

Statistics 206

Homework 1 Solution

Due : October 2, 2019, In Class

1. Review Appendix A3 and the matrix notes on canvas under “Files”.
2. \mathbf{Z} is an n -dimensional random vector with expectation $\mathbf{E}(\mathbf{Z})$ and variance-covariance matrix:

$$\mathbf{Var}(\mathbf{Z}) = \mathbf{Cov}(\mathbf{Z}, \mathbf{Z}) = \Sigma.$$

A is an $s \times n$ nonrandom matrix and B is a $t \times n$ nonrandom matrix. Show the following:

(a) $\mathbf{E}(A\mathbf{Z}) = A\mathbf{E}(\mathbf{Z})$.

Proof.

$$(A\mathbf{Z})_j = \sum_k a_{jk} \mathbf{Z}_k, j = 1, \dots, s$$

$$(E(A\mathbf{Z}))_j = E((A\mathbf{Z})_j) = E\left(\sum_k a_{jk} \mathbf{Z}_k\right) = \sum_k a_{jk} E(\mathbf{Z}_k) = (AE(\mathbf{Z}))_j, j = 1, \dots, s$$

□

(b) $\mathbf{Cov}(A\mathbf{Z}, B\mathbf{Z}) = A\Sigma B^T$. So in particular, $\mathbf{Var}(A\mathbf{Z}) = A\Sigma A^T$.

Proof. Define

$$W = A\mathbf{Z}, U = B\mathbf{Z}, C = \mathbf{Cov}(W, U), D = A\Sigma B^T$$

$$\begin{aligned} C_{ij} &= \mathbf{Cov}(W_i, U_j) \\ &= \mathbf{Cov}\left(\sum_k a_{ik} \mathbf{Z}_k, \sum_l b_{jl} \mathbf{Z}_l\right) \\ &= \sum_k \sum_l a_{ik} b_{jl} \mathbf{Cov}(\mathbf{Z}_k, \mathbf{Z}_l) \\ &= \sum_k \sum_l a_{ik} b_{jl} \Sigma_{kl} \\ &= D_{ij}, i = 1, \dots, s, j = 1, \dots, t \end{aligned}$$

□

3. Derive the following.

(a) $\sum_{i=1}^n (X_i - \bar{X}) = 0, \quad \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n (X_i - \bar{X})X_i = \sum_{i=1}^n X_i^2 - n(\bar{X})^2.$

Proof.

$$\sum_{i=1}^n (X_i - \bar{X}) = \sum_{i=1}^n X_i - \sum_{i=1}^n \bar{X} = n\bar{X} - n\bar{X} = 0.$$

$$\begin{aligned} \sum_{i=1}^n (X_i - \bar{X})^2 &= \sum_{i=1}^n (X_i^2 - 2X_i\bar{X} + (\bar{X})^2) \\ &= \sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + \sum_{i=1}^n (\bar{X})^2 \\ &= \sum_{i=1}^n X_i^2 - 2n(\bar{X})^2 + n(\bar{X})^2 \\ &= \sum_{i=1}^n X_i^2 - n(\bar{X})^2. \end{aligned}$$

□

$$(b) \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \sum_{i=1}^n (X_i - \bar{X})Y_i = \sum_{i=1}^n X_i Y_i - n\bar{X} \bar{Y}.$$

Proof.

$$\begin{aligned} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) &= \sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i \bar{Y} - \sum_{i=1}^n \bar{X} Y_i + \sum_{i=1}^n \bar{X} \bar{Y} \\ &= \sum_{i=1}^n X_i Y_i - n\bar{X} \bar{Y} - n\bar{X} \bar{Y} + n\bar{X} \bar{Y} \\ &= \sum_{i=1}^n X_i Y_i - n\bar{X} \bar{Y}. \end{aligned}$$

□

4. Least-squares principle.

(a) State the least-squares principle.

For a given line: $y = b_0 + b_1 x$, the *sum of squared vertical deviations* of the observations $\{(X_i, Y_i)\}_{i=1}^n$ from the corresponding points on the line is:

$$Q(b_0, b_1) = \sum_{i=1}^n (Y_i - (b_0 + b_1 X_i))^2.$$

The *least squares (LS) principle* is to fit the observed data by minimizing the sum of squared vertical deviations.

(b) Derive the LS estimators for simple linear regression model.

Proof. From the lecture notes, the LS estimators can be derived by finding b_0 and b_1 which satisfy the normal equations.

$$nb_0 + b_1 \sum_{i=1}^n X_i = \sum_{i=1}^n Y_i \dots (1)$$

$$b_0 \sum_{i=1}^n X_i + b_1 \sum_{i=1}^n X_i^2 = \sum_{i=1}^n X_i Y_i \dots (2)$$

From equation (1), $b_0 = \bar{Y} - b_1 \bar{X}$.

Using this in equation (2) we have

$$\begin{aligned} b_0 n \bar{X} + b_1 \sum_{i=1}^n X_i^2 &= \sum_{i=1}^n X_i Y_i \\ \Rightarrow n \bar{X} \bar{Y} + b_1 [\sum_{i=1}^n X_i^2 - n \bar{X}^2] &= \sum_{i=1}^n X_i Y_i \\ \Rightarrow b_1 &= \frac{\sum_{i=1}^n X_i Y_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2} \end{aligned}$$

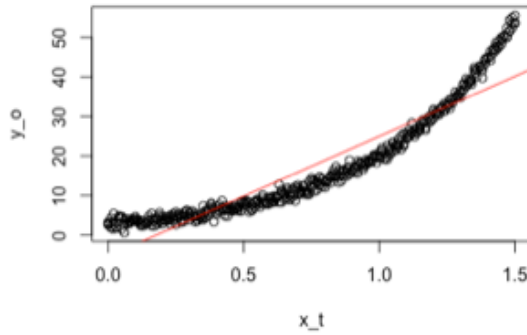
Now from 2(a) and (b), $\hat{\beta}_1 = \frac{\sum_{i=1}^n X_i Y_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$.

From equation (1), $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$. □

(c) Assume the observations follow:

$$Y_i = \exp(a + bX_i) + \epsilon_i, \quad i = 1, \dots, n,$$

where $a, b \in \mathbb{R}$ are unknown parameters and ϵ_i s are uncorrelated random variables with $E(\epsilon_i) = 0, Var(\epsilon_i) = \sigma^2$. Describe how to estimate the regression function (equivalently, a, b) by least-squares principle. (Notes: You only need to provide a description. This is an example of a nonlinear regression model.)



Plot implies linear regression model is not accurate for non-linear dataset.

$$Q(a, b) = \sum_{i=1}^n (Y_i - \exp(a + bX_i))^2.$$

The *least squares (LS) principle* is to fit the observed data by minimizing the sum of squared vertical deviations.

$$(\hat{a}, \hat{b}) = \operatorname{argmin}_{a,b} \sum_{i=1}^n (Y_i - \exp(a + bX_i))^2.$$

5. Tell true or false (with a brief explanation) of the following statements with regard to simple linear regression.

(a) The least squares line always passes the center of the data (\bar{X}, \bar{Y}) .

ANS. True. since $y = \bar{Y} + \hat{\beta}_1(x - \bar{X})$.

(b) The true regression line fits the data the best.

ANS. False. LS line fits the data better.

(c) If $\bar{X} = 0$, $\bar{Y} = 0$, then $\hat{\beta}_0 = 0$ no matter what is $\hat{\beta}_1$.

ANS. True since $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1\bar{X}$

(d) Given the sample size, the larger the range of X_i s, the smaller the standard errors of $\hat{\beta}_0, \hat{\beta}_1$ tend to be.

ANS. True. since $s\{\hat{\beta}_0\}$ and $s\{\hat{\beta}_1\}$ has $\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}$ in the denominator.

(e) If the correlation coefficient r_{XY} between X_i s and Y_i s is such that $|r_{XY}| < 1$, then we will observe regression effect.

ANS. TRUE. When $|r_{XY}| < 1$, then one standard deviation change in X would amount to less than one standard deviation change in Y (on average) and such a phenomena is called the regression effect (a.k.a. regression towards the mean).

6. Properties of the residuals under simple linear regression model. Recall that

$$\begin{aligned} e_i &= Y_i - \hat{Y}_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i), \quad i = 1, \dots, n. \\ &= (Y_i - \bar{Y}) - \hat{\beta}_1 (X_i - \bar{X}). \end{aligned}$$

Show that

(a) $\sum_{i=1}^n e_i = 0$.

Proof.

$$\begin{aligned}
\sum_{i=1}^n e_i &= \sum_{i=1}^n (Y_i - \bar{Y}) - \sum_{i=1}^n \hat{\beta}_1 (X_i - \bar{X}) \\
&= \left(\sum_{i=1}^n Y_i - \sum_{i=1}^n \bar{Y} \right) - \hat{\beta}_1 \left(\sum_{i=1}^n X_i - \sum_{i=1}^n \bar{X} \right) \\
&= (n\bar{Y} - n\bar{Y}) - \hat{\beta}_1 (n\bar{X} - n\bar{X}) \\
&= 0.
\end{aligned}$$

□

(b) $\sum_{i=1}^n X_i e_i = 0.$

Proof.

$$\begin{aligned}
\sum_{i=1}^n X_i e_i &= \sum_{i=1}^n X_i ((Y_i - \bar{Y}) - \hat{\beta}_1 (X_i - \bar{X})) \\
&= \left(\sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i \bar{Y} \right) - \hat{\beta}_1 \left(\sum_{i=1}^n X_i^2 - \sum_{i=1}^n X_i \bar{X} \right) \\
&= \left(\sum_{i=1}^n X_i Y_i - n\bar{X} \bar{Y} \right) - \hat{\beta}_1 \left(\sum_{i=1}^n X_i^2 - n(\bar{X})^2 \right) \\
&= \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) - \hat{\beta}_1 \sum_{i=1}^n (X_i - \bar{X})^2 \\
&= \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) - \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \sum_{i=1}^n (X_i - \bar{X})^2 \\
&= 0
\end{aligned}$$

□

(c) $\sum_{i=1}^n \hat{Y}_i e_i = 0.$

Proof. By parts (a) and (b), and $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i.$

□

7. Under the simple linear regression model, show the following.

- (a) The LS estimator $\hat{\beta}_0$ is an unbiased estimator of β_0 and derive the formula for its variance.

Proof.

$$\begin{aligned}
E(\hat{\beta}_0) &= E(\bar{Y}) - E(\hat{\beta}_1 \bar{X}) = \frac{1}{n} \sum (\beta_0 + \beta_1 X_i) - \bar{X} E(\hat{\beta}_1) \\
&= \beta_0 + \beta_1 \bar{X} - \bar{X} \beta_1 \quad \text{use the fact that } E(\hat{\beta}_1) = \beta_1 \\
&= \beta_0.
\end{aligned}$$

$$\begin{aligned}
\text{Var}(\hat{\beta}_1) &= \text{Var} \left(\frac{\sum (X_i - \bar{X}) Y_i}{\sum (X_i - \bar{X})^2} \right) \quad \text{use 3(b)} \\
&= \frac{\sigma^2}{\sum (X_i - \bar{X})^2}, \quad \text{as } Y_i \text{s are uncorrelated}
\end{aligned}$$

$$\begin{aligned}
\text{Var}(\hat{\beta}_0) &= \text{Var}(\bar{Y} - \hat{\beta}_1 \bar{X}) \quad \text{use 4(b)} \\
&= \text{Var}(\bar{Y}) + \bar{X}^2 \text{Var}(\hat{\beta}_1) \quad \text{use 7(b)} \\
&= \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum (X_i - \bar{X})^2} \right)
\end{aligned}$$

□

(b) \bar{Y} and $\hat{\beta}_1$ are uncorrelated. (Hint: Write them as linear combinations of Y_i s.)

Proof.

$$\begin{aligned}
\bar{Y} &= \frac{1}{n} \sum_{i=1}^n Y_i, \quad \hat{\beta}_1 = \sum \frac{X_i - \bar{X}}{\sum_{j=1}^n (X_j - \bar{X})^2} Y_i \\
\text{Cov}(\bar{Y}, \hat{\beta}_1) &= \sum_{i=1}^n \frac{1}{n} \frac{X_i - \bar{X}}{\sum_{j=1}^n (X_j - \bar{X})^2} \text{Var}(Y_i) \quad (\text{by } Y_i \text{s are uncorrelated}) \\
&= \frac{1}{n} \frac{\sum_{i=1}^n X_i - n\bar{X}}{\sum_{j=1}^n (X_j - \bar{X})^2} \sigma^2 = 0 \quad \text{since } \sum_{i=1}^n X_i - n\bar{X} = 0.
\end{aligned}$$

□

8. Under the simple linear regression model:

(a) Show that

$$SSE = \sum_{i=1}^n (Y_i - \bar{Y})^2 - \hat{\beta}_1^2 \sum_{i=1}^n (X_i - \bar{X})^2.$$

Proof.

$$\begin{aligned}
\hat{Y}_i &= \hat{\beta}_0 + \hat{\beta}_1 X_i = \bar{Y} + \hat{\beta}_1 (X_i - \bar{X}) \\
SSE &= \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - \bar{Y} - \hat{\beta}_1 (X_i - \bar{X}))^2 \\
&= \sum_{i=1}^n (Y_i - \bar{Y})^2 + \hat{\beta}_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 - 2\hat{\beta}_1 \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \\
&= \sum_{i=1}^n (Y_i - \bar{Y})^2 + \hat{\beta}_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 - 2\hat{\beta}_1 \cdot \hat{\beta}_1 \sum_{i=1}^n (X_i - \bar{X})^2 \\
&\quad [\text{since } \hat{\beta}_1 \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})] \\
&= \sum_{i=1}^n (Y_i - \bar{Y})^2 - \hat{\beta}_1^2 \sum_{i=1}^n (X_i - \bar{X})^2
\end{aligned}$$

□

(b) Derive $E(\hat{\beta}_1^2)$.

Proof.

$$\begin{aligned}
E(\hat{\beta}_1^2) &= \text{Var}(\hat{\beta}_1) + E(\hat{\beta}_1)^2 \\
&= \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} + \beta_1^2 \\
\text{since } \text{Var}(\hat{\beta}_1) &= \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \\
\text{and } E(\hat{\beta}_1) &= \beta_1 \text{ from lecture notes.}
\end{aligned}$$

□

(c) Show that

$$E((Y_i - \bar{Y})^2) = \beta_1^2 (X_i - \bar{X})^2 + E((\epsilon_i - \bar{\epsilon})^2),$$

where $\bar{\epsilon} = \frac{1}{n} \sum_{i=1}^n \epsilon_i$.

Proof.

$$\begin{aligned}
E((Y_i - \bar{Y})^2) &= E((\beta_0 + \beta_1 X_i + \epsilon_i - \beta_0 - \beta_1 \bar{X} - \bar{\epsilon}))^2 \\
&= E(\beta_1 (X_i - \bar{X}) + \epsilon_i - \bar{\epsilon})^2 \\
&= \beta_1^2 (X_i - \bar{X})^2 + E((\epsilon_i - \bar{\epsilon})^2) + 2\beta_1 (X_i - \bar{X}) E(\epsilon_i - \bar{\epsilon}) \\
&= \beta_1^2 (X_i - \bar{X})^2 + E((\epsilon_i - \bar{\epsilon})^2) \text{ as } E(\epsilon_i - \bar{\epsilon}) = E(\epsilon_i) - E(\bar{\epsilon}) = 0
\end{aligned}$$

□

(d) Show that $E(SSE) = (n-2)\sigma^2$. (Hint: What is $E(\sum_{i=1}^n (\epsilon_i - \bar{\epsilon})^2)$?)

Proof.

$$\begin{aligned}
\text{Var}(\epsilon_i - \bar{\epsilon}) &= \text{Var}(\epsilon_i) + \text{Var}(\bar{\epsilon}) - 2\text{Cov}(\epsilon_i, \bar{\epsilon}) \\
&= \sigma^2 + \frac{\sigma^2}{n} - 2\frac{\sigma^2}{n} = \frac{(n-1)\sigma^2}{n} \\
\text{as } \text{Cov}(\epsilon_i, \bar{\epsilon}) &= \text{Cov}(\epsilon_i, \frac{\sum_{i=1}^n \epsilon_i}{n}) = \frac{\text{Var}(\epsilon_i)}{n} = \frac{\sigma^2}{n} \\
\therefore E\left(\sum_{i=1}^n (\epsilon_i - \bar{\epsilon})^2\right) &= \sum_{i=1}^n E(\epsilon_i - \bar{\epsilon})^2 = \sum_{i=1}^n \text{Var}(\epsilon_i - \bar{\epsilon}) = \frac{n(n-1)\sigma^2}{n} = (n-1)\sigma^2
\end{aligned}$$

From 8(a),

$$\begin{aligned}
E(SSE) &= \sum_{i=1}^n E(Y_i - \bar{Y})^2 - E(\hat{\beta}_1^2) \sum_{i=1}^n (X_i - \bar{X})^2, \text{ then from 8(b) and 8(c)} \\
&= \beta_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n E((\epsilon_i - \bar{\epsilon})^2) - \left(\frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} + \beta_1^2 \right) \sum_{i=1}^n (X_i - \bar{X})^2 \\
&= \beta_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 + (n-1)\sigma^2 - \sigma^2 - \beta_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 \\
&= (n-2)\sigma^2
\end{aligned}$$

□