Online Matrix Completion with Side Information

Mark Herbster, Stephen Pasteris, Lisa Tse

Abstract

We present online transductive and inductive algorithms for binary matrix completion with side information. For these algorithms, we prove novel mistake and expected regret bounds. In the case of no side information, the bounds scale with the dimensionality of the matrix. In the case of ideal side information, these scale with the number of underlying latent factors.

The Model

Online Matrix Completion

On a trial t = 1, ..., T:

- 1. the learner is queried by the *environment* to predict matrix entry (i_t, j_t)
- 2. the learner predicts a label $\hat{y}_t \in \{-1, 1\}$
- 3. the learner receives a label $y_t \in \{-1, 1\}$ from the environment and
- 4. a mistake is incurred if $y_t \neq \hat{y}_t$.

	<i>m</i> Movies								
n Users		1	1		-1				-1
						1		1	
		-1		1				1	
					1		-1		
	-1		1						
			1					-1	

Side Information

Assume that we are given additional side information about each row and column. For instance with movies, we might have "genre" information and for the users, we might have "demographics". For our purposes, we assume that this can be summarised as $m \times m$ positive definite matrix M for the row side information and $n \times n$ positive definite matrix N for the column side information.

Previous results in online learning

Matrix completion on real-valued matrices [1]

Binary matrix completion [2]

Mistake bound:

$$\sum_{t \in [T]} [y_t \neq \hat{y}_t] \le \tilde{\mathcal{O}} \left((m+n) \operatorname{mc}^2(U) \right).$$

No side information

Matrix completion with graph side information on specific matrix classes with only mistake bounds [3,4]

Main Results

General regret and mistake bounds

Mistake bound:

$$\sum_{t \in [T]} [y_t \neq \hat{y}_t] \le \tilde{\mathcal{O}} \left(\mathsf{mc}^2(\mathbf{U}) \mathcal{D} \right).$$

Expected regret bound:

$$\sum_{t \in [T]} \mathbb{E}[y_t \neq \hat{y}_t] - \sum_{t \in [T]} [y_t \neq U_{i_t j_t}] \le \tilde{\mathcal{O}}\left(\sqrt{\mathcal{D}} ||\boldsymbol{U}||_{\max}^2 T\right)$$

for all $U \in \{-1, 1\}^{m \times n}$.

The *max-norm*:

$$egin{aligned} egin{aligned} egin{aligne} egin{aligned} egin{aligned} egin{aligned} egin$$

where the minimum is over all matrices P, Q and every integer d.

The margin complexity:

$$\mathsf{mc}(\boldsymbol{U}) := \min_{\boldsymbol{V} \in \mathsf{SP}^1(\boldsymbol{U})} \|\boldsymbol{V}\|_{\mathsf{max}} = \min_{\boldsymbol{P}\boldsymbol{Q}^\top \in \mathsf{SP}(\boldsymbol{U})} \max_{ij} \frac{\|\boldsymbol{P}_i\| \|\boldsymbol{Q}_j\|}{|\langle \boldsymbol{P}_i, \boldsymbol{Q}_j \rangle|}$$

where the minimum is over all matrices P, Q and every integer d. Here, $SP(U) = \{V \in \Re^{m \times n} : \forall_{ij} V_{ij} U_{ij} > 0\}$ and $SP^1(U) = \{V \in \Re^{m \times n} : \forall_{ij} V_{ij} U_{ij} \ge 1\}$

The quasi-dimension:

$$\mathcal{D}^{\gamma}_{\boldsymbol{M},\boldsymbol{N}}(\boldsymbol{U}) := \min_{\hat{\boldsymbol{P}}\hat{\boldsymbol{Q}}^{\top} = \gamma\boldsymbol{U}} \mathcal{R}_{\boldsymbol{M}} \operatorname{tr} \left(\hat{\boldsymbol{P}}^{\top} \boldsymbol{M} \hat{\boldsymbol{P}} \right) + \mathcal{R}_{\boldsymbol{N}} \operatorname{tr} \left(\hat{\boldsymbol{Q}}^{\top} \boldsymbol{N} \hat{\boldsymbol{Q}} \right) \,,$$

where the minimum is over all row-normalized matrices \hat{P} , \hat{Q} and every integer d and where $\mathcal{R}_{M} := \max_{i \in [m]} M_{ii}^{+}$. If the infimum does not exist then $\mathcal{D}_{M,N}^{\gamma}(U) := +\infty$.

Biclustered matrices

The class of (k, ℓ) -binary-biclustered matrices is defined as

 $\mathbb{B}_{k,\ell}^{m,n} = \{ \boldsymbol{U} \in \{-1,1\}^{m \times n} : \boldsymbol{r} \in [k]^m, \boldsymbol{c} \in [\ell]^n, \boldsymbol{U}^* \in \{-1,1\}^{k \times \ell}, U_{ij} = U_{r_i c_j}^*, i \in [m], j \in [n] \}$







Feature vectors in well-separated boxes with the min kernel

Graph Laplacians



 $\mathcal{D} \le \mathcal{O}(k+\ell)$

 $\mathcal{D} \le \mathcal{O}(k^2 + \ell^2)$

Algorithms

Transductive algorithm: instance of MEG [5]



Inductive algorithm

For the inductive algorithm, we are given kernel functions instead of the matrices *M* and *N*. On each trial, we are also given feature vectors for the row and column.

We can show that the transductive and inductive algorithms are predictionequivalent, up to the value of \Re_M and \Re_N .

References

[1] E. Hazan, S. Kale, and S. Shalev-Shwartz (2012)
[2] M. Herbster, S. Pasteris, and M. Pontil (2008)
[3] M. Herbster, S. Pasteris, and M. Pontil (2015)
[4] C. Gentile, M. Herbster, and S. Pasteris (2013)
[5] K. Tsuda, G. Rätsch, and M.K. Warmuth (2005)



