

# Sparse covariance and precision matrix estimation under matched design

Yukun Li, Raji Balasubramanian

University of Massachusetts Amherst

March 2021

# Covariance Matrix and Precision Matrix

- Covariance matrix and precision matrix are two essential parts in multivariate analysis since they represent marginal and conditional dependence structures respectively.
- Let  $\mathbf{X}$  be multivariate normal with covariance matrix  $\Sigma$ , and the precision matrix  $\Omega$  is defined to be the inverse of the covariance matrix  $\Omega = \Sigma^{-1}$ . For  $i \neq j$ ,

$\Sigma_{ij} = 0$  if and only if  $X_i$  and  $X_j$  are marginally independent

$\Omega_{ij} = 0$  if and only if  $X_i$  and  $X_j$  are conditionally independent  
given all other variables

- The estimation of high dimensional covariance/precision matrix based on few sample observations is a difficult problem. Sample covariance matrix is singular and noninvertible.

# Matched case-control Study

- Case-control study: One of the observational study designs for performing clinical research. Economical, quick to perform, and easy to implement.
- Matched case-control study: Match the cases and controls for confounding factors.
- Matched data are not randomly collected.
- Problem: Recover the covariance/precision matrix of variables under matched case-control design

# Inverse probability weighting

- A well-known technique used in controlling for selection biases in non-experimental studies and in many missing data problems.
- Create a pseudo population by giving a weight to each observation in the case control sample.
- Weight: inverse of sampling probability. Up-weight the observations which have low probability of being in the case control sample, and down-weight those that have high probability.

- Covariance estimation: Sample covariance with inverse probability weighting

$$\mathbf{S}_{ipw} = \frac{1}{\sum_{i=1}^n w_i - 1} \sum_{i=1}^n w_i (\mathbf{X}_i - \hat{\boldsymbol{\mu}}^*) (\mathbf{X}_i - \hat{\boldsymbol{\mu}}^*)^T \quad (1)$$

where  $w_i$  is weight for the  $i^{th}$  subject,  $\hat{\boldsymbol{\mu}}^* = \frac{1}{\sum_{i=1}^n w_i} \sum_{i=1}^n w_i \mathbf{X}_i$  is the weighted sample mean.

- Precision estimation: Inverse of estimated covariance

$$\hat{\boldsymbol{\Omega}} = \mathbf{S}_{ipw}^{-1} \quad (2)$$

# Method: High Dimension

- Given a set  $\mathbf{X}$  of  $n$  i.i.d. vectors following  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , we have the weighted log likelihood function

$$\log f(\mathbf{X}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{\sum_{i=1}^n w_i}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^n w_i (\mathbf{X}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{X}_i - \boldsymbol{\mu}) + c$$

Adopting the idea from Bien & Tibshirani (2010) and Friedman et al. (2007), we can add lasso penalty to the likelihood function.

- Covariance matrix:

$$\hat{\boldsymbol{\Sigma}} = \operatorname{argmax}_{\boldsymbol{\Sigma}} \{-\log \det \boldsymbol{\Sigma} - \operatorname{tr}(\mathbf{S}_{ipw}^* \boldsymbol{\Sigma}^{-1}) - \lambda \|\boldsymbol{\Sigma}\|_1\} \quad (3)$$

- Precision matrix:

$$\hat{\boldsymbol{\Omega}} = \operatorname{argmax}_{\boldsymbol{\Omega}} \{\log \det \boldsymbol{\Omega} - \operatorname{tr}(\mathbf{S}_{ipw}^* \boldsymbol{\Omega}) - \lambda \|\boldsymbol{\Omega}\|_1\} \quad (4)$$

where  $\mathbf{S}_{ipw}^* = \frac{1}{\sum_{i=1}^n w_i} \sum_{i=1}^n w_i (\mathbf{X}_i - \boldsymbol{\mu})^T (\mathbf{X}_i - \boldsymbol{\mu})$ , and  $\lambda$  is the tuning parameter

# Method: High Dimensional Covariance Matrix Estimation

- Problem (3)

$$\hat{\Sigma} = \operatorname{argmax}_{\Sigma} \{-\log \det \Sigma - \operatorname{tr}(\mathbf{S}_{ipw}^* \Sigma^{-1}) - \lambda \|\Sigma\|_1\}$$

- This problem is not convex. According to Bien & Tibshirani (2010), it can be solved by majorize-minimize iteration

$$\hat{\Sigma}^{(t)} = \operatorname{argmin}_{\Sigma} \left[ \operatorname{tr} \left\{ \left( \hat{\Sigma}^{(t-1)} \right)^{-1} \Sigma \right\} + \operatorname{tr} \left( \Sigma^{-1} \mathbf{S}_{ipw}^* \right) + \lambda \|\Sigma\|_1 \right]$$

- Then generalized gradient descent can be applied to solve this problem

$$\Sigma \leftarrow \mathcal{S} \left\{ \Sigma - t \left( \Sigma_0^{-1} - \Sigma^{-1} \mathbf{S}_{ipw}^* \Sigma^{-1} \right), \lambda t \right\} \quad (5)$$

where  $\mathcal{S}$  is the elementwise soft-thresholding operator defined by  $\mathcal{S}(A, B)_{ij} = \operatorname{sign}(A_{ij}) (A_{ij} - B_{ij})_+$ .

- Problem (4)

$$\hat{\Omega} = \operatorname{argmax}_{\Omega} \{ \log \det \Omega - \operatorname{tr}(\mathbf{S}_{ipw}^* \Omega) - \lambda \|\Omega\|_1 \}$$

- Graphical Lasso (Friedman et al.(2007)) can be applied to solve this problem by optimizing over each row and corresponding column of the covariance matrix in a block coordinate descent fashion.



## Method: Weight estimation

- Suppose the subjects in the cohort are divided into  $J$  strata based on  $Z$
- Randomly select  $n_1$  cases from the  $N_1$  cases in the cohort
- Select controls by matching for strata  $J$
- The weight can be estimated by

$$\hat{w}_{i1} = N_1/n_1$$

$$\hat{w}_{i0} = N_{0j}/n_{0j} \text{ for control subject in stratum } j$$

where

$N_1$ : number of case subjects in the cohort

$n_1$ : number of case subjects in the matched data

$N_{0j}$ : number of control subjects in stratum  $j$  in the cohort

$n_{0j}$ : number of control subjects in stratum  $j$  in the matched sample

# Simulation Study: Low Dimension Scenario

- Cohort sample size  $N = 20000$
- Matching variable:  $\mathbf{Z} = (Z_1, Z_2)$ , where  $Z_1 \sim N(0, 1)$ ,  $Z_2 \sim Ber(0.5)$
- Variables of interest:  $\mathbf{X} = (X_1, X_2, \dots, X_{20})^T \sim N(\boldsymbol{\mu}_{\mathbf{X}|\mathbf{Z}}, Cov(\mathbf{X}|\mathbf{Z}))$  where

$$\boldsymbol{\mu}_{\mathbf{X}|\mathbf{Z}} = \begin{cases} 0.1 + \alpha_Z * Z_1 + 0.1 * Z_2 & \text{for } X_1, \dots, X_{10} \\ 0 & \text{Others} \end{cases}$$

and  $Cov(\mathbf{X}|\mathbf{Z})$  has 20% non-zero elements randomly selected from  $U(-2, -1)$  or  $U(1, 2)$ .

- Case-control status  $Y$

$$\text{logit}(Y) = \beta_0 + \sum_{i=1}^{10} 0.5X_i + \beta_Z * Z_1 + 0.2 * Z_2$$

- Set  $\beta_Z = 1.2$  and  $\alpha_Z = \{0.3, 0.5, 0.7\}$  to investigate the effect of  $\alpha_Z$ .
- $\beta_0$ : chosen to make the proportion of cases in the cohort is 20%

# Simulation Study: Low Dimension Scenario

- Matching procedure
  - Cohort subjects are divided into 10 strata based on their propensity scores  $P(Y = 1|Z_1, Z_2)$ , which can be estimated by logistic regression;
  - Randomly select 500 cases from the cohort;
  - Match the cases with 500 controls according to the strata defined by propensity scores;
  - Bin the adjacency strata if their  $\frac{\# \text{ of controls}}{\# \text{ of cases}}$  are greater than 50
- We compare the following 4 covariance estimators to the true covariance matrix of  $\mathbf{X}$ :
  - Sample covariance of  $\mathbf{X}$  from random samples :  $\mathbf{S}_{random}$
  - Sample covariance of  $\mathbf{X}$  from case-control samples:  $\mathbf{S}_{CC}$
  - Sample covariance of  $\mathbf{X}$  from matched case-control samples:  $\mathbf{S}_{MCC}$
  - Proposed weighted covariance of  $\mathbf{X}$  from matched case-control samples:  $\mathbf{S}_{ipw}$

# Simulation Study: Low Dimension Scenario

- Compare the four estimates with the true covariance matrix of  $\mathbf{X}$

$$\begin{aligned} \text{Cov}(\mathbf{X}) &= E(\text{Cov}(\mathbf{X}|\mathbf{Z})) + \text{Cov}(E(\mathbf{X}|\mathbf{Z})) \\ &= \text{Cov}(\mathbf{X}|\mathbf{Z}) + \text{Cov}(\boldsymbol{\mu}_{\mathbf{X}|\mathbf{Z}}) \end{aligned}$$

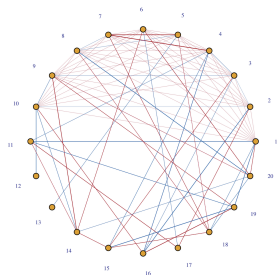


Figure: Visualization of true  $\text{Cov}(\mathbf{X})$

- Error metric:

$$\text{Average Relative Bias} = \frac{1}{p^2} \sum_i^p \sum_j^p \frac{\hat{\Sigma}_{ij} - \Sigma_{ij}}{1 + \text{abs}(\Sigma_{ij})} \quad (6)$$

- For low dimensional setting, comparison of precision matrix estimators can be directly obtained by inverting the covariance estimates.

# Simulation Result: Low Dimension Scenario

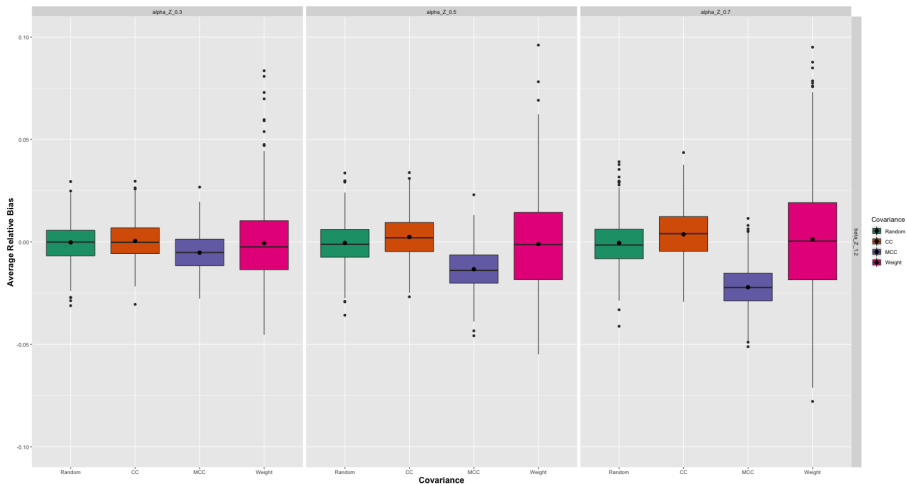


Figure: Covariance matrix estimation for low dimension scenario

# Effect of Binning

- Bin the adjacency strata if their  $\frac{\# \text{ of controls}}{\# \text{ of cases}}$  are greater than 50, 25 and 12.5

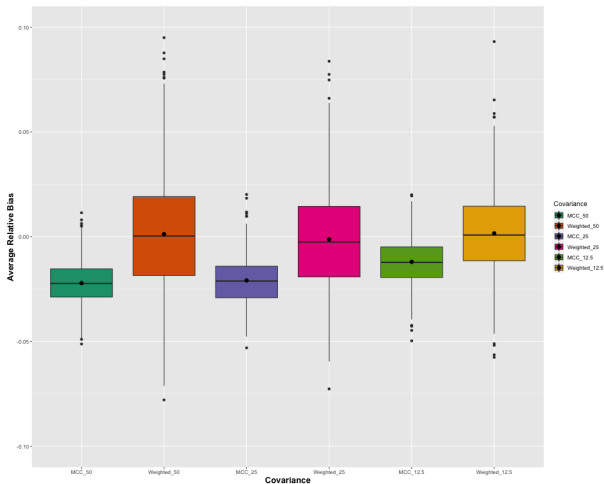
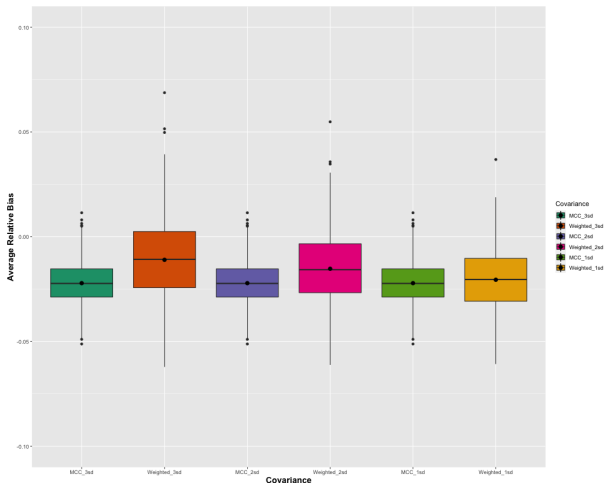


Figure: Effect of binning when  $\alpha_Z = 0.7$  and  $\langle \beta_Z \rangle = 1.2$

# Effect of Thresholding

- Thresholding can help to reduce the effect of extreme weights by truncating them at a maximum allowable weight.
- Choice of threshold:  $\text{mean} + 3/2/1 \cdot \text{SD}$





- Proposed weighted estimator in low dimensional setting and a likelihood-based estimation procedure in high dimensional setting
- Simulation for low dimensional setting shows simple covariance estimation from matched case control sample has significant bias, and estimation from our proposed method has almost no bias for all simulation scenarios.
- Variance reduction methods are needed for our proposed method.
- Simulation for high dimensional setting will be conducted.

- Bien, Jacob, and Robert J. Tibshirani. "Sparse estimation of a covariance matrix." *Biometrika* 98.4 (2011): 807-820.
- Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. "Sparse inverse covariance estimation with the graphical lasso." *Biostatistics* 9.3 (2008): 432-441.