

Supplementary Materials:

A Unified Algorithm for One-class Structured Matrix Factorization with Side Information

Supp-1 Detailed Derivations of Algorithms

Supp-1.1 Detailed Derivation of $L^-(\text{vec}(W))$

From (4), (6), and (7), we have

$$\begin{aligned}
L^-(\text{vec}(W)) &= \left\| \sqrt{\text{diag}(\mathbf{p})} (\bar{a} \mathbf{1}_m \mathbf{1}_n^\top - XWH^\top) \sqrt{\text{diag}(\mathbf{q})} \right\|_F^2 \\
&= \text{tr} \left(\sqrt{\text{diag}(\mathbf{q})} (\bar{a} \mathbf{1}_n \mathbf{1}_m^\top - HW^\top X^\top) \text{diag}(\mathbf{p}) \times \right. \\
&\quad \left. (\bar{a} \mathbf{1}_m \mathbf{1}_n^\top - XWH^\top) \sqrt{\text{diag}(\mathbf{q})} \right) \\
&= \text{tr} \left((\bar{a} \mathbf{1}_n \mathbf{1}_m^\top - HW^\top X^\top) \text{diag}(\mathbf{p}) \times \right. \\
&\quad \left. (\bar{a} \mathbf{1}_m \mathbf{1}_n^\top - XWH^\top) \text{diag}(\mathbf{q}) \right) \\
&= \text{tr} (\bar{a}^2 \mathbf{1}_n \mathbf{1}_m^\top \text{diag}(\mathbf{p}) \mathbf{1}_m \mathbf{1}_n^\top \text{diag}(\mathbf{q})) \\
&\quad + \text{tr} (HW^\top X^\top \text{diag}(\mathbf{p}) XWH^\top \text{diag}(\mathbf{q})) \\
&\quad - 2 \text{tr} (\bar{a} \mathbf{1}_n \mathbf{1}_m^\top \text{diag}(\mathbf{p}) XWH^\top \text{diag}(\mathbf{q})),
\end{aligned}$$

where $\mathbf{1}_m$ and $\mathbf{1}_n$ are vectors of ones with order m and n , respectively, and $\text{tr}(\cdot)$ denotes the trace of a matrix (i.e., the sum of diagonal entries).

Supp-1.2 Gradient Calculation $\nabla g(\tilde{\mathbf{w}})$

The following property is useful in subsequent derivations.

Property 1. Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_m]^\top$, $H = [\mathbf{h}_1, \dots, \mathbf{h}_n]^\top$, and $D \in \mathbb{R}^{m \times n}$. We have

$$X^\top DH = \sum_{ij} \mathbf{x}_i D_{ij} \mathbf{h}_j^\top \quad (1.1)$$

From (7), the gradient of $g(\tilde{\mathbf{w}})$ is

$$\nabla g(\tilde{\mathbf{w}}) = \nabla L^+(\tilde{\mathbf{w}}) + \nabla L^-(\tilde{\mathbf{w}}) + \lambda_w \nabla \mathcal{R}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}}),$$

where

$$\begin{aligned}
\nabla L^+(\tilde{\mathbf{w}}) &= \sum_{(i,j) \in \Omega^+} \ell_{ij}^{+'} \tilde{\mathbf{x}}_{ij} = \sum_{(i,j) \in \Omega^+} \text{vec} \left(\ell_{ij}^{+'} \mathbf{x}_i \mathbf{h}_j^\top \right) \\
&= \text{vec} (X^\top D^+ H). \quad (1.2)
\end{aligned}$$

The last equality follows from (1.1), $\ell_{ij}^{+'}$ is the abbreviation of $\ell_{ij}^{+'}(Y_{ij}, \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij})$ and,

$$D_{ij}^+ = \begin{cases} \ell_{ij}^{+'}, & \forall (i, j) \in \Omega^+, \\ 0, & \text{otherwise.} \end{cases}$$

Similar to the situation of computing $L^+(\tilde{\mathbf{w}})$ in Section 3.1, if $B = XW$ has been calculated, then constructing a sparse matrix D^+ requires only $O(|\Omega^+|k)$ operations and $O(|\Omega^+|)$ space. With a similar derivation in (1.2), we have

$$\nabla L^-(\tilde{\mathbf{w}}) = \text{vec} (X^\top D^- H),$$

where D^- is a dense $m \times n$ matrix with $D_{ij}^- = \ell_{ij}^{-'}$, $\forall (i, j) \in [m] \times [n]$. From the definition of $\ell^-(\cdot, \cdot)$ in (6),

$$\begin{aligned}
\frac{\partial}{\partial b} \ell^-(\bar{a}, b) &= 2p_i q_j (b - \bar{a}) \\
\Rightarrow \ell_{ij}^{-'} &= 2p_i q_j (\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij} - \bar{a}) = 2p_i q_j \left((XWH^\top)_{ij} - \bar{a} \right).
\end{aligned}$$

Thus, we have

$$D^- = 2 \text{diag}(\mathbf{p}) (XWH^\top - \bar{a} \mathbf{1}_m \mathbf{1}_n^\top) \text{diag}(\mathbf{q}).$$

Straightforward computation of entire D^- requires $O(mnk + \text{nnz}(X)k)$ time and $O(mn)$ space; both requirements are infeasible for large values of m and n . Fortunately, $\nabla L^-(\tilde{\mathbf{w}}) = \text{vec}(X^\top D^- H)$ can be computed without explicitly forming D^- as follows:

$$\begin{aligned}
&X^\top D^- H \\
&= 2X^\top \text{diag}(\mathbf{p}) XWH^\top \text{diag}(\mathbf{q}) H \\
&\quad - 2\bar{a} X^\top \text{diag}(\mathbf{p}) \mathbf{1}_m \mathbf{1}_n^\top \text{diag}(\mathbf{q}) H \quad (1.3) \\
&= 2X^\top \text{diag}(\mathbf{p}) \underbrace{(XW)}_B \underbrace{(H^\top \text{diag}(\mathbf{q}) H)}_M - 2\bar{a} \underbrace{(X^\top \mathbf{p})}_d \underbrace{(\mathbf{q}^\top H)}_{\mathbf{k}^\top},
\end{aligned}$$

where B , M , \mathbf{d} , and \mathbf{k} are the same matrices/vectors used in the objective value evaluation introduced in Algorithm 1. Combining (1.2), (1.3), and $\tilde{\mathbf{w}} = \text{vec}(W)$, $\nabla g(\tilde{\mathbf{w}})$ can be computed using the following sequence of matrix-matrix products:

$$\begin{aligned}
\nabla g(\tilde{\mathbf{w}}) &= \text{vec} (X^\top D^+ H + X^\top D^- H + \lambda_w \nabla \mathcal{R}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}})) \\
&= \text{vec} (X^\top [D^+ H + 2 \text{diag}(\mathbf{p}) (BM) + 2\lambda_w \lambda_g LB] \\
&\quad + 2\lambda_w W - 2\bar{a} \mathbf{d} \mathbf{k}^\top),
\end{aligned}$$

where the last equality follows from (1.3) and

$$\begin{aligned}
\nabla \mathcal{R}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}}) &= 2 \text{vec} ((I + \lambda_g X^\top L X) W) \\
&= 2 \text{vec} (W + \lambda_g X^\top L B).
\end{aligned}$$

Thus, the cost is

$$O(\text{nnz}(X)k + |\Omega^+|k + \text{nnz}(L)k + mk^2 + dk) \quad (1.4)$$

time. A detailed procedure is shown in Algorithm 2. Note that if the computation of $\nabla g(\tilde{\mathbf{w}})$ follows immediately after the objective value evaluation of the same $\tilde{\mathbf{w}}$, then M , \mathbf{d} , \mathbf{q} , \mathbf{k} , and B can be re-used to save additional computation.²

Note that we may use a different sequence of matrix-matrix products for (1.3) if the number of features d is very small such that

$$d^2 + dk < mk. \quad (1.5)$$

By pre-computing $\bar{M} = X^\top \text{diag}(\mathbf{p}) X$ in $O(\text{nnz}(X)d)$ time and storing it in $O(d^2)$ space, the following sequence to compute $\nabla g(\tilde{\mathbf{w}})$,

$$\begin{aligned}
\nabla g(\tilde{\mathbf{w}}) &= \text{vec} (X^\top D^+ H + X^\top D^- H + \lambda_w \nabla \mathcal{R}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}})) \\
&= \text{vec} (X^\top (D^+ H + 2\lambda_w \lambda_g LB) + 2\bar{M} W \\
&\quad + 2\lambda_w W - 2\bar{a} \mathbf{d} \mathbf{k}^\top),
\end{aligned}$$

costs

$$O(\text{nnz}(X)k + |\Omega^+|k + \text{nnz}(L)k + d^2k + dk^2)$$

time, which is smaller than that in (1.4) because of (1.5).

² $\text{diag}(\mathbf{p})(BM)$ can also be re-used.

Supp-1.3 Hessian-vector Multiplication $\nabla^2 g(\tilde{\mathbf{w}})\tilde{\mathbf{s}}$

Given a current point $\tilde{\mathbf{w}}$ and a vector $\tilde{\mathbf{s}} \in \mathbb{R}^{dk}$, computing $\nabla^2 g(\tilde{\mathbf{w}})\tilde{\mathbf{s}}$ is a common operation required in many iterative optimization algorithms such as conjugate gradient. Let $S \in \mathbb{R}^{d \times k}$ be the matrix such that $\tilde{\mathbf{s}} = \text{vec}(S)$. From (7), we have

$$\nabla^2 g(\tilde{\mathbf{w}})\tilde{\mathbf{s}} = \nabla^2 L^+(\tilde{\mathbf{w}})\tilde{\mathbf{s}} + \nabla^2 L^-(\tilde{\mathbf{w}})\tilde{\mathbf{s}} + \lambda_w \nabla^2 \mathcal{R}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}})\tilde{\mathbf{s}},$$

where

$$\begin{aligned} \nabla^2 L^+(\tilde{\mathbf{w}})\tilde{\mathbf{s}} &= \sum_{(i,j) \in \Omega^+} \ell_{ij}^{+''} \tilde{\mathbf{x}}_{ij} \tilde{\mathbf{x}}_{ij}^\top \tilde{\mathbf{s}} \\ &= \sum_{(i,j) \in \Omega^+} \ell_{ij}^{+''} ((\mathbf{h}_j \mathbf{h}_j^\top) \otimes (\mathbf{x}_i \mathbf{x}_i^\top)) \text{vec}(S) \\ &= \sum_{(i,j) \in \Omega^+} \text{vec}(\ell_{ij}^{+''} \mathbf{x}_i \mathbf{x}_i^\top S \mathbf{h}_j \mathbf{h}_j^\top) \quad (1.6) \\ &= \sum_{(i,j) \in \Omega^+} \text{vec}(\mathbf{x}_i \underbrace{\ell_{ij}^{+''} \mathbf{x}_i^\top S \mathbf{h}_j}_{U_{ij}^+} \mathbf{h}_j^\top) \\ &= \text{vec}(X^\top U^+ H). \quad (1.7) \end{aligned}$$

Note that (1.6) is from (4), (1.7) is from (1.1), U^+ is a sparse matrix with

$$U_{ij}^+ = \begin{cases} \ell_{ij}^{+''} \mathbf{x}_i^\top S \mathbf{h}_j, & \forall (i,j) \in \Omega^+, \\ 0 & \text{otherwise,} \end{cases}$$

and $\ell_{ij}^{+''} = \ell_{ij}^{+''}(Y_{ij}, \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij})$, $\forall (i,j) \in \Omega^+$. Similar to the gradient computation, $\ell_{ij}^{+''}$ can be computed in $O(k)$ time if $B = XW$ is pre-computed. To efficiently compute U_{ij}^+ , $\forall (i,j) \in \Omega^+$, we first compute a matrix-matrix product XS in $O(\text{nnz}(X)k)$ time and store it in an $m \times k$ matrix $N = [\dots, \mathbf{n}_i, \dots]^\top$ such that $\mathbf{n}_i = S^\top \mathbf{x}_i$. Then the entire sparse matrix U^+ can be constructed in $O(|\Omega^+|k + \text{nnz}(X)k)$ time.

With a similar derivation as above, we have

$$\nabla^2 L^-(\tilde{\mathbf{w}})\tilde{\mathbf{s}} = \text{vec}(X^\top U^- H),$$

where U^- is a dense $m \times n$ matrix with

$$U_{ij}^- = \ell_{ij}^{-''} \mathbf{x}_i^\top S \mathbf{h}_j, \quad \forall (i,j) \in [m] \times [n].$$

From the definition of $\ell^-(\cdot, \cdot)$ in (6),

$$\frac{\partial^2}{\partial b^2} \ell_{ij}^-(\bar{a}, b) = 2p_i q_j, \quad \forall b \quad \Rightarrow \quad \ell_{ij}^{-''} = 2p_i q_j.$$

Thus, we have

$$U^- = 2 \text{diag}(\mathbf{p}) X S H^\top \text{diag}(\mathbf{q}).$$

Similar to the situation in the gradient computation, computing $\nabla^2 L^-(\tilde{\mathbf{w}})\tilde{\mathbf{s}} = \text{vec}(X^\top U^- H)$ can be done without forming the entire U^- explicitly as follows:

$$X^\top U^- H = 2X^\top \text{diag}(\mathbf{p}) \underbrace{(XS)}_N \underbrace{H^\top \text{diag}(\mathbf{q})H}_M, \quad (1.8)$$

where N is available when we construct the sparse matrix U^+ , and M can be pre-computed and stored (same M in the objective value evaluation and gradient computation). Combining (1.7), (1.8), and $\tilde{\mathbf{w}} = \text{vec}(W)$, $\nabla^2 g(\tilde{\mathbf{w}})\tilde{\mathbf{s}}$ can be computed using the following sequence of matrix-matrix products:

$$\begin{aligned} \nabla^2 g(\tilde{\mathbf{w}})\tilde{\mathbf{s}} &= \text{vec}(X^\top U^+ H + X^\top U^- H + \lambda_w \nabla^2 \mathcal{R}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}})\tilde{\mathbf{s}}) \\ &= \text{vec}(X^\top [U^+ H + 2 \text{diag}(\mathbf{p}) N M] + \lambda_w \nabla^2 \mathcal{R}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}})\tilde{\mathbf{s}}) \\ &= \text{vec}(X^\top [U^+ H + 2 \text{diag}(\mathbf{p}) N M + 2\lambda_w \lambda_g L N] + 2\lambda_w S), \end{aligned}$$

where the last equality follows from

$$\begin{aligned} \nabla^2 \mathcal{R}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}})\tilde{\mathbf{s}} &= 2 \text{vec}(S + \lambda_g X^\top L X S) \\ &= 2 \text{vec}(S + \lambda_g X^\top L N). \end{aligned}$$

Thus, the cost is

$$O(\text{nnz}(X)k + |\Omega^+|k + \text{nnz}(L)k + mk^2 + dk)$$

time. Note that we can move the computation of $\ell_{ij}^{+''}$, $\forall (i,j) \in \Omega^+$ into the pre-processing phase to save some repeated operations when the Hessian-vector products are performed many times with the same $\tilde{\mathbf{w}}$, a common situation in iterative optimization algorithms such as conjugate gradient. Further, when the squared loss $\ell(a, b) = (a-b)^2$ is used, for entries in Ω^+ , $\ell_{ij}^{+''} = 2(1-p_i q_j)$ is independent of the choice of $\tilde{\mathbf{w}}$ and can be pre-computed in the beginning of the entire optimization procedure. A detailed procedure is presented in Algorithm 3.

Similar to the gradient calculation, when $d^2 + dk < mk$, we can reduce the cost by using another sequence of matrix-matrix products to compute $\nabla^2 g(\tilde{\mathbf{w}})\tilde{\mathbf{s}}$ in

$$O(|\Omega^+|k + \text{nnz}(X)k + \text{nnz}(L)k + d^2k + dk^2)$$

time with $\bar{M} = X^\top \text{diag}(\mathbf{p})X$ as follows:

$$X^\top (U^+ H + 2\lambda_w \lambda_g L N) + 2\bar{M} S M + 2\lambda_w S.$$

Supp-1.4 Line Search Procedure and Trust Region Method

Line search procedures and trust region methods are two common techniques to determine a step size and guarantee the asymptotic convergence for an optimization algorithm. In this section, we briefly introduce how to apply these two techniques to (7).

Line Search Procedure A line search procedure is used to ensure the sufficient decrease of the objective function value after a search direction ΔW is given. Common choices for the search direction include the negative gradient direction and the Newton direction. However, even if ΔW is a descent direction, $g(W + \Delta W)$ is not necessary smaller than $g(W)$. Thus, in a backward line search procedure, we try a sequence of step size $\alpha = 1, \beta, \beta^2, \dots$, with $\beta < 1$ such that

$$g(W + \alpha \Delta W) < g(W) + \sigma \alpha \langle \nabla g(W), \Delta W \rangle, \quad (1.9)$$

where $\sigma < 1/2$. Note that even with our proposed efficient function value evaluation procedure (Algorithm 1), recalculating function value at $W + \alpha \Delta W$ for each α can still be expensive.

We propose an effective trick to reduce the cost for line search by taking that $\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij}$ is a linear function of $\tilde{\mathbf{w}}$, $\ell_{ij}^-(\tilde{\mathbf{a}}, \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij})$ is a quadratic function of $\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij}$, and $\mathcal{R}(\tilde{\mathbf{w}})$ is a quadratic function of $\tilde{\mathbf{w}}$. The idea is to have

$$\begin{aligned} & g(W + \alpha\Delta W) \\ &= L^+(W + \alpha\Delta W) + L^-(W + \alpha\Delta W) + \lambda_w \mathcal{R}(W + \alpha\Delta W) \\ &= g(W) - L^+(W) + L^+(W + \alpha\Delta W) + \alpha \text{Val.1} + \alpha^2 \text{Val.2}, \end{aligned}$$

where Val.1 and Val.2 can be calculated with the cost of one function evaluation. If these two values are available, then the sufficient decrease condition in (1.9) can be easily checked as follows.

$$\begin{aligned} & g(W + \alpha\Delta W) - g(W) \\ &= -L^+(W) + L^+(W + \alpha\Delta W) + \alpha \text{Val.1} + \alpha^2 \text{Val.2} \\ &< \sigma \alpha \langle \nabla g(W), \Delta W \rangle. \end{aligned}$$

Moreover, when the line-search procedure terminates, we have the next iterate $W + \alpha\Delta W$ and the new function value. From (8)-(13) in the main text,

$$\begin{aligned} & L^-(W + \alpha\Delta W) \\ &= \tilde{\mathbf{a}}^2 (\mathbf{p}^\top \mathbf{1}_m) (\mathbf{q}^\top \mathbf{1}_n) + \langle B + \alpha\Delta B, \text{diag}(\mathbf{p})(B + \alpha\Delta B)M \rangle \\ &\quad - 2\tilde{\mathbf{a}} \mathbf{d}^\top (W + \alpha\Delta W) \mathbf{k} \\ &= L^-(W) + \alpha^2 \langle \Delta B, \text{diag}(\mathbf{p})\Delta B M \rangle + \\ &\quad \alpha \langle \langle \Delta B, \text{diag}(\mathbf{p})BM \rangle + \langle B, \text{diag}(\mathbf{p})\Delta B M \rangle - 2\tilde{\mathbf{a}} \mathbf{d}^\top \Delta W \mathbf{k} \rangle \end{aligned}$$

and

$$\begin{aligned} & \mathcal{R}(W + \alpha\Delta W) \\ &= \text{tr} \left((W + \alpha\Delta W)^\top (W + \alpha\Delta W) \right) \\ &\quad + \lambda_g \text{tr} \left((B + \alpha\Delta B)^\top L(B + \alpha\Delta B) \right) \\ &= \mathcal{R}(W) + \alpha (2 \text{tr}(W^\top \Delta W) + \lambda_g B^\top L \Delta B) \\ &\quad + \alpha^2 \text{tr}(\Delta W^\top \Delta W + \lambda_g \Delta B^\top L \Delta B), \\ &= \mathcal{R}(W) + \alpha (2 \langle W, \Delta W \rangle + 2\lambda_g \langle B, L \Delta B \rangle) \\ &\quad + \alpha^2 \left(\|\Delta W\|_F^2 + \lambda_g \langle \Delta B, L \Delta B \rangle \right), \end{aligned}$$

where

$$\Delta B = X \Delta W.$$

The above calculation shows that in addition to the matrix B , we can maintain

$$\hat{M} = \text{diag}(\mathbf{p})BM, \text{ and } \hat{L} = LB.$$

Before the line-search procedure,

$$\Delta B = X \Delta W, \quad (1.10)$$

$$\Delta \hat{M} = \text{diag}(\mathbf{p})\Delta B M, \quad (1.11)$$

$$\Delta \delta = \mathbf{d}^\top \Delta W \mathbf{k}, \text{ and} \quad (1.12)$$

$$\Delta \hat{L} = L \Delta B \quad (1.13)$$

are calculated. Then

$$\begin{aligned} \text{Val.1} &= \langle \Delta B, \hat{M} \rangle + \langle B, \Delta \hat{M} \rangle - 2\tilde{\mathbf{a}} \Delta \delta \\ &\quad + 2\lambda_w \left(\langle W, \Delta W \rangle + \lambda_g \langle B, \Delta \hat{L} \rangle \right), \end{aligned}$$

$$\begin{aligned} \text{Val.2} &= \langle \Delta B, \Delta \hat{M} \rangle \\ &\quad + \lambda_w \left(\|\Delta W\|_F^2 + \lambda_g \langle \Delta B, \Delta \hat{L} \rangle \right). \end{aligned}$$

Clearly, the computation for (1.10)-(1.13) costs

$$O(\text{nnz}(X)k + mk^2 + dk + \text{nnz}(L)k),$$

which is the same as the cost of one function evaluation. For Val.1 and Val.2, the cost $O(dk + mk)$ is much smaller. To check the sufficient decrease condition (1.9), we pre-calculate $\nabla g(W)^\top \Delta W$ before the line-search procedure, so the calculation under an α is $O(|\Omega^+|k)$ for calculating $L^+(W + \alpha\Delta W)$.

If the squared loss is considered on entries in Ω^+ , the same technique for calculating $L^-(W + \alpha\Delta W)$ can be applied for calculating $L^+(W + \alpha\Delta W)$. Then the $O(|\Omega^+|k)$ cost is needed only before the line-search procedure rather than at each step of checking an α value.

Trust Region Methods Trust region method is an alternative technique of line search to guarantee asymptotic convergence, which is also widely used in many machine learning applications (Lin, Weng, and Keerthi 2008). At the t -th iteration of a trust region method, we have the current iterate $\tilde{\mathbf{w}}^t$, a size Δ^t of trust region, and a quadratic model $q_t(\tilde{\mathbf{s}})$ as the approximation of the value $g(\tilde{\mathbf{w}}^t + \tilde{\mathbf{s}}) - g(\tilde{\mathbf{w}}^t)$. Two common choices for $q_t(\tilde{\mathbf{s}})$ are the first-order approximation

$$q_t(\tilde{\mathbf{s}}) = \nabla g(\tilde{\mathbf{w}}^t)^\top \tilde{\mathbf{s}} + \frac{1}{2} \tilde{\mathbf{s}}^\top \tilde{\mathbf{s}} \quad (1.14)$$

and the second-order approximation

$$q_t(\tilde{\mathbf{s}}) = \nabla g(\tilde{\mathbf{w}}^t)^\top \tilde{\mathbf{s}} + \frac{1}{2} \tilde{\mathbf{s}}^\top \nabla^2 g(\tilde{\mathbf{w}}^t) \tilde{\mathbf{s}}. \quad (1.15)$$

Next we find $\tilde{\mathbf{s}}^t$ to approximately solve

$$\min_{\|\tilde{\mathbf{s}}\| \leq \Delta^t} q_t(\tilde{\mathbf{s}}). \quad (1.16)$$

If (1.16) is the choice, we can apply conjugate gradient to obtain $\tilde{\mathbf{s}}^t$ with our efficient Hessian-vector product procedure in Algorithm 3. In Algorithm 4, we present a detailed conjugate gradient procedure to solve (1.16). See Lin, Weng, and Keerthi (2008) for more details about the convergence property and the implementation issues of Algorithm 4.

We then update $\tilde{\mathbf{w}}^t$ and Δ^t as follows:

$$\tilde{\mathbf{w}}^{t+1} = \begin{cases} \tilde{\mathbf{w}}^t + \tilde{\mathbf{s}} & \text{if } \rho^t > \eta_0, \\ \tilde{\mathbf{w}}^t & \text{otherwise,} \end{cases} \quad (1.17)$$

$$\Delta^{t+1} \in \begin{cases} [\sigma_1 \min(\|\tilde{\mathbf{s}}^t\|, \Delta^t), \sigma_2 \Delta^t] & \text{if } \rho^t \leq \eta_1, \\ [\sigma_1 \Delta^t, \sigma_3 \Delta^t] & \text{if } \rho^t \in (\eta_1, \eta_2), \\ [\Delta^t, \sigma_3 \Delta^t] & \text{if } \rho^t \geq \eta_2, \end{cases} \quad (1.18)$$

where ρ^t is the ratio of the actual reduction in the function to the predicted reduction in the quadratic model:

$$\rho^t = \frac{g(\tilde{\mathbf{w}}^t + \tilde{\mathbf{s}}^t) - g(\tilde{\mathbf{w}}^t)}{q_t(\tilde{\mathbf{s}}^t)},$$

and η_0, η_1, η_2 , and $\sigma_1, \sigma_2, \sigma_3$ are pre-specified positive constants such that $\eta_1 < \eta_2 < 1$ and $\sigma_1 < \sigma_2 < 1 < \sigma_3$. The idea behind a trust region method is that if the ratio ρ^t is large enough, we can accept the current step $\tilde{\mathbf{s}}^t$ and enlarge the size of trust region Δ^t ; otherwise, we reject the current step and shrink the size of trust region. More details can be found in Lin, Weng, and Keerthi (2008).

We can see that the computation of the ratio ρ^t requires a function value evaluation at $g(\tilde{\mathbf{w}}^t + \tilde{\mathbf{s}}^t)$ and $q_t(\tilde{\mathbf{s}}^t)$. To compute $g(\tilde{\mathbf{w}}^t + \tilde{\mathbf{s}}^t)$, we can apply Algorithm 1. To compute $q_t(\tilde{\mathbf{s}}^t)$, we can utilize our proposed algorithms as subroutines to perform a straightforward computation. For example, assume the second-order approximation in (1.15) is considered and the conjugate gradient described in Algorithm 4 is applied to solve (1.16). In the procedure in Algorithm 4, $\nabla g(\tilde{\mathbf{w}}^t)$ is required as an input. Further, we have the final residual vector

$$\mathbf{r}^{i^*} = -\nabla q_t(\tilde{\mathbf{s}}^{i^*}) = -\nabla g(\tilde{\mathbf{w}}^t) - \nabla^2 g(\tilde{\mathbf{w}}^t) \tilde{\mathbf{s}}^{i^*}$$

available when Algorithm 4 terminates, where i^* is the number of CG iterations performed, and $\tilde{\mathbf{s}}^t = \tilde{\mathbf{s}}^{i^*}$. Then, we can compute

$$q_t(\tilde{\mathbf{s}}^t) = -\frac{1}{2} \left((-\nabla g(\tilde{\mathbf{w}}^t))^\top \tilde{\mathbf{s}}^t - (\tilde{\mathbf{s}}^t)^\top \mathbf{r}^{i^*} \right)$$

in $O(dk)$ time.

Supp-1.5 Complete Optimization Procedures

With the efficient procedures proposed in Section 3, we are able to implement many optimization algorithms to solve (7) such as gradient descent with line search, nonlinear conjugate gradient, truncated Newton methods with line search, and trust region Newton methods. To illustrate how Algorithms 1-3 are used in practice, in Algorithm 5, we use a trust region Newton method (TRON) (Lin, Weng, and Keerthi 2008) as an example to show the complete optimization procedure.

We give another Algorithm 6 to demonstrate the use of line search in a gradient descent method for solving (7). Note that we use a Newton direction in Algorithm 5 and a negative gradient direction in Algorithm 6, respectively, though in either algorithm any descent direction can be considered.

Algorithm 4 Conjugate gradient procedure to approximately solve the trust-region subproblem (1.16).

Input: $\nabla g(\tilde{\mathbf{w}}^t)$, a relative stopping parameter $\xi < 1$, and a size of trust region $\Delta^t > 0$

Output: $\tilde{\mathbf{s}}^t = \tilde{\mathbf{s}}^{i^*}$ and $\mathbf{r}^{i^*} = -\nabla g(\tilde{\mathbf{w}}^t) - \nabla^2 g(\tilde{\mathbf{w}}^t) \tilde{\mathbf{s}}^t$

- 1: // Initialization
- 2: $\tilde{\mathbf{s}}^0 = \mathbf{0}$ // initial point
- 3: $\mathbf{r}^0 = -\nabla g(\tilde{\mathbf{w}}^t)$ // initial residual
- 4: $\mathbf{d}^0 = \mathbf{r}^0$ // initial direction
- 5: **for** $i = 0, 1, 2, 3, \dots$ (CG iterations) **do**
- 6: **if** $\|\mathbf{r}^i\| \leq \xi \|\mathbf{r}^0\|$ **then**
- 7: Stop the **for** loop and return $\tilde{\mathbf{s}}^t = \tilde{\mathbf{s}}^i$ and \mathbf{r}^i
- 8: Compute and store $\nabla^2 g(\tilde{\mathbf{w}}^t) \mathbf{d}^i$ by Algorithm 3
- 9: $\alpha^i = \|\mathbf{r}^i\|^2 / \langle \mathbf{d}^i, \nabla^2 g(\tilde{\mathbf{w}}^t) \mathbf{d}^i \rangle$
- 10: $\tilde{\mathbf{s}}^{i+1} = \tilde{\mathbf{s}}^i + \alpha^i \mathbf{d}^i$
- 11: **if** $\|\tilde{\mathbf{s}}^{i+1}\| \geq \Delta^t$ **then**
- 12: Compute τ such that $\|\tilde{\mathbf{s}}^i + \tau \mathbf{d}^i\| = \Delta^t$
- 13: $\tilde{\mathbf{s}}^{i+1} = \tilde{\mathbf{s}}^i + \tau \mathbf{d}^i$
- 14: $\mathbf{r}^{i+1} = \mathbf{r}^i - \tau \nabla^2 g(\tilde{\mathbf{w}}^t) \mathbf{d}^i$
- 15: Stop the **for** loop and return $\tilde{\mathbf{s}}^t = \tilde{\mathbf{s}}^{i+1}$ and \mathbf{r}^{i+1}
- 16: $\mathbf{r}^{i+1} = \mathbf{r}^i - \alpha^i \nabla^2 g(\tilde{\mathbf{w}}^t) \mathbf{d}^i$
- 17: $\beta^i = \|\mathbf{r}^{i+1}\|^2 / \|\mathbf{r}^i\|^2$
- 18: $\mathbf{d}^{i+1} = \mathbf{r}^{i+1} + \beta^i \mathbf{d}^i$

Algorithm 5 A trust-region Newton method for (7).

Input: feature matrix $X = [\mathbf{x}_1, \dots, \mathbf{x}_m]^\top \in \mathbb{R}^{m \times d}$, label matrix $Y = [\mathbf{y}_1, \dots, \mathbf{y}_m]^\top \in \mathbb{R}^{m \times n}$, $H \in \mathbb{R}^{n \times k}$, and parameters $k, \lambda_w, \lambda_g, \mathbf{p}, \mathbf{q}, \bar{a}$, and an initial size $\Delta > 0$ of trust region. Choose positive constants $\eta_0, \eta_1 < \eta_2 < 1$ and $\sigma_1 < \sigma_2 < 1 < \sigma_3$.

```
1:  $M = H^\top \text{diag}(\mathbf{q})H$   $\dots O(nk^2)$ 
2:  $\mathbf{d} = X^\top \mathbf{p}$   $\dots O(\text{nnz}(X))$ 
3:  $\mathbf{k} = H^\top \mathbf{q}$   $\dots O(nk)$ 
4: Compute  $\mathbf{p}^\top \mathbf{1}_m$  and  $\mathbf{q}^\top \mathbf{1}_n$   $\dots O(m+n)$ 
5:  $B = XW$ , where  $B = [\dots, \mathbf{b}_i, \dots]^\top$   $\dots O(\text{nnz}(X)k)$ 
6:  $\hat{M} = \text{diag}(\mathbf{p})BM$   $\dots O(mk^2)$ 
7:  $\delta = \mathbf{d}^\top W\mathbf{k}$   $\dots O(dk)$ 
8:  $\hat{L} = LB$   $\dots O(\text{nnz}(L)k)$ 
9:  $L^+ = \sum_{(i,j) \in \Omega^+} \ell_{ij}(Y_{ij}, \mathbf{b}_i^\top \mathbf{h}_j)$   $O(|\Omega^+|k)$ 
10:  $L^- = \bar{a}^2(\mathbf{p}^\top \mathbf{1}_m)(\mathbf{q}^\top \mathbf{1}_n) + \langle B, \hat{M} \rangle - 2\bar{a}\delta$   $\dots O(mk)$ 
11:  $\text{obj} = L^+ + L^- + \lambda_w(\|W\|_F^2 + \lambda_g \langle B, \hat{L} \rangle)$ 
12:  $\dots O(dk + mk)$ 
13: // Compute gradient
14:  $D_{ij}^+ = \ell_{ij}^{+'}(Y_{ij}, \mathbf{b}_i^\top \mathbf{h}_j), \forall (i, j) \in \Omega^+$   $\dots O(|\Omega^+|k)$ 
15:  $G = X^\top(D^+H + 2\hat{M})$   $\dots O(|\Omega^+|k + \text{nnz}(X)k + mk^2)$ 
16:  $G = G - 2\bar{a}\mathbf{d}\mathbf{k}^\top + 2\lambda_w\lambda_g X^\top \hat{L}$ 
17:  $\dots O(dk + \text{nnz}(X)k)$ 
18: // Pre-processing for Hessian-vector products
19:  $\ell_{ij}^{+''} = \ell_{ij}^{+''}(Y_{ij}, \mathbf{b}_i^\top \mathbf{h}_j), \forall (i, j) \in \Omega^+$   $\dots O(|\Omega^+|k)$ 
20: for  $t = 1, 2, 3, \dots$  do
21: // Solve (1.16) approximately
22: Get  $\tilde{\mathbf{s}} = \text{vec}(S)$  and  $\mathbf{r}^{i*} = \text{vec}(R)$  by Algorithm 4 with  $\Delta$ 
23:
24: // Compute  $q_t(\tilde{\mathbf{s}})$ 
25:  $q_t(\tilde{\mathbf{s}}) = -\frac{1}{2}(-\langle G, S \rangle - \langle R, S \rangle)$   $\dots O(dk)$ 
26:
27: // Compute  $g(\tilde{\mathbf{w}}^t + \tilde{\mathbf{s}})$ 
28:  $W_{\text{new}} = W + S$   $\dots O(dk)$ 
29:  $B_{\text{new}} = XW_{\text{new}}$   $\dots O(\text{nnz}(X)k)$ 
30:  $\hat{M} = \text{diag}(\mathbf{p})B_{\text{new}}M$   $O(mk^2)$ 
31:  $\delta = \mathbf{d}^\top W_{\text{new}}\mathbf{k}$   $\dots O(dk)$ 
32:  $\hat{L} = LB_{\text{new}}$   $\dots O(\text{nnz}(L)k)$ 
33:  $L^+ = \sum_{(i,j) \in \Omega^+} \ell_{ij}(Y_{ij}, (\mathbf{w}_{\text{new}})_i^\top \mathbf{h}_j)$   $\dots O(|\Omega^+|k)$ 
34:  $L^- = \bar{a}^2(\mathbf{p}^\top \mathbf{1}_m)(\mathbf{q}^\top \mathbf{1}_n) + \langle B_{\text{new}}, \hat{M} \rangle - 2\bar{a}\delta$ 
35:  $\dots O(mk)$ 
36:  $\text{obj}_{\text{new}} = L^+ + L^- + \lambda_w(\|W_{\text{new}}\|_F^2 + \lambda_g \langle B_{\text{new}}, \hat{L} \rangle)$ 
37:  $\dots O(mk + dk)$ 
38:  $\rho = (\text{obj}_{\text{new}} - \text{obj})/q_t(\tilde{\mathbf{s}})$   $\dots O(1)$ 
39: Update  $\Delta$  based on  $\rho$  and (1.18)  $\dots O(1)$ 
40: if  $\rho > \eta_0$  then
41:  $\text{obj} = \text{obj}_{\text{new}}$   $\dots O(1)$ 
42:  $W = W_{\text{new}}$   $\dots O(dk)$ 
43:  $B = B_{\text{new}}$   $\dots O(mk)$ 
44: // Compute gradient
45:  $D_{ij}^+ = \ell_{ij}^{+'}(Y_{ij}, \mathbf{b}_i^\top \mathbf{h}_j), \forall (i, j) \in \Omega^+$   $\dots O(|\Omega^+|k)$ 
46:  $G = X^\top(D^+H + 2\hat{M})$ 
47:  $\dots O(|\Omega^+|k + \text{nnz}(X)k + mk^2)$ 
48:  $G = G - 2\bar{a}\mathbf{d}\mathbf{k}^\top + 2\lambda_w\lambda_g X^\top \hat{L}$ 
49:  $\dots O(dk + \text{nnz}(X)k)$ 
50: // Pre-processing for Hessian-vector products
51:  $\ell_{ij}^{+''} = \ell_{ij}^{+''}(Y_{ij}, \mathbf{b}_i^\top \mathbf{h}_j), \forall (i, j) \in \Omega^+$   $\dots O(|\Omega^+|k)$ 
```

Algorithm 6 A gradient-descent method with backtracking line search to minimize (7).

Input: feature matrix $X = [\mathbf{x}_1, \dots, \mathbf{x}_m]^\top \in \mathbb{R}^{m \times d}$, label matrix $Y = [\mathbf{y}_1, \dots, \mathbf{y}_m]^\top \in \mathbb{R}^{m \times n}$, $H \in \mathbb{R}^{n \times k}$, and parameters $k, \lambda_w, \lambda_g, \mathbf{p}, \mathbf{q}$, and \bar{a} . Choose $\beta < 1$ and $\sigma < 1/2$.

```
1:  $M = H^\top \text{diag}(\mathbf{q})H$   $\dots O(nk^2)$ 
2:  $\mathbf{d} = X^\top \mathbf{p}$   $\dots O(\text{nnz}(X))$ 
3:  $\mathbf{k} = H^\top \mathbf{q}$   $\dots O(nk)$ 
4: Compute  $\mathbf{p}^\top \mathbf{1}_m$  and  $\mathbf{q}^\top \mathbf{1}_n$   $\dots O(m+n)$ 
5:  $B = XW$ , where  $B = [\dots, \mathbf{b}_i, \dots]^\top$   $\dots O(mk^2)$ 
6:  $\hat{M} = \text{diag}(\mathbf{p})BM$   $\dots O(mk^2)$ 
7:  $\hat{L} = LB$   $\dots O(\text{nnz}(L)k)$ 
8:  $L^+ = \sum_{(i,j) \in \Omega^+} \ell_{ij}(Y_{ij}, \mathbf{b}_i^\top \mathbf{h}_j)$   $O(|\Omega^+|k)$ 
9: for  $t = 1, 2, 3, \dots$  do
10: // Compute gradient
11:  $D_{ij}^+ = \ell_{ij}^{+'}(Y_{ij}, \mathbf{b}_i^\top \mathbf{h}_j), \forall (i, j) \in \Omega^+$   $\dots O(|\Omega^+|k)$ 
12:  $G = X^\top(D^+H + 2\hat{M})$   $\dots O(|\Omega^+|k + \text{nnz}(X)k + mk^2)$ 
13:  $G = G - 2\bar{a}\mathbf{d}\mathbf{k}^\top + 2\lambda_w\lambda_g X^\top \hat{L}$   $\dots O(dk + \text{nnz}(X)k)$ 
14:
15: // Obtain a descent direction  $\Delta W$ 
16:  $\Delta W = -G$ 
17:
18: // Perform backtracking line search
19:  $\Delta B = X\Delta W$   $\dots O(\text{nnz}(X)k)$ 
20:  $\Delta \hat{M} = \text{diag}(\mathbf{p})\Delta B M$   $\dots O(mk^2)$ 
21:  $\Delta \delta = \mathbf{d}^\top \Delta W \mathbf{k}$   $\dots O(dk)$ 
22:  $\Delta \hat{L} = L\Delta B$   $\dots O(\text{nnz}(L)k)$ 
23:  $\text{Val.1} = \langle \Delta B, \hat{M} \rangle + \langle B, \Delta \hat{M} \rangle - 2\bar{a}\Delta \delta$ 
24:  $+ 2\lambda_w(\langle W, \Delta W \rangle + \lambda_g \langle B, \Delta \hat{L} \rangle) \dots O(dk + mk)$ 
25:  $\text{Val.2} = \langle \Delta B, \Delta \hat{M} \rangle$ 
26:  $+ \lambda_w(\|\Delta W\|_F^2 + \lambda_g \langle \Delta B, \Delta \hat{L} \rangle) \dots O(dk + mk)$ 
27:  $c = \langle G, \Delta W \rangle$   $\dots O(dk)$ 
28: for  $\alpha = 1, \beta, \beta^2, \dots$  do
29:  $B_{\text{new}} = B + \alpha \Delta B$   $\dots O(mk)$ 
30:  $L_{\text{new}}^+ = \sum_{(i,j) \in \Omega^+} \ell_{ij}(Y_{ij}, \mathbf{b}_i^{\text{new}\top} \mathbf{h}_j)$   $O(|\Omega^+|k)$ 
31: if  $-L^+ + L_{\text{new}}^+ + \alpha \text{Val.1} + \alpha^2 \text{Val.2} < \alpha \sigma c$  then
32:  $W = W + \alpha \Delta W$   $\dots O(dk)$ 
33:  $B = B_{\text{new}}$   $\dots O(mk)$ 
34:  $\hat{M} = \hat{M} + \alpha \Delta \hat{M}$   $\dots O(mk)$ 
35:  $\hat{L} = \hat{L} + \alpha \Delta \hat{L}$   $\dots O(mk)$ 
36:  $L^+ = L_{\text{new}}^+$   $\dots O(1)$ 
37: Break the line-search for-loop
```

Supp-2 Experimental Details

Supp-2.1 Evaluation Criteria and Data Sources

We used Precision@ k and nDCG@ k as evaluation criteria. For a predicted score vector $\hat{\mathbf{y}} \in \mathbb{R}^L$ and the true label vector $\mathbf{y} \in \{0, 1\}^L$, Precision@ k is defined as

$$p@k \equiv \frac{1}{k} \sum_{l \in \text{toprank}_k(\hat{\mathbf{y}})} y_l,$$

and nDCG@ k is defined as

$$\text{DCG}@k \equiv \sum_{l \in \text{toprank}_k(\hat{\mathbf{y}})} \frac{y_l}{\log(l+1)},$$

$$\text{nDCG}@k \equiv \frac{\text{DCG}@k}{\sum_{l=1}^{\min(k, \|\mathbf{y}\|_0)} \frac{1}{\log(l+1)}}.$$

For multi-label learning, the data sets are downloaded from *Mulan: A Java Library for Multi-Label Learning*³ and *The Extreme Classification Repository*.⁴ To avoid the 0/0 situation in the calculation of nDCG, data instances in the test set without any active labels (i.e., $\|\mathbf{y}\|_0 = 0$) are removed. For recommender systems with graph information, the data sources have been mentioned in the main text.

Supp-2.2 Parameter Selection

To select a suitable set of parameters, we hold out 1/5 of the training instances for multi-label problems. Similarly, we hold out 1/5 of the observed entries for graph structured one-class MF. The chosen parameters are those that give the highest validation performance on the hold out set. Since the best parameter for different measures ($p@1$ to $p@5$) may slightly vary, we used the one corresponding to the highest $p@5$. However, if the best parameter for $p@5$ has a very bad performance for other measures on the hold out set, we will choose the second best parameter for $p@5$.

For our methods, we fix \bar{a} in (3) to -1 because of the following reasons. First, for LR-wLR, we have

$$\ell_{ij}^+(Y_{ij}, \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij}) = \ell_{ij}^-(Y_{ij}, \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij})$$

$$= \log\left(1 + e^{-Y_{ij} \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij}}\right),$$

where $Y_{ij} = \pm 1$. Therefore, for $(i, j) \in \Omega^-$, its $Y_{ij} = \bar{a}$ should be -1 . For other formulations such as LR-wSQ, we have

$$\ell^-(Y_{ij}, \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij}) = (\bar{a} - \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij})^2.$$

The logistic loss $\ell^+(Y_{ij}, \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij})$ intends to have

$$(i, j) \in \Omega^+ \quad \text{if } \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij} \geq 0,$$

but for the squared loss $\ell^-(Y_{ij}, \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij})$, it intends to have

$$(i, j) \in \Omega^- \quad \text{if } \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij} \approx 0.$$

Clearly, a conflict occurs when $\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_{ij} > 0$ but ≈ 0 . On the other hand, for LEML (Yu et al. 2014), which is SQ-SQ, we use $\bar{a} = 0$ in order to be consistent with their experiment setting. For our methods that include a weighted

loss (SQ-wSQ, LR-wSQ and LR-wLR), $p_i q_j$ in assumption (3) is set to a constant ρ . We perform a grid search for $\log_2(\rho) \in \{-9, -7, -5, -3, -1, 0\}$. For Subsampled approaches, we fix $\rho = 1$ and perform the grid search on the ratio $|\Omega^-|/|\Omega^+| \in \{0.25, 0.5, 1, 2\}$.

For regularization parameters λ_w and λ_h , we set $\lambda_w = \lambda_h$ and perform a grid search on $\log_2(\lambda_w) \in \{-6, -4, -2, 0, 2, 4, 6\}$. For problems without graph information, we set $\lambda_g = 0$. For problems with graph information, λ_g is non-zero and is obtained by a grid search on $\log_2(\lambda_g \lambda_w) \in \{-6, -3, 0, 3, 6, 8, 10\}$. We consider the same setting for LEML.

For the kernel Nyström method in the last experiment, we use the same regularization and weight-loss parameters selected earlier for the setting without the Nyström method. For the Gaussian kernel $K(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|^2)$, we check $\gamma \in \{0.5/\sigma^2, 1/\sigma^2, 2/\sigma^2\}$, where σ^2 is the variance of all feature values in the entire training data set. We also conduct a grid search on the number of columns using $D \in \{500, 1000\}$.

For all one-class MF formulations, we run 15 alternating iterations in our validation process and obtain the validation performance at each iteration. We then use the iteration index yielding the best validation performance as the number of iterations to be run in training a model for predicting the test set. Note that the held-out validation set is included so we use the whole training set for generating the final model.

For the non-linear method SLEEC, since it includes 8 hyper-parameters, a grid search on this high dimensional space is not easy. Therefore we first fix the number of clusters to 1. Then the number of learners also becomes 1 because the use of several learners is due to the randomness in clustering. Next we conduct a grid search on the two most sensitive parameters: the k in kNN and the number of neighbors considered during the singular value projection method as suggested by the paper. For all other parameters we use the suggested values given by the authors.

Supp-2.3 More Experimental Results

The results are shown in both tables and figures.

- In Figure Supp-1, we show the results of Full versus Subsampled for multi-label learning in bar charts.
- In Figure Supp-2 and Figure Supp-3, we show the results of various loss functions for three types of one-class MF problems in bar charts.
- In Figure Supp-4, we show the results of the comparison of non-linear multi-label classifiers in bar charts.
- In Table Supp-1, we show the detailed results of recommender systems with and without graph information in terms of precision@1 to precisions@5 and nDCG@1 to nDCG@5.
- In Table Supp-2, we show the detailed results of the comparison among non-linear multi-label classifiers in terms of precision@1 to precisions@5 and nDCG@1 to nDCG@5.
- In Table Supp-3, we show the detailed results of the multi-label learning in terms of precision@1 to precisions@5 and nDCG@1 to nDCG@5.

³<http://mulan.sourceforge.net/datasets-mlc.html>

⁴<https://manikvarma.github.io/downloads/XC/XMLRepository.html>

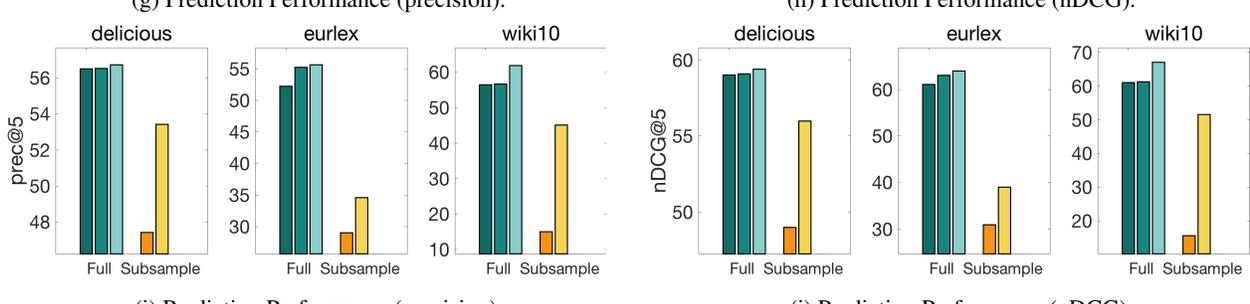
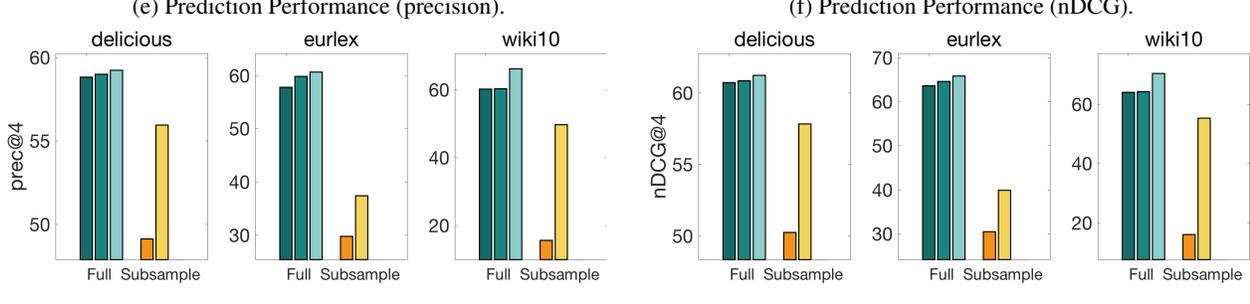
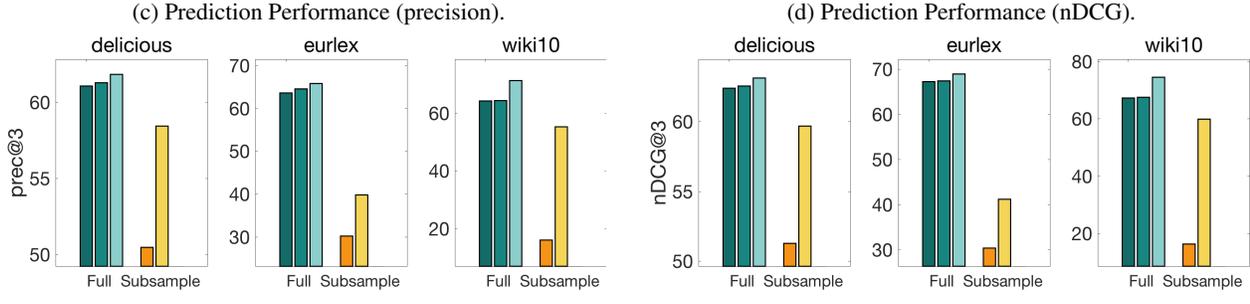
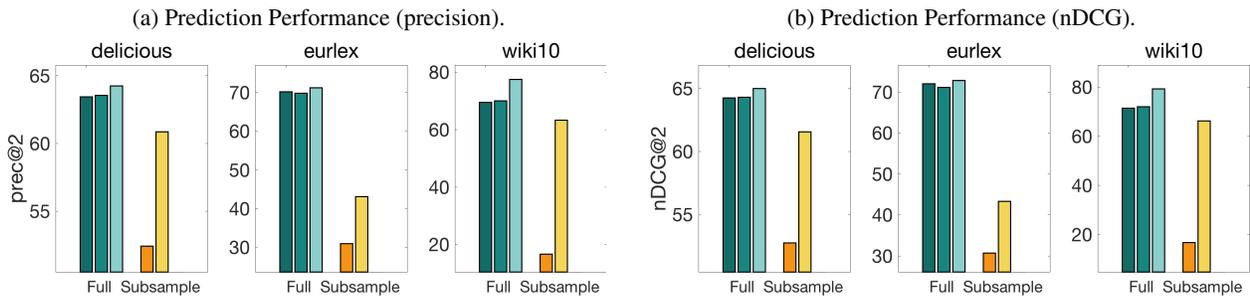
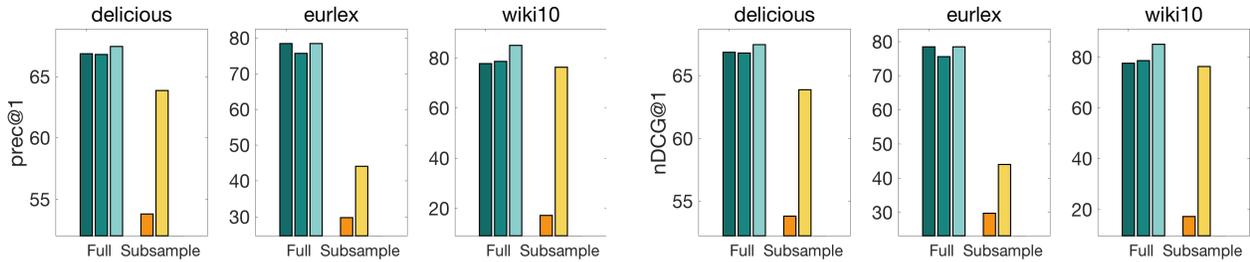


Figure Supp-1: Full v.s. Subsampled on multi-label learning with various loss functions. See Figure 1 for the legend information.

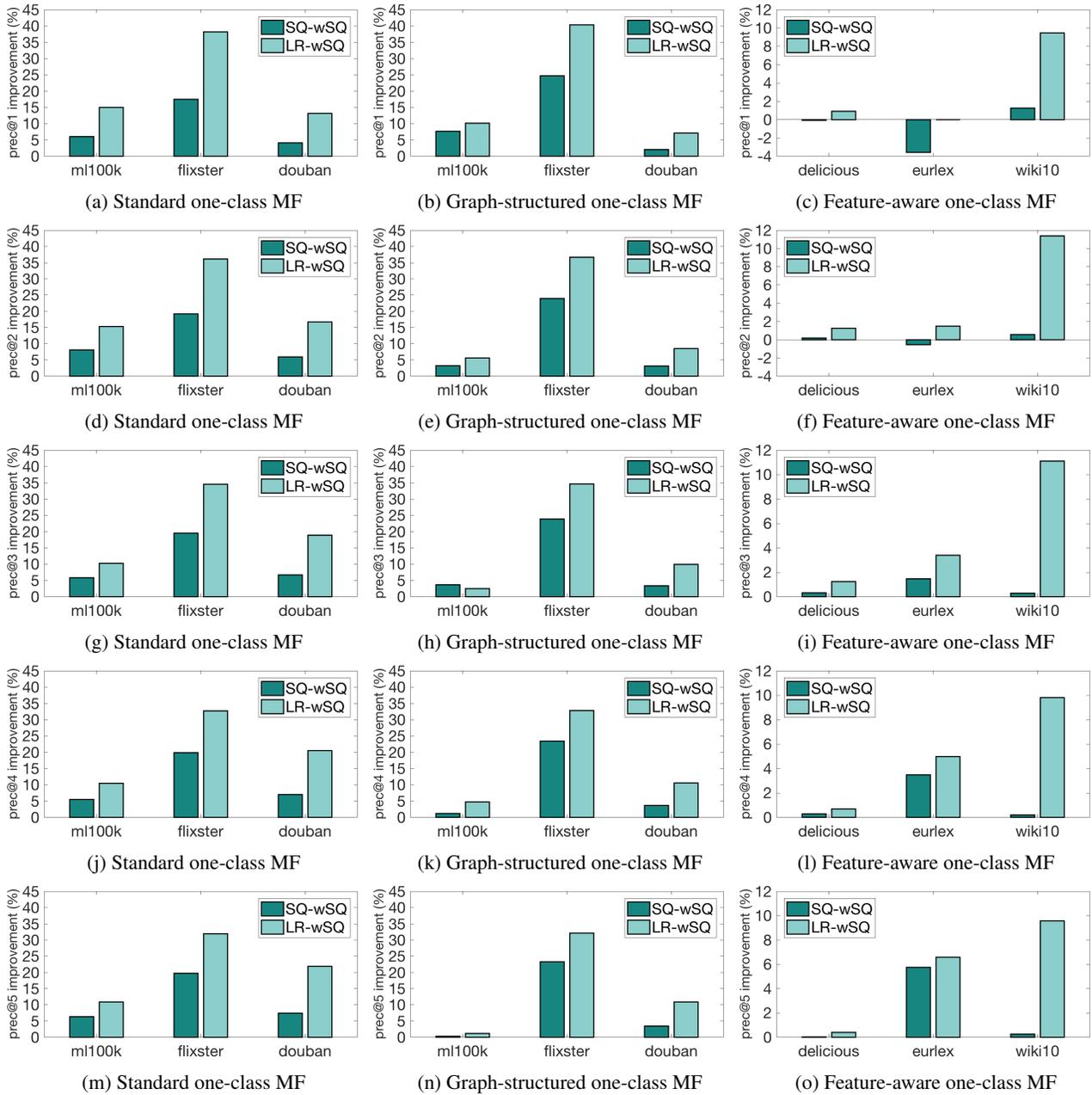


Figure Supp-2: Comparison on various loss functions. Y-axis is the improvement of Precision@k over the SQ-SQ formulation in percentage.

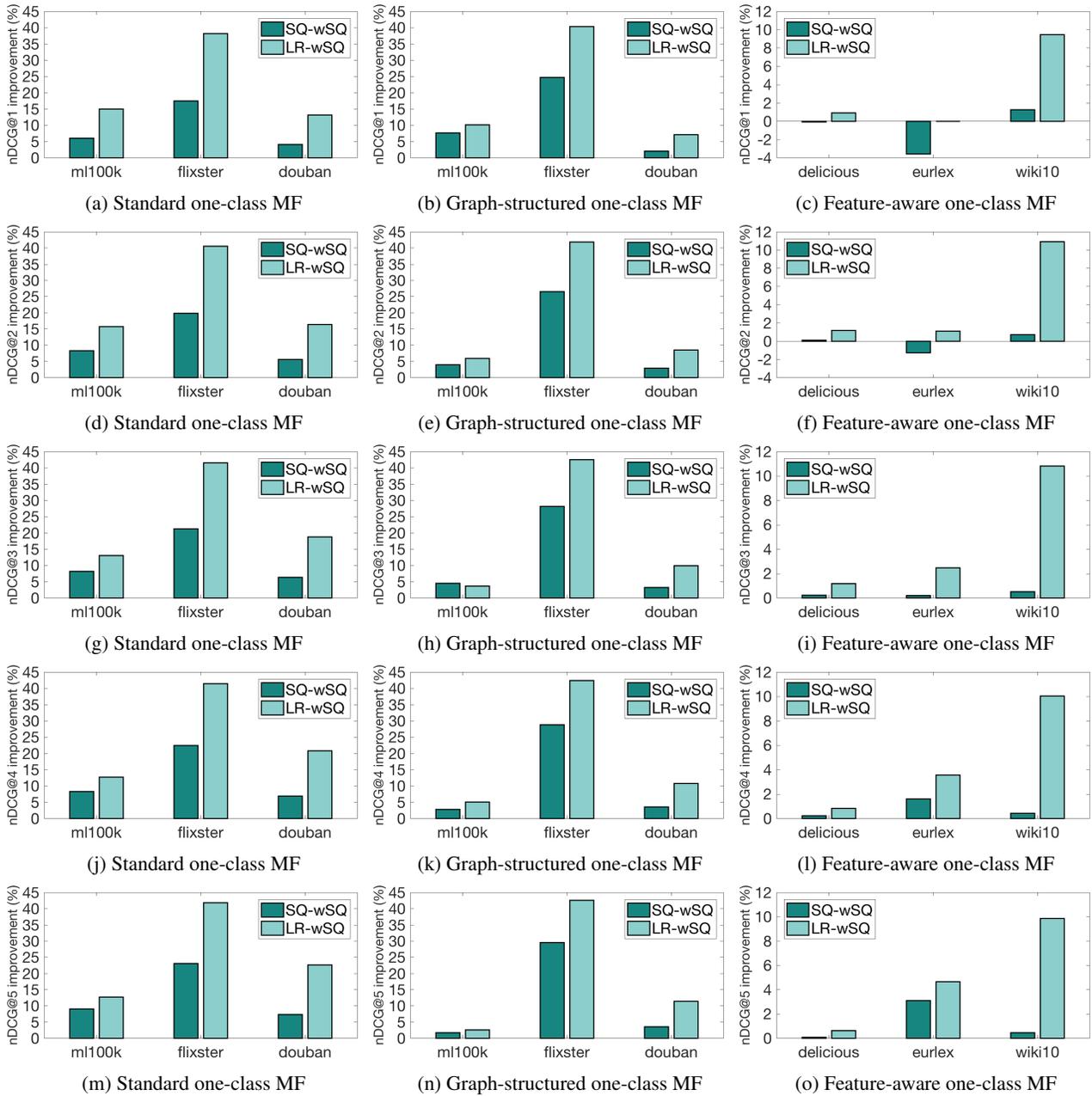
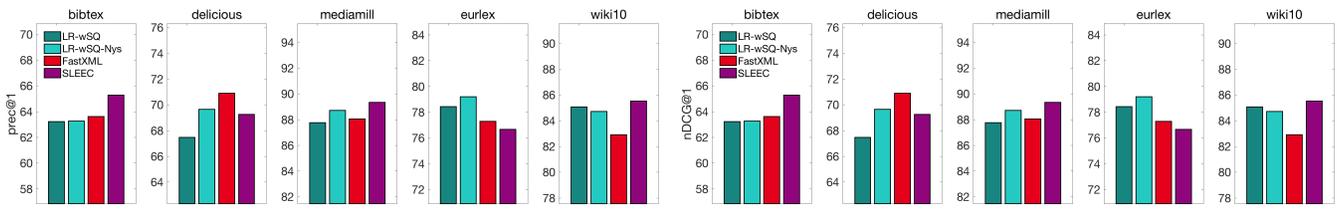
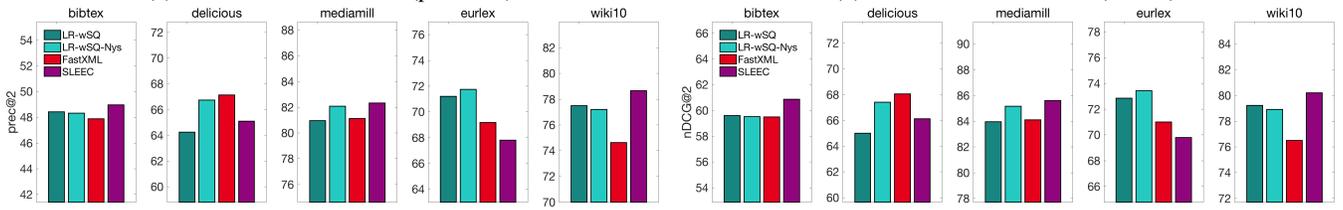


Figure Supp-3: Comparison on various loss functions. Y-axis is the improvement of nDCG@k over the SQ-SQ formulation in percentage.



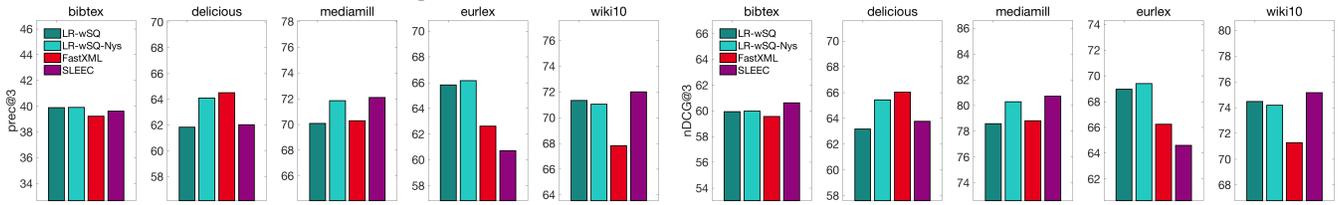
(a) Prediction Performance (precision).

(b) Prediction Performance (nDCG).



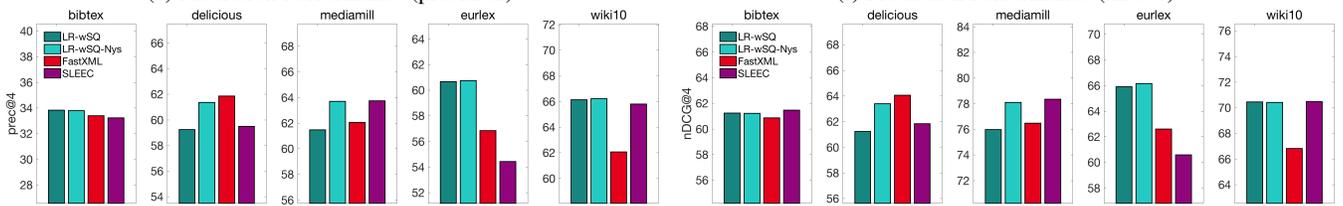
(c) Prediction Performance (precision).

(d) Prediction Performance (nDCG).



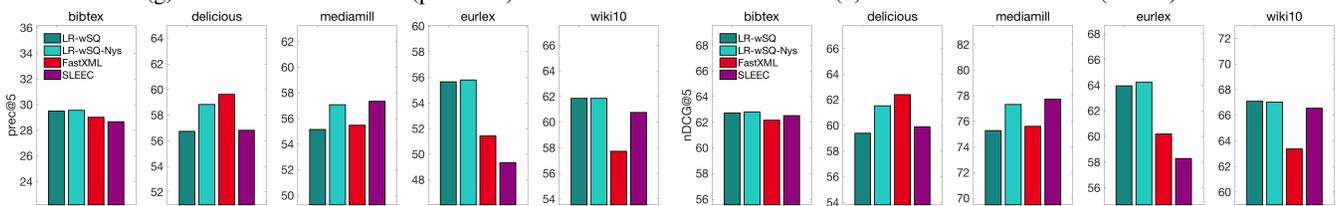
(e) Prediction Performance (precision).

(f) Prediction Performance (nDCG).



(g) Prediction Performance (precision).

(h) Prediction Performance (nDCG).



(i) Prediction Performance (precision).

(j) Prediction Performance (nDCG).

Figure Supp-4: Comparison on non-linear multi-label classifiers.

Table Supp-1: Comparison on graph structured one-class MF. Loss1-Loss2 denotes the formulation with the Loss1 on entries in Ω^+ and Loss2 on entries in Ω^- . wSQ/wLR denote the weighted square/logistic loss functions respectively. Note that the SQ-SQ formulation for the *graph structured one-class MF* part is equivalent to the formulation considered in Rao et al. (2015).

		(a) Precision@k					(b) nDCG@k					
		time	p@1	p@2	p@3	p@4	p@5	n@1	n@2	n@3	n@4	n@5
		<i>Standard one-class MF</i>					<i>Standard one-class MF</i>					
ml100k	SQ-SQ	1.2	26.93	22.20	20.00	18.47	16.95	26.93	23.99	23.00	22.84	22.54
	SQ-wSQ	0.9	28.56	23.99	21.16	19.48	18.04	28.56	25.98	24.88	24.72	24.59
	LR-wSQ	1.8	30.98	25.61	22.04	20.41	18.80	30.98	27.76	26.00	25.75	25.42
		<i>Graph structured one-class MF</i>					<i>Graph structured one-class MF</i>					
ml100k	SQ-SQ	1.3	28.55	24.62	22.04	20.26	19.05	28.55	26.54	25.57	25.28	25.38
	SQ-wSQ	0.9	30.75	25.43	22.85	20.49	19.12	30.75	27.59	26.73	25.99	25.82
	LR-wSQ	7.0	31.45	26.01	22.58	21.21	19.28	31.45	28.12	26.50	26.55	26.04
		<i>Standard one-class MF</i>					<i>Standard one-class MF</i>					
flixster	SQ-SQ	23.7	14.87	12.41	10.96	10.06	9.35	14.87	14.46	14.44	14.64	14.84
	SQ-wSQ	51.2	17.46	14.79	13.09	12.05	11.19	17.46	17.33	17.51	17.94	18.26
	LR-wSQ	161.3	20.54	16.89	14.74	13.34	12.34	20.54	20.33	20.44	20.72	21.04
		<i>Graph structured one-class MF</i>					<i>Graph structured one-class MF</i>					
flixster	SQ-SQ	27.4	14.66	12.37	10.97	10.06	9.35	14.66	14.35	14.38	14.59	14.78
	SQ-wSQ	54.0	18.28	15.34	13.60	12.42	11.53	18.28	18.16	18.43	18.80	19.15
	LR-wSQ	180.6	20.58	16.91	14.77	13.36	12.36	20.58	20.36	20.50	20.77	21.08
		<i>Standard one-class MF</i>					<i>Standard one-class MF</i>					
douban	SQ-SQ	64.8	17.42	15.07	13.56	12.47	11.62	17.42	15.84	14.97	14.44	14.09
	SQ-wSQ	100.3	18.13	15.97	14.47	13.35	12.48	18.13	16.74	15.93	15.43	15.13
	LR-wSQ	514.3	19.71	17.59	16.12	15.03	14.16	19.71	18.44	17.78	17.45	17.29
		<i>Graph structured one-class MF</i>					<i>Graph structured one-class MF</i>					
douban	SQ-SQ	58.3	18.89	16.61	15.07	13.97	13.11	18.89	17.41	16.60	16.14	15.89
	SQ-wSQ	75.0	19.27	17.12	15.58	14.47	13.56	19.27	17.93	17.14	16.71	16.45
	LR-wSQ	571.4	20.23	18.02	16.57	15.44	14.53	20.23	18.89	18.24	17.88	17.70

Table Supp-2: Comparison of state-of-the-art approaches on multi-label learning.

		(a) Precision@k					(b) nDCG@k					
		time	p@1	p@2	p@3	p@4	p@5	n@1	n@2	n@3	n@4	n@5
bibtex	LR-wSQ	22	63.22	48.43	39.89	33.83	29.50	63.22	59.59	59.93	61.24	62.73
	LR-wSQ-Nys	54	63.26	48.33	39.91	33.80	29.55	63.26	59.53	59.97	61.21	62.81
	FastXML	17	63.62	47.89	39.22	33.39	29.01	63.62	59.49	59.55	60.87	62.16
	SLEEC	249	65.29	48.97	39.63	33.24	28.76	65.29	60.87	60.62	61.46	62.52
delicious	LR-wSQ	23	67.47	64.24	61.85	59.25	56.73	67.47	64.99	63.16	61.23	59.40
	LR-wSQ-Nys	94	69.64	66.74	64.11	61.35	58.84	69.64	67.41	65.43	63.39	61.54
	FastXML	34	70.90	67.14	64.52	61.87	59.63	70.90	68.05	66.04	64.06	62.40
	SLEEC	1,124	69.27	65.10	62.03	59.49	56.82	69.27	66.11	63.78	61.84	59.89
mediamill	LR-wSQ	61	87.77	80.96	70.08	61.48	55.15	87.77	83.96	78.60	75.97	75.28
	LR-wSQ-Nys	144	88.72	82.08	71.85	63.70	57.08	88.72	85.17	80.29	78.10	77.32
	FastXML	256	88.04	81.12	70.29	62.06	55.48	88.04	84.10	78.81	76.48	75.61
	SLEEC	1,764	89.34	82.33	72.09	63.73	57.34	89.34	85.60	80.73	78.36	77.73
eurlex	LR-wSQ	404	78.43	71.21	65.82	60.67	55.64	78.43	72.84	68.97	65.91	63.96
	LR-wSQ-Nys	1,134	79.20	71.75	66.15	60.75	55.80	79.20	73.44	69.41	66.15	64.25
	FastXML	164	77.29	69.17	62.63	56.84	51.44	77.29	71.00	66.25	62.60	60.19
	SLEEC	1,497	76.67	67.80	60.71	54.45	49.34	76.67	69.80	64.61	60.58	58.28
wiki10	LR-wSQ	449	85.07	77.50	71.33	66.14	61.85	85.07	79.21	74.47	70.46	67.10
	LR-wSQ-Nys	666	84.72	77.20	71.05	66.22	61.86	84.72	78.90	74.18	70.41	67.02
	FastXML	744	82.91	74.63	67.82	62.06	57.73	82.91	76.50	71.27	66.83	63.37
	SLEEC	183	85.52	78.65	72.00	65.81	60.75	85.52	80.20	75.16	70.48	66.55

Table Supp-3: Comparison of various Full and Subsampled formulations on multi-label learning. Loss1-Loss2 denotes the formulation with the Loss1 on entries in Ω^+ and Loss2 on entries in Ω^- . wSQ/wLR denote the weighted square/logistic loss functions respectively. Note that the Full approach with SQ-SQ is equivalent to the LEML (Yu et al. 2014) formulation.

		(a) Precision@k					(b) nDCG@k					
		time	p@1	p@2	p@3	p@4	p@5	n@1	n@2	n@3	n@4	n@5
		<i>the Subsampled approach</i>					<i>the Subsampled approach</i>					
bibtex	SQ-SQ	12	46.04	23.16	31.31	27.59	24.63	46.04	45.47	46.71	48.87	50.66
	LR-LR	9	52.92	39.42	32.63	28.27	25.08	52.92	49.17	49.59	51.32	52.94
			<i>the Full approach</i>					<i>the Full approach</i>				
	SQ-SQ	13	63.14	47.58	38.77	33.09	28.81	63.14	58.99	58.95	60.34	61.71
	SQ-wSQ	12	63.30	48.35	39.78	33.74	29.37	63.30	59.61	59.96	61.26	62.65
	LR-wSQ	22	63.22	48.43	39.89	33.83	29.50	63.22	59.59	59.93	61.24	62.73
	LR-wLR	85	61.91	46.66	38.21	32.66	28.43	61.91	57.75	57.91	59.40	60.76
		<i>the Subsampled approach</i>					<i>the Subsampled approach</i>					
delicious	SQ-SQ	8	53.80	52.42	50.46	49.13	47.41	53.80	52.73	51.28	50.26	49.02
	LR-LR	19	63.89	60.86	58.44	55.95	53.43	63.89	61.56	59.71	57.84	55.99
			<i>the Full approach</i>					<i>the Full approach</i>				
	SQ-SQ	3	66.88	63.44	61.09	58.84	56.51	66.88	64.24	62.42	60.72	59.03
	SQ-wSQ	8	66.81	63.55	61.29	59.01	56.52	66.81	64.31	62.57	60.86	59.07
	LR-wSQ	23	67.47	64.24	61.85	59.25	56.73	67.47	64.99	63.16	61.23	59.40
	LR-wLR	446	67.91	65.04	62.27	59.81	57.27	67.91	65.74	63.68	61.85	59.99
		<i>the Subsampled approach</i>					<i>the Subsampled approach</i>					
mediamill	SQ-SQ	83	83.20	79.21	69.62	61.66	55.50	83.20	81.59	77.17	75.10	74.60
	LR-LR	93	87.62	80.72	69.87	61.37	55.08	87.62	83.74	78.39	75.82	75.15
			<i>the Full approach</i>					<i>the Full approach</i>				
	SQ-SQ	8	87.70	81.05	69.89	61.23	54.72	87.70	83.99	78.43	75.75	74.88
	SQ-wSQ	41	84.85	79.82	69.92	61.79	55.47	84.85	82.41	77.75	75.57	74.92
	LR-wSQ	61	87.77	80.96	70.08	61.48	55.15	87.77	83.96	78.60	75.97	75.28
	LR-wLR	437	87.78	81.17	70.38	61.85	55.33	87.78	84.17	78.87	76.33	75.50
		<i>the Subsampled approach</i>					<i>the Subsampled approach</i>					
eurlex	SQ-SQ	93	29.75	30.94	30.23	29.75	29.04	29.75	30.67	30.38	30.47	30.92
	LR-LR	166	44.08	43.09	39.83	37.39	34.58	44.08	43.32	41.17	39.90	38.99
			<i>the Full approach</i>					<i>the Full approach</i>				
	SQ-SQ	335	78.43	70.18	63.67	57.80	52.21	78.43	72.04	67.30	63.63	61.12
	SQ-wSQ	672	75.63	69.79	64.60	59.82	55.21	75.63	71.11	67.45	64.66	63.02
	LR-wSQ	404	78.43	71.21	65.82	60.67	55.64	78.43	72.84	68.97	65.91	63.96
		<i>the Subsampled approach</i>					<i>the Subsampled approach</i>					
wiki10	SQ-SQ	121	17.23	16.58	16.14	15.59	15.04	17.23	16.73	16.38	15.98	15.56
	LR-LR	764	76.30	63.29	55.28	49.72	45.09	76.30	66.23	59.91	55.40	51.62
			<i>the Full approach</i>					<i>the Full approach</i>				
	SQ-SQ	247	77.71	69.58	64.20	60.23	56.45	77.71	71.42	67.20	64.04	61.07
	SQ-wSQ	872	78.69	69.97	64.39	60.35	56.59	78.69	71.94	67.55	64.31	61.34
	LR-wSQ	449	85.07	77.50	71.33	66.14	61.85	85.07	79.21	74.47	70.46	67.10