

Graphic Lasso: Clustering accuracy for precision matrix model

Jiaxin Hu

January 30, 2021

1 Accuracy

Consider the optimization problem.

$$\begin{aligned}
 \min_{\mathbf{U}, \{\Theta^l\}} \quad & Q(\mathbf{U}, \{\Theta^l\}) = \sum_{k=1}^K \text{tr}(S^k \Omega^k) - \log \det(\Omega^k), \\
 \text{s.t.} \quad & \Omega^k = \sum_{l=1}^r u_{kl} \Theta^l, \\
 & |\Theta^l|_1 < b, \quad \text{for } l \in [r], \\
 & \mathbf{U} \in \{0, 1\}^{K \times r} \text{ is a membership matrix,} \\
 & \{\Theta^l\} \text{ are irreducible and invertible.}
 \end{aligned} \tag{1}$$

Notations

Note that the definitiona of confusion matrix and MCR differ with the definitions in tensor block model (TBM) by a factor K (d in TBM). Therefore, the final result includes a K in the dominator in this model.

1. Σ^l for $l = 1, \dots, r$: the true covariance matrices, where $(\Theta^l)^{-1} = \Sigma^l$
2. $D = \llbracket D_{ij} \rrbracket \in [0, 1]^{r \times r}$: the confusion matrix between the true membership matrix \mathbf{U} and the estimation $\hat{\mathbf{U}}$, where $D_{ij} = \sum_{k=1}^K \mathbf{I}(u_{ki} = \hat{u}_{kj} = 1)$.
3. $I_l = \{k : u_{kl} = 1\}$ for $l = 1, \dots, r$: the index set of the categories belong to the group l . The sets I_l rely on the membership \mathbf{U} . Note that $\sum_{l=1}^r |I_l| = K$. Let \hat{I}_l denote the index sets according to the estimation $\hat{\mathbf{U}}$.
4. $MCR(\hat{\mathbf{U}}, \mathbf{U}) = \max_{l, a \neq a' \in [r]} \min\{D_{al}, D_{a'l}\}$: the misclassification rate.
5. $\delta = \min_{k \neq l \in [r]} \min_{(i,j)} (\Sigma_{ij}^k - \Sigma_{ij}^l)^2$: the minimal gap between the true covariance matrices.
6. $\varphi \min(\cdot)$: the minimal singular value of the matrix.
7. $\varphi \max(\cdot)$: the maximal singular value of the matrix,
8. $n_k > 0$ for $k \in [K]$: the n_k is the sample size for the k -th category.

Theorem 1.1. Suppose the singular value for the true precision matrices are bounded, i.e., $0 < s < \min_{l \in [r]} \varphi_{\min}(\Theta^l) \leq \max_{l \in [r]} \varphi_{\max}(\Theta^l) < \tau < \infty$, where $s < \tau$ are positive constants. The minimal gap between the precision matrices $\delta = \min_{k \neq l \in [r]} \min_{(i,j)} (\Sigma_{ij}^k - \Sigma_{ij}^l)^2$ is larger than 0. Let $\hat{\mathbf{U}}$ denote the minimizer of the objective function (1). Then, for any $\epsilon \in [0, 1]$, we have

$$\mathbb{P}(MCR(\hat{\mathbf{U}}, \mathbf{U}) \geq \epsilon) \leq C_1 \exp \left(-\frac{C_2 \min_{k \in [K]} n_k \epsilon^2 \delta^2}{s^2 K^2 r^2 b^2 \max_{l \in [r]} \|\Theta^l\|_{\max}^2} \right),$$

where C_1, C_2 are two constants only depend on τ .

Proof. Recall the objective function

$$Q(\mathbf{U}, \{\Theta^l\}) = \sum_{k=1}^K \text{tr}(S^k \Omega^k) - \log \det(\Omega^k) = \sum_{k=1}^K \langle S^k, \Omega^k \rangle - \log \det(\Omega^k).$$

The deviation between the true parameters $\{\mathbf{U}, \{\Theta^l\}\}$ and estimations $\{\hat{\mathbf{U}}, \{\hat{\Theta}^l\}\}$ comes from two aspects: the estimation of $\{\Theta^l\}$ and the misclassification (estimation of \mathbf{U}). We tease apart these two parts.

1. First, we suppose the membership \mathbf{U} is given. We now assess the stochastic error due to the estimation of $\{\Theta^l\}$, conditional on \mathbf{U} . Note that the objective function is convex. The optimal $\{\Theta^l\}$ satisfies the first order condition.

$$\frac{\partial Q}{\partial \Theta_{ij}^l} = \sum_{k \in I_l} S_{ij}^k - |I_l| \frac{C(\Theta^l)_{ij}}{\det(\Theta^l)} = 0,$$

where $C(A)_{ij}$ is the cofactor of matrix A corresponding to the element A_{ij} . Note that Θ^l should be a symmetric matrix. The cofactor matrix $C(\Theta^l) = C^T(\Theta^l)$. Then, the derivation of the matrix Θ^l is equal to

$$\frac{\partial Q}{\partial \Theta^l} = \sum_{k \in I_l} S^k - |I_l| \frac{C^T(\Theta^l)}{\det(\Theta^l)} = 0,$$

which implies that

$$\hat{\Theta}^l = \left(\frac{\sum_{k \in I_l} S^k}{|I_l|} \right)^{-1}, \quad \text{for } l \in [r].$$

Therefore, the estimation of $\{\Theta^l\}$ is a function of \mathbf{U} . Consider the function $F(\mathbf{U}) = Q(\mathbf{U}, \{\hat{\Theta}^l\})$. By a straightforward calculation, we have

$$\begin{aligned} F(\mathbf{U}) &= \sum_{l=1}^r \sum_{k \in I_l} \langle S^k, \left(\frac{\sum_{k \in I_l} S^k}{|I_l|} \right)^{-1} \rangle - |I_l| \log \det \left(\left(\frac{\sum_{k \in I_l} S^k}{|I_l|} \right)^{-1} \right) \\ &= \sum_{l=1}^r \langle \sum_{k \in I_l} S^k, \left(\frac{\sum_{k \in I_l} S^k}{|I_l|} \right)^{-1} \rangle - |I_l| \log \det \left(\left(\frac{\sum_{k \in I_l} S^k}{|I_l|} \right)^{-1} \right) \\ &= \sum_{l=1}^r |I_l| p - |I_l| \log \det \left(\frac{\sum_{k \in I_l} S^k}{|I_l|} \right)^{-1}. \end{aligned}$$

Note that $\sum_{l=1}^r |I_l|p = Kp$ is independent with the membership. We only need to consider the second term. For simplicity, we define

$$F(\mathbf{U}) = - \sum_{l=1}^r |I_l| \log \det \left(\left(\frac{\sum_{k \in I_l} S^k}{|I_l|} \right)^{-1} \right).$$

Note that $\mathbb{E} \left[\frac{\sum_{k \in I_l} S^k}{|I_l|} \right] = \frac{\sum_{a=1}^r D_{al} \Sigma^l}{|I_l|}$. Correspondingly, we define the population version of $F(\mathbf{U})$ as following.

$$G(\mathbf{U}) = - \sum_{l=1}^r |I_l| \log \det \left(\left(\frac{\sum_{a=1}^r D_{al} \Sigma^a}{|I_l|} \right)^{-1} \right),$$

where $\Sigma^k = \mathbb{E}[S^k]$ is the true covariance matrices. Therefore, the deviation $F(\mathbf{U}) - G(\mathbf{U})$ quantifies the stochastic error due to the estimation of $\{\Theta^l\}$.

2. Next, we free \mathbf{U} and quantify the total deviation. Considering the maximizer,

$$\hat{\mathbf{U}} = \arg \min_{\mathbf{U}} F(\mathbf{U}).$$

The corresponding $G(\hat{\mathbf{U}})$ is

$$G(\hat{\mathbf{U}}) = - \sum_{l=1}^r |\hat{I}_l| \log \det \left(\left(\frac{\sum_{a=1}^r D_{al} \Sigma^a}{|\hat{I}_l|} \right)^{-1} \right),$$

and the function $G(\mathbf{U})$ with true membership is

$$G(\mathbf{U}) = - \sum_{l=1}^r |I_l| \log \det \left(\left(\frac{\sum_{k \in I_l} \Sigma^l}{|I_l|} \right)^{-1} \right) = \sum_{l=1}^r |I_l| \log \det(\Sigma^l).$$

Then, the deviation $G(\hat{\mathbf{U}}) - G(\mathbf{U})$ measures the stochastic error of the misclassification.

Now back to the probability of misclassification rate. By Lemma 1, we have

$$\mathbb{P}(MCR(\hat{\mathbf{U}}, \mathbf{U}) \geq \epsilon) \leq \mathbb{P}(G(\hat{\mathbf{U}}) - G(\mathbf{U}) \leq -\frac{1}{4s} \epsilon \delta).$$

Notice that the total deviation between $\hat{\mathbf{U}}$ and \mathbf{U} is able to be decomposed into three parts.

$$\begin{aligned} F(\hat{\mathbf{U}}) - F(\mathbf{U}) &= [F(\hat{\mathbf{U}}) - G(\hat{\mathbf{U}})] + [G(\hat{\mathbf{U}}) - G(\mathbf{U})] + [G(\mathbf{U}) - F(\mathbf{U})] \\ &\leq 2m - \frac{1}{4s} \epsilon \delta, \end{aligned}$$

where $m = \sup_{\mathbf{U}} |F(\mathbf{U}) - G(\mathbf{U})|$. Since $\hat{\mathbf{U}}$ is the minimizer of the objective function, we know that $F(\hat{\mathbf{U}}) - F(\mathbf{U}) \leq 0$. Therefore, we obtain the accuracy of misclassification rate

$$\begin{aligned} \mathbb{P}(MCR(\hat{\mathbf{U}}, \mathbf{U}) \geq \epsilon) &\leq \mathbb{P}(F(\hat{\mathbf{U}}) - F(\mathbf{U}) \leq 2m - \frac{1}{4s}\epsilon\delta) \\ &\leq \mathbb{P}(m \geq \frac{1}{8s}\epsilon\delta) \\ &\leq \mathbb{P}\left(\max_{k, (i,j)} |S_{(i,j)}^k - \mathbb{E}[S_{(i,j)}^k]| \geq \frac{\epsilon\delta}{8sKr b \max_{l \in [r]} \|\Theta^l\|_{\max}}\right), \\ &\leq C_1 \exp\left(-\frac{C_2 \min_{k \in [K]} n_k \epsilon^2 \delta^2}{s^2 K^2 r^2 b^2 \max_{l \in [r]} \|\Theta^l\|_{\max}^2}\right) \end{aligned}$$

where the third inequality follow by Lemma 2, the last inequality follows by the Lemma 3, and C_1, C_2 are two constants. \square

Lemma 1. Assume the minimal singular-value of the true precision matrices is lower bounded $\min_{l \in [r]} \varphi_{\min}(\Theta^l) > s$, where s is a positive constant, and the minimal gap between covariance matrices $\delta > 0$. For any fixed $\epsilon > 0$, suppose the misclassification rate $MCR(\hat{\mathbf{U}}, \mathbf{U}) \geq \epsilon$, we have

$$G(\hat{\mathbf{U}}) - G(\mathbf{U}) \leq -\frac{1}{4s}\epsilon\delta.$$

Proof. Note that for an invertible matrix A , $\det(A^{-1}) = \frac{1}{\det(A)}$. Recall the formula of $G(\mathbf{U})$ and $G(\hat{\mathbf{U}})$. We have

$$G(\mathbf{U}) = \sum_{l=1}^r |I_l| \log \det(\Sigma^l), \quad \text{and} \quad G(\hat{\mathbf{U}}) = \sum_{l=1}^r |\hat{I}_l| \log \det\left(\frac{\sum_{a=1}^r D_{al} \Sigma^a}{|\hat{I}_l|}\right).$$

Note that

$$\sum_{l=1}^r |\hat{I}_l| \left(\frac{\sum_{a=1}^r D_{al} \log \det(\Sigma^a)}{|\hat{I}_l|}\right) = \sum_{a=1}^r \sum_{l=1}^r D_{al} \log \det(\Sigma^a) = G(\mathbf{U}),$$

where the second equality follows by the fact that $\sum_{l=1}^r D_{al} = |I_a|$. Since $MCR(\hat{\mathbf{U}}, \mathbf{U}) \geq \epsilon$, there exist $l, k \neq k' \in [r]$ such that $\min\{D_{kl}, D_{k'l}\} \geq \epsilon$. Let $\tilde{\Sigma} = \frac{\sum_{a=1}^r D_{al} \Sigma^a}{|\hat{I}_l|}$, and $\Delta = \Sigma - \tilde{\Sigma}$ for some matrix Σ . Consider the function $f(t) = \log \det(\tilde{\Sigma} + t\Delta)$. By Taylor Expansion, we have

$$\log \det(\Sigma) - \log \det(\tilde{\Sigma}) = f(1) - f(0) = f'(0) + \frac{f''(\xi)}{2}, \quad \text{for some } \xi \in [0, 1], \quad (2)$$

where

$$f'(0) = \langle \tilde{\Sigma}, \Delta \rangle, \quad \text{and} \quad f''(\xi) = \text{vec}(\Delta)^T (\tilde{\Sigma} + \xi\Delta)^{-1} \otimes (\tilde{\Sigma} + \xi\Delta)^{-1} \text{vec}(\Delta). \quad (3)$$

Particularly, by the definition of singular value, we have the lower bound of the second derivative

$$f''(\xi) = \text{vec}(\Delta)^T (\tilde{\Sigma} + \xi\Delta)^{-1} \otimes (\tilde{\Sigma} + \xi\Delta)^{-1} \text{vec}(\Delta) \geq \|\Delta\|_F^2 s, \quad (4)$$

where $\|\cdot\|_F$ is the matrix Frobenius norm.

Let $\Delta^l = \Sigma^l - \tilde{\Sigma}$, $l \in [r]$. Combining the Taylor Expansion (2) with the lower bound (4), we have

$$\begin{aligned} \left(\frac{\sum_{a=1}^r D_{al} \log \det(\Sigma^a)}{|\hat{I}_l|} \right) - \log \det(\tilde{\Sigma}) &= \sum_{a=1}^l \frac{D_{al}}{|\hat{I}_l|} \left[\log \det(\Sigma^a) - \log \det(\tilde{\Sigma}) \right] \\ &\geq \sum_{a=1}^r \frac{D_{al}}{|\hat{I}_l|} \left(\langle \tilde{\Sigma}, \Delta^a \rangle + \frac{1}{2} s \|\Delta^a\|_F^2 \right) \\ &\geq \frac{D_{kl}}{2|\hat{I}_l|} s \|\Delta^k\|_F^2 + \frac{D_{k'l}}{2|\hat{I}_l|} s \|\Delta^{k'}\|_F^2, \end{aligned}$$

where the last inequality follows by the fact that $\sum_{a=1}^r \frac{D_{al}}{|\hat{I}_l|} \langle \tilde{\Sigma}, \Delta^a \rangle = 0$. By the inequality $\frac{1}{2} \|A + B\|_F^2 \leq \|A\|_F^2 + \|B\|_F^2$, we obtain that

$$\left(\frac{\sum_{a=1}^r D_{al} \log \det(\Sigma^a)}{|\hat{I}_l|} \right) - \log \det(\tilde{\Sigma}) \geq \frac{\min\{D_{kl}, D_{k'l}\} s}{|\hat{I}_l|} \|\Sigma^k - \Sigma^{k'}\|_F^2 \geq \frac{\epsilon}{4s|\hat{I}_l|} \delta. \quad (5)$$

For other $l' \in [r]/l$, since $\log \det(\cdot)$ is a convex function, by Jensen's inequality, we have

$$\left(\frac{\sum_{a=1}^r D_{al'} \log \det(\Sigma^a)}{|\hat{I}_{l'}|} \right) - \log \det \left(\frac{\sum_{a=1}^r D_{al'} \Sigma^a}{|\hat{I}_{l'}|} \right) \geq 0. \quad (6)$$

Combining the the inequality (5) and (6), we obtain the misclassification error

$$G(\hat{\mathbf{U}}) - G(\mathbf{U}) = \sum_{l=1}^r |\hat{I}_l| \log \det \left(\frac{\sum_{a=1}^r D_{al} \Sigma^a}{|\hat{I}_l|} \right) - \sum_{l=1}^r |\hat{I}_l| \left(\frac{\sum_{a=1}^r D_{al} \log \det(\Sigma^a)}{|\hat{I}_l|} \right) \leq \frac{1}{4s} \epsilon \delta.$$

□

Lemma 2. Suppose we have $|\Theta^l|_1 < b$ for all $l \in [r]$, where $|A|_1$ is the number of nonzero elements in matrix A . Then, we have

$$|F(\mathbf{U}) - G(\mathbf{U})| \leq Krb \max_{l \in [r]} \|\Theta^l\|_{\max} \max_{k, (i,j)} |S_{(i,j)}^k - \mathbb{E}[S_{(i,j)}^k]|$$

Proof. Recall the formula of $F(\mathbf{U})$ and $G(\mathbf{U})$, where \mathbf{U} may not be the true membership matrix. We have

$$|F(\mathbf{U}) - G(\mathbf{U})| \leq \sum_{l=1}^r |I_l| \left| \log \det \left(\frac{\sum_{k \in I_l} S^k}{|I_l|} \right) - \log \det \left(\mathbb{E} \left[\frac{\sum_{k \in I_l} S^k}{|I_l|} \right] \right) \right|.$$

Consider the function $f(t) = \log \det \left(\frac{\sum_{k \in I_l} S^k}{|I_l|} + t\Delta \right)$, where $\Delta = \mathbb{E} \left[\frac{\sum_{k \in I_l} S^k}{|I_l|} \right] - \frac{\sum_{k \in I_l} S^k}{|I_l|}$. By the previous calculation (3), we know that $f(t)$ is a convex function. Then, the function is locally Lipschitz with $L = \sup_t |f'(t)|$. Therefore, we have

$$\begin{aligned} |F(\mathbf{U}) - G(\mathbf{U})| &\leq \sum_{l=1}^r |I_l| |f(1) - f(0)| \\ &\leq \sum_{l=1}^r |I_l| |f'(1)| \\ &\leq K \sup \left| \left\langle \left(\mathbb{E} \left[\frac{\sum_{k \in I_l} S^k}{|I_l|} \right] \right)^{-1}, \frac{\sum_{k \in I_l} S^k}{|I_l|} - \mathbb{E} \left[\frac{\sum_{k \in I_l} S^k}{|I_l|} \right] \right\rangle \right|. \end{aligned}$$

Since $(A)^{-1}$ is convex function of A , we have

$$\left\| \left(\mathbb{E} \left[\frac{\sum_{k \in I_l} S^k}{|I_l|} \right] \right)^{-1} \right\|_{\max} \leq \left\| \left(\frac{\sum_{k \in I_l} \mathbb{E}[S^k]^{-1}}{|I_l|} \right) \right\|_{\max} \leq \max_{l \in [r]} \left\| \Theta^l \right\|_{\max}.$$

We also have the sparsity

$$\left| \left(\frac{\sum_{k \in I_l} \mathbb{E}[S^k]^{-1}}{|I_l|} \right) \right|_1 \leq rb.$$

Therefore, we obtain the upper bound

$$|F(\mathbf{U}) - G(\mathbf{U})| \leq Krb \max_{l \in [r]} \left\| \Theta^l \right\|_{\max} \max_{k, (i,j)} |S_{(i,j)}^k - \mathbb{E}[S_{(i,j)}^k]|.$$

□

Lemma 3. Let $Z_i \sim_{i.i.d.} \mathcal{N}(0, \Sigma)$ and $\varphi_{\max}(\Sigma) \leq \tau < \infty$. Let $\Sigma = \llbracket \Sigma_{ij} \rrbracket$, then

$$P \left(\left| \sum_{i=1}^n Z_{ij} Z_{ik} - n \Sigma_{jk} \right| \geq n\nu \right) \leq c_1 e^{-c_2 n \nu^2}, \text{ for } |\nu| \leq \delta,$$

where c_1, c_2, δ depends on τ only.

Proof. See Lemma 1 of Rothman et.al.

□