $n_l \leq 2P \cdot n$  for all  $l \in [L]$  i.e. every label will have at most  $2P \cdot n$  training points for which its value is seen. The following result shows that this is a high probability event:

**Lemma 1.** For any  $\delta > 0$ , if  $n \geq \frac{1}{2p^2} \log \frac{L}{\delta}$ , then with probability  $1 - \delta$ , we have

$$\mathbb{P}[\mathcal{E}_{\max}] \geq 1 - \delta$$

*Proof.* For any  $l \in [L]$ , an application of Chernoff's bound for Boolean random variables tells us that with probability at least  $1 - \exp(-2np_l^2)$ , we have  $n_l \leq 2p_l \cdot n \leq 2P \cdot n$ . Taking a union bound and using  $p_l \geq p$  finishes the proof.  $\square$

Conditioning on the event  $\mathcal{E}_{\max}$  shall allow us to get a control over the spectral norm of the matrix  $X_\epsilon^1$  by getting a bound on the sub-Gaussian norm of the individual columns of  $X_\epsilon^1$ . We show below, that conditioning on this event does not affect the Rademacher average calculations. A simple calculation shows that  $\mathbb{E}_{X,\epsilon} \left[ \left\| \|X_\epsilon^1\|_2^2 \mathbf{1} \right\| \right] \leq n\Sigma$ . If we have  $n > \frac{1}{2p^2} \log \frac{L\Sigma}{Pd(\eta^2 + \sigma_1)}$ , we have  $\mathbb{P}[\neg \mathcal{E}_{\max}] < \frac{Pd(\eta^2 + \sigma_1)}{\Sigma}$ . This gives us the following bound:

$$\begin{aligned} \mathbb{E}_{X,1,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \right] &= \mathbb{E}_{X,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \mid \mathcal{E}_{\max} \right] \mathbb{P}[\mathcal{E}_{\max}] + \mathbb{E}_{X,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \mid \neg \mathcal{E}_{\max} \right] (1 - \mathbb{P}[\mathcal{E}_{\max}]) \\ &= \mathbb{E}_{X,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \mid \mathcal{E}_{\max} \right] (1 - \delta) + \mathbb{E}_{X,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \mid \neg \mathcal{E}_{\max} \right] \delta \\ &\leq \mathbb{E}_{X,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \mid \mathcal{E}_{\max} \right] + n\Sigma \left( \frac{Pd(\eta^2 + \sigma_1)}{\Sigma} \right) \\ &\leq \mathcal{O} \left( \mathbb{E}_{X,\epsilon} \left[ \|X_\epsilon^1\|_2^2 \mid \mathcal{E}_{\max} \right] \right) \end{aligned}$$

where the last step follows since our subsequent calculations will show that  $\mathbb{E}_{X,\epsilon} \left[ \left\| X_\epsilon^1 \right\|_2^2 \middle| \mathcal{E}_{\max} \right] = \mathcal{O}(nPd(\eta^2 + \sigma_1))$ . Thus, it suffices to bound  $\mathbb{E}_{X,\epsilon} \left[ \left\| X_\epsilon^1 \right\|_2^2 \middle| \mathcal{E}_{\max} \right] = \mathbb{E}_{X,\epsilon} \left[ \left\| V \right\|_2^2 \middle| \mathcal{E}_{\max} \right]$ . For sake of brevity we will omit the conditioning term from now on.

For simplicity let  $A_l = \frac{V_l}{c}$  where  $c = \eta \cdot \sqrt{2P \cdot n}$  and  $A = [A_1 A_2 \dots A_L]$ . Thus

$$\mathbb{E}_{X,1,\epsilon} \left[ \left\| X_\epsilon^1 \right\|_2^2 \right] = c^2 \cdot \mathbb{E}_{X,1,\epsilon} \left[ \left\| A \right\|_2^2 \right]$$

We first bound the sub-Gaussian norm of the column vectors  $A_l$ . For any vector  $\mathbf{v} \in \mathbb{R}^d$ , we have:

$$\begin{aligned} \mathbb{E} \left[ \exp \left( A_l^\top \mathbf{v} \right) \right] &= \mathbb{E} \left[ \exp \left( \frac{1}{c} \sum_{i \in I_l} \epsilon_i \langle \mathbf{x}_i, \mathbf{v} \rangle \right) \right] \\ &= \left( \mathbb{E} \left[ \exp \left( \langle \mathbf{x}, \frac{1}{c} \epsilon \mathbf{v} \rangle \right) \right] \right)^{n_l} \\ &\leq \left( \exp \left( \left\| \frac{1}{c} \epsilon \mathbf{v} \right\|_2^2 \eta^2 / 2 \right) \right)^{n_l} \\ &= \exp \left( \frac{n_l}{2\eta^2 P \cdot n} \left\| \mathbf{v} \right\|_2^2 \eta^2 / 2 \right) \\ &\leq \exp \left( \left\| \mathbf{v} \right\|_2^2 / 2 \right) \end{aligned}$$

where, in the second step, we have used the fact that  $\mathbf{x}_i, \mathbf{x}_j$  and  $\epsilon_i, \epsilon_j$  are independent for  $i \neq j$ , in the third step we have used the sub-Gaussian properties of  $\mathbf{x}$  and in the fourth step, we have use the fact that the event  $\mathcal{E}_{\max}$  holds. This shows us that the sub-Gaussian norm of the column vector  $A_l$  is bounded i.e.  $\|A_l\|_{\psi_2} \leq 1$ .

We now proceed to bound  $\mathbb{E}_{X,\epsilon} \left[ \left\| A \right\|_2^2 \right] = \mathbb{E}_{X,\epsilon} \left[ \left\| A^\top \right\|_2^2 \right]$ . Our proof proceeds by an application of a Bernstein-type inequality followed by a covering number argument and finishing off by bounding the expectation in terms of the cumulative distribution function. The first two parts of the proof proceed on the lines of the proof of Theorem 5.39 in (Vershynin, 2012) For any fixed vector  $\mathbf{v} \in \mathcal{S}^{d-1}$ , the set of unit norm vectors in  $d$  dimensions, we have:

$$\|A\mathbf{v}\|_2^2 = \sum_{l=1}^L \langle A_l, \mathbf{v} \rangle^2 =: \sum_{l=1}^L Z_l^2$$

Now observe that conditioned on  $\mathbf{l}$ ,  $I_t \cap I_{t'} = \varphi$  if  $t \neq t'$  and thus, conditioned on  $\mathbf{l}$ , the variables  $Z_t, Z_{t'}$  are independent for  $t \neq t'$ . This will allow us to apply the following Bernstein-type inequality

**Theorem 5** ((Vershynin, 2012), Corollary 5.17). *Let  $X_1, \dots, X_N$  be independent centered sub-exponential variables with bounded sub-exponential norm i.e. for all  $i$ , we have  $\|X_i\|_{\psi_1} \leq B$  for some  $B > 0$ . Then for some absolute constant  $c_1 > 0$ , we have for any  $\epsilon > 0$ ,*

$$\mathbb{P} \left[ \sum_{i=1}^N X_i \geq \epsilon N \right] \leq \exp \left( -c_1 \min \left\{ \frac{\epsilon^2}{B^2}, \frac{\epsilon}{B} \right\} N \right).$$

To apply the above result, we will first bound expectation of the random variables  $Z_l^2$ .

$$\mathbb{E} \left[ Z_l^2 \right] = \mathbb{E} \left[ \langle A_l, \mathbf{v} \rangle^2 \right] = \mathbb{E} \left[ \left( \frac{1}{c} \sum_{i \in I_l} \epsilon_i \langle \mathbf{x}_i, \mathbf{v} \rangle \right)^2 \right] = \frac{n_l}{c^2} \mathbb{E} \left[ \langle \mathbf{x}, \mathbf{v} \rangle^2 \right] \leq \frac{n_l \sigma_1}{c^2} \leq \frac{\sigma_1}{\eta^2}$$

where the fourth inequality follows from definition of the top singular norm  $\sigma_1$  of  $X := \mathbb{E} \left[ \mathbf{x} \mathbf{x}^\top \right]$  and the last inequality follows from the event  $\mathcal{E}_{\max}$ . The above calculation gives us a bound on the expectation of  $Z_l^2$  which will be used to center it. Since we have already established  $\|A_l\|_{\psi_2} \leq 1$ , we automatically get  $\|Z_l\|_{\psi_2} \leq 1$ . Using standard inequalities between



the sub-exponential norm  $\|\cdot\|_{\psi_1}$  and the sub-Gaussian norm  $\|\cdot\|_{\psi_2}$  of random variables (for instance, see Vershynin, 2012, Lemma 5.14) we also have

$$\|Z_l^2 - \mathbb{E} \llbracket Z_l^2 \rrbracket\|_{\psi_1} \leq 2 \|Z_l^2\|_{\psi_1} \leq 4 \|Z_l\|_{\psi_2}^2 \leq 4.$$

Applying Theorem 5 to the variables  $X_l = Z_l^2 - \mathbb{E} \llbracket Z_l^2 \rrbracket$ , we get

$$\mathbb{P} \left[ \sum_{l=1}^L Z_l^2 - L \frac{\sigma_1}{\eta^2} \geq \epsilon L \right] \leq \exp(-c_1 L \min\{\epsilon^2, \epsilon\})$$

where  $c_1 > 0$  is an absolute constant. Thus with probability at least  $1 - \exp(-c_1 L \min\{\epsilon^2, \epsilon\})$ , for a fixed vector  $\mathbf{v} \in \mathcal{S}^{d-1}$ , we have the inequality

$$\|A\mathbf{v}\|_2^2 \leq \left( \frac{\sigma_1}{\eta^2} + \epsilon \right) L$$

Applying a union bound over a  $\frac{1}{4}$ -net  $\mathcal{N}_{1/4}$  over  $\mathcal{S}^{d-1}$  (which can be of size at most  $9^d$ ), we get that with probability at most  $1 - 9^d \exp(-c_1 L \min\{\epsilon^2, \epsilon\})$ , we have the above inequality for every vector  $\mathbf{v} \in \mathcal{N}_{1/4}$  as well. We note that this implies a bound on the spectral norm of the matrix  $A$  (see Vershynin, 2012, Lemma 5.4) and get the following bound

$$\|A\|_2^2 \leq 2 \left( \frac{\sigma_1}{\eta^2} + \epsilon \right) L$$

Put  $\epsilon = c_2 \cdot \frac{d}{L} + \frac{\epsilon'}{L}$  where  $c_2 = \max\left\{1, \frac{\ln 9}{c_1}\right\}$  and suppose  $d \geq L$ . Since  $c_2 \geq 1$ , we have  $\epsilon \geq 1$  which gives  $\min\{\epsilon, \epsilon^2\} = \epsilon$ . This gives us with probability at least  $1 - \exp(-c_1 \epsilon')$ ,

$$\|A\|_2^2 \leq 2 \left( L \frac{\sigma_1}{\eta^2} + c_2 d + \epsilon' \right)$$

Consider the random variable  $Y = \frac{\|A\|_2^2}{2} - L \frac{\sigma_1}{\eta^2} - c_2 d$ . Then we have  $\mathbb{P}[Y > \epsilon] \leq \exp(-c_1 \epsilon)$ . Thus we have

$$\mathbb{E} \llbracket Y \rrbracket = \int_0^\infty \mathbb{P}[Y > \epsilon] d\epsilon \leq \int_0^\infty \exp(-c_1 \epsilon) d\epsilon = \frac{1}{c_1}$$

This gives us

$$\mathbb{E} \llbracket \|A\|_2^2 \rrbracket \leq 2 \left( L \frac{\sigma_1}{\eta^2} + c_2 d + \frac{1}{c_1} \right)$$

and consequently,

$$\mathbb{E}_{X, \mathbf{1}, \epsilon} \llbracket \|X_\epsilon^1\|_2^2 \rrbracket = c^2 \cdot \mathbb{E}_{X, \mathbf{1}, \epsilon} \llbracket \|A\|_2^2 \rrbracket \leq 4\eta^2 P \cdot n \left( L \frac{\sigma_1}{\eta^2} + c_2 d + \frac{1}{c_1} \right) \leq \mathcal{O} \left( n\eta^2 P \left( d + L \frac{\sigma_1}{\eta^2} \right) \right) \leq \mathcal{O} \left( nPd(\eta^2 + \sigma_1) \right)$$

where the last step holds when  $d \geq L$ . Thus, we are able to bound the Rademacher averages, for some absolute constant  $c_3$  as

$$\mathcal{R}_n(\mathcal{Z}) \leq \frac{\lambda}{n\sqrt{\Sigma}} \sqrt{\left\| \mathbb{E}_{X, \mathbf{1}, \epsilon} \llbracket \|X_\epsilon^1\|_2^2 \rrbracket \right\|} \leq c_3 \lambda \sqrt{\frac{Pd(\eta^2 + \sigma_1)}{n\Sigma}},$$

which allows us to make the following claim:

**Theorem 6.** *Suppose we learn a predictor using the trace norm regularized formulation  $\hat{Z} = \arg \inf_{\|Z\|_{\text{tr}} \leq \lambda} \hat{\mathcal{L}}(Z)$  over a set of  $n$  training points. Further suppose that, for any  $l \in [L]$ , the probability of observing the value of label  $l$  is given by  $p_l$  and let  $P = \max_{l \in [L]} p_l$ . Then with probability at least  $1 - \delta$ , we have*

$$\mathcal{L}(\hat{Z}) \leq \arg \inf_{\|Z\|_{\text{tr}} \leq \lambda} \mathcal{L}(Z) + \mathcal{O} \left( s\lambda \sqrt{\frac{dP(\eta^2 + \sigma_1)}{n\Sigma}} \right) + \mathcal{O} \left( s \sqrt{\frac{\log \frac{1}{\delta}}{n}} \right),$$

where the terms  $\eta, \sigma_1, \Sigma$  are defined by the data distribution as before.

Essentially, the above result indicates that if some label is observed too often, as would be the case when  $P = \Omega(1)$ , we get no benefit from trace norm regularization since this is akin to a situation with fully observed labels. However, if the distribution on the labels is close to uniform i.e.  $P = \mathcal{O}(\frac{1}{L})$ , the above calculation lets us bound the Rademacher average, and consequently, the excess risk as

$$\mathcal{R}_n(\mathcal{Z}) \leq c_3 \lambda \sqrt{\frac{d(\eta^2 + \sigma_1)}{nL\Sigma}},$$

thus proving the first part of Theorem 4.

We now notice that However, in case our data distribution is near isotropic, i.e.  $\Sigma \gg \sigma_1$ , then this result gives us superior bounds. For instance, if the data points are generated from a standard normal distribution, then we have  $\sigma_1 = 1$ ,  $\Sigma = d$  and  $\eta = 1$  using which we can bound the Rademacher average term as

$$\mathcal{R}_n(\mathcal{Z}) \leq c_3 \lambda \sqrt{\frac{2}{nL}},$$

which gives us the second part of Theorem 4.

## D. Lower Bounds for Uniform Convergence-based Proofs

In this section, we show that our analysis for Theorems 3 and 4 are essentially tight. In particular, we show for each case, a data distribution such that the deviation of the empirical losses from the population risks is, up to a constant factor, the same as predicted by the results. We state these lower bounds in two separate subsections below:

### D.1. Lower Bound for Trace Norm Regularization

In this section we shall show that for general distribution, Theorem 3 is tight. Recall that Theorem 3 predicts that for a predictor  $\hat{Z}$  learned using a trace norm regularized formulation satisfies, with constant probability (i.e.  $\delta = \Omega(1)$ ),

$$\mathcal{L}(\hat{Z}) \leq \hat{\mathcal{L}}(\hat{Z}) + \mathcal{O}\left(\lambda \sqrt{\frac{1}{n}}\right),$$

where, for simplicity as well as w.l.o.g., we have assumed  $s = 1$ . We shall show that this result is tight by demonstrating the following lower bound:

**Claim 7.** *There exists a data-label distribution and a loss function such that the empirical risk minimizer learned as  $\hat{Z} = \arg \inf_{\|Z\|_{\text{tr}} \leq \lambda} \hat{\mathcal{L}}(Z)$  has, with constant probability, its population risk lower bounded by*

$$\mathcal{L}(\hat{Z}) \geq \hat{\mathcal{L}}(\hat{Z}) + \Omega\left(\lambda \sqrt{\frac{1}{n}}\right),$$

thus establishing the tightness claim. Our proof will essentially demonstrate this by considering a non-isotropic data distribution (since, for isotropic distributions, Theorem 4 shows that a tighter upper bound is actually possible). For simplicity, and w.l.o.g., we will prove the result for  $\lambda = 1$ . Let  $\boldsymbol{\mu} \in \mathbb{R}^d$  be a fixed unit vector and consider the following data distribution

$$\mathbf{x}_i = \zeta_i \boldsymbol{\mu},$$

where  $\zeta_i$  are independent Rademacher variables and a trivial label distribution

$$\mathbf{y}_i = \mathbf{1},$$

where  $\mathbf{1} \in \mathbb{R}^L$  is the all-ones vector. Note that the data distribution satisfies  $\mathbb{E}[\|\mathbf{x}\|_2^2] = 1$  and thus, satisfies the assumptions of Theorem 3. Let  $\omega_i^l = 1$  iff the label  $l$  is observed for the  $i^{\text{th}}$  training point. Note that for any  $i$ , we have  $\sum_{l=1}^L \omega_i^l = 1$  and that for any  $l \in [L]$ ,  $\omega_i^l = 1$  with probability  $1/L$ . Also consider the following loss function

$$\ell(\mathbf{y}^l, f^l(\mathbf{x}; Z)) = \langle Z_l, \mathbf{y}^l \mathbf{x} \rangle$$

Let

$$\hat{Z} = \arg \inf_{\|Z\|_{\text{tr}} \leq 1} \hat{\mathcal{L}}(Z) = \arg \inf_{\|Z\|_{\text{tr}} \leq 1} \frac{1}{n} \langle Z, \boldsymbol{\mu} \mathbf{v}^\top \rangle$$

where  $\mathbf{v}$  is the vector

$$\mathbf{v} = \left[ \sum_{i=1}^n \zeta_i \omega_i^1 \quad \sum_{i=1}^n \zeta_i \omega_i^2 \quad \dots \quad \sum_{i=1}^n \zeta_i \omega_i^L \right]$$

Clearly, since  $x$  is a centered distribution and  $\ell$  is a linear loss function,  $\mathcal{L}(\hat{Z}) = 0$ . However, by Hölder's inequality, we also have

$$\hat{Z} = -\frac{\boldsymbol{\mu} \mathbf{v}^\top}{\|\mathbf{v}\|_2},$$

and thus,  $\hat{\mathcal{L}}(\hat{Z}) = -\frac{1}{n} \|\mathbf{v}\|_2$  since  $\|\boldsymbol{\mu}\|_2 = 1$ . The following lemma shows that with constant probability,  $\|\mathbf{v}\|_2 \geq \sqrt{n/2}$  which shows that  $\mathcal{L}(\hat{Z}) \geq \hat{\mathcal{L}}(\hat{Z}) + \Omega\left(\sqrt{\frac{1}{n}}\right)$ , thus proving the lower bound.

**Lemma 2.** *With probability at least 3/4, we have  $\|\mathbf{v}\|_2^2 \geq n/2$ .*

*Proof.* We have

$$\begin{aligned} \|\mathbf{v}\|_2^2 &= \sum_{l=1}^L \left( \sum_{i=1}^n \zeta_i \omega_i^l \right)^2 = \sum_{l=1}^L \sum_{i=1}^n \omega_i^l + \sum_{l=1}^L \sum_{i \neq j} \zeta_i \omega_i^l \zeta_j \omega_j^l \\ &= n + \sum_{i \neq j} \zeta_i \zeta_j \langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_j \rangle = n + W, \end{aligned}$$

where  $\boldsymbol{\omega}_i = [\omega_i^1, \omega_i^2, \dots, \omega_i^L]$ . Now clearly  $\mathbb{E}[W] = 0$  and as the following calculation shows,  $\mathbb{E}[W^2] \leq 2n^2/L$  which, by an application of Tchebysheff's inequality, gives us, for  $L > 32$ , with probability at least 3/4,  $|W| \leq n/2$  and consequently  $\|\mathbf{v}\|_2^2 \geq n/2$ . We give an estimation of the variance of  $Z$  below.

$$\begin{aligned} \mathbb{E}[W^2] &= \mathbb{E} \left[ \sum_{i_1 \neq j_1, i_2 \neq j_2} \zeta_{i_1} \zeta_{j_1} \langle \boldsymbol{\omega}_{i_1}, \boldsymbol{\omega}_{j_1} \rangle \zeta_{i_2} \zeta_{j_2} \langle \boldsymbol{\omega}_{i_2}, \boldsymbol{\omega}_{j_2} \rangle \right] \\ &= 2\mathbb{E} \left[ \sum_{i \neq j} \langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_j \rangle^2 \right] = 2n(n-1)\mathbb{E}[\langle \boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \rangle] \leq \frac{2n^2}{L}, \end{aligned}$$

where we have used the fact that  $\langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_j \rangle^2 = \langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_j \rangle$  since  $\langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_j \rangle = 0$  or 1, and that  $\mathbb{E}[\langle \boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \rangle] = \frac{1}{L}$  since that is the probability of the same label getting observed for  $\mathbf{x}_1$  and  $\mathbf{x}_2$ .  $\square$

## D.2. Lower Bound for Frobenius Norm Regularization

In this section, we shall prove that even for isotropic distributions, Frobenius norm regularization cannot offer  $\mathcal{O}\left(\frac{1}{\sqrt{nL}}\right)$ -style bounds as offered by trace norm regularization.

**Claim 8.** *There exists an isotropic, sub-Gaussian data distribution and a loss function such that the empirical risk minimizer learned as  $\hat{Z} = \arg \inf_{\|Z\|_F \leq \lambda} \hat{\mathcal{L}}(Z)$  has, with constant probability, its population risk lower bounded by*

$$\mathcal{L}(\hat{Z}) \geq \hat{\mathcal{L}}(\hat{Z}) + \Omega\left(\lambda \sqrt{\frac{1}{n}}\right),$$

whereas an empirical risk minimizer learned as  $\hat{Z} = \arg \inf_{\|Z\|_{\text{tr}} \leq \lambda} \hat{\mathcal{L}}(Z)$  over the same distribution has, with probability at least  $1 - \delta$ , its population risk bounded by

$$\mathcal{L}(\hat{Z}) \leq \hat{\mathcal{L}}(\hat{Z}) + \mathcal{O}\left(\lambda \sqrt{\frac{1}{nL}}\right) + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right).$$

We shall again prove this result for  $\lambda = 1$ . We shall retain the distribution over labels as well as the loss function from our previous discussion in Appendix D.1. We shall also reuse  $\omega_i^l$  to denote the label observation pattern. We shall however use Rademacher vectors to define the data distribution i.e. each of the  $d$  coordinates of the vector  $\mathbf{x}$  obeys the law

$$r \sim \frac{1}{2}(\mathbb{1}_{\{r=1\}} + \mathbb{1}_{\{r=-1\}}).$$

Thus we sample  $\mathbf{x}_i$  as

$$\mathbf{x}_i = \frac{1}{\sqrt{d}} [r_i^1, r_i^2, \dots, r_i^d],$$

where each coordinate is independently sampled. We now show that this distribution satisfies the assumptions of Theorem 4. We have  $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \frac{1}{d} \cdot \mathbb{I}$  where  $\mathbb{I}$  is the  $d \times d$  identity matrix. Thus  $\sigma_1 = \frac{1}{d}$  and  $\Sigma = 1$ . We also have, for any  $\mathbf{v} \in \mathbb{R}^d$ ,

$$\begin{aligned} \mathbb{E}[\exp(\mathbf{x}^\top \mathbf{v})] &= \mathbb{E}\left[\exp\left(\sum_{j=1}^d \mathbf{x}^j \mathbf{v}^j\right)\right] = \prod_{j=1}^d \mathbb{E}[\exp(\mathbf{x}^j \mathbf{v}^j)] \\ &= \prod_{j=1}^d \frac{1}{2} \left( \exp\left(\frac{1}{\sqrt{d}} \mathbf{v}^j\right) + \exp\left(-\frac{1}{\sqrt{d}} \mathbf{v}^j\right) \right) \\ &= \prod_{j=1}^d \cosh\left(\frac{1}{\sqrt{d}} \mathbf{v}^j\right) \leq \prod_{j=1}^d \exp\left(\frac{1}{d} (\mathbf{v}^j)^2\right) \\ &= \exp\left(\sum_{j=1}^d \frac{1}{d} (\mathbf{v}^j)^2\right) = \exp\left(\frac{1}{d} \|\mathbf{v}\|_2^2\right), \end{aligned}$$

where the second equality uses the independence of the coordinates of  $\mathbf{x}$ . Thus we have  $\eta^2 = \frac{2}{d}$ . Thus, this distribution fulfills all the preconditions of Theorem 4. Note that had trace norm regularization been applied, then by applying Theorem 4, we would have gotten an excess error of

$$\mathcal{O}\left(\sqrt{\frac{d(\eta^2 + \sigma_1)}{nL\Sigma}}\right) = \mathcal{O}\left(\sqrt{\frac{d(2/d + 1/d)}{nL \cdot 1}}\right) = \mathcal{O}\left(\sqrt{\frac{1}{nL}}\right)$$

whereas, as the calculation given below shows, Frobenius norm regularization cannot guarantee an excess risk better than  $\mathcal{O}\left(\sqrt{\frac{1}{n}}\right)$ . Suppose we do perform Frobenius norm regularization in this case. Then we have

$$\hat{Z} = \arg \inf_{\|Z\|_F \leq 1} \hat{\mathcal{L}}(Z) = \arg \inf_{\|Z\|_F \leq 1} \frac{1}{n} \langle Z, X \rangle,$$

where  $X$  is the matrix

$$X = \begin{bmatrix} \sum_{i=1}^L \omega_i^1 \mathbf{x}_i & \sum_{i=1}^L \omega_i^2 \mathbf{x}_i & \dots & \sum_{i=1}^L \omega_i^L \mathbf{x}_i \end{bmatrix}.$$

As before,  $\mathcal{L}(\hat{Z}) = 0$  since the data distribution is centered and the loss function is linear. By a similar application of Hölder's inequality, we can also get

$$\hat{Z} = -\frac{X}{\|X\|_F},$$

and thus,  $\hat{\mathcal{L}}(\hat{Z}) = -\frac{1}{n} \|X\|_F$ . The following lemma shows that with constant probability,  $\|X\|_F \geq \sqrt{n/2}$  which shows that  $\mathcal{L}(\hat{Z}) \geq \hat{\mathcal{L}}(\hat{Z}) + \Omega\left(\sqrt{\frac{1}{n}}\right)$ , thus proving the claimed inability of Frobenius norm regularization to give  $\mathcal{O}\left(\frac{1}{\sqrt{nL}}\right)$ -style bounds even for isotropic distributions.

**Lemma 3.** *With probability at least 3/4, we have  $\|X\|_F^2 \geq n/2$ .*

*Proof.* We have

$$\begin{aligned}\|X\|_F^2 &= \sum_{l=1}^L \left\| \sum_{i=1}^n \omega_i^l \mathbf{x}_i \right\|_2^2 = \sum_{l=1}^L \sum_{i=1}^n \omega_i^l \|\mathbf{x}_i\|_2^2 + \sum_{l=1}^L \sum_{i \neq j} \omega_i^l \omega_j^l \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ &= \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 + \sum_{i \neq j} \langle \mathbf{x}_i, \mathbf{x}_j \rangle \langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_j \rangle = n + W\end{aligned}$$

where as before,  $\boldsymbol{\omega}_i = [\omega_i^1, \omega_i^2, \dots, \omega_i^L]$ . We will, in the sequel prove that  $|W| \leq n/2$ , thus establishing the claim. Clearly  $\mathbb{E}[W] = 0$  and as the following calculation shows,  $\mathbb{E}[W^2] \leq 2n^2/Ld$  which, by an application of Tchebysheff's inequality, gives us, for  $Ld > 32$ , with probability at least  $3/4$ ,  $|W| \leq n/2$  and consequently  $\|X\|_F^2 \geq n/2$ . We give an estimation of the variance of  $W$  below.

$$\begin{aligned}\mathbb{E}[W^2] &= \mathbb{E} \left[ \sum_{i_1 \neq j_1, i_2 \neq j_2} \langle \mathbf{x}_{i_1}, \mathbf{x}_{j_1} \rangle \langle \boldsymbol{\omega}_{i_1}, \boldsymbol{\omega}_{j_1} \rangle \langle \mathbf{x}_{i_2}, \mathbf{x}_{j_2} \rangle \langle \boldsymbol{\omega}_{i_2}, \boldsymbol{\omega}_{j_2} \rangle \right] \\ &= 2\mathbb{E} \left[ \sum_{i \neq j} \langle \mathbf{x}_i, \mathbf{x}_j \rangle^2 \langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_j \rangle^2 \right] = 2n(n-1)\mathbb{E}[\langle \mathbf{x}_1, \mathbf{x}_2 \rangle^2 \langle \boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \rangle^2] \\ &= 2n(n-1)\mathbb{E}[\langle \mathbf{x}_1, \mathbf{x}_2 \rangle^2] \mathbb{E}[\langle \boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \rangle^2] \leq \frac{2n^2}{Ld},\end{aligned}$$

where we have used the fact that data points and label patterns are sampled independently.  $\square$

## E. More Experimental Results

### E.1. Evaluation Criteria

Given a test set  $\{\mathbf{x}_i, \mathbf{y}_i : i = 1, \dots, n\}$ , three criteria are used to evaluate the performance for an real-valued predictor  $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$ :

- **Top- $K$  accuracy:** for each instance, we select the  $K$  labels with the largest decision values for prediction. The average accuracy among all instances are reported as the top- $K$  accuracy.
- **Hamming-loss:** for each pair of instance  $\mathbf{x}$  and label index  $j$ , we round the decision value  $f^j(\mathbf{x})$  to 0 or 1.

$$\text{Hamming Loss} = \frac{1}{nL} \sum_{i=1}^n \sum_{j=1}^L I[\text{round}(f^j(\mathbf{x}_i)) \neq \mathbf{y}_i^j]$$

- **Average AUC:** we follow (Bucak et al., 2009) to calculate area under ROC curve for each instance and report the average AUC among all test instances.

## E.2. Speedup Results Due to Multi-core Computation

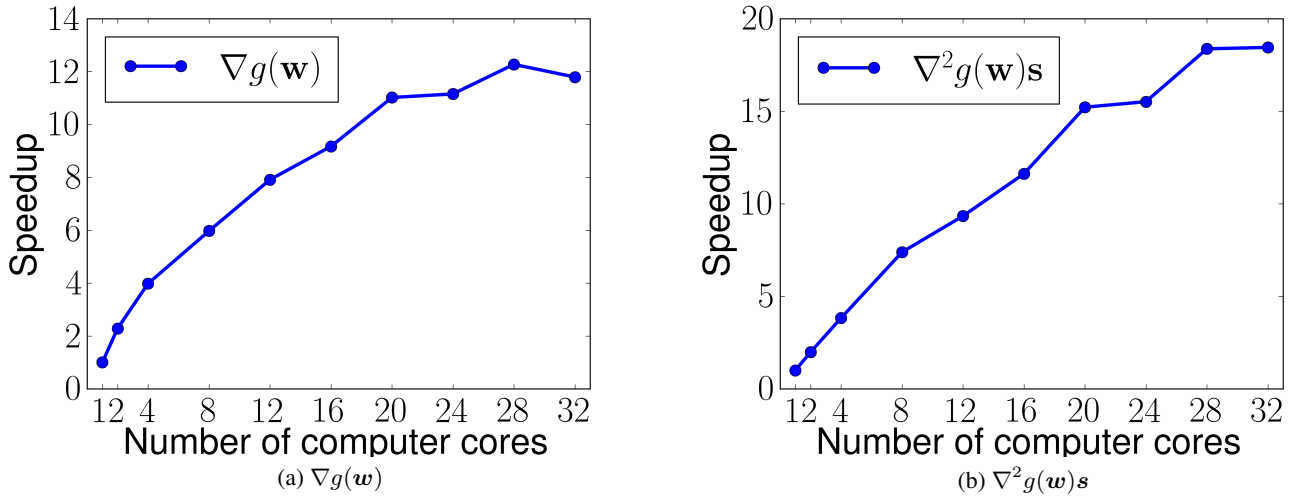


Figure 3. Speedup results for our proposed fast gradient calculation and Hessian-vector multiplication.

## E.3. Detailed Results with Full Labels

- Table 6 shows the top-1 accuracy results for the case with fully observed labels.
- Table 7 shows the top-3 accuracy results for the case with fully observed labels.
- Table 8 shows the top-5 accuracy results for the case with fully observed labels.
- Table 9 shows the Hamming loss results for the case with fully observed labels.
- Table 10 shows the average AUC results for the case with fully observed labels.

## E.4. Detailed Results with Missing Labels

- Table 11 shows the top-1 accuracy results for the case with various missing ratios and dimension reduction rates.
- Table 12 shows the top-3 accuracy results for the case with various missing ratios and dimension reduction rates.
- Table 13 shows the top-5 accuracy results for the case with various missing ratios and dimension reduction rates.
- Table 14 shows the Hamming loss results for the case with various missing ratios and dimension reduction rates.
- Table 15 shows the average AUC results for the case with various missing ratios and dimension reduction rates.

Table 6. Comparison for dimensionality reductions approach on fully observed  $Y$  with various rank. SQ for squared loss, LR for logistic loss, SH for squared hinge loss, and WAR for weighted approximated-rank loss

	$k/L$	Top-1 Accuracy					
		SQ	LEML LR	SH	BCS SQ	CPLST SQ	WSABIE WAR
bibtex	20%	<b>58.33</b>	46.20	46.52	41.43	55.55	48.51
	40%	<b>60.99</b>	50.78	40.68	54.63	58.73	52.37
	60%	<b>61.99</b>	51.37	39.24	57.53	60.36	51.45
	80%	<b>63.38</b>	52.64	39.96	59.76	62.31	53.04
	100%	<b>63.94</b>	53.76	38.41	60.24	63.02	53.24
autofood	20%	86.84	84.21	<b>89.47</b>	68.42	52.63	47.37
	40%	<b>92.11</b>	89.47	<b>92.11</b>	28.95	55.26	86.84
	60%	73.68	<b>89.47</b>	86.84	71.05	52.63	65.79
	80%	<b>94.74</b>	89.47	89.47	81.58	57.89	78.95
	100%	81.58	<b>89.47</b>	86.84	84.21	57.89	60.53
compphys	20%	92.50	87.50	<b>97.50</b>	70.00	52.50	65.00
	40%	<b>95.00</b>	92.50	<b>95.00</b>	65.00	50.00	47.50
	60%	<b>95.00</b>	92.50	<b>95.00</b>	72.50	47.50	70.00
	80%	95.00	87.50	<b>97.50</b>	75.00	50.00	45.00
	100%	95.00	<b>97.50</b>	<b>97.50</b>	67.50	50.00	52.50
delicious	20%	<b>67.16</b>	57.39	61.07	59.50	66.53	48.35
	40%	<b>66.66</b>	51.62	56.20	61.16	66.25	47.25
	60%	<b>66.28</b>	50.96	51.59	63.08	66.22	47.38
	80%	<b>66.25</b>	51.55	49.11	62.10	66.22	45.59
	100%	<b>66.28</b>	50.83	46.53	63.45	66.22	46.25

Table 7. Comparison for dimensionality reductions approach on fully observed  $Y$  with various rank. SQ for squared loss, LR for logistic loss, SH for squared hinge loss, and WAR for weighted approximated-rank loss

	$k/L$	Top-3 Accuracy					
		SQ	LEML LR	SH	BCS SQ	CPLST SQ	WSABIE WAR
bibtex	20%	<b>34.16</b>	25.65	27.37	21.74	31.99	28.77
	40%	<b>36.53</b>	28.20	24.81	28.95	34.53	30.05
	60%	<b>38.00</b>	28.68	23.26	32.25	36.01	31.11
	80%	<b>38.58</b>	29.42	23.04	34.09	36.75	31.21
	100%	<b>38.41</b>	30.25	22.36	34.87	36.91	31.24
autofood	20%	<b>81.58</b>	80.70	<b>81.58</b>	53.51	42.98	66.67
	40%	76.32	<b>80.70</b>	78.95	50.88	42.11	70.18
	60%	70.18	80.70	<b>81.58</b>	64.91	41.23	60.53
	80%	80.70	80.70	<b>85.09</b>	73.68	42.98	72.81
	100%	75.44	80.70	<b>82.46</b>	65.79	42.98	64.04
compphys	20%	<b>80.00</b>	<b>80.00</b>	<b>80.00</b>	42.50	40.83	49.17
	40%	<b>80.00</b>	78.33	79.17	60.00	37.50	39.17
	60%	<b>80.00</b>	<b>80.00</b>	<b>80.00</b>	51.67	39.17	49.17
	80%	80.00	78.33	<b>80.83</b>	53.33	39.17	52.50
	100%	80.00	79.17	<b>81.67</b>	62.50	39.17	56.67
delicious	20%	<b>61.20</b>	53.68	57.27	53.01	61.13	42.87
	40%	<b>61.23</b>	49.13	52.95	56.20	61.08	42.05
	60%	<b>61.15</b>	46.76	49.58	57.07	61.09	42.22
	80%	<b>61.13</b>	48.06	47.34	57.09	61.09	42.01
	100%	<b>61.12</b>	46.11	45.92	57.91	61.09	41.34

Table 8. Comparison for dimensionality reductions approach on fully observed  $Y$  with various rank. SQ for squared loss, LR for logistic loss, SH for squared hinge loss, and WAR for weighted approximated-rank loss

	$k/L$	Top-5 Accuracy					
		LEML			BCS	CPLST	WSABIE
		SQ	LR	SH	SQ	SQ	WAR
bibtex	20%	<b>24.49</b>	19.24	20.33	15.39	23.11	21.92
	40%	<b>26.84</b>	20.61	18.54	19.95	24.96	22.47
	60%	<b>27.66</b>	20.99	17.61	22.43	26.07	23.33
	80%	<b>28.20</b>	21.48	17.46	24.07	26.47	23.44
	100%	<b>28.01</b>	22.03	16.83	24.48	26.47	23.44
autofood	20%	<b>81.05</b>	80.00	75.79	44.21	36.84	66.32
	40%	73.68	<b>78.42</b>	76.84	51.05	36.32	66.84
	60%	69.47	<b>78.95</b>	78.42	57.37	36.32	60.53
	80%	74.74	78.95	<b>80.53</b>	68.95	36.84	66.84
	100%	72.63	78.42	<b>83.16</b>	62.11	36.84	61.58
compphys	20%	72.00	<b>73.50</b>	72.50	32.50	37.50	46.00
	40%	73.00	74.00	<b>74.50</b>	54.50	35.50	41.00
	60%	73.00	<b>74.00</b>	<b>74.00</b>	43.50	34.50	44.00
	80%	73.00	73.00	<b>74.00</b>	47.50	36.00	46.50
	100%	72.50	72.50	<b>73.00</b>	54.50	36.00	49.50
delicious	20%	<b>56.46</b>	49.46	52.94	47.91	56.30	39.79
	40%	<b>56.39</b>	45.66	49.54	51.61	56.28	39.27
	60%	<b>56.28</b>	43.22	46.93	52.85	56.23	38.97
	80%	<b>56.27</b>	44.03	45.43	52.92	56.23	39.27
	100%	<b>56.27</b>	42.11	44.24	53.28	56.23	38.41

Table 9. Comparison for dimensionality reductions approach on fully observed  $Y$  with various rank. SQ for squared loss, LR for logistic loss, SH for squared hinge loss, and WAR for weighted approximated-rank loss

	$k/L$	Hamming Loss				
		LEML			BCS	CPLST
		SQ	LR	SH	SQ	SQ
bibtex	20%	<b>0.0126</b>	0.0211	0.0231	0.0150	<b>0.0127</b>
	40%	<b>0.0124</b>	0.0240	0.0285	0.0140	0.0126
	60%	<b>0.0123</b>	0.0233	0.0320	0.0132	0.0126
	80%	<b>0.0123</b>	0.0242	0.0343	0.0130	0.0125
	100%	<b>0.0122</b>	0.0236	0.0375	0.0129	0.0125
autofood	20%	<b>0.0547</b>	0.0621	0.0588	0.0846	0.0996
	40%	0.0590	0.0608	<b>0.0578</b>	0.0846	0.0975
	60%	0.0593	0.0611	<b>0.0586</b>	0.0838	0.0945
	80%	0.0572	0.0611	<b>0.0569</b>	1.0000	0.0944
	100%	0.0603	0.0617	<b>0.0586</b>	1.0000	0.0944
compphys	20%	0.0457	0.0470	<b>0.0456</b>	0.0569	0.0530
	40%	<b>0.0454</b>	0.0466	0.0456	0.0569	0.0526
	60%	<b>0.0454</b>	0.0469	0.0460	0.0569	0.0530
	80%	0.0464	0.0484	<b>0.0456</b>	0.0569	0.0755
	100%	0.0453	0.0469	<b>0.0450</b>	0.0569	0.0755
delicious	20%	<b>0.0181</b>	0.0196	0.0187	0.0189	<b>0.0182</b>
	40%	<b>0.0181</b>	0.0221	0.0198	0.0186	<b>0.0182</b>
	60%	<b>0.0182</b>	0.0239	0.0207	0.0187	<b>0.0182</b>
	80%	<b>0.0182</b>	0.0253	0.0212	0.0186	<b>0.0182</b>
	100%	<b>0.0182</b>	0.0260	0.0216	0.0186	<b>0.0182</b>



Table 10. Comparison for dimensionality reductions approach on fully observed  $Y$  with various rank. SQ for squared loss, LR for logistic loss, SH for squared hinge loss, and WAR for weighted approximated-rank loss

	$k/L$	Average AUC					
		SQ	LEML LR	SH	BCS SQ	CPLST SQ	WSABIE WAR
bibtex	20%	0.8910	0.8677	0.8541	0.7875	0.8657	<b>0.9055</b>
	40%	0.9015	0.8809	0.8467	0.8263	0.8802	<b>0.9092</b>
	60%	0.9040	0.8861	0.8505	0.8468	0.8854	<b>0.9089</b>
	80%	0.9035	0.8875	0.8491	0.8560	0.8882	<b>0.9164</b>
	100%	0.9024	0.8915	0.8419	0.8614	0.8878	<b>0.9182</b>
autofood	20%	0.9565	<b>0.9598</b>	0.9424	0.7599	0.7599	0.8779
	40%	0.9277	<b>0.9590</b>	0.9485	0.7994	0.7501	0.8806
	60%	0.8815	<b>0.9582</b>	0.9513	0.8282	0.7552	0.8518
	80%	0.9280	<b>0.9588</b>	0.9573	0.8611	0.7538	0.8520
	100%	0.9361	<b>0.9581</b>	0.9561	0.8718	0.7539	0.8471
compphys	20%	0.9163	0.9223	<b>0.9274</b>	0.6972	0.7692	0.8212
	40%	<b>0.9199</b>	0.9157	0.9191	0.7881	0.7742	0.8066
	60%	<b>0.9179</b>	0.9143	0.9098	0.7705	0.7705	0.8040
	80%	0.9187	0.9003	<b>0.9220</b>	0.7820	0.7806	0.7742
	100%	<b>0.9205</b>	0.9040	0.8977	0.7884	0.7804	0.7951
delicious	20%	0.8854	0.8588	<b>0.8894</b>	0.7308	0.8833	0.8561
	40%	0.8827	0.8534	<b>0.8868</b>	0.7635	0.8814	0.8553
	60%	0.8814	0.8517	<b>0.8852</b>	0.7842	0.8834	0.8523
	80%	0.8814	0.8468	<b>0.8845</b>	0.7941	0.8834	0.8558
	100%	0.8814	0.8404	<b>0.8836</b>	0.8000	0.8834	0.8557

Table 11. Comparison for  $Y$  with missing labels

dataset	$\frac{k}{L}$	$\frac{ \Omega }{nL}$	Top-1 Accuracy						
			Squared			Logistic		Squared Hinge	
			LEML	BCS	BR	LEML	BR	LEML	BR
bibtex	20%	5%	30.30	30.22	<b>42.90</b>	41.51	<b>46.68</b>	30.42	<b>44.97</b>
		10%	39.84	33.56	<b>44.53</b>	41.99	<b>51.09</b>	33.44	<b>48.55</b>
		20%	<b>48.35</b>	40.12	46.08	43.06	<b>55.94</b>	37.22	<b>52.84</b>
		40%	<b>52.37</b>	41.79	43.82	42.27	<b>58.57</b>	40.24	<b>55.39</b>
	40%	5%	34.35	39.17	<b>42.90</b>	43.42	<b>46.68</b>	31.13	<b>44.97</b>
		10%	42.11	39.96	<b>44.53</b>	46.00	<b>51.09</b>	29.03	<b>48.55</b>
		20%	<b>51.97</b>	45.49	46.08	47.40	<b>55.94</b>	32.05	<b>52.84</b>
		40%	<b>56.38</b>	50.10	43.82	49.70	<b>58.57</b>	38.17	<b>55.39</b>
	60%	5%	36.58	41.87	<b>42.90</b>	43.54	<b>46.68</b>	42.54	<b>44.97</b>
		10%	<b>45.53</b>	45.13	44.53	39.36	<b>51.09</b>	31.37	<b>48.55</b>
		20%	<b>53.52</b>	49.54	46.08	46.12	<b>55.94</b>	33.28	<b>52.84</b>
		40%	<b>57.18</b>	54.19	43.82	48.83	<b>58.57</b>	32.13	<b>55.39</b>
autofood	20%	5%	<b>7.89</b>	0.00	<b>7.89</b>	<b>7.89</b>	<b>7.89</b>	<b>7.89</b>	<b>7.89</b>
		10%	44.74	2.63	<b>50.00</b>	<b>55.26</b>	44.74	<b>50.00</b>	<b>50.00</b>
		20%	<b>63.16</b>	0.00	57.89	<b>73.68</b>	47.37	<b>68.42</b>	57.89
		40%	60.53	15.79	<b>78.95</b>	<b>81.58</b>	68.42	<b>86.84</b>	78.95
	40%	5%	<b>10.53</b>	<b>10.53</b>	7.89	<b>7.89</b>	<b>7.89</b>	<b>13.16</b>	7.89
		10%	<b>57.89</b>	7.89	50.00	<b>60.53</b>	44.74	<b>55.26</b>	50.00
		20%	<b>76.32</b>	31.58	57.89	<b>78.95</b>	47.37	<b>76.32</b>	57.89
		40%	60.53	5.26	<b>78.95</b>	<b>84.21</b>	68.42	<b>84.21</b>	78.95
	60%	5%	7.89	<b>10.53</b>	7.89	<b>7.89</b>	<b>7.89</b>	<b>7.89</b>	<b>7.89</b>
		10%	<b>57.89</b>	23.68	50.00	<b>57.89</b>	44.74	<b>55.26</b>	50.00
		20%	<b>73.68</b>	57.89	57.89	<b>78.95</b>	47.37	<b>76.32</b>	57.89
		40%	63.16	36.84	<b>78.95</b>	<b>81.58</b>	68.42	<b>89.47</b>	78.95
compphys	20%	5%	<b>62.50</b>	35.00	42.50	<b>45.00</b>	<b>45.00</b>	<b>67.50</b>	42.50
		10%	<b>75.00</b>	10.00	52.50	<b>67.50</b>	52.50	<b>55.00</b>	52.50
		20%	<b>72.50</b>	7.50	52.50	<b>72.50</b>	52.50	<b>70.00</b>	52.50
		40%	<b>87.50</b>	5.00	52.50	<b>77.50</b>	52.50	<b>80.00</b>	52.50
	40%	5%	<b>65.00</b>	60.00	42.50	<b>45.00</b>	<b>45.00</b>	<b>65.00</b>	42.50
		10%	<b>70.00</b>	17.50	52.50	<b>65.00</b>	52.50	<b>72.50</b>	52.50
		20%	<b>72.50</b>	52.50	52.50	<b>70.00</b>	52.50	<b>75.00</b>	52.50
		40%	<b>80.00</b>	42.50	52.50	<b>80.00</b>	52.50	<b>80.00</b>	52.50
	60%	5%	<b>67.50</b>	52.50	42.50	<b>45.00</b>	<b>45.00</b>	<b>65.00</b>	42.50
		10%	<b>70.00</b>	52.50	52.50	<b>67.50</b>	52.50	<b>67.50</b>	52.50
		20%	<b>77.50</b>	52.50	52.50	<b>80.00</b>	52.50	<b>80.00</b>	52.50
		40%	<b>82.50</b>	52.50	52.50	<b>80.00</b>	52.50	<b>80.00</b>	52.50

Table 12. Comparison for Y with missing labels

dataset	$\frac{k}{L}$	$\frac{ \Omega }{nL}$	Top-3 Accuracy							
			Squared			Logistic		Squared Hinge		
			LEML	BCS	BR	LEML	BR	LEML	BR	
bibtex	20%	5%	16.06	14.29	<b>22.19</b>	21.74	<b>24.47</b>	16.10	<b>23.29</b>	
		10%	20.95	16.29	<b>24.10</b>	22.88	<b>28.43</b>	17.64	<b>26.69</b>	
		20%	<b>26.34</b>	18.78	25.78	23.21	<b>31.92</b>	21.06	<b>29.56</b>	
	40%	5%	<b>30.17</b>	21.55	26.26	23.61	<b>34.50</b>	23.05	<b>31.99</b>	
		10%	18.73	18.99	<b>22.19</b>	22.84	<b>24.47</b>	17.03	<b>23.29</b>	
		20%	22.49	20.16	<b>24.10</b>	25.18	<b>28.43</b>	16.62	<b>26.69</b>	
	60%	5%	<b>28.50</b>	23.84	25.78	25.79	<b>31.92</b>	18.97	<b>29.56</b>	
		10%	<b>32.74</b>	27.58	26.26	27.18	<b>34.50</b>	21.18	<b>31.99</b>	
		20%	18.81	21.09	<b>22.19</b>	22.62	<b>24.47</b>	22.48	<b>23.29</b>	
	autofood	20%	5%	23.96	24.06	<b>24.10</b>	19.84	<b>28.43</b>	17.28	<b>26.69</b>
			10%	<b>29.07</b>	27.05	25.78	25.13	<b>31.92</b>	19.14	<b>29.56</b>
			20%	<b>33.55</b>	31.13	26.26	27.66	<b>34.50</b>	19.46	<b>31.99</b>
40%		5%	<b>30.70</b>	11.40	19.30	<b>29.82</b>	17.54	<b>38.60</b>	19.30	
		10%	<b>52.63</b>	5.26	33.33	<b>50.88</b>	23.68	<b>57.02</b>	33.33	
		20%	59.65	10.53	<b>62.28</b>	<b>70.18</b>	53.51	<b>66.67</b>	61.40	
60%		5%	57.89	20.18	<b>71.93</b>	<b>76.32</b>	63.16	<b>75.44</b>	71.93	
		10%	<b>26.32</b>	15.79	19.30	<b>29.82</b>	17.54	<b>31.58</b>	19.30	
		20%	<b>59.65</b>	12.28	33.33	<b>51.75</b>	23.68	<b>53.51</b>	33.33	
80%		5%	<b>67.54</b>	35.09	62.28	<b>71.05</b>	53.51	<b>64.04</b>	61.40	
		10%	55.26	33.33	<b>71.93</b>	<b>78.07</b>	63.16	<b>77.19</b>	71.93	
		20%	<b>25.44</b>	8.77	19.30	<b>28.95</b>	17.54	<b>22.81</b>	19.30	
compphys	20%	5%	<b>52.63</b>	35.09	33.33	<b>50.00</b>	23.68	<b>61.40</b>	33.33	
		10%	<b>68.42</b>	35.09	62.28	<b>73.68</b>	53.51	<b>71.05</b>	61.40	
		20%	57.02	23.68	<b>71.93</b>	<b>75.44</b>	63.16	<b>74.56</b>	71.93	
	40%	5%	<b>46.67</b>	32.50	28.33	<b>40.00</b>	28.33	<b>40.00</b>	28.33	
		10%	<b>53.33</b>	9.17	37.50	<b>59.17</b>	29.17	<b>40.83</b>	37.50	
		20%	<b>62.50</b>	10.83	31.67	<b>60.83</b>	28.33	<b>61.67</b>	31.67	
	60%	5%	<b>69.17</b>	26.67	43.33	<b>73.33</b>	33.33	<b>70.83</b>	43.33	
		10%	<b>45.83</b>	27.50	28.33	<b>37.50</b>	28.33	<b>41.67</b>	28.33	
		20%	<b>57.50</b>	20.83	37.50	<b>60.00</b>	29.17	<b>55.83</b>	37.50	
	80%	5%	<b>65.00</b>	35.83	31.67	<b>60.00</b>	28.33	<b>61.67</b>	31.67	
		10%	<b>68.33</b>	32.50	43.33	<b>70.83</b>	33.33	<b>73.33</b>	43.33	
		20%	<b>45.00</b>	30.83	28.33	<b>35.83</b>	28.33	<b>45.00</b>	28.33	
100%	5%	<b>59.17</b>	26.67	37.50	<b>61.67</b>	29.17	<b>56.67</b>	37.50		
	10%	<b>65.00</b>	29.17	31.67	<b>60.83</b>	28.33	<b>64.17</b>	31.67		
	20%	<b>71.67</b>	30.00	43.33	<b>65.83</b>	33.33	<b>70.83</b>	43.33		

Table 13. Comparison for Y with missing labels

dataset	$\frac{k}{L}$	$\frac{ \Omega }{nL}$	Top-5 Accuracy							
			Squared			Logistic		Squared Hinge		
			LEML	BCS	BR	LEML	BR	LEML	BR	
bibtex	20%	5%	11.71	10.32	<b>16.14</b>	16.34	<b>17.74</b>	12.07	<b>17.32</b>	
		10%	15.42	11.55	<b>17.77</b>	16.91	<b>20.80</b>	13.11	<b>19.65</b>	
		20%	<b>19.51</b>	13.26	18.81	17.07	<b>23.95</b>	15.52	<b>22.12</b>	
	40%	5%	<b>22.05</b>	15.32	19.13	17.55	<b>25.57</b>	17.57	<b>23.30</b>	
		10%	13.53	13.25	<b>16.14</b>	17.02	<b>17.74</b>	12.70	<b>17.32</b>	
		20%	16.25	14.30	<b>17.77</b>	18.78	<b>20.80</b>	12.24	<b>19.65</b>	
	60%	5%	<b>20.56</b>	17.36	18.81	19.05	<b>23.95</b>	14.46	<b>22.12</b>	
		10%	<b>23.75</b>	19.73	19.13	19.67	<b>25.57</b>	15.86	<b>23.30</b>	
		20%	13.61	14.78	<b>16.14</b>	16.62	<b>17.74</b>	16.56	<b>17.32</b>	
	autofood	20%	5%	16.99	17.31	<b>17.77</b>	14.41	<b>20.80</b>	12.91	<b>19.65</b>
			10%	<b>21.10</b>	19.51	18.81	18.23	<b>23.95</b>	14.17	<b>22.12</b>
			20%	<b>24.50</b>	22.31	19.13	20.38	<b>25.57</b>	14.95	<b>23.30</b>
40%		5%	<b>35.26</b>	8.42	25.26	<b>34.21</b>	21.58	<b>36.84</b>	25.26	
		10%	<b>46.84</b>	6.84	35.79	<b>51.05</b>	32.11	<b>48.95</b>	35.79	
		20%	50.53	10.53	<b>57.89</b>	<b>66.84</b>	52.11	<b>60.53</b>	57.89	
60%		5%	52.11	16.84	<b>68.42</b>	<b>73.16</b>	56.32	<b>72.11</b>	68.42	
		10%	<b>32.11</b>	17.89	25.26	<b>31.58</b>	21.58	<b>30.00</b>	25.26	
		20%	<b>49.47</b>	10.00	35.79	<b>50.53</b>	32.11	<b>45.26</b>	35.79	
80%		5%	<b>64.74</b>	32.11	57.89	<b>66.32</b>	52.11	<b>60.53</b>	57.89	
		10%	50.53	28.95	<b>68.42</b>	<b>73.16</b>	56.32	<b>74.74</b>	68.42	
		20%	<b>31.58</b>	17.37	25.26	<b>31.05</b>	21.58	<b>30.00</b>	25.26	
compphys	20%	5%	<b>50.53</b>	31.58	35.79	<b>52.63</b>	32.11	<b>53.68</b>	35.79	
		10%	<b>64.74</b>	28.95	57.89	<b>68.42</b>	52.11	<b>67.89</b>	57.89	
		20%	58.95	20.00	<b>68.42</b>	<b>71.58</b>	56.32	<b>69.47</b>	68.42	
	40%	5%	<b>34.50</b>	23.00	25.00	<b>28.50</b>	26.00	<b>34.50</b>	25.00	
		10%	<b>50.50</b>	13.50	28.50	<b>51.50</b>	24.00	<b>41.50</b>	29.00	
		20%	<b>52.00</b>	11.50	36.50	<b>55.00</b>	30.00	<b>53.00</b>	36.50	
	60%	5%	<b>60.50</b>	24.00	38.00	<b>64.50</b>	31.00	<b>64.00</b>	38.50	
		10%	<b>34.50</b>	22.00	25.00	<b>29.50</b>	26.00	<b>33.50</b>	25.00	
		20%	<b>53.50</b>	29.00	28.50	<b>51.50</b>	24.00	<b>46.00</b>	29.00	
	80%	5%	<b>56.50</b>	31.00	36.50	<b>55.50</b>	30.00	<b>52.50</b>	36.50	
		10%	<b>59.50</b>	26.00	38.00	<b>61.50</b>	31.00	<b>62.50</b>	38.50	
		20%	<b>36.00</b>	22.00	25.00	<b>27.50</b>	26.00	<b>33.50</b>	25.00	
100%	5%	<b>53.00</b>	24.50	28.50	<b>50.50</b>	24.00	<b>50.50</b>	29.00		
	10%	<b>56.50</b>	29.00	36.50	<b>54.00</b>	30.00	<b>55.50</b>	36.50		
	20%	<b>61.00</b>	32.00	38.00	<b>61.50</b>	31.00	<b>63.50</b>	38.50		

Table 14. Comparison for Y with missing labels

dataset	$\frac{k}{L}$	$\frac{ \Omega }{nL}$	Hamming Loss							
			Squared			Logistic		Squared Hinge		
			LEML	BCS	BR	LEML	BR	LEML	BR	
bibtex	20%	5%	0.0158	0.1480	<b>0.0144</b>	0.0143	<b>0.0138</b>	0.0180	<b>0.0137</b>	
		10%	<b>0.0146</b>	0.1360	0.0156	0.0144	<b>0.0134</b>	0.0187	<b>0.0135</b>	
		20%	<b>0.0136</b>	0.1179	0.0193	0.0156	<b>0.0132</b>	0.0210	<b>0.0136</b>	
	40%	5%	<b>0.0131</b>	0.0994	0.0251	0.0174	<b>0.0128</b>	0.0242	<b>0.0141</b>	
		10%	0.0152	0.2837	<b>0.0144</b>	0.0141	<b>0.0138</b>	0.0175	<b>0.0137</b>	
		20%	<b>0.0149</b>	0.2716	0.0156	0.0141	<b>0.0134</b>	0.0211	<b>0.0135</b>	
	60%	5%	<b>0.0136</b>	0.2496	0.0193	0.0150	<b>0.0132</b>	0.0226	<b>0.0136</b>	
		10%	<b>0.0128</b>	0.2271	0.0251	0.0160	<b>0.0128</b>	0.0269	<b>0.0141</b>	
		20%	0.0154	0.4082	<b>0.0144</b>	0.0145	<b>0.0138</b>	0.0154	<b>0.0137</b>	
	autofood	20%	5%	<b>0.0147</b>	0.3978	0.0156	0.0163	<b>0.0134</b>	0.0215	<b>0.0135</b>
			10%	<b>0.0138</b>	0.3726	0.0193	0.0157	<b>0.0132</b>	0.0252	<b>0.0136</b>
			20%	<b>0.0129</b>	0.3638	0.0251	0.0172	<b>0.0128</b>	0.0312	<b>0.0141</b>
40%		5%	<b>0.0924</b>	0.1727	0.0942	<b>0.0918</b>	0.0991	<b>0.0884</b>	0.0942	
		10%	<b>0.0807</b>	0.1449	0.0837	<b>0.0832</b>	0.0854	<b>0.0811</b>	0.0837	
		20%	<b>0.0750</b>	0.1436	0.0760	<b>0.0686</b>	0.0843	<b>0.0697</b>	0.0760	
60%		5%	0.0780	0.1399	<b>0.0752</b>	<b>0.0655</b>	0.0838	<b>0.0629</b>	0.0750	
		10%	<b>0.0919</b>	0.2887	0.0942	<b>0.0919</b>	0.0991	<b>0.0941</b>	0.0942	
		20%	<b>0.0801</b>	0.2264	0.0837	<b>0.0812</b>	0.0854	<b>0.0814</b>	0.0837	
compphys		20%	5%	<b>0.0671</b>	0.2445	0.0760	<b>0.0681</b>	0.0843	<b>0.0697</b>	0.0760
			10%	<b>0.0671</b>	0.2445	0.0760	<b>0.0681</b>	0.0843	<b>0.0697</b>	0.0760
			20%	0.0903	0.2042	<b>0.0752</b>	<b>0.0647</b>	0.0838	<b>0.0648</b>	0.0750
	40%	5%	<b>0.0932</b>	0.4189	0.0942	<b>0.0921</b>	0.0991	<b>0.0937</b>	0.0942	
		10%	0.0840	0.4144	<b>0.0837</b>	<b>0.0817</b>	0.0854	<b>0.0817</b>	0.0837	
		20%	<b>0.0689</b>	0.3596	0.0760	<b>0.0676</b>	0.0843	<b>0.0692</b>	0.0760	
	60%	5%	<b>0.0724</b>	0.3384	0.0752	<b>0.0650</b>	0.0838	<b>0.0645</b>	0.0750	
		10%	<b>0.0555</b>	0.1391	0.0556	<b>0.0554</b>	0.0555	0.0567	<b>0.0556</b>	
		20%	<b>0.0536</b>	0.1446	0.0565	<b>0.0542</b>	0.0569	<b>0.0554</b>	0.0565	
	compphys	20%	5%	<b>0.0524</b>	0.1431	0.0566	<b>0.0518</b>	0.0566	<b>0.0518</b>	0.0566
			10%	<b>0.0484</b>	0.1048	0.0543	<b>0.0489</b>	0.0561	<b>0.0488</b>	0.0543
			20%	0.0567	0.2924	<b>0.0556</b>	<b>0.0555</b>	0.0555	0.0566	<b>0.0556</b>
40%		5%	<b>0.0532</b>	0.2532	0.0565	<b>0.0535</b>	0.0569	<b>0.0532</b>	0.0565	
		10%	<b>0.0518</b>	0.2569	0.0566	<b>0.0513</b>	0.0566	<b>0.0518</b>	0.0566	
		20%	<b>0.0505</b>	0.1766	0.0543	<b>0.0495</b>	0.0561	<b>0.0484</b>	0.0543	
60%		5%	0.0558	0.4394	<b>0.0556</b>	0.0556	<b>0.0555</b>	<b>0.0555</b>	0.0556	
		10%	<b>0.0532</b>	0.4148	0.0565	<b>0.0532</b>	0.0569	<b>0.0544</b>	0.0565	
		20%	<b>0.0516</b>	0.3797	0.0566	<b>0.0519</b>	0.0566	<b>0.0517</b>	0.0566	
40%		<b>0.0486</b>	0.3563	0.0543	<b>0.0495</b>	0.0561	<b>0.0480</b>	0.0543		

Table 15. Comparison for Y with missing labels

dataset	$\frac{k}{L}$	$\frac{ \Omega }{nL}$	Average AUC							
			Squared			Logistic		Squared Hinge		
			LEML	BCS	BR	LEML	BR	LEML	BR	
bibtex	20%	5%	0.7115	0.6529	<b>0.7789</b>	0.8066	<b>0.8123</b>	0.7363	<b>0.7998</b>	
		10%	0.7665	0.6756	<b>0.7954</b>	0.8208	<b>0.8561</b>	0.7371	<b>0.8210</b>	
		20%	<b>0.8269</b>	0.7111	0.8087	0.8205	<b>0.8941</b>	0.7859	<b>0.8378</b>	
	40%	5%	<b>0.8674</b>	0.7375	0.8104	0.8347	<b>0.9153</b>	0.8167	<b>0.8530</b>	
		10%	0.7379	0.7182	<b>0.7789</b>	<b>0.8164</b>	0.8123	0.7396	<b>0.7998</b>	
		20%	0.7730	0.7353	<b>0.7954</b>	0.8370	<b>0.8561</b>	0.7351	<b>0.8210</b>	
	60%	5%	<b>0.8332</b>	0.7817	0.8087	0.8392	<b>0.8941</b>	0.7813	<b>0.8378</b>	
		10%	<b>0.8724</b>	0.8097	0.8104	0.8639	<b>0.9153</b>	0.8038	<b>0.8530</b>	
		20%	0.7376	0.7445	<b>0.7789</b>	<b>0.8132</b>	0.8123	<b>0.8051</b>	0.7998	
	autofood	20%	5%	0.7778	0.7831	<b>0.7954</b>	0.7639	<b>0.8561</b>	0.7444	<b>0.8210</b>
			10%	<b>0.8367</b>	0.8264	0.8087	0.8251	<b>0.8941</b>	0.7755	<b>0.8378</b>
			20%	<b>0.8753</b>	0.8504	0.8104	0.8716	<b>0.9153</b>	0.7899	<b>0.8530</b>
40%		5%	<b>0.7170</b>	0.5198	0.6451	<b>0.7070</b>	0.6356	<b>0.7235</b>	0.6450	
		10%	<b>0.8083</b>	0.5578	0.7576	<b>0.8194</b>	0.7259	<b>0.8131</b>	0.7576	
		20%	0.8043	0.5804	<b>0.8178</b>	<b>0.8797</b>	0.7712	<b>0.8665</b>	0.8178	
60%		5%	0.8007	0.5807	<b>0.8860</b>	<b>0.9317</b>	0.8087	<b>0.9237</b>	0.8857	
		10%	<b>0.7129</b>	0.6299	0.6451	<b>0.7029</b>	0.6356	<b>0.7157</b>	0.6450	
		20%	<b>0.8218</b>	0.6517	0.7576	<b>0.8198</b>	0.7259	<b>0.8175</b>	0.7576	
compphys		20%	5%	<b>0.8634</b>	0.6322	0.8178	<b>0.8796</b>	0.7712	<b>0.8644</b>	0.8178
			10%	0.8131	0.6848	<b>0.8860</b>	<b>0.9319</b>	0.8087	<b>0.9260</b>	0.8857
			20%	<b>0.7175</b>	0.6013	0.6451	<b>0.7045</b>	0.6356	<b>0.7128</b>	0.6450
	40%	5%	<b>0.8206</b>	0.6316	0.7576	<b>0.8196</b>	0.7259	<b>0.8213</b>	0.7576	
		10%	<b>0.8725</b>	0.6758	0.8178	<b>0.8800</b>	0.7712	<b>0.8781</b>	0.8178	
		20%	0.8141	0.6351	<b>0.8860</b>	<b>0.9315</b>	0.8087	<b>0.9255</b>	0.8857	
	60%	5%	<b>0.6486</b>	0.5727	0.6457	<b>0.6479</b>	0.6424	<b>0.6488</b>	0.6457	
		10%	<b>0.7478</b>	0.5691	0.7235	<b>0.7473</b>	0.7147	<b>0.7556</b>	0.7235	
		20%	<b>0.7908</b>	0.5729	0.7459	<b>0.7921</b>	0.7297	<b>0.8101</b>	0.7459	
	compphys	20%	5%	<b>0.8172</b>	0.6788	0.7728	<b>0.8416</b>	0.7413	<b>0.8718</b>	0.7730
			10%	<b>0.6474</b>	0.6049	0.6457	<b>0.6478</b>	0.6424	<b>0.6480</b>	0.6457
			20%	<b>0.7509</b>	0.6295	0.7235	<b>0.7481</b>	0.7147	<b>0.7437</b>	0.7235
40%		5%	<b>0.7964</b>	0.6442	0.7459	<b>0.7913</b>	0.7297	<b>0.7849</b>	0.7459	
		10%	<b>0.8192</b>	0.6651	0.7728	<b>0.8371</b>	0.7413	<b>0.8561</b>	0.7730	
		20%	0.6443	0.6089	<b>0.6457</b>	<b>0.6468</b>	0.6424	<b>0.6601</b>	0.6457	
60%		5%	<b>0.7504</b>	0.6505	0.7235	<b>0.7489</b>	0.7147	<b>0.7421</b>	0.7235	
		10%	<b>0.7991</b>	0.6687	0.7459	<b>0.7854</b>	0.7297	<b>0.8064</b>	0.7459	
		20%	<b>0.8269</b>	0.7240	0.7728	<b>0.8378</b>	0.7413	<b>0.8659</b>	0.7730	