

Debiased Lasso

Lecturer: Chao Gao

Scribe: Xinze Li

In previous lectures, we obtain estimators of precision matrix. Now it's time for us to do inference, which is a harder problem.

1 Debiased Lasso

First suppose that

$$y = X\beta + z, \quad z \sim \mathcal{N}(0, I_n) \quad (1)$$

Also suppose that $X = (X_1, \dots, X_n)^\top$, $X_i \in \mathbb{R}^p$ and that X_i are i.i.d. drawn from multivariate normal distribution $\mathcal{N}(0, \Omega^{-1})$, where $\Omega = \Sigma$ is the precision matrix. And so we could write the ordinary least square estimator if $n > p$

$$\hat{\beta}_{OLS} = (X^\top X)^{-1} X^\top y \sim \mathcal{N}(\beta, (X^\top X)^{-1}) \quad (2)$$

Also suppose that $\tilde{\beta}$ is the Lasso estimator

$$\tilde{\beta} = \arg \max_{\beta} (\|y - X\beta\|^2 - \lambda \|\beta\|_1) \quad (3)$$

We know that $\mathbf{E}y = X\beta$ and thus $\mathbf{E}X^\top y = X^\top X\beta$. And so

$$\tilde{\beta} - \hat{\beta}_{OLS} = (X^\top X)^{-1} (X^\top X\tilde{\beta} - X^\top y) \quad (4)$$

is approximately the bias. A tautology from this equation is

$$\begin{aligned} \hat{\beta} &= \tilde{\beta} + (X^\top X)^{-1} (X^\top y - X^\top X\tilde{\beta}) \\ &= \tilde{\beta} + \left(\frac{1}{n} X^\top X\right)^{-1} \cdot \frac{1}{n} (X^\top y - X^\top X\tilde{\beta}) \end{aligned} \quad (5)$$

So now our idea is to replace $(\frac{1}{n} X^\top X)^{-1}$ by a matrix close to the precision matrix. Let's call this approximation matrix M , and now in general we could write

$$\hat{\beta} = \tilde{\beta} + M \cdot \frac{1}{n} (X^\top y - X^\top X\tilde{\beta}) \quad (6)$$

This idea first appeared in the article [Zhang and Zhang, 2011]. It is also rediscovered in article [Javanmard and Montanari, 2015] and [van de Geer et al., 2014]. The results in the latter two papers are a subset of [Zhang and Zhang, 2011], which is hard to read.

1.1 Known Covariance

If Σ the covariance matrix is known, then a natural attempt is to plug in $M = \Omega = \Sigma^{-1}$. Suppose $\widehat{\Sigma} = \frac{1}{n}X^\top X$ is the sample covariance matrix, we have

$$\begin{aligned}
\widehat{\beta} &= \widetilde{\beta} + \frac{1}{n}\Omega X^\top (y - X\widetilde{\beta}) \\
&= \widetilde{\beta} + \frac{1}{n}\Omega X^\top (X\beta + z - X\widetilde{\beta}) \\
&= \widetilde{\beta} + \frac{1}{n}\Omega X^\top X\beta + \frac{1}{n}\Omega X^\top z - \frac{1}{n}\Omega X^\top X\widetilde{\beta} \\
&= \beta + \frac{1}{n}\Omega X^\top z + (I_p - \Omega\widehat{\Sigma}) (\widetilde{\beta} - \beta)
\end{aligned} \tag{7}$$

Rearranging the order gives

$$\sqrt{n} (\widehat{\beta} - \beta) = \frac{1}{\sqrt{n}}\Omega X^\top z + \sqrt{n} (I_p - \Omega\widehat{\Sigma}) (\widetilde{\beta} - \beta) \tag{8}$$

Since inference is applied on each component, let's analyze each component of this vector.

$$\sqrt{n} (\widehat{\beta}_j - \beta_j) = \frac{1}{\sqrt{n}}\Omega_j^\top X^\top z + \sqrt{n} (e_j - \widehat{\Sigma}\Omega_j)^\top \cdot (\widetilde{\beta} - \beta) \tag{9}$$

Clearly the first term of the RHS is Gaussian given X .

$$\frac{1}{\sqrt{n}}\Omega_j^\top X^\top z|X \sim \mathcal{N}(0, \Omega_j^\top \widehat{\Sigma}\Omega_j) \tag{10}$$

And for the second term, we use Holder's inequality

$$\left| \sqrt{n} (e_j - \widehat{\Sigma}\Omega_j)^\top (\widetilde{\beta} - \beta) \right| \leq \sqrt{n} \|e_j - \widehat{\Sigma}\Omega_j\|_\infty \cdot \|\widetilde{\beta} - \beta\|_1 \tag{11}$$

For the latter part of the RHS $\|\widetilde{\beta} - \beta\|_1$, we know from properties of lasso estimator that w.h.p.

$$\|\widetilde{\beta} - \beta\|_1 \lesssim s \sqrt{\frac{\log p}{n}} \tag{12}$$

And for the former part of RHS, we have

$$\begin{aligned}
&\|e_j - \widehat{\Sigma}\Omega_j\|_\infty \\
&= \max_{1 \leq l \leq p} \frac{1}{n} \left| \sum_{i=1}^n [(\Omega_j^\top X_i) (e_l^\top X_i) - \mathbf{E}(\Omega_j^\top X_i) (e_l^\top X_i)] \right| \\
&\leq \sqrt{\frac{\log p}{n}}
\end{aligned} \tag{13}$$

Details see HW5 Problem2 solution. So combining equation 13 and 12 and plugging the bound into 11 gives the following upper bound

$$\left| \sqrt{n} \left(e_j - \Omega_j^\top \hat{\Sigma} \right) \left(\tilde{\beta} - \beta \right) \right| \lesssim s \frac{\log p}{\sqrt{n}} \quad (14)$$

Thus, if $\frac{s \log p}{\sqrt{n}} \rightarrow 0$ and $\Omega_j^\top \hat{\Sigma} \Omega_j \gtrsim 1$ not degenerate w.h.p. (to be proved), then

$$\frac{\sqrt{n} \left(\hat{\beta}_j - \beta_j \right)}{\sqrt{\Omega_j^\top \hat{\Sigma} \Omega_j}} \rightsquigarrow \mathcal{N}(0, 1) \quad (15)$$

1.2 Unknown Covariance

Following the process as in eqs. (7) to (11), suppose $M = (m_1, m_2, \dots, m_p)^\top$ and that M depends only on X , we deduce that

$$\begin{aligned} \hat{\beta} &= \beta + \frac{1}{n} M X^\top z + \left(I_p - M \hat{\Sigma} \right) \left(\tilde{\beta} - \beta \right) \\ \sqrt{n} \left(\hat{\beta} - \beta \right) &= \frac{1}{\sqrt{n}} M X^\top z + \sqrt{n} \left(I_p - M \hat{\Sigma} \right) \left(\tilde{\beta} - \beta \right) \\ \sqrt{n} \left(\hat{\beta}_j - \beta_j \right) &= \frac{1}{\sqrt{n}} m_j^\top X^\top z + \sqrt{n} \left(e_j - \hat{\Sigma} m_j \right)^\top \left(\tilde{\beta} - \beta \right) \\ &\quad \frac{1}{\sqrt{n}} m_j^\top X^\top z | X \sim \mathcal{N} \left(0, m_j^\top \hat{\Sigma} m_j \right) \\ \left| \sqrt{n} \left(e_j - \hat{\Sigma} m_j \right)^\top \left(\tilde{\beta} - \beta \right) \right| &\leq \sqrt{n} \left\| e_j - \hat{\Sigma} m_j \right\|_\infty \cdot \left\| \tilde{\beta} - \beta \right\|_1 \end{aligned} \quad (16)$$

In order to obtain the normality, note that we should restrict $\left\| e_j - \hat{\Sigma} m_j \right\|_\infty$ to a scale of $\sqrt{\frac{\log p}{n}}$ such that $\left| \sqrt{n} \left(e_j - \hat{\Sigma} m_j \right)^\top \left(\tilde{\beta} - \beta \right) \right| \lesssim s \frac{\log p}{\sqrt{n}}$ as before. In the same time, we would like the confidence interval be as short as possible. In order to do this, we introduce the following convex program (j is fixed)

$$\begin{aligned} \min_{m \in \mathbb{R}^p} \quad & m^\top \hat{\Sigma} m \\ \text{s.t.} \quad & \left\| e_j - \hat{\Sigma} m_j \right\|_\infty \leq C \sqrt{\frac{\log p}{n}} \end{aligned} \quad (17)$$

Note that this problem is certainly feasible since we have just proved that Ω_j is a feasible point in 13. Let m_j be the solution to 17, then as we previously stated, if $\frac{s \log p}{\sqrt{n}} \rightarrow 0$ and $m_j^\top \hat{\Sigma} m_j \gtrsim 1$ not degenerate w.h.p. (to be proved), then

$$\frac{\sqrt{n} \left(\hat{\beta}_j - \beta_j \right)}{\sqrt{m_j^\top \hat{\Sigma} m_j}} \rightsquigarrow \mathcal{N}(0, 1) \quad (18)$$

Now let's prove that $m_j^\top \widehat{\Sigma} m_j$ is indeed bounded away from zero. Suppose $\mu = C\sqrt{\frac{\log p}{n}}$. Since $\|e_j - \widehat{\Sigma} m_j\|_\infty \leq \mu$, we have

$$\left| e_j^\top \widehat{\Sigma} m_j - 1 \right| \leq \mu \quad (19)$$

And so

$$\begin{aligned} m^\top \widehat{\Sigma} m &\geq m^\top \widehat{\Sigma} m + ce_j^\top \widehat{\Sigma} m - ce_j^\top \widehat{\Sigma} m \\ &\geq m^\top \widehat{\Sigma} m + c(1 - \mu) - ce_j^\top \widehat{\Sigma} m \end{aligned} \quad (20)$$

Taking minimum on the both sides gives

$$\begin{aligned} \min_{m \text{ feasible}} m^\top \widehat{\Sigma} m &\geq c(1 - \mu) + \min_{m \in \mathbb{R}^p} (m^\top \widehat{\Sigma} m - ce_j^\top \widehat{\Sigma} m) \\ &\geq c(1 - \mu) - \frac{c^2}{4} e_j^\top \widehat{\Sigma} e_j \\ &\geq \frac{(1 - \mu)^2}{\widehat{\Sigma}_{jj}} \end{aligned} \quad (21)$$

Now note that $\mu = C\sqrt{\frac{\log p}{n}}$ is small. $\widehat{\Sigma}_{jj} = \Sigma_{jj} - (\widehat{\Sigma}_{jj} - \Sigma_{jj})$, Σ_{jj} is bounded away from zero, and $\widehat{\Sigma}_{jj} - \Sigma_{jj}$ is of order $O\left(\frac{1}{\sqrt{n}}\right)$ goes to zero. We conclude with the following theorem

Theorem 1.1. *If $\frac{s \log p}{\sqrt{n}} \rightsquigarrow 0$, then*

$$\frac{\sqrt{n}(\widehat{\beta}_j - \beta_j)}{\sqrt{m_j^\top \widehat{\Sigma} m_j}} \rightsquigarrow \mathcal{N}(0, 1) \quad (22)$$

1.3 Estimate Noise Level

If $y = X\beta + \sigma z$, then we could use lasso residual (not debiased lasso) to estimate the noise level σ

$$\begin{aligned} \widehat{\sigma}^2 &= \frac{1}{n} \|y - X\widetilde{\beta}\|^2 \\ &= \frac{1}{n} \|X\beta - X\widetilde{\beta} + \sigma z\|^2 \\ &= \frac{\sigma^2}{n} \|z\|^2 + \frac{1}{n} \|X(\beta - \widetilde{\beta})\|^2 + \frac{2\sigma}{n} z^\top X(\beta - \widetilde{\beta}) \end{aligned} \quad (23)$$

And from bound deduced in last lecture, we deduce that

$$\sqrt{n}(\widehat{\sigma}^2 - \sigma^2) = \sigma^2 \frac{1}{\sqrt{n}} \sum_{i=1}^n (z_i^2 - 1) + O_p\left(\frac{s \log p}{\sqrt{n}}\right) \quad (24)$$

Note that $\sigma^2 \frac{1}{\sqrt{n}} \sum_{i=1}^n (z_i^2 - 1) \rightsquigarrow \mathcal{N}(0, 2\sigma^4)$. Thus we have

$$\frac{\sqrt{n}(\widehat{\sigma}^2 - \sigma^2)}{\sqrt{2}\sigma^2} \rightsquigarrow \mathcal{N}(0, 1) \quad (25)$$

References

- [Javanmard and Montanari, 2015] Javanmard, A. and Montanari, A. (2015). De-biasing the lasso: Optimal sample size for gaussian designs.
- [van de Geer et al., 2014] van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):11661202.
- [Zhang and Zhang, 2011] Zhang, C.-H. and Zhang, S. S. (2011). Confidence intervals for low-dimensional parameters in high-dimensional linear models.