

VE414 Lecture 21

Jing Liu

UM-SJTU Joint Institute

November 19, 2019

- So far we have largely used data to only estimate **un**observable,

$$Y$$

- **Linear regression model** is a way to study the relationship of an observable

$$Y$$

in terms of a set of other observable variables

$$X_1, X_2, \dots, X_k$$

specifically, it is a type of smoothly changing model for

$$f_{Y|\{X_1, X_2, \dots\}}$$

in which the conditional expectation $\mathbb{E}[Y | \{X_1, \dots, X_k\}]$ has a form that is linear in a set of **un**observable β_i , which are often known as the parameters

$$\mathbb{E}[Y | \{X_1, \dots, X_k\}] = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k = \mathbf{x}^T \boldsymbol{\beta}$$

- In addition to being linear,

$$\mathbb{E}[Y \mid \{X_1, \dots, X_k\}] = \mathbf{x}^T \boldsymbol{\beta}$$

the variability around the mean, i.e. the error,

$$Y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i$$

is often assumed to be normal

$$\varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \text{Normal}(0, \sigma^2)$$

- Under the above specification, we have the following density function

$$\begin{aligned} f_{\{Y_1, Y_2, \dots, Y_n\} | \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \boldsymbol{\beta}, \sigma^2\}} &= \prod_{i=1}^n f_{Y_i | \{\mathbf{x}_i, \boldsymbol{\beta}, \sigma^2\}} \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2\right) \end{aligned}$$

- We can put the density function into a vector form,

$$\begin{aligned} f_{\{Y_1, Y_2, \dots, Y_n\} | \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \boldsymbol{\beta}, \sigma^2\}} &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \text{RSS}\right) \end{aligned}$$

where residual sum of squares is given by

$$\text{RSS} = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

- Thus our model in vector form is $\mathbf{Y} | \{\mathbf{X}, \boldsymbol{\beta}, \sigma^2\} \sim \text{Normal}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$.

Q: What would frequentists do next?

- Frequentists would maximise the likelihood by treating the density function as a function of the unknown parameters, which is equivalent to minimise

$$\text{RSS}(\mathbf{b}) = (\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b})$$

- Recall to minimise a function,

$$\text{RSS}(\mathbf{b}) = (\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{b}^T \mathbf{X}^T \mathbf{y} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b}$$

we set the gradient to zero,

$$\nabla \text{RSS} = 0 - 2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{b}$$

Setting this to zero, we have

$$\hat{\beta}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Hence, the fitted value is given by

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{P} \mathbf{y}$$

and the residual can be found using

$$\hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{P}) \mathbf{y}$$

- With more linear algebra, we have

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta + \mathbf{e}) = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{e}$$

which means it is unbiased as expected,

$$\mathbb{E} [\hat{\beta} \mid \mathbf{X}] = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E} [\varepsilon \mid \mathbf{X}] = \beta$$

- The variance is given by

$$\begin{aligned} \text{Var} [\hat{\beta} \mid \mathbf{X}] &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Var} [\varepsilon \mid \mathbf{X}] \left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right)^T \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned}$$

- With the normal assumption, we see

$$\hat{\beta} \sim \text{Normal} \left(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- To estimate σ^2 , frequentists typically use the following

$$\hat{\sigma}^2 = \frac{1}{n - k - 1} \hat{\mathbf{e}}^T \hat{\mathbf{e}} \quad \text{where} \quad \hat{\mathbf{e}} = (\mathbf{I} - \mathbf{P}) \mathbf{y}$$

which is unbiased as well as being consistent.

- It can be shown the residual

$$\hat{\mathbf{e}} = (\mathbf{I} - \mathbf{P}) \mathbf{y} = (\mathbf{I} - \mathbf{P}) (\mathbf{X}\boldsymbol{\beta} + \mathbf{e})$$

is an unbiased and consistent estimator of the error \mathbf{e} , and the variance is

$$\begin{aligned} \text{Var} [\hat{\mathbf{e}} \mid \mathbf{X}] &= \text{Var} [(\mathbf{I} - \mathbf{P}) (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) \mid \mathbf{X}] \\ &= (\mathbf{I} - \mathbf{P}) \text{Var} [\boldsymbol{\varepsilon} \mid \mathbf{X}] (\mathbf{I} - \mathbf{P})^T \\ &= (\mathbf{I} - \mathbf{P}) \sigma^2 \mathbf{I} (\mathbf{I} - \mathbf{P})^T = \sigma^2 (\mathbf{I} - \mathbf{P}) \end{aligned}$$

- Thus with the normal assumption, we have

$$\hat{\mathbf{e}} \sim \text{Normal} (\mathbf{0}, \sigma^2 (\mathbf{I} - \mathbf{P}))$$

Q: How would Bayesian approach the same problem?

$$f_{\mathbf{Y}|\{\mathbf{X}, \beta, \sigma^2\}} = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \text{RSS}(\beta)\right)$$

where

$$\text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) = \mathbf{y}^T \mathbf{y} - 2\beta^T \mathbf{X}^T \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{X} \beta$$

- Using a normal prior for $\beta \sim \text{Normal}(\beta_0, \Sigma_0)$, we have

$$\begin{aligned} f_{\beta}(\beta) &= \frac{1}{\sqrt{(2\pi)^k \det(\Sigma_0)}} \exp\left(-\frac{1}{2} (\beta - \beta_0)^T \Sigma_0^{-1} (\beta - \beta_0)\right) \\ &\propto \exp\left(-\frac{1}{2} \beta^T \Sigma_0^{-1} \beta + \beta^T \Sigma_0^{-1} \beta_0\right) \end{aligned}$$

Q: What is the conditional posterior of β ?

$$f_{\beta|\{\sigma^2, \mathbf{Y}, \mathbf{X}\}}$$

- Using the **precision parameter** in the likelihood instead of σ^2 , that is

$$\tau = \frac{1}{\sigma^2}$$

and using a gamma prior for $\tau \sim \text{Gamma}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$,

$$\begin{aligned} f_\tau &= \frac{(\nu_0 \sigma_0^2 / 2)^{\nu_0 / 2}}{\Gamma(\nu_0 / 2)} \tau^{\nu_0 / 2 - 1} \exp\left(-\frac{\nu_0 \sigma_0^2}{2} \tau\right) \\ &\propto \tau^{\nu_0 / 2 - 1} \exp\left(-\frac{\nu_0 \sigma_0^2}{2} \tau\right) \end{aligned}$$

Q: What is the conditional posterior of τ ?

$$f_{\sigma^2 | \{\beta, \mathbf{Y}, \mathbf{X}\}}$$

Q: How can we sample from the Joint posterior?

$$f_{\{\beta, \sigma^2\} | \{\mathbf{Y}, \mathbf{X}\}}$$

- Since both conditionals are readily available, and both are pretty standard,

$$\boldsymbol{\beta} \mid \{\sigma^2, \mathbf{Y}, \mathbf{X}\} \sim \text{Normal}(\mathbf{m}, \mathbf{V})$$

$$\sigma^2 \mid \{\boldsymbol{\beta}, \mathbf{Y}, \mathbf{X}\} \sim \text{Inverse-Gamma}(\alpha, \beta)$$

where

$$\mathbf{m} = (\boldsymbol{\Sigma}_0^{-1} + \mathbf{X}^T \mathbf{X} / \sigma^2)^{-1} (\boldsymbol{\Sigma}_0^{-1} \boldsymbol{\beta}_0 + \mathbf{X}^T \mathbf{y} / \sigma^2)$$

$$\mathbf{V} = (\boldsymbol{\Sigma}_0^{-1} + \mathbf{X}^T \mathbf{X} / \sigma^2)^{-1}$$

$$\alpha = \frac{\nu_0 + n}{2}; \quad \beta = \frac{\nu_0 \sigma_0^2 + \text{RSS}(\boldsymbol{\beta})}{2}$$

and positivity is satisfied, using Gibbs sampling is then straightforward

$$(\boldsymbol{\beta}, \sigma^2) \in \mathbb{R}^k \times (0, \infty)$$

- If other priors are used, we will have a different joint and a different sampling scheme, but the essences of Bayesian linear regression are the same.