

STAT 525 Lecture 16

October 26, 2017

1 Topics

1. Linear regression
2. Multiple testing
3. Prediction

Question: How does y vary as a function of X ?

Interested in $p(\mathbf{y} \mid \mathbf{X}, \theta)$

Assume $(y_i, \mathbf{x}_i)_{i=1}^n$ are exchangeable.

$n = \#$ of observations, $p = \#$ of predictors

2 Linear model

$$\begin{aligned}\mathbb{E}[y_i \mid \mathbf{x}, \boldsymbol{\beta}] &= \beta_1 x_1 + \cdots + \beta_p x_p \\ &= \mathbf{x}^T \boldsymbol{\beta}\end{aligned}$$

Ordinary linear model

$$\begin{aligned}\text{Var}[y_i \mid \mathbf{x}, \boldsymbol{\theta}] &= \sigma^2 \\ \boldsymbol{\theta} &= \begin{bmatrix} \boldsymbol{\beta} & \sigma \end{bmatrix}\end{aligned}$$

Normal linear model

$$\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\epsilon}_{n \times 1}$$

where

$$\begin{aligned}\epsilon_i &\stackrel{\text{iid}}{\sim} N(0, \sigma^2) \\ \boldsymbol{\epsilon} &\sim MVN(0, \sigma^2 \mathbb{I})\end{aligned}$$

Key points

1. Flexible choices for \mathbf{X}
2. Transformation of \mathbf{X} and / or \mathbf{y} for linearity and Gaussian of $\boldsymbol{\epsilon}$.
3. Generation for other variance model.

3 Bayesian Conditional Modeling

Data (\mathbf{y}, \mathbf{X})

Likelihood:

$$p(\mathbf{y}, \mathbf{X} \mid \boldsymbol{\theta}, \psi)$$

Priors

$$p(\boldsymbol{\theta}, \psi)$$

Posterior

$$p(\boldsymbol{\theta}, \psi \mid \mathbf{y}, \mathbf{X}) = p(\psi \mid \boldsymbol{\theta}, \mathbf{y}, \mathbf{X}) p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X})$$

Assume

$$p(\psi \mid \boldsymbol{\theta}, \mathbf{y}, \mathbf{X}) = p(\psi \mid \mathbf{X})$$

- Prior independence: $p(\boldsymbol{\theta}, \psi) = p(\boldsymbol{\theta}) p(\psi)$
- ψ conditional independent of \mathbf{y} given \mathbf{X}

$$p(\boldsymbol{\theta}, \psi \mid \mathbf{y}, \mathbf{X}) = p(\psi \mid \mathbf{X}) p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X})$$

Posterior inference for $\boldsymbol{\theta}$ only need the density $p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X})$ or

$$p(\boldsymbol{\theta} \mid \mathbf{y}, \mathbf{X}) \propto p(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{X}) p(\boldsymbol{\theta} \mid \mathbf{X})$$

Commentary on the intercept:

1. if $\mathbf{x}_{i1} = 1$ for all i , usually assume $p(\beta_1) \propto 1$
2. Usually center and scale \mathbf{y} and columns of \mathbf{X}

$$y_i = \frac{y_i - \text{mean}(\mathbf{y})}{\text{sd}(\mathbf{y})}$$

4 General Case

The prior on $\boldsymbol{\beta}$ is

$$\boldsymbol{\beta} \sim \text{MVN}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta)$$

Then the likelihood is

$$\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}, \sigma^2 \sim \text{MVN}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbb{I})$$

The posterior is

$$\begin{aligned} p(\boldsymbol{\beta} \mid \mathbf{y}, \mathbf{X}, \sigma^2) &\propto p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}, \sigma^2) p(\boldsymbol{\beta}) \\ &\propto |2\pi\sigma^2\mathbb{I}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\sigma^2\mathbb{I})^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \times \\ &\quad |2\pi\boldsymbol{\Sigma}_\beta|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)^T (\boldsymbol{\Sigma}_\beta)^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)\right) \\ &\propto \exp\left(-\frac{1}{2}\sigma^{-2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + (\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)^T \boldsymbol{\Sigma}_\beta^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_\beta)\right) \\ &\propto \exp\left(-\frac{1}{2}\left(\boldsymbol{\beta}^T \mathbf{Q}_\beta \boldsymbol{\beta} - 2\boldsymbol{\beta}^T \mathbf{l}_\beta\right)\right) \end{aligned}$$

where

$$\begin{aligned}\mathbf{Q}_\beta &= \sigma^{-2} \mathbf{X}^T \mathbf{X} + \Sigma_\beta^{-1} \\ \mathbf{l}_\beta &= \sigma^{-2} \mathbf{X} \mathbf{y} + \Sigma_\beta^{-1} \boldsymbol{\mu}_\beta\end{aligned}$$

Therefore, the posterior for β is

$$[\beta \mid \mathbf{y}, \mathbf{X}, \sigma^2] \sim MVN(\mathbf{Q}_\beta^{-1} \mathbf{l}_\beta, \mathbf{Q}_\beta^{-1})$$

5 g-prior (Zeller 1986)

$$[\beta \mid \sigma^2, \mathbf{X}, g] \sim MVN\left(\boldsymbol{\mu}_\beta, \underbrace{g\sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1}}_{\Sigma_\beta}\right)$$

We plug the new $\boldsymbol{\mu}_\beta$ and Σ_β into the section 4, then

$$\begin{aligned}\hat{\mathbf{Q}}_\beta &= \sigma^{-2} \mathbf{X}^T \mathbf{X} + g^{-1} \sigma^{-2} \mathbf{X}^T \mathbf{X} \\ &= \sigma^{-2} \mathbf{X}^T \mathbf{X} [1 + g^{-1}] \\ \hat{\mathbf{l}}_\beta &= \sigma^{-2} \mathbf{X}^T \mathbf{y} + g^{-1} \sigma^{-2} \mathbf{X}^T \mathbf{X} \boldsymbol{\mu}_\beta \\ &= \sigma^{-2} \mathbf{X}^T (\mathbf{y} + g^{-1} \mathbf{X} \boldsymbol{\mu}_\beta)\end{aligned}$$

Therefore we have

$$\begin{aligned}\mathbf{Q}_\beta^{-1} &= \frac{g}{g+1} \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1} \\ \mathbf{Q}_\beta^{-1} \mathbf{l}_\beta &= \frac{g}{g+1} \cancel{\sigma^2} [\mathbf{X}^T \mathbf{X}]^{-1} \cancel{\sigma^2} \mathbf{X}^T (\mathbf{y} + g^{-1} \mathbf{X} \boldsymbol{\mu}_\beta) \\ &= \frac{g}{g+1} \underbrace{[\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T \mathbf{y}}_{\hat{\boldsymbol{\beta}}_{\text{ols}}} + \frac{1}{g+1} \boldsymbol{\mu}_\beta\end{aligned}$$

The posterior expectation for β is

$$\begin{aligned}\mathbb{E}[\beta \mid \mathbf{X}, \mathbf{y}, \sigma^2, g] &= (1 - \kappa) \hat{\boldsymbol{\beta}} + \kappa \boldsymbol{\mu}_\beta \\ \kappa &= \frac{1}{g+1}\end{aligned}$$

The posterior for β is

$$[\beta \mid \mathbf{X}, \mathbf{y}, \sigma^2, g] \sim MVN\left((1 - \kappa) \beta + \kappa \boldsymbol{\mu}_\beta, \frac{g}{g+1} \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1}\right)$$

Invariance of g prior

Let $\mathbf{X}^* = \mathbf{X} \mathbf{H}$ for some $p \times p$ \mathbf{H}

Posterior for $\beta \mid \mathbf{X}, \mathbf{y}$ shall be the same as $\mathbf{H} \beta^* \mid \mathbf{y}, \mathbf{X}^*$

stratify when $\boldsymbol{\mu}_\beta \rightarrow 0$ and $\Sigma_\beta = \kappa (\mathbf{X}^T \mathbf{X})^{-1}$ for $\kappa > 0$.

Flat prior

A common non-informative choice is

$$p(\boldsymbol{\beta}, \sigma^2) = p(\boldsymbol{\beta} | \sigma^2) p(\sigma^2) \propto \frac{1}{\sigma^2}$$

Special case of g prior (in the posterior) when $g \rightarrow \infty$ and $p(\sigma^2) \propto \frac{1}{\sigma^2}$

when $g \rightarrow \infty$ then

$$\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}, \sigma^2 \sim \left(\hat{\boldsymbol{\beta}}_{\text{ols}}, \sigma^2 [\mathbf{X}^T \mathbf{X}]^{-1} \right)$$

6 Penalized regression

6.1 Ridge regression

Suppose $\beta_j | \sigma_\beta^2 \stackrel{\text{iid}}{\sim} N(0, \sigma_\beta^2)$ for $j = 1, \dots, p$

The log-posterior (up to an additive constant) is

$$\begin{aligned} \log p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}, \sigma^2, \sigma_\beta^2) &\propto \log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \sigma^2, \sigma_\beta^2) + \log p(\boldsymbol{\beta}) \\ &\propto -\frac{1}{2\sigma^2} \sum_{i=1}^n \|y_i - \mathbf{x}_i^T \boldsymbol{\beta}\|^2 + \sum_{j=1}^p \log p(\beta_j) \\ &\propto -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 - \frac{1}{2\sigma_\beta^2} \sum_{j=1}^p \beta_j^2 \\ &\propto -\frac{1}{2\sigma^2} \underbrace{\left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}}_{\text{Ridge Regression}} \end{aligned}$$

where

$$\lambda = \frac{\sigma^2}{\sigma_\beta^2}$$

The posterior expectation for $\boldsymbol{\beta}$ is

$$\mathbb{E}[\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}, \sigma^2, \lambda] = \underbrace{\left[\mathbf{X}^T \mathbf{X} + \lambda \mathbb{I} \right]^{-1} \mathbf{X}^T \mathbf{y}}_{\text{Ridge Estimator}}$$

λ is the tuning parameter for freq. For Bayesian $\lambda = \frac{\sigma^2}{\sigma_\beta^2}$. We just need to estimate σ_β^2 and put a prior on σ_β^2 .

Key points

1. Posterior mode (and here posterior mean) solves the optimization problem

$$\begin{aligned} \hat{\boldsymbol{\beta}}_R &= \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \sum_{j=1}^p \beta_j^2 \\ &= \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \text{ s.t. } \|\boldsymbol{\beta}\|^2 \leq s \end{aligned}$$

- (a) for $\lambda > 0$, allowed to have $p \geq n$
- (b) Ridge regression is useful for $p > n$ and / or correlated predictions.

2. Lasso regression

$$\hat{\beta}_L = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- (a) often have $\hat{\beta}_{L,j} = 0$ for many j provides sparse solution.
- (b) overshrinkage large β_j