

# Stats 401 Lab 7

Naomi Giertych

10/19/2018

# Announcements

- ▶ Midterm Monday in class

# Outline

- ▶ Probability Model
- ▶ Covariance Review
- ▶ Exam Practice and Questions

# Probability Model

(Looking ahead to HW 6)

- ▶ Recall: Probability Model is an assignment of probabilities to possible outcomes. We don't observe these probabilities, but we observe a random sample of them, e.g. a response variable and  $p$  predictor variables.
- ▶ This means there exists a probability model of the form  $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ . Note:  $\boldsymbol{\beta}$  and  $\boldsymbol{\epsilon}$  are unknown and  $\mathbb{X}$  is the **observed** explanatory matrix.
- ▶ However, when we observe data, we only have the sample version of the linear model  $\mathbf{y} = \mathbb{X}\mathbf{b} + \mathbf{e}$

# Examples of Writing these models in different forms

## Subscript form:

$$Y_i = \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_p x_{i,p} + \epsilon_i, \text{ where}$$

- ▶  $Y_i$  is the  $i$ th observation of the response variable  $Y$
- ▶  $\beta_1, \beta_2, \dots, \beta_p$  are the true coefficients of the explanatory variables
- ▶  $x_{i,1}, x_{i,2}, \dots, x_{i,p}$  are the observed  $p$  predictor variables for observation  $i$ ; note:  $x_{i,p}$  is set to 1 so our model includes an intercept
- ▶  $\epsilon_i$  is the true error of the  $i$ th observation

# Examples of Writing these models in different forms

Full matrix form:

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

where

- ▶  $Y_i$  is the  $i$ th observation of the response variable  $Y$
- ▶  $\beta_1, \beta_2, \dots, \beta_p$  are the true coefficients of the explanatory variables
- ▶  $x_{i,1}, x_{i,2}, \dots, x_{i,p}$  are the observed  $p$  predictor variables for observation  $i$ ; note:  $x_{i,p}$  is set to 1 so our model includes an intercept
- ▶  $\epsilon_i$  is the true error of the  $i$ th observation

# Examples of Writing these models in different forms

## Matrix Notation:

$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where

- ▶  $\mathbf{Y}$  is the  $(n \times 1)$  vector of the response variable
- ▶  $\boldsymbol{\beta}$  is the  $(p \times 1)$  vector of the true coefficients of the explanatory variables
- ▶  $\mathbb{X}$  is the  $(n \times p)$  design matrix of the  $p$  explanatory variables
- ▶  $\boldsymbol{\epsilon}$  is the  $(n \times 1)$  vector of the true errors

$$\hat{\beta}$$

Recall:  $\mathbf{b} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X} \mathbf{y}$

Similar to the probability model for  $\mathbf{y} = \mathbb{X} \mathbf{b} + \mathbf{e}$ ,  $\mathbf{b}$  has a probability model.

$$\hat{\beta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X} \mathbf{Y}$$

Key idea: We have a population model for the response variables ( $\mathbf{Y}$ ). We can estimate the true coefficients using  $\hat{\beta}$ . However, since we only get a random draw of  $\mathbf{Y}$ , known as  $\mathbf{y}$ , we can only use an estimate of  $\hat{\beta}$  known as  $\mathbf{b}$ .



# Multivariate Random Variables

Recall in last lab, we discussed bivariate random variables and the bivariate normal distributions, and we extended these concepts to multivariate random variables.

- ▶ For example, we might have the random vector  $\mathbf{X} = (X_1, X_2, \dots, X_p)$
- ▶ Another natural example of this random vector variable would be the vector of the  $p$  predictor variables of my probability model.

# Multivariate Random Variables

- ▶ Summary statistics for a multivariate random variable include the expected value vector and the variance-covariance matrix
- ▶ The expected value vector  $E(\mathbf{X}) = (E(X_1), \dots, E(X_p))$  tells us the means for each component of  $\mathbf{X}$
- ▶ The variance-covariance matrix gives the variances for each component along the diagonal and the pairwise covariances in the other entries:

$$\mathbb{V} = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{Cov}(X_2, X_p) \\ \vdots & \vdots & & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \dots & \text{Var}(X_p) \end{bmatrix}$$

## Lab Activity (Part 1)

### Fitting a probability model.

The director of the CDC wants to assess how well rates of hospital-acquired infections (`Infection.risk`) can be predicted using properties of a hospital. She expects to use the average length of stay (`Length.of.stay`), the estimated percentage of patients acquiring an infection in hospital (`Infection.risk`), the average number of cultures for each patient without signs or symptoms of hospital-acquired infection, times 100 (`Culture`), and the number of X-ray procedures divided by number of patients without signs or symptoms of pneumonia, times 100 (`X.ray`).

- Write the probability model in subscript form, in full matrix form, and using matrix notation.

## Lab activity (Part 2)

Fitting a sample linear model.

She collects a dataset for 113 hospitals with the variables mentioned above and fits the linear model below. Explain how this linear model is different from the one in Part 1. Write this model in subscript form, using the numbers below.

```
senic <- read.table(  
  "https://ionides.github.io/401w18/hw/hw09/senic.txt",  
  header = T)  
senic <- senic[,-1]  
  
lm1 <- lm(Infection.risk ~ Length.of.stay +  
          Infection.risk + Culture + X.ray,  
          data = senic)  
summary(lm1)
```

##

## Lab Activity (Part 3)

Let  $\mathbf{Y} = (Y_1, Y_2, Y_3)$  be a random vector with mean vector  $(2, 4, 6)$  and variance/covariance matrix

$$\mathbb{V} = \begin{bmatrix} 6 & 2 & 3 \\ 2 & 4 & 5 \\ 3 & 5 & 1 \end{bmatrix}$$

$W_1$  is the sum of  $Y_1$ ,  $Y_2$ , and  $Y_3$ .

$W_2$  is the sum of  $Y_1$ , ( $Y_2$  multiplied by 2), and ( $Y_3$  multiplied by -1).

$W_3$  is the sum of  $Y_1$ , and ( $Y_2$  and  $Y_3$  both multiplied by -1).

1. State the above in matrix notation.
2. Find the expectation of the random vector  $\mathbf{W}$ .
3. Find the variance/covariance matrix of  $\mathbf{W}$ .

# Exam Questions

- ▶ What questions do you have about concepts or from the practice midterm or hw?

## Lab Ticket

Using the same random variables as Lab Activity Part 3.

- ▶ Find the probability that  $Y_1$  is bigger than  $Y_2$
- ▶ Find the probability that  $W_1$  is bigger than  $W_2$ .