# 36-617: Applied Linear Models

Multivariate Regression Brian Junker 132E Baker Hall brian@stat.cmu.edu

# Reading, HW, Quiz

#### Reading

- □ This week: Sheather Ch 5 (Skip Ch 4 for now)
- For next week: Sheather Ch 6
- HW 03 due tonight 1159pm
- HW 04 will be on Canvas later today Due next Monday

#### Quiz on Ch 5 in class – today!

## Outline

- Matrix Form of Multiple Regression Model
- Multivariate Normal Distribution
- ML/LS Estimates
- Two Interpretations
- Distributional Properties
- SS Decompositions and F Statistics
- Some Comments

# Matrix Form of Multiple Regression Model

• Let 
$$X_i = (x_{i0}, ..., x_{ip})$$
 and  $\beta = (\beta_0, ..., \beta_p)^T$ ; then  
 $y_i = \beta_0 X_{i0} + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i$   
 $= X_i \beta + \epsilon_i$ 

• If we also stack  $Y = (y_1, ..., y_n)^T$ ,  $X = (X_1^T, ..., X_n^T)^T$ , and  $\epsilon = (\epsilon_1, ..., \epsilon_n)^T$ , we can write

$$Y_{n \times 1} = X_{n \times k} \beta_{k \times 1} + \epsilon_{n \times 1} \quad (k = p + 1)$$

# Matrix Form of Multiple Regression Model

$$Y = X\beta + \epsilon$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{10} & x_{11} & \cdots & x_{1p} \\ x_{20} & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n0} & x_{n1} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Usually  $x_{i0} \equiv 1$ , so we get

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

#### **Multivariate Normal Distribution**

In the model

$$y_i = X_i\beta + \epsilon_i, i = 1, \dots, n$$

it is usual to assume  $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ 

Recall V ~ N(0,1) iff  $f(y) = \frac{1}{\sqrt{2\pi}}e^{-y^2/2}$ 

• Y ~ N(
$$\mu$$
, $\sigma$ <sup>2</sup>) iff  $\frac{y-\mu}{\sigma} \sim N(0,1)$ 

$$f(y) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

# Multivariate Normal Distribution

• 
$$\mathbf{Y} = (\mathbf{Y}_1, ..., \mathbf{Y}_n)^{\mathsf{T}} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}) = N \left( \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \cdots & 1 \end{bmatrix} \right)$$

$$f(y_1, \dots, y_n) = \prod_{i=1}^n f(y_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-y_i^2/2}$$
  
• Y ~ N( $\mu$ , $\Sigma$ ) iff  
 $\Sigma = 1/2 (X_{i-1}) = N(0, I)$  with some using formula to the same using formula

$$\Sigma^{-1/2}(Y-\mu) \sim N(0,I) \qquad \Big|_{\rm fo}^{\dots}$$

...and some ugly formula for  $f(y_1, ..., y_n)$ ...

#### Multivariate Normal Distribution

• When Y ~ N(
$$\mu$$
, $\Sigma$ ), then

$$\mu = (\mu_1, \mu_2, \dots, \mu_n)^T = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_n^2 \end{bmatrix}$$

is the *variance-covariance matrix*:

is the *mean vector*.

$$E[y_i] = \mu_i, \ i = 1, \dots, n$$

$$Var(y_i) = \sigma_i^2, i = 1, \dots, n$$
  
$$Cov(y_i, y_j) = \sigma_{ij}, i, j = 1, \dots, n$$

### Many Equivalent Model Formulations

• Y 
$$\sim$$
 N(X $eta$ ,  $\sigma^2$  I), i.e.:

- □  $E[y_i] = X_i\beta$ , i = 1, ...,n
- □ Var(y<sub>i</sub>) =  $\sigma^2$ , i=1, ..., n
- □  $Cov(y_i, y_j) = 0, \forall i \neq j$
- We could also write

■ Y = X
$$\beta$$
 +  $\epsilon$ ,  $\epsilon \sim N(0,\sigma^2 I)$   
■ y<sub>i</sub> = X<sub>i</sub> $\beta$  +  $\epsilon_i$ ,  $\epsilon_i \sim N(0,\sigma^2)$ , iid  
■ y<sub>i</sub> =  $\beta_1 X_{i1}$  + ... +  $\beta_k X_{iK}$  +  $\epsilon$ ,  $\epsilon_i \sim N(0,\sigma^2)$ , iid

# ML/LS Estimates

• 
$$\mathbf{Y}_{i} = \mathbf{X}_{i}\beta + \epsilon_{i}, \ \epsilon_{i} \sim \mathbf{N}(\mathbf{0}, \sigma^{2})$$

- So  $Var(Y_i) = E[(Y_i X_i\beta)^2] = E[\epsilon_i^2] = Var(\epsilon_i) = \sigma^2$
- Then we can estimate (MoM!):

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - X_i \beta)^2$$

 Fitting the model is basically just finding values β to minimize

$$\frac{1}{n} \sum_{i=1}^{n} (y_i - X_i \beta)^2 = \frac{1}{n} (Y - X \beta)^T (Y - X \beta)$$

• It turns out that  $\hat{\beta} = (X^T X)^{-1} X^T y$ 

### ML/LS Estimates

$$y = X\beta + \epsilon, \ \epsilon \sim N(0, \sigma^2 I)$$

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

The "residual SD" is the square root of

$$\hat{\sigma}^2 = \frac{1}{n-k} \sum_{i=1}^n (y_i - X_i \hat{\beta})^2 = \frac{1}{n-k} (y - X \hat{\beta})^T (y - X \hat{\beta})$$

We will see below that  

$$Var(\hat{\beta}) = (X^T X)^{-1} \sigma^2$$
  
 $Var(\hat{y}) = X(X^T X)^{-1} X^T \sigma^2$ 

# ML/LS Estimates - Example

```
> str(kidiq)
> y <- cbind(kidig$kid.score)</pre>
> X <- cbind(1, with(kidiq,</pre>
   cbind(mom.hs,mom.iq)))
> \dim(X)
[1] 434 3
> n <- dim(X)[1]
> k < - dim(X)[2]
> V <- solve(t(X) %*% X)
> beta.hat <- V %*% t(X) %*% y
> res.var <- t(y - X%*%beta.hat)</pre>
   %*% (y - X%*%beta.hat) / (n-k)
> res.sd <- sqrt(res.var)</pre>
>
> var.beta <- V * c(res.var)</pre>
> beta0.sd <- sqrt(var.beta[1,1])</pre>
> beta1.sd <- sqrt(var.beta[2,2])</pre>
> beta2.sd <- sqrt(var.beta[3,3])</pre>
```

```
> round(cbind(beta.hat,
   c(beta0.sd, beta1.sd,
   beta2.sd)), 2)
        [,1] [,2]
       25.73 5.88
mom.hs 5.95 2.21
mom.iq 0.56 0.06
> round(res.sd,2)
      [, 1]
[1,] 18.14
> summary(lm(kid.score ~ mom.hs
   + mom.iq, data=kidiq))
           coef.est coef.se
(Intercept) 25.73
                     5.88
mom.hs
           5.95
                     2.21
mom.iq
          0.56
                     0.06
n = 434, k = 3
res sd = 18.14, R-Squared = 0.21
```

estimates-by-hand.r

#### Two Interpretations

- In the model y<sub>i</sub> = β<sub>0</sub> X<sub>i0</sub> + ... + β<sub>p</sub> X<sub>ip</sub> + ε<sub>i</sub>, β<sub>j</sub> is the change in y for a unit change in X<sub>j</sub>, holding the other X's fixed
   Since β<sub>j</sub> estimates β<sub>j</sub>, β<sub>j</sub> inherits this interpretation also.
- When we look at "added variable plots", we will see that  $\hat{\beta}_j$  measures the variation in y, left after controlling for the other X's, that is uniquely attributable to  $X_j$ .

Distributional Properties:  $\hat{eta}$ 

$$\mathsf{Fact:}\ Y \sim N(\mu, \Sigma) \Rightarrow AY \sim N(A\mu, A\Sigma A^T)$$

$$y \sim N(X\beta, \sigma^2 I)$$

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y \\ E[\hat{\beta}] &= E[(X^T X)^{-1} X^T y] = (X^T X)^{-1} X^T E[y] \\ &= (X^T X)^{-1} X^T X \beta = \beta \\ \text{Var} (\hat{\beta}) &= \text{Var} ((X^T X)^{-1} X^T y) \\ &= (X^T X)^{-1} X^T \text{Var} (y) X (X^T X)^{-1} \\ &= (X^T X)^{-1} X^T X (X^T X)^{-1} \sigma^2 = (X^T X)^{-1} \sigma^2 \\ \Rightarrow \hat{\beta} &\sim N(\beta, (X^T X)^{-1} \sigma^2) \end{aligned}$$

# Hat Matrix H, & Distribution of $\hat{y}$ and $\hat{e}$

$$\begin{split} \hat{y} &= X\hat{\beta} &= X\left[(X^TX)^{-1}X^Ty\right] = \left[X(X^TX)^{-1}X^T\right]y = Hy\\ HX &= X(X^TX)^{-1}X^TX = X\\ &\Rightarrow \forall \beta^*, HX\beta^* = X\beta^* \\ H^T &= H \text{ and } (I-H)^T = (I-H)\\ HH &= H \text{ and } (I-H)(I-H) = (I-H)\\ E[\hat{y}] &= E[Hy] = HE[y] = HX\beta = X\beta\\ \text{Var}(\hat{y}) &= \text{Var}(Hy) = H\text{Var}(y)H^T = HH\sigma^2 = H\sigma^2\\ &\Rightarrow \hat{y} \equiv Hy \sim N(X\beta, H\sigma^2)\\ \text{Similarly,} \hat{e} = y - \hat{y} \equiv (I-H)y \sim N(0, (I-H)\sigma^2) \end{split}$$

#### SS Decompositions and F Statistics

Fact: 
$$y^T A y + y^T B y = (y^T A + y^T B) y = y^T (A + B) y$$

$$SST = \sum_{i=1}^{n} (y_i - \bar{y})^2 = (y - \bar{y})^T (y - \bar{y})$$
  
$$= y^T (I - H_1)^T (I - H_1) y = y^T (I - H_1) y$$
  
$$SS_{reg} = \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 = (\hat{y} - \bar{y})^T (\hat{y} - \bar{y}) = y^T (H - H_1) y$$

$$RSS = \sum_{i=1}^{n} (y_i - \hat{y})^2 = (y - \hat{y})^T (y - \hat{y}) = y^T (I - H) y$$

$$\Rightarrow SS_{reg} + RSS = y^T (H - H_1)y + y^T (I - H)y = y^T (I - H_1)y = SST$$

- Cochran's Theorem implies  $SS_{reg}/\sigma^2$  and  $RSS/\sigma^2$  are indep.  $\chi^2$ 's under  $H_0: \beta_1 = \cdots = \beta_p = 0$
- The df for each  $\chi^2$  will be the rank, or equivalently the trace, of each defining matrix. Using tr(AB) = tr(BA): tr(H) = p + 1,  $tr(H_1) = 1$ , tr(I) = n, so  $df(SS_{reg}/\sigma^2) = p$ ,  $df(RSS/\sigma^2) = n p 1$ ,  $df(SST/\sigma^2) = n 1$

$$(H_1 = hat matrix for y = \beta_0 + \varepsilon)$$
 16

# SS Decompositions and F Statistics

#### The foregoing lead to the traditional Analysis of Variance Table

Source of variation	Degrees of freedom (df)	Sums of squares (SS)	Mean square (MS)	F
Regression	p	SSreg	SSreg/p	$F = \frac{SSreg/p}{RSS/(n-p-1)}$
Residual	n-p-1	RSS	RSS/(n-p-1)	, ( 1 )
Total	n-1	SST		

■ As before we can define "multiple R<sup>2</sup>":

$$R^{2} = \frac{SS_{reg}}{SST} = 1 - \frac{RSS}{SST} \quad (= \operatorname{Corr}(y, \hat{y})^{2})$$
  
"Adjusted R<sup>2</sup>": mean-squares instead of sums of squares, to account for capitalization on chance  
$$R_{adj}^{2} = 1 - \frac{RSS/(n-p-1)}{SST/(n-1)}$$

#### SS Decompositions and F Statistics

Since  $SST = y^T (I - H_1)y = RSS_{H_1}$ , the residual sum-of-square from the smallest model (intercept-only), the F statistic from the Anova table can be written as

$$F = \frac{SS_{reg}/p}{RSS/n - p - 1} = \frac{(RSS_{H_1} - RSS_H)/(df_{H_1} - df_H)}{RSS_H/df_H} \quad (*)$$

This idea, and the sum-of-square calculations we did earlier, can be generalized so that, if  $H_{full}$  and  $H_{reduced}$  are hat matrices from a "full" model, and from a "reduced" model obtained by linear restrictions on the "full" model, then the *partial* F statistic

$$F = \frac{(RSS_{H_{reduced}} - RSS_{H_{full}})/(df_{H_{reduced}} - df_{H_{full}})}{RSS_{H_{full}}/df_{H_{full}}}$$

will have an F distribution under the null hypothesis that the linear restrictions are true.

(\*) The df throughout are *residual* df, that is, tr(I-H)

#### Some Comments

- It's good to know the "canonical" theory of the linear model and the Analysis of Variance table
  - Distribution assumptions and multiple testing matters
  - We will more fully discuss later in the course
- The "linear restrictions" for the partial F statistic usually amount to just setting some β's = 0. This is especially useful when a regressor is categorical, since a categorical X is recoded as a set of dummy variables, one for each level of X
- The partial F test brings us into "variable selection"
   We will more fully discuss variable selection later as well!

### Summary

- Matrix Form of Multiple Regression Model
- Multivariate Normal Distribution
- ML/LS Estimates
- Two Interpretations
- Distributional Properties
- SS Decompositions and F Statistics
- Some Comments