

# Nonparametric approach for binary matrix completion

Miaoyan Wang, Sep 1, 2020

## 1 Problem

Suppose that we observe a subset of entries from a binary matrix,  $\{y_{ij} \in \{-1, 1\} : (i, j) \in \Omega\}$ , where  $\Omega \subset [d_1] \times [d_2]$  is the index set of observed entries. How to predict the unobserved entries  $\{y_{ij} \in \{-1, 1\} : (i, j) \in \Omega^c\}$ ?

$$\begin{bmatrix} -1 & ? & ? & -1 & ? \\ ? & 1 & ? & ? & ? \\ -1 & ? & ? & -1 & ? \\ ? & ? & -1 & ? & 1 \end{bmatrix} \quad (1)$$

## 2 Earlier solution

First, we perform probability estimation based on parametric models. Assume  $y_{ij}$  are independent Bernoulli random variables with success probabilities  $P(y_{ij} = 1)$  for all  $(i, j) \in [d_1] \times [d_2]$ . We model the probability matrix using the GLM logistic model,

$$\mathbb{P}(y_{ij} = 1) = \frac{e^{\theta_{ij}}}{1 + e^{\theta_{ij}}}, \quad \text{where } \Theta = [\theta_{ij}] \in \mathbb{R}^{d_1 \times d_2} \text{ is a rank-}r \text{ matrix.}$$

Define the rank- $r$  maximum log-likelihood estimator  $\hat{\Theta}^{\text{MLE}} = \llbracket \hat{\theta}_{ij}^{\text{MLE}} \rrbracket = \arg \min_{\Theta \in \mathcal{P}(r, \alpha)} L(\Theta)$ , where

$$L(\Theta) = - \sum_{(i,j) \in \Omega} \log(e^{y_{ij}\theta_{ij}} + 1), \quad \text{and} \quad (2)$$

$$P(r, \alpha) = \{\Theta \in \mathbb{R}^{d_1 \times d_2} : \text{rank}(\Theta) \leq r \text{ and } \|\Theta\|_{\infty} \leq \alpha\}.$$

Second, we perform prediction using plug-in estimates,

$$\hat{y}_{ij} = \text{sign } \hat{\theta}_{ij}^{\text{MLE}}, \quad \text{for all } (i, j) \in \Omega^c.$$

## 3 New proposal

If our goal is to predict the unobserved entries by two labels  $\{-1, 1\}$ , there is no need to estimate the probability. We could directly perform the prediction in a nonparametric fashion. This scenario reduces to a special case of our matrix-valued classification problem.

1. Feature space:

$$\begin{aligned}\mathcal{X} &= \{\mathbf{X} \in \{0, 1\}^{d_1 \times d_2} \mid \text{only one entry of } \mathbf{X} \text{ is one, and others are zero}\} \\ &= \{\mathbf{e}_i \otimes \mathbf{e}_j : (i, j) \in [d_1] \times [d_2]\}.\end{aligned}$$

2. Outcome space:  $\mathcal{Y} \in \{0, 1\}$ .

3. Uniform marginal distribution  $\mathcal{P}(\mathbf{X})$  over  $\mathcal{X}$ . No other joint distribution assumptions on  $P(\mathbf{X}, y)$ ;

4. i.i.d. training set:  $\{(\mathbf{X}_{ij}, y_{ij}) : (i, j) \in \Omega\}$ , where  $\mathbf{X}_{ij} = \mathbf{e}_i \otimes \mathbf{e}_j \in \{0, 1\}^{d_1 \times d_2}$  is an indicator matrix specifying the observed index, and  $y_{ij} \in \{-1, 1\}$  is the observed label at index  $(i, j)$ . For example, the features in the training sample for problem (1) are

$$\mathbf{X}_1 = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix}, \quad \mathbf{X}_2 = \begin{bmatrix} 0 & \cdots & 1 & 0 \\ 0 & \ddots & \vdots & \vdots \\ 0 & \cdots & 0 & 0 \end{bmatrix}, \quad \cdots, \quad \mathbf{X}_7 = \begin{bmatrix} 0 & \cdots & 0 & 0 \\ 0 & \ddots & \vdots & \vdots \\ 0 & \cdots & 0 & 1 \end{bmatrix}.$$

5. Define the rank- $r$  large-margin estimator  $\hat{\Theta}^{\text{margin}} = \llbracket \hat{\theta}_{ij}^{\text{margin}} \rrbracket = \arg \min_{\Theta \in \mathcal{P}(r, \alpha)} L(\Theta)$ , where

$$L(\Theta) = \sum_{(i,j) \in \Omega} [1 - y_{ij} \langle \mathbf{X}_{ij}, \Theta \rangle]_+, \text{ and} \quad (3)$$

$$\mathcal{P}(r, \alpha) = \{\Theta \in \mathbb{R}^{d_1 \times d_2} : \text{rank}(\Theta) \leq r \text{ and } \|\Theta\|_\infty \leq \alpha\}.$$

Here, we have omitted the intercept for simplicity.

6. Predict unobserved entries using  $\hat{y}_{ij} = \text{sign } \hat{\theta}_{ij}^{\text{margin}}$ .

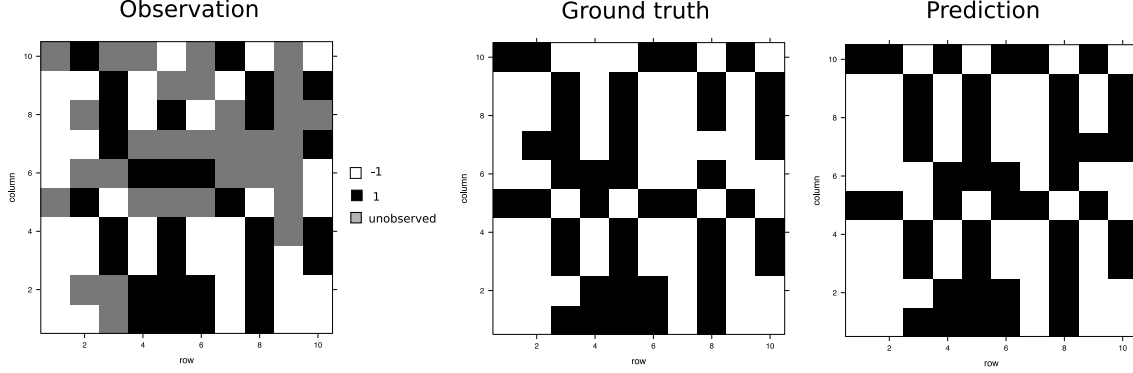
7. Nonparametric probability estimation  $\hat{\mathbb{P}}(y_{ij} = 1 | \mathbf{X}_{ij})$  is also possible using a sequence of weighted low-rank classifications (3).

## 4 Numerical experiments

### 4.1 Missing data imputation

dimension  $d_1 = d_2 = 10$ ; rank = 2; cost = 1; observation probability  $p = 0.6$ .

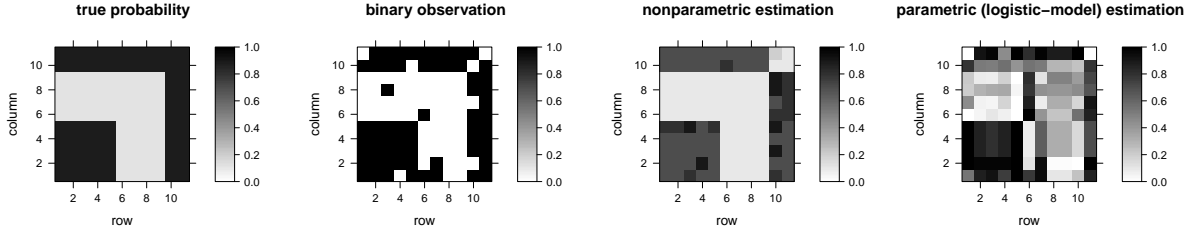
	Unobserved		Observed	
	pred = 1	pred = -1	pred = 1	pred = -1
true = 1	16	3	36	1
true = -1	1	12	1	30



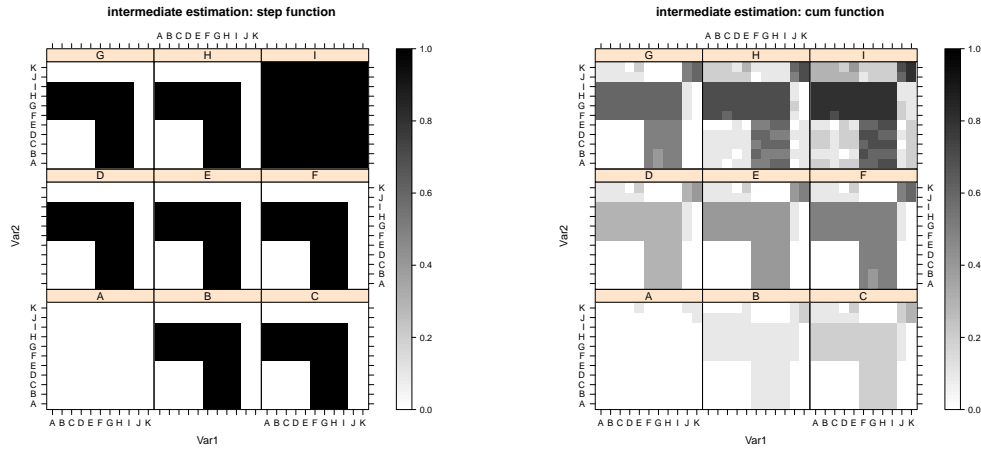
## 4.2 Probability estimation

dimension  $d_1 = d_2 = 11$ ; cost = 1; observation probability  $p = 1$  (no missing data).

Goal: estimate probability matrix  $P \in [0, 1]^{d_1 \times d_2}$  from binary observations  $\mathbf{Y} = \{0, 1\}^{d_1 \times d_2}$ .



Intermediate steps.



Define cumulative probability matrices  $P_h = \mathbf{1}(P \leq \frac{h}{10}) \in \{0, 1\}^{d_1 \times d_2}$  for  $h = 1, \dots, 10$ . Panel A:  $\frac{1}{10} \sum_{h \leq 1} P_h$ , Panel B:  $\frac{1}{10} \sum_{h \leq 2} P_h, \dots$ , Panel I:  $\frac{1}{10} \sum_{h \leq 9} P_h$ .

## 5 Theory

**Definition 1** (Misclassification error). Let  $\mathbf{Y} = \llbracket y_{ij} \rrbracket$ ,  $\mathbf{Z} = \llbracket z_{ij} \rrbracket \in \{0, 1\}^{d_1 \times d_2}$  be two binary matrices. We define the misclassification error (MCE),

$$\text{MCE}(\mathbf{Y}, \mathbf{Z}) = \frac{1}{d_1 d_2} \sum_{(i,j) \in [d_1] \times [d_2]} \mathbb{1}\{y_{ij} \neq z_{ij}\}.$$

**Theorem 5.1** (Generalization error bounds). *Consider a binary target matrix  $\mathbf{Y} = \llbracket y_{ij} \rrbracket \in \{-1, 1\}^{d_1 \times d_2}$  whose entries are independent realizations from some unknown distributions  $\text{Ber}(p_{ij})$ , for all  $(i, j) \in [d_1] \times [d_2]$ . Suppose that we observe a subset of entries,  $\mathbf{Y}_\Omega := \{y_{ij}\}_{(i,j) \in \Omega}$ , where  $\Omega \subset [d_1] \times [d_2]$  is a random set with  $|\Omega|$  entries, and each entry in  $\Omega$  is an i.i.d. drawn uniformly from  $[d_1] \times [d_2]$ . Let  $\hat{\Theta} = \llbracket \hat{\theta}_{ij} \rrbracket \in \mathcal{P}(r, \alpha)$  denote any estimator based on the observations  $\mathbf{Y}_\Omega$ , where  $r$  is the rank bound and  $\alpha$  is the infinity norm bound. Then, with probability at least  $1 - \delta$  over  $\mathbf{Y}$  and the sample selection  $\Omega$ , the following bound holds uniformly for all  $\hat{\Theta} \in \mathcal{P}(r, \alpha)$ ,*

$$\underbrace{\text{MCE}(\mathbf{Y}, \text{sign } \hat{\Theta})}_{\text{misclassification error in targeted matrix}} \leq \underbrace{\frac{L}{|\Omega|} \sum_{(i,j) \in \Omega} \text{surrogate-loss}(y_{ij} \hat{\theta}_{ij})}_{\text{surrogate loss in sample}} + C_1 \alpha \sqrt{\frac{(d_1 + d_2)r}{|\Omega|}} + C_2 \sqrt{\frac{\log(3/\delta)}{2|\Omega|}},$$

where  $C_1, C_2 > 0$  are two universal constants, and  $L > 0$  is the Lipschitz constant of the surrogate loss,

$$L = \begin{cases} 1, & \text{for hinge loss } S(t) = (1 - t)_+, \\ \frac{1}{\log 2}, & \text{for logistic loss } S(t) = \log_2(e^t + 1), \end{cases}$$

In particular, the generalization error of  $\hat{\Theta}$  converges to zero as long as the sample size  $|\Omega| \geq \tilde{\mathcal{O}}(d_{\max} r)$ .

**Corollary 1** (Large-margin estimator). *Consider the same set-up as in Theorem 5.1. Let  $\Theta^* \in \mathbb{R}^{d_1 \times d_2}$  be the optimal estimator in  $\mathcal{P}(r, \alpha)$  that minimizes the MCE, i.e.,*

$$\Theta^* = \arg \min_{\Theta \in \mathcal{P}(r, \alpha)} \text{MCE}(\mathbf{Y}, \text{sign } \Theta).$$

Then, for the constrained MLE defined in (2) and large-margin estimator defined in (3), we have

$$\begin{aligned} \text{MCE}(\mathbf{Y}, \text{sign } \hat{\Theta}^{\text{margin}}) - \text{MCE}(\mathbf{Y}, \text{sign } \Theta^*) &\leq C_1 \alpha \sqrt{\frac{(d_1 + d_2)r}{|\Omega|}} + C_2 \sqrt{\frac{\log(1/\delta)}{2|\Omega|}}, \\ \text{MCE}(\mathbf{Y}, \text{sign } \hat{\Theta}^{\text{MLE}}) - \text{MCE}(\mathbf{Y}, \text{sign } \Theta^*) &\leq C_1 \alpha \sqrt{\frac{(d_1 + d_2)r}{|\Omega|}} + C_2 \sqrt{\frac{\log(1/\delta)}{2|\Omega|}}, \end{aligned}$$

with probability at least  $1 - \delta$  over  $\mathbf{Y}$  and the sample selection  $\Omega$ .

**Remark 1** (Approximation error). What is the total estimation error of  $\text{sign } \hat{\Theta}$  from the bayes rule? Two sources of error: generalization error + approximation error. The approximation error reaches zero when the “bayes rule” binary matrix is included in the set of candidate sign matrices. Namely, there exists a low-rank, entrywise bounded matrix  $\Theta^* \in \mathcal{P}(r, \alpha)$  such that

$$\Theta^* \stackrel{\text{equal in sign}}{=} \llbracket p_{ij} - 0.5 \rrbracket, \text{ or equivalently, } \text{MCE}(\Theta^*, \underbrace{\text{sign}(p_{ij} - 0.5)}_{\text{“bayes rule” binary matrix}}) = 0.$$

**Remark 2.** Given a bayes rule binary matrix, how can we tell whether it is the sign matrix for some low-rank matrix in  $\mathbb{P}(r, \alpha)$ ? For matrix completion problem, the sample size  $|\Omega|$  is always smaller than the feature dimension  $d_1 d_2$ . **What does “decision boundary” mean when the feature space is discrete?**

**Remark 3.** If full rank, then  $\theta_{ss'} = \text{intercept} = \text{sample average} = \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} y_{ij}^{\text{test}}$  for all  $(s, s') \in \Omega^c$ . Non-vanishing MCE unless  $|\Omega| \approx d_1 d_2$ .

**Remark 4.** **MCE vs. MSE. sharpness compared to earlier paper?**

**Remark 5** (Nonlinear extension).

## 6 Proofs

*Proof of Theorem 5.1.* Because the desired conclusion is a uniform bound over all  $\hat{\Theta} \in \mathcal{P}(r, \alpha)$ , we write  $\Theta$  in place of  $\hat{\Theta}$  for notational convenience. One should note that  $\Theta$  is a random variable depending on the realizations of the training set  $\{y_{ij}\}_{(i,j) \in \Omega}$ .

Define the function class  $\mathcal{F} = \{f(\mathbf{X}): \mathbf{X} \mapsto \langle \mathbf{X}, \Theta \rangle \mid \Theta \in \mathcal{P}(r, M)\}$ . Given the features in the training set,  $\{\mathbf{X}_{ij} = \mathbf{e}_i \otimes \mathbf{e}_j : (i, j) \in \Omega\}$ , we consider the empirical Rademacher complexity of  $\mathcal{F}$  conditional on  $\Omega$  and the training set. Let  $\{\xi_{ij}\}$  a set of i.i.d. Rademacher random variables with equal probability on  $\pm 1$ , then

$$\mathcal{R}_\Omega(\mathcal{F}) = \frac{2}{|\Omega|} \mathbb{E}_{\xi_{ij}} \left\{ \sup_{\Theta \in \mathcal{P}(r, \alpha)} \sum_{(i,j) \in \Omega} \xi_{ij} \Theta_{ij} \right\}. \quad (4)$$

Note that  $\Theta \in \mathcal{P}(r, M)$  implies that  $\|\Theta\|_{\max} \leq \sqrt{r} \alpha$ , where  $\|\Theta\|_{\max} = \min_{\Theta = \mathbf{U}^T \mathbf{V}} \{\|\mathbf{U}\|_{2, \infty} \|\mathbf{V}\|_{2, \infty}\}$  denotes the matrix max-qnorm. Therefore, the inequality (4) is upper bounded,

$$\begin{aligned} \frac{2}{|\Omega|} \mathbb{E}_{\xi_{ij}} \left\{ \sup_{\Theta \in \mathcal{P}(r, M)} \sum_{(i,j) \in \Omega} \xi_{ij} \Theta_{ij} \right\} &\leq \frac{2}{|\Omega|} \mathbb{E}_{\xi_{ij}} \left\{ \sup_{\|\Theta\|_{\max} \leq \sqrt{r} \alpha} \sum_{(i,j) \in \Omega} \xi_{ij} \Theta_{ij} \right\} \\ &\leq c \alpha \sqrt{\frac{r(d_1 + d_2)}{|\Omega|}}, \end{aligned}$$

where the last inequality follows from Ghadermarzy et al. [2019, Lemma 31].

Using the generalization error inequality in the earlier notes, we have that, with probability at least  $1 - \delta$  over the sample selection  $\Omega$  and training data  $\{y_{ij}^{\text{train}}\}_{(i,j) \in \Omega}$ , the following bound holds uniformly over  $\Theta = \llbracket \theta_{ij} \rrbracket \in \mathcal{P}(r, M)$ ,

$$\mathbb{P} [y^{\text{test}} \neq \text{sign} \langle \mathbf{X}^{\text{test}}, \Theta \rangle] \leq \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} \text{hinge-loss}(y_{ij}^{\text{train}}, \theta_{ij}) + C_1 \alpha \sqrt{\frac{r(d_1 + d_2)}{|\Omega|}} + \sqrt{\frac{\log(1/\delta)}{2|\Omega|}}, \quad (5)$$

for some constant  $C_1 > 0$ , where the probability at the left hand side is taken with respect to  $(\mathbf{X}^{\text{test}}, y^{\text{test}}) \in \{\mathbf{e}_s \otimes \mathbf{e}'_{s'} : (s, s') \in [d_1] \times [d_2]\} \times \{0, 1\}$ , independent of training data  $\{y_{ij}^{\text{train}}\}_{(i,j) \in \Omega}$ .

Now, the i.i.d. uniform sampling assumption implies the mutual independence of the events  $\mathbb{1}\{y_{ss'} \neq \text{sign} \theta_{ss'}\}$  and marginal uniform distribution  $\mathbb{P}(\mathbf{X}^{\text{test}} = \mathbf{e}_s \otimes \mathbf{e}'_{s'}) = \frac{1}{d_1 d_2}$  for all  $(s, s') \in [d_1] \times [d_2]$ . By properties of conditional expectation and concentration inequality, we have, for any  $\alpha > 0$ ,

$$\begin{aligned} \mathbb{P} [y^{\text{test}} \neq \text{sign} \langle \mathbf{X}^{\text{test}}, \Theta \rangle] &= \frac{1}{d_1 d_2} \sum_{(s,s') \in [d_1] \times [d_2]} \mathbb{E}_{y_{ss'}^{\text{test}}} \mathbb{1}\{y_{ss'}^{\text{test}} \neq \text{sign} \theta_{ss'}\} \\ &\geq \frac{1}{d_1 d_2} \sum_{(s,s') \in [d_1] \times [d_2]} \mathbb{1}\{y_{ss'}^{\text{test}} \neq \text{sign} \theta_{ss'}\} - C\alpha \sqrt{\frac{1}{d_1 d_2}}, \end{aligned} \quad (6)$$

where the last statement holds with probability at least  $1 - \exp(-\alpha^2)$  over the test matrix  $\mathbf{Y}^{\text{test}} = \llbracket y_{ss'}^{\text{test}} \rrbracket \in \{0, 1\}^{d_1 \times d_2}$ .

Combining (5) and (6) with  $\alpha = \sqrt{\log(1/\delta)}$  yields the uniform bound for all  $\Theta \in \mathcal{P}(r, \alpha)$ ,

$$\text{MCE}(\mathbf{Y}^{\text{test}}, \text{sign} \Theta) \leq \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} \text{hinge-loss}(y_{ij}^{\text{train}}, \theta_{ij}) + C_1 \alpha \sqrt{\frac{r(d_1 + d_2)}{|\Omega|}} + C_2 \sqrt{\frac{\log(1/\delta)}{2|\Omega|}}, \quad (7)$$

with probability at least  $1 - 2\alpha$  taken jointly over the test data  $\mathbf{Y}^{\text{test}}$ , sample selection  $\Omega$ , and training data  $\{y_{ij}^{\text{train}}\}_{(i,j) \in \Omega}$ . Note that in the bound (7), the test data  $\mathbf{Y}^{\text{test}}$  at the left hand side and training data  $\{y_{ij}^{\text{train}}\}_{(i,j) \in \Omega}$  at the right hand side are independent of each other.

We write the target binary matrix  $\mathbf{Y} = \llbracket y_{ij} \rrbracket \in \{0, 1\}^{d_1 \times d_2}$  as

$$y_{ij} = \begin{cases} y_{ij}^{\text{train}}, & (i, j) \in \Omega, \\ y_{ij}^{\text{test}}, & (i, j) \in \Omega^c. \end{cases}$$

Then, the classification error satisfies

$$\text{MCE}(\mathbf{Y}, \text{sign} \Theta) = \frac{1}{d_1 d_2} \left\{ \sum_{(s,s') \in \Omega^c} \mathbb{1}\{y_{ss'}^{\text{test}} \neq \text{sign} \theta_{ss'}\} + \sum_{(s,s') \in \Omega} \mathbb{1}\{y_{ss'}^{\text{test}} \neq \text{sign} \theta_{ss'}\} \right\}$$

$$\begin{aligned}
& + \frac{1}{d_1 d_2} \left\{ \sum_{(s,s') \in \Omega} \mathbb{1} \{y_{ss'}^{\text{train}} \neq \text{sign } \theta_{ss'}\} - \sum_{(s,s') \in \Omega} \mathbb{1} \{y_{ss'}^{\text{test}} \neq \text{sign } \theta_{ss'}\} \right\} \\
& \leq \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} \text{hinge-loss}(y_{ij}^{\text{train}}, \theta_{ij}) + C_1 \alpha \sqrt{\frac{r(d_1 + d_2)}{|\Omega|}} + C'_2 \sqrt{\frac{\log(1/\delta)}{2|\Omega|}},
\end{aligned}$$

with probability at least  $1 - 3\delta$ , where the last line follows from (7) and the concentration inequality for  $\sum_{(s,s') \in \Omega} [\mathbb{1} \{y_{ss'}^{\text{train}} \neq \text{sign } \theta_{ss'}\} - \mathbb{1} \{y_{ss'}^{\text{test}} \neq \text{sign } \theta_{ss'}\}]$ .  $\square$

## References

Navid Ghadermarzy, Yaniv Plan, and Özgür Yilmaz. Near-optimal sample complexity for convex tensor completion. *Information and Inference: A Journal of the IMA*, 8(3):577–619, 2019.