

1. (a) let $y = \text{gpa}$
 $x_1 = \text{sat}$
 $x_2 = \text{hsavg.}$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 (x_{i2})^2 + \epsilon_i$$

$$(b) y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & x_{12}^2 \\ 1 & x_{21} & x_{22} & x_{22}^2 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n2}^2 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Z_c β_c

(c) Reduced model of excluding hsavg.
 is same as model where $\beta_2 = \beta_3 = 0$

so test $H_0: \beta_2 = \beta_3 = 0$ (Reduced model)

is that $y_i = \beta_0 + \beta_1 x_{i1} + \tilde{\epsilon}_i$

HA: not H_0
 Complete Model is model of part (a), above.

Test statistic:

$$\frac{SSE_R - SSE_c}{2} \sim F_{\text{numdf}=2, \text{denof}=n-(3+1)}$$

$$\text{where } SSE_c = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

where \hat{y}_i is estimate of y_i under the complete model, namely

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \hat{\beta}_3 x_{i2}^2$$

$$\text{where } \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{pmatrix} = (Z_c^T Z_c)^{-1} Z_c^T y$$

$$\text{similarly, } SSE_R = \sum_{i=1}^n (y_i - \hat{y}_{iR})^2$$

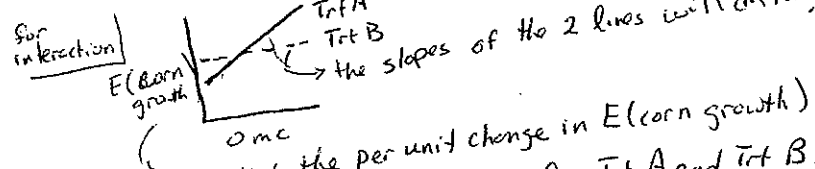
where \hat{y}_{iR} is estimate of y_i under the reduced model, namely

$$\hat{y}_{iR} = \tilde{\beta}_0 + \tilde{\beta}_1 x_{i1}, \text{ where}$$

$$\begin{pmatrix} \tilde{\beta}_0 \\ \tilde{\beta}_1 \end{pmatrix} = (Z_R^T Z_R)^{-1} Z_R^T y, \text{ where } Z_R = \begin{pmatrix} 1 & x_{11} \\ 1 & x_{21} \\ \vdots & \vdots \\ 1 & x_{n1} \end{pmatrix}$$

(2) Set up dummy f.v. $x_{i1} = \begin{cases} 1 & \text{if plot } i \text{ got Trt A} \\ 0 & \text{otherwise} \end{cases}$ (no interaction)
 for $i=1, 2, \dots, 18+15$, $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 \text{omc}_i + \epsilon_i$

(b) if interaction: $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 \text{omc}_i + \beta_3 (x_{i1} \cdot \text{omc}_i) + \epsilon_i$



This means that the per unit change in E(corn growth) as omc changes, would differ for Trt A and Trt B. (Thus the estimated effect of Applying Trt A vs Trt B, depends on the omc of the plot.)

Chapter 12 and Time Series

1. Consider the regression of (college) gpa by predicting variables of sat and high school grade point average (hs gpa). The proposed model consists of a linear trend over sat and a quadratic trend over hsavg. A random set of n students are obtained and for each the gpa, sat, hsavg is recorded.

- Present this model for a generic student i of the sample.
- Over all $i = 1, 2, \dots, n$, show the above model in terms of vectors and matrices (e.g. one component would show the terms that would be in the "Z" matrix).
- How would you test the null hypothesis of a reduced model that excluded hsavg as a predicting variable? (You should clearly state the null and alternative hypothesis, present the test statistic, state its distribution under the null hypothesis, explain how each component of the test statistic would be computed, and explain how you would decide whether or not to reject the null hypothesis, at the $\alpha = 0.05$ level.

(ALSO SEE 1. d. ON NEXT PAGE)

2. Fertilizer treatment A is given to $n_A = 18$ randomly selected plots from a large region; also, Fertilizer treatment B is given to $n_B = 15$ independently selected random plots from the same region. Prior to treatment, the organic matter content (omc) of each sample plot was measured. It is desired to model corn growth (y) in the plots as being potentially associated with which treatment is given and also a possible linear effect of a plot's organic matter content, with the assumption that there is no interaction between treatment and organic matter effects.

- Write this model in multiple regression form: with β 's, variables and error, being sure to clearly describe the variables.
- Now, if interaction between treatment and organic matter is added, then how would the model (in regression form) change and what would be the practical result in terms of the plot of E(corn growth) vs omc, shown separately for each respective treatment.

3. Consider that the monthly unemployment rate (y) can be modeled as the sum of a linear long term trend, a linear influence of the previous month's gdp, a continuous (cos and sin) seasonal adjustment, and first-order autocorrelated errors. Data on these are recorded for the last 48 months ($t=1, 2, \dots, 48$)

- Write the above model in terms of y_t, β 's, variables and ϕ .
- How would you predict next month's unemployment rate?

(3) (a) long-term linear trend

$$y_t = \beta_0 + \beta_1 t + \beta_2 \text{gdp}_{t-1} + \beta_3 \left(\cos \frac{2\pi t}{12} \right) + \beta_4 \left(\sin \frac{2\pi t}{12} \right) + R_t$$

$$\text{and } R_t = \phi R_{t-1} + \epsilon_t, \text{ for } t=2, 3, \dots, 48.$$

(b) Now $R_{49} = \phi R_{48} + \epsilon$, so $\hat{R}_{49} = \phi \hat{R}_{48}$ and where

$$\hat{R}_{48} = y_{48} - (\hat{\beta}_0 + \hat{\beta}_1 \cdot 48 + \hat{\beta}_2 \text{gdp}_{47} + \hat{\beta}_3 \left(\cos \frac{2\pi(48)}{12} \right) + \hat{\beta}_4 \left(\sin \frac{2\pi(48)}{12} \right))$$

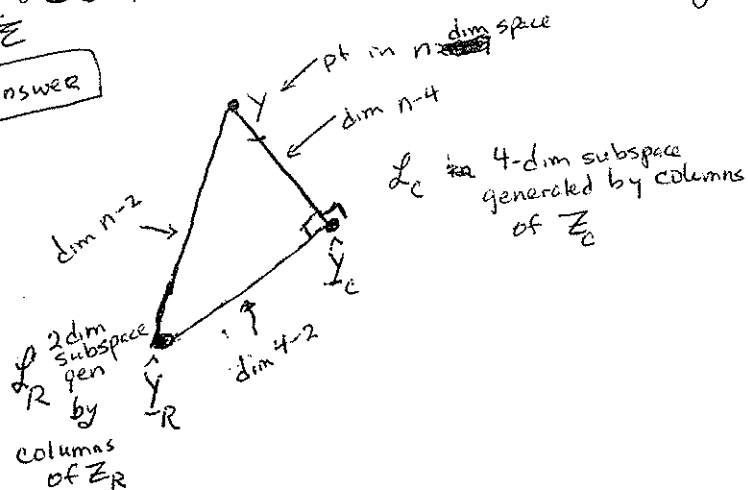
So, predict next month's unemployment rate by:

$$\hat{y}_{49} = \hat{\beta}_0 + \hat{\beta}_1 \cdot 49 + \hat{\beta}_2 \text{gdp}_{48} + \hat{\beta}_3 \left(\cos \frac{2\pi(49)}{12} \right) + \hat{\beta}_4 \left(\sin \frac{2\pi(49)}{12} \right) + \hat{R}_{49}$$

see above for \hat{R}_{49} .

1. d. Present the n -dimensional geometric representation of this test.

Answer



$$\text{Test statistic} = \frac{[(\hat{y}_c - \hat{y}_R)'(\hat{y}_c - \hat{y}_R)] / (4-2)}{[(y - \hat{y}_c)'(y - \hat{y}_c)] / (n-4)} \sim F_{H_0} \begin{matrix} \text{numdf} = 4-2 \\ \text{denof} = n-4 \end{matrix}$$

FOR ADDITIONAL EXPLANATION,
SEE 1-PAGE SUMMARY OF
EFRON ARTICLE ON NEXT PAGE

Other suggested problems from Chapter 12.

In book, Chpt. 12 #1, 2, 5, 12, 16, 28, 37, 32, 39,
44, 54, 60, 68, 70, 75, 88, 90,
101, 109

COMPLETE MODEL (2) quadratic: $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i$
REDUCED MODEL (1) linear: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$
 $H_0: \beta_2 = 0$

The spaces $\mathcal{S}(1)$ and $\mathcal{S}(2)$ are indicated in Fig. 3, along with the corresponding least-squares points $\hat{\mu}(1)$ and $\hat{\mu}(2)$. Because $\mathcal{S}(1) \subset \mathcal{S}(2)$ we must have

$$(3.5) \quad \sum_{i=1}^n (y_i - \hat{y}_{i(2)})^2 = SSE_C \Rightarrow \|y - \hat{\mu}(2)\|^2 \leq \|y - \hat{\mu}(1)\|^2 = SSE_R = \sum_{i=1}^n (y_i - \hat{y}_{i(1)})^2$$

In other words, increasing the explanatory space \mathcal{S} decreases the residual vector $r = y - \hat{\mu}$. As we shall see this does *not* mean that big models are always better than small ones.

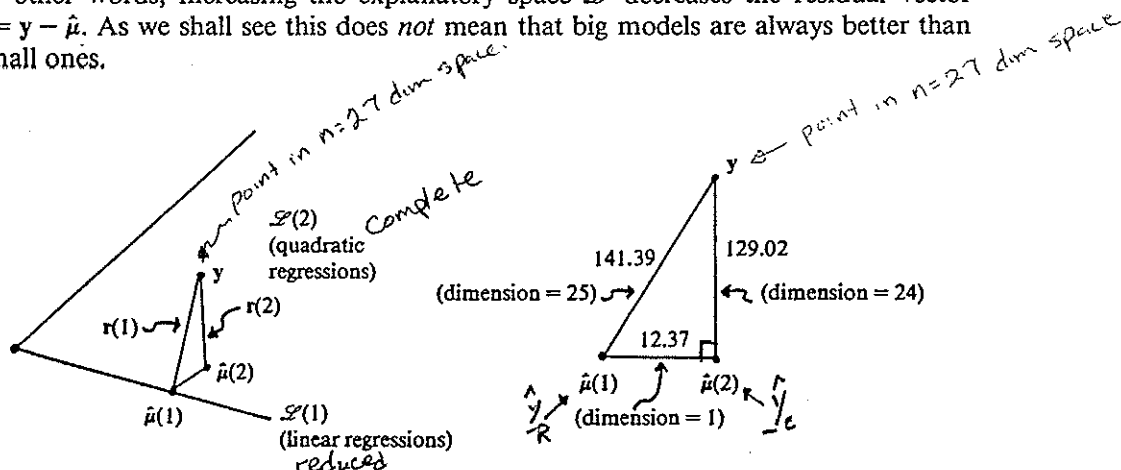


FIG. 3. On the left is a schematic diagram comparing linear and quadratic models for the hormone data of §1. Because $\mathcal{S}(1) \subset \mathcal{S}(2)$, the residual vector must be larger for $\mathcal{S}(1)$. The squared lengths $\|y - \hat{\mu}(1)\|^2$ and $\|y - \hat{\mu}(2)\|^2$ are shown at right. Traditional theory indicates that there is no strong reason to believe that quadratic regression is superior to linear regression for the hormone data.

Figure 3 indicates the squared residual lengths $\|y - \hat{\mu}(1)\|^2 = 141.39$ and $\|y - \hat{\mu}(2)\|^2 = 129.02$ for the hormone data, from which Pythagoras' theorem gives $\|\hat{\mu}(2) - \hat{\mu}(1)\|^2 = 12.37$. Does a decrease of 12.37 in squared residual length indicate a genuine advantage to the bigger model? Here is the traditional answer to that question:

(i) $y - \hat{\mu}(2)$ takes its value in the 24-dimensional subspace of \mathbb{R}^{27} orthogonal to the three-dimensional space $\mathcal{S}(2)$.

(ii) $\hat{\mu}(2) - \hat{\mu}(1)$ takes its value in a one-dimensional space, the portion of $\mathcal{S}(2)$ orthogonal to $\mathcal{S}(1)$.

(iii) If there is no true quadratic effect, that is, if the true mean vector μ lies in $\mathcal{S}(1)$, then we expect the ratio

$$(3.6) \quad F = \frac{(\hat{y}_{(2)} - \hat{y}_{(1)})^2 / 1}{\|y - \hat{\mu}(2)\|^2 / 24}$$

H_0 : no quadratic effect
 $F \sim F_{1,24}$
 $F_{num df=1} = k-q$
 $F_{den df=24} = n - (k+1)$

to approximately equal 1 (since then there is nothing special about $\hat{\mu}(2) - \hat{\mu}(1)$ compared to the other 24 components of $y - \hat{\mu}(1) = [\hat{\mu}(2) - \hat{\mu}(1)] + [y - \hat{\mu}(2)]$).

(iv) The observed value of F in this case is $12.37 / (129.02 / 24) = 2.30$, so we have to decide whether or not 2.30 "approximately equals 1."

(v) The theoretical distribution of F can be computed for model (2.6), (2.7); assuming that the true quadratic effect is zero,

$$(3.7) \quad \text{Prob}\{F > 2.30\} = .14 \quad \text{P-value of test}$$

(vi) An "achieved significance level" of .14 is not considered significant evidence for the existence of a genuine quadratic component in μ . We expect bigger values of F , and smaller achieved significance levels, if μ is genuinely quadratic. The conventional borderline for significance is .05, in this case $F \geq 4.26$. To state things in

$F_{(.05)}$
 $num df = 1$
 $den df = 24$