

# ORF 525: Statistical Foundations of Data Science

Jianqing Fan — Frederick L. Moore'18 Professor of Finance

Problem Set #1

Spring 2021

Due Friday, February 12, 2021.

1. Consider the linear model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where  $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{W})$  with known positive definite matrix  $\mathbf{W}$ , and  $\mathbf{X}$  is of full rank.

(a) Show that the general least-squares estimator, which minimizes  $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{W}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ , is the best linear unbiased estimator.

(b) Give explicitly the least-squares estimators for  $\boldsymbol{\beta}$  and  $\sigma^2$ .

(c) If  $\mathbf{W} = (1 - \rho)\mathbf{I} + \rho\mathbf{1}\mathbf{1}^T$  is a equi-correlation matrix, what is the least-square estimate  $\hat{\boldsymbol{\beta}}$ ?

*Hint: Use Sherman-Morrison formula if necessary.*

## Answers:

(a) Let's denote by  $\hat{\boldsymbol{\beta}}$  the general least-squares estimator of  $\boldsymbol{\beta}$ :

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{W}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

Multiplying  $\mathbf{W}^{-\frac{1}{2}}$  on both sides of the linear model, and by denoting

$$\tilde{\mathbf{X}} = \mathbf{W}^{-\frac{1}{2}} \mathbf{X}, \quad \tilde{\mathbf{y}} = \mathbf{W}^{-\frac{1}{2}} \mathbf{y}, \quad \text{and} \quad \tilde{\boldsymbol{\varepsilon}} = \mathbf{W}^{-\frac{1}{2}} \boldsymbol{\varepsilon},$$

we can have  $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\boldsymbol{\beta} + \tilde{\boldsymbol{\varepsilon}}$  where  $\tilde{\boldsymbol{\varepsilon}} \sim N(0, \sigma^2 \mathbf{I})$  and  $\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2$ . The result then follows from the course.

(b) We clearly have

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2 = \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}} = \left( \mathbf{X}^T \mathbf{W}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{W}^{-1} \mathbf{y}$$

and

$$\hat{\sigma}^2 = \frac{\|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\hat{\boldsymbol{\beta}}\|_2^2}{n - p} = \frac{\|\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}\|_2^2}{n - p} = \frac{\|\tilde{\mathbf{y}} - \mathbf{P}\tilde{\mathbf{y}}\|_2^2}{n - p},$$

where  $\mathbf{P} = \tilde{\mathbf{X}} \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T = \mathbf{W}^{-1/2} \mathbf{X} \left( \mathbf{X}^T \mathbf{W}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{W}^{-1/2}$  is a projector of rank  $p$ . Then,  $\hat{\sigma}^2$  can be rewritten as

$$\hat{\sigma}^2 = \frac{\tilde{\mathbf{y}}^T (\mathbf{I} - \mathbf{P})^2 \tilde{\mathbf{y}}}{n - p} = \frac{\tilde{\mathbf{y}}^T (\mathbf{I} - \mathbf{P}) \tilde{\mathbf{y}}}{n - p} = \frac{\mathbf{y}^T \mathbf{W}^{-1/2} (\mathbf{I} - \mathbf{P}) \mathbf{W}^{-1/2} \mathbf{y}}{n - p}.$$

(c) Suppose that  $\mathbf{W} \in \mathbb{R}^{n \times n}$ . If  $\mathbf{W}$  is the equi-correlation matrix with correlation  $\rho$ , namely

$$\mathbf{W} = (1 - \rho)\mathbf{I} + \rho\mathbf{1}\mathbf{1}^T$$

where  $\mathbf{1}$  stands for a vector in  $\mathbb{R}^n$  with all elements equal to 1. Then, since  $\rho \neq 1$  (as  $\mathbf{W}$  is positive definite), by Sherman-Morrison formula, we have

$$\mathbf{W}^{-1} = a\mathbf{I} + b\mathbf{1}\mathbf{1}^T, \quad a = \frac{1}{1 - \rho}, \quad b = \frac{1}{1 - \rho} \frac{-\rho}{1 + (n - 1)\rho}.$$

We denote by

$$\mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, \text{ and } \alpha = \mathbb{1}^T \mathbf{P} \mathbb{1},$$

then

$$\begin{aligned} & (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \\ &= (\mathbf{X}^T (a\mathbf{I} + b\mathbb{1}\mathbb{1}^T) \mathbf{X})^{-1} \mathbf{X}^T (a\mathbf{I} + b\mathbb{1}\mathbb{1}^T) \\ &= (a\mathbf{X}^T \mathbf{X} + b(\mathbf{X}^T \mathbb{1})(\mathbf{X}^T \mathbb{1})^T)^{-1} (a\mathbf{X}^T + b\mathbf{X}^T \mathbb{1}\mathbb{1}^T) \\ &= \left( \frac{1}{a} (\mathbf{X}^T \mathbf{X})^{-1} - \frac{b}{a^2} \frac{[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}] \cdot [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}]^T}{1 + \frac{b}{a} \mathbb{1}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}} \right) \cdot (a\mathbf{X}^T + b\mathbf{X}^T \mathbb{1}\mathbb{1}^T) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \frac{b}{a} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}] \cdot \mathbb{1}^T - \frac{\frac{b}{a}}{1 + \frac{b}{a} \alpha} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}] \cdot [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}]^T \cdot \mathbf{X}^T \\ &\quad - \frac{b}{a} \cdot \frac{\frac{b}{a}}{1 + \frac{b}{a} \alpha} \cdot [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}] \cdot \alpha \cdot \mathbb{1}^T \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \frac{b}{a} \left( 1 - \frac{\frac{b}{a} \alpha}{1 + \frac{b}{a} \alpha} \right) \cdot [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}] \cdot \mathbb{1}^T - \frac{\frac{b}{a}}{1 + \frac{b}{a} \alpha} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{1}] \cdot \mathbb{1}^T \cdot \mathbf{P} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{I} - \gamma \mathbb{1}\mathbb{1}^T (\mathbf{I} - \mathbf{P})) \end{aligned}$$

where  $\gamma = \frac{-b}{a + b\mathbb{1}^T \mathbf{P} \mathbb{1}}$ . Hence, the least-squares estimator for  $\boldsymbol{\beta}$  with an equi-correlation variance matrix is

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{I} - \gamma \mathbb{1}\mathbb{1}^T (\mathbf{I} - \mathbf{P})) \mathbf{y}.$$

2. Consider the multiple regression model  $Y_i = \mathbf{X}_i^T \boldsymbol{\beta} + \varepsilon_i$ , where  $\varepsilon_i \sim N(0, \sigma^2)$  for  $i = 1, \dots, n$ .

- (a) Show that the maximum likelihood estimator is equivalent to the least-squares estimator, which finds  $\boldsymbol{\beta}$  to minimize

$$\sum_{i=1}^n (Y_i - \mathbf{X}_i^T \boldsymbol{\beta})^2$$

and

$$\hat{\sigma} = \sqrt{\frac{\text{RSS}}{n}},$$

where  $\text{RSS} = \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \hat{\boldsymbol{\beta}})^2$  and  $\hat{\boldsymbol{\beta}}$  is the least-squares solution.

- (b) Show that  $\text{RSS} \sim \sigma^2 \chi_{n-p}^2$ , where  $p$  is the rank of  $\mathbf{X}$  (full rank for simplicity).  
(c) Prove that  $1 - \alpha$  CI for  $\beta_j$  is  $\hat{\beta}_j \pm t_{n-p}(1 - \alpha/2) \sqrt{v_j \text{RSS}/(n-p)}$ , where  $v_j$  is the  $j^{\text{th}}$  diagonal element of  $(\mathbf{X}^T \mathbf{X})^{-1}$ .  
(d) Dropping the normality assumption, if  $\{\mathbf{X}_i\}$  are i.i.d. from a population with  $E\mathbf{X}\mathbf{X}^T = \boldsymbol{\Sigma}$  and independent of  $\{\varepsilon_i\}_{i=1}^n$ , which are i.i.d. from a population with  $E\varepsilon = 0$  and  $\text{var}(\varepsilon) = \sigma^2$ , show that

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N(0, \sigma^2 \boldsymbol{\Sigma}^{-1}).$$

**Answers:**

(a) A straight forward calculation yields

$$\log(P(\mathbf{Y}|\beta, \sigma, \mathbf{X})) \propto -n \cdot \log(\sigma) - \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|^2,$$

Thus, the MLE becomes (maximization and minimization are unique here):

$$\hat{\beta} = \arg \max_{\beta} P(\mathbf{Y}|\beta, \sigma, \mathbf{X}) = \arg \min_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|^2 = \hat{\beta}^{OLS}.$$

which is identical to the OLS estimator. And by taking the derivative w.r.t.  $\sigma$ , we can also have

$$\hat{\sigma} = \sqrt{\frac{\|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2}{n}}$$

(b) Let us denote  $\mathbf{P} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}$ , then  $\mathbf{P}$  is an idempotent matrix so that  $\mathbf{I} - \mathbf{P}$  and

$$Tr(\mathbf{I} - \mathbf{P}) = n - Tr((\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}) = n - tr(\mathbf{I}_p) = n - p.$$

Since an idempotent matrix can only have eigenvalues 0 or 1, the spectral decomposition gives us

$$\mathbf{I} - \mathbf{P} = \mathbf{U}\mathbf{D}\mathbf{U}^T, \quad \mathbf{D} = \begin{bmatrix} \mathbf{I}_{n-p} & 0 \\ 0 & 0 \end{bmatrix}$$

for some orthogonal matrix  $\mathbf{U} \in \mathbb{R}^{n \times n}$  (i.e.  $\mathbf{U}^T\mathbf{U} = \mathbf{I}_n$ ). We then have

$$\|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 = \mathbf{Y}^T(\mathbf{I} - \mathbf{P})\mathbf{Y} = (\mathbf{U}^T\mathbf{Y})^T\mathbf{D}(\mathbf{U}^T\mathbf{Y}) = \sum_{i=1}^{n-p} (\mathbf{U}^T\mathbf{Y})_i^2$$

We know that  $\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I}_n)$ , and thus  $\mathbf{U}^T\mathbf{Y} \sim N(\mathbf{U}^T\mathbf{X}\beta, \sigma^2\mathbf{I}_n)$  which implies that the components  $\{(\mathbf{U}^T\mathbf{Y})_i\}_{i=1}^n$  are independent with each other. Hence, the quantity  $\|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2$  is a sum of  $(n-p)$  iid normal distributed random variable  $N(0, \sigma^2)$ . By definition

$$RSS = n\hat{\sigma}^2 = \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2 \sim \sigma^2\chi_{n-p}^2.$$

(c) It is easy to show that  $\hat{\beta} \sim N(\beta, \sigma^2(\mathbf{X}^T\mathbf{X})^{-1})$ . So, for any vector  $u \in \mathbb{R}^d$  where  $d$  is the dimension of vector  $X_i$ , we must have

$$\frac{u^T(\hat{\beta} - \beta)}{\sigma \cdot \sqrt{u^T(\mathbf{X}^T\mathbf{X})^{-1}u}} \sim N(0, 1)$$

and thus by independence of  $\hat{\beta}$  and  $\hat{\sigma}^2$ , we derive

$$\frac{u^T(\hat{\beta} - \beta)}{\sigma \cdot \sqrt{u^T(\mathbf{X}^T\mathbf{X})^{-1}u}} \cdot \left( \frac{n\hat{\sigma}^2}{\sigma^2} \frac{1}{n-p} \right)^{-\frac{1}{2}} \sim t_{n-p}.$$

If we choose  $u$  that has all zeros but 1 on the  $j$ th coordinate, then we get

$$\beta_j \sim \hat{\beta}_j - \sqrt{\frac{n\hat{\sigma}^2}{n-p}} \cdot (\mathbf{X}^T\mathbf{X})_{jj}^{-1} \cdot t_{n-p}.$$

Then the result for CI follows.

(d) We can write

$$\sqrt{n}(\hat{\beta} - \beta) = \left(\frac{1}{n}\mathbf{X}^T\mathbf{X}\right)^{-1} \cdot \frac{1}{\sqrt{n}}\mathbf{X}^T\boldsymbol{\varepsilon}.$$

We also observe by LLN that

$$\frac{1}{n}\mathbf{X}^T\mathbf{X} = \frac{1}{n}\sum_{i=1}^n \mathbf{X}_i\mathbf{X}_i^T \longrightarrow \mathbb{E}[\mathbf{X}_1\mathbf{X}_1^T] = \boldsymbol{\Sigma} \quad \text{in probability,}$$

Moreover, by CLT, we have a convergence in distribution for

$$\frac{1}{\sqrt{n}}\mathbf{X}^T\boldsymbol{\varepsilon} = \frac{1}{\sqrt{n}}\sum_{i=1}^n \mathbf{X}_i\epsilon_i \longrightarrow N(0, \text{Var}(\mathbf{X}_1\epsilon_1)) = N(0, \mathbb{E}[\epsilon_1^2\mathbf{X}_1\mathbf{X}_1^T]) = N(0, \sigma^2\boldsymbol{\Sigma})$$

where the last equality comes from  $\mathbf{X}_1$  and  $\epsilon_1$  are independent and  $\epsilon_1 \in \mathbb{R}$ . By Slutsky's theorem and the continuous mapping theorem, we can drive

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} \boldsymbol{\Sigma}^{-1}N(0, \sigma^2\boldsymbol{\Sigma}) = N(0, \sigma^2\boldsymbol{\Sigma}^{-1}).$$

3. Let us consider the 129 macroeconomic time series again as described in the lecture notes

<http://orfe.princeton.edu/%7Ejqfan/fan/classes/245/chap11.pdf>

Let  $Y_t = \Delta \log(\text{PCE}_t)$  be the changes in the personal consumption expenditure. Let us take

$$\begin{aligned} X_{t,1} &= \text{Unrate}_{t-1}, \quad X_{t,2} = \Delta \log(\text{IndPro}_{t-1}), \quad X_{t,3} = \Delta \log(\text{M2Real}_{t-1}), \\ X_{t,4} &= \Delta \log(\text{CPI}_{t-1}), \quad X_{t,5} = \Delta \log(\text{SPY}_{t-1}), \quad X_{t,6} = \text{HouSta}_{t-1}, \quad X_{t,7} = \text{FedFund}_{t-1} \end{aligned}$$

Let us again take the last 10 years data as testing set and remaining as training set. Conduct a similar analysis as those in the lecture notes. Answer in particular the following questions.

- What are  $\hat{\sigma}^2$ , adjusted  $R^2$  and insignificant variables?
- Now perform the stepwise deletion, eliminating one least significant variable at a time (by iteratively looking at the smallest  $|t|$ -statistic) until all variables are statistically significant. Let us call this model as model  $\widehat{\mathcal{M}}$ . (For the stepwise deletion by AIC, the function `step` can do the job automatically)
- Using model  $\widehat{\mathcal{M}}$ , what are root mean-square prediction error and mean absolute deviation prediction error for the test sample?

*Note: It is possible that the resulting model is not as good as the vanilla least squares; It's also possible that the resulting out of sample  $R^2$  is negative, since predicting the difference of log PCE is not an easy task in general. If we take the predictors at time  $t$ , then it is the association study and the prediction is somewhat better.*

- Compute the standardized residuals based on model  $\widehat{\mathcal{M}}$ . Present the time series plot of the residuals, fitted values versus the standardized residuals, and QQ plot for the standardized residuals.

**Answers:** cf R codes.

4. Zillow is an online real estate database company that was founded in 2006. The most important task for Zillow is to predict the house price. However, their accuracy has been criticized a lot. According to Fortune, "Zillow has Zestimated the value of 57 percent of U.S. housing stock, but only 65 percent of that could be considered 'accurate' – by its definition, within 10 percent of the actual selling price. And even that accuracy isn't equally distributed". Therefore, Zillow needs your help to build a housing pricing model to improve their accuracy. Download and read the data (traing data: 15129 cases, testing data: 6484 cases)

```
train.data <- read.csv('train.data.csv', header=TRUE)
test.data <- read.csv('test.data.csv', header=TRUE)
train.data$zipcode <- as.factor(train.data$zipcode)
test.data$zipcode <- as.factor(test.data$zipcode)
```

where the last two lines make sure that zip code is treated as factor. Letting  $\mathcal{T}$  as a test set, define out-of-sample  $R^2$  as of a prediction method  $\{\hat{y}_i^{pred}\}$  as

$$R^2 = 1 - \frac{\sum_{i \in \mathcal{T}} (y_i - \hat{y}_i^{pred})^2}{\sum_{i \in \mathcal{T}} (y_i - \bar{y}^{pred})^2},$$

where  $\bar{y}^{pred} = \text{ave}(\{y_i\}_{i \in \mathcal{T}_0})$  and  $\mathcal{T}_0$  is the training set.

- Calculate out-of-sample  $R^2$  using variables "bathrooms", "bedrooms", "sqft\_living", and "sqft\_lot".
- Calculate out-of-sample  $R^2$  using the 4 variables above along with interaction terms.
- Compare the result with the nonparametric model using Gaussian kernel (standardize predictors first) with  $\gamma = 0.1^2/2$  and  $\lambda = 0.1$ .
- Add the factor zipcode to (b) and compute out-of-sample  $R^2$ .
- Add the following additional variables to (d):  $X_{12} = I(\text{view} == 0)$ ,  $X_{13} = L^2$ ,  $X_{13+i} = (L - \tau_i)_+^2$ ,  $i = 1, \dots, 9$ , where  $\tau_i$  is  $10 * i^{th}$  percentile and  $L$  is the size of living area ("sqft\_living"). Compute out-of-sample  $R^2$ .

**Answers:** cf R codes.

5. Prove the representer theorem in the lecture note (Theorem 1.4). You are allowed to consult the book, but not allow to have verbatim copy. You need to write the solution of your own with at least some changes of notation.

**Answer:** cf Theorem 2.6 in the book.