



JOHNS HOPKINS  
BLOOMBERG SCHOOL  
*of* PUBLIC HEALTH

## Lecture 7

---

Vector representation of MLR continued,  
assessing the impact of Gaussian residuals assumption

# MLR model expressed in vector notation

We have for each subject  $i, i=1, \dots, n$

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_p X_{pi} + \varepsilon_i$$

$$\begin{array}{l} i \\ 1 \\ 2 \\ \vdots \\ n \end{array} \quad \begin{array}{l} Y_1 = 1 \cdot \beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \dots + \beta_p X_{p1} + \varepsilon_1 \\ Y_2 = 1 \cdot \beta_0 + \beta_1 X_{12} + \beta_2 X_{22} + \dots + \beta_p X_{p2} + \varepsilon_2 \\ \vdots \\ Y_n = 1 \cdot \beta_0 + \beta_1 X_{1n} + \beta_2 X_{2n} + \dots + \beta_p X_{pn} + \varepsilon_p \end{array}$$

$$\underset{\sim}{Y}_{n \times 1} = \underset{\sim}{1}_{n \times 1} \beta_0 + \beta_1 \underset{\sim}{X}_{1 \times 1} + \beta_2 \underset{\sim}{X}_{2 \times 1} + \dots + \beta_p \underset{\sim}{X}_{p \times 1} + \underset{\sim}{\varepsilon}_{n \times 1}$$

$$\underset{\sim}{Y}_{n \times 1} = \underset{\sim}{X}_{n \times (p+1)} \underset{\text{design matrix}}{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}} + \underset{\sim}{\varepsilon}_{n \times 1}$$

# MLR model expressed in vector notation

$$\Rightarrow \underset{\sim}{Y}_{n \times 1} = \underset{n \times (p+1)}{X} \underset{(p+1) \times 1}{\beta} + \underset{\sim}{\varepsilon}_{n \times 1}$$

Say  $n=3, p=2 \Rightarrow Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \varepsilon_i$

$$\underset{\sim}{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \end{pmatrix} \quad \underset{\sim}{X} = \begin{bmatrix} 1 & X_{11} & X_{21} \\ 1 & X_{12} & X_{22} \\ 1 & X_{13} & X_{23} \end{bmatrix} \quad \underset{\sim}{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} \quad \underset{\sim}{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \end{pmatrix}$$

$$\underset{3 \times 3}{X} \underset{3 \times 1}{\beta} = \left( \begin{array}{c} \\ \\ \end{array} \right)_{3 \times 1} = \underset{\sim}{\mu}_{3 \times 1}$$

Defines the mean of  $Y_i$  for each  $i$

## MLR model expressed in vector notation

What about distribution of  $\underline{\varepsilon}$ ? and  $\underline{y}$ ?

In general, we can define the multivariate normal distribution as:  $\underline{y} \sim \text{MVN}(\underline{\mu}, V)$   
where  $\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$   $\underline{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}$   $V = \begin{bmatrix} v_{11} & v_{12} & \dots & v_{1n} \\ v_{21} & v_{22} & \dots & v_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ v_{n1} & v_{n2} & \dots & v_{nn} \end{bmatrix}$   
 $v_{ii} = \text{Var}(y_i)$   
 $v_{ij} = v_{ji} = \text{Cov}(y_i, y_j)$

If  $\varepsilon_i \sim N(0, \sigma^2)$ , independent

$$\underline{\varepsilon} \sim \text{MVN}(\underline{0}, \sigma^2 \underline{I}) \quad \underline{I} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} \text{ identity matrix}$$

$$\underline{y} \sim \text{MVN}(\underline{X}\beta, \sigma^2 \underline{I})$$



## MLE or LS solution expressed in vector notation

$$\text{MLR model: } \underset{\sim}{y} = \underset{\sim}{X} \underset{\sim}{\beta} + \underset{\sim}{\varepsilon}, \quad \underset{\sim}{\varepsilon} \sim \text{MVN}(\underset{\sim}{0}, \sigma^2 \underset{\sim}{I})$$

MLE or least squares: Going to drop the " $\sim$ "  
Choose  $\hat{\beta}$  and  $\hat{\sigma}^2$  to minimize  $\sum_{i=1}^n (y_i - x_i \beta)^2$   
 $x_i = i^{\text{th row of } X}$

$$\sum_{i=1}^n (y_i - x_i \beta)^2 = (\underset{1 \times n}{Y} - \underset{n \times 1}{X} \beta)' (\underset{n \times 1}{Y} - \underset{n \times 1}{X} \beta)$$

$$\begin{aligned} U_{\beta}(\beta) &= \frac{d}{d\beta} (\underset{1 \times n}{Y} - \underset{n \times 1}{X} \beta)' (\underset{n \times 1}{Y} - \underset{n \times 1}{X} \beta) \\ &= X' (\underset{n \times 1}{Y} - \underset{n \times 1}{X} \beta) \end{aligned}$$

Set to 0, solve for  $\hat{\beta} = (X'X)^{-1} X'Y$



## Predicted values and residuals in vector notation

$$\hat{\beta} = (X'X)^{-1}X'y$$

$$\hat{Y} = X\hat{\beta} =$$

$$\hat{R} = y - \hat{Y} =$$



## Distribution of $\hat{\beta}$

Note that if  $\underline{Y} \sim \text{MVN}(\underline{\mu}, V)$ , then  
 $A\underline{Y} \sim \text{MVN}(A\underline{\mu}, AVA')$

$$\hat{\beta} = \underbrace{(X'X)^{-1}X'}_A Y$$

$$E(\hat{\beta}) = E(AY) =$$

$$\text{Var}(\hat{\beta}) = \text{Var}(AY) =$$



## Distribution of $\hat{Y}$

$$\hat{Y} = X\hat{\beta} = X(X'X)^{-1}X'Y = HY$$





## Properties of the Hat matrix

1) hat matrix is symmetric

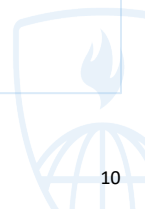
$$H' =$$

2) hat matrix is idempotent :  $H \cdot H = H$



## Distribution of $\hat{R}$

$$\hat{R} = Y - \hat{Y} = Y - HY =$$



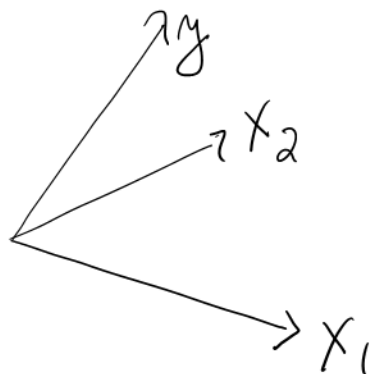
## Relationship between $\hat{Y}$ and $\hat{R}$

$$\text{Cov}(\hat{Y}, \hat{R}) = E \left[ H Y \{ (I-H) Y \}' \right]$$

$HY \quad (I-H)Y$

# Geometry of least squares

Consider  $y \sim_{n \times 1}$ ,  $X_1 \sim_{n \times 1}$ ,  $X_2 \sim_{n \times 1}$



$H$  projects  $y$  onto the plane spanned by  $X_1, X_2$   
 $\Rightarrow \hat{y} = Hy = X\hat{\beta}$

1) minimize the distance between  $y$  and  $\hat{y} = X\hat{\beta}$

2) Shortest distance is the one that has a right angle between the predicted value and residual

3) residual is orthogonal to the plane spanned by  $X$

4) Score equations:  
 $X'(y - X\hat{\beta}) = 0$

# Simulation study

- ▶ We derived the distribution of the estimated regression coefficients assuming the residuals were Gaussian.
- ▶ Does approximate normality of the estimated regression coefficients hold even when the residuals are non-Gaussian?



## Next time....

- ▶ Deriving the distribution of linear combinations of regression coefficients
- ▶ Deriving the distribution of non-linear combinations of regression coefficients using the Delta method
- ▶ LAB: You will generate the distribution of combinations of regression coefficients using bootstrap!

