# CHAPTER 2
## MULTIPLE REGRESSION

### 2.1 REVIEW OF VECTOR AND MATRIX

* Matrix Cookbook (Petersen and Pedersen) Buch

Differentiation: Let $x$ be a column vector in $\mathbb{R}^d$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} = [x_1, x_2 \quad x_d]^T$$

Then for a function $F: \mathbb{R}^d \to \mathbb{R}$,

$$\frac{df(x)}{dx} = \begin{bmatrix} dF/\partial(x_1) \\ dF/\partial(x_d) \end{bmatrix} \quad \begin{array}{l} \leftarrow \text{Differentiate each individually} \\ = \nabla_x f(x) \\ \text{Gradient of } x \end{array}$$

① IF $f(x) = c$, a constant, then $\frac{df}{dx} = 0$ is a $d|x|$ vector of 0's

linear combination of x's

② IF $f(x) = a^T x = a_1 x_1 + a_2 x_2 + \cdots + a_d x_d$

$$df/dx = \begin{bmatrix} a_1 \\ \vdots \\ a_d \end{bmatrix} = a$$

③ Similarly if $f(x) = x^T a = a_1 x_1 + a_2 x_2 + \cdots + a_d x_d = a^T x$

$$\frac{dF}{dx} = a$$

④ Consider $f(x) = x^T A x$ ← a $d \times d$ matrix   $\frac{dF}{dx} = ?$

Look at $Ax$: $d \times 1$ vector with $j^{th}$ entry $[Ax]_j = \sum_{k=1}^{d} A_{jk} x_k$

$j^{th}$ row of A by corresponding row of x

Then $x^T A x = [x_1, \quad x_d] A x$

$$= \sum_{j=1}^{d} x_j [Ax]_j$$

Sum over x?

$$= \sum_{j=1}^{d} x_j \sum_{k=1}^{d} A_{jk} x_k$$

$$= \sum_{j=1}^{d} \sum_{k=1}^{d} x_j A_{jk} x_k$$

$$= f(x)$$

$$\frac{dF}{dx_L} = \frac{d}{dx}\left[\sum_{j=1}^{d}\sum_{H=1}^{d} x_T A_{JH} X_H\right] \qquad x^T A x = F(x)$$

$$= \frac{d}{dx_L}\left[\sum_{J=1}^{d}\left(A_{JJ} x_J{}^2 + x_J \sum A_{JH} X_H\right)\right]$$

$$= \frac{d}{dx_L}\left[A_{LL} X_L{}^2 + X_L \sum A_{LH} X_H\right] + \frac{d}{dx_L}\left[\sum_{J\neq L} x_J A_{JL} X_L\right]$$

$$= \left[2 A_{LL} X_L{}^2 + \sum_{H\neq L} A_{LH} X_H + \sum_{J\neq L} A_{JL} X_J\right]$$

$$= 2 A_{LL} X_L + \sum_{H\neq L}\left[A_{LH} + A_{HL}\right] X_H$$

Mostly we are concerned with a Symetric $\Rightarrow A_{LH} = A_{HL}$  $(A = A^T)$

$\rightsquigarrow 2 A_L X_L + \sum_{H\neq L} (2 A_{LH}) X_H$

$= 2 \sum_{K=1}^{d} A_{LH} X_K \qquad = dF/dx_L \qquad \text{w.r.t} \quad L^{th} x$

$$\frac{dF}{dx} = 2 A x$$

## Moments of Random Variables

$x = \begin{bmatrix} x_i \\ x_d \end{bmatrix}$ is a d-dimensional random vector (or d-dimensional random variable)

$$\mathbb{E}[x] = \begin{bmatrix} \mathbb{E}[x_1] \\ \mathbb{E}[x_d] \end{bmatrix}$$

If $a$ and $b$ are constants and $x, y$ are random vectors (of same dimension) the

$$\mathbb{E}[ax + by] = a\,\mathbb{E}[x] + b\,\mathbb{E}[y]$$

Covariance matrix of $x$

$\text{Var}[x]$

$\text{Cov}[x]$ $\begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1d} \\ & & & \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_{dd} \end{bmatrix}$  How they vary with each other.

Symetric matrix

- Where $\sigma_{ii} = \text{Var}[x_i]$  $\sigma_{ij} = \text{Cov}[x_i, x_j]$  when $i \neq j$
- The Covariance matrix is Symetric $\Rightarrow \Sigma = \Sigma^T$

IS    ALSM1

If the off diagonal entries are all zero, then the elements of $x$ are uncorrelated

$$\begin{bmatrix} \sigma_1^2 & & 0 \\ & \sigma_2^2 & \\ 0 & & \sigma_d^2 \end{bmatrix}$$

$Var[a^T x] = a^T \Sigma a$

$Cov[a^T x, b^T x] = a^T \Sigma b$

## 2.2 MATRIX FORMULATION OF THE SLR MODEL

- The SLR model is written $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$  $i = 1, \ldots, n$
- We can write this in matrix terms using random vectors.

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

$$\underline{y} = X\beta + \underline{\varepsilon}$$

$y$ : Response Variable

$X$ : Design Matrix

$\beta$ : Parameter / Coefficient Vector

$\varepsilon$ : Error Vector

$E[y] = X\beta$

$E[\varepsilon] = 0$    zero vector → $n \times 1$ vector of zeros

$Var[y] = Var[\varepsilon] = \begin{bmatrix} \sigma^2 & & 0 \\ & \sigma^2 & \\ 0 & & \sigma^2 \end{bmatrix} = \sigma^2 I$   independent ⇒ 0 in off diagonals

↖ Identity matrix $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

## 2.3 MULTIPLE REGRESSION

Consider extending the SLR model such that the dependent variable (response) has mean depending on a number of predictors/independent variables $X_1 \ldots X_p$

$$y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{2i} + \ldots + \beta_p X_{ip} + \varepsilon_i$$

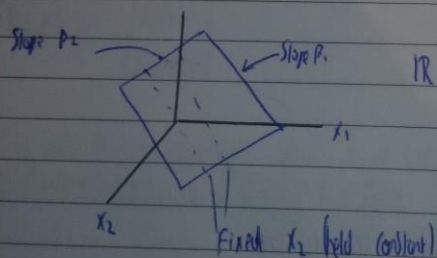$$\mathbb{E}[y_i \mid X_{i1}, X_{i2}, \ldots, X_{pi}]$$

The model can be written in matrix notation $y = X\beta + \varepsilon$

$$
\underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{n \times 1}
\quad
\underbrace{\begin{bmatrix} 1 & X_{11} & X_{21} & \cdots & X_{p1} \\ 1 & X_{12} & X_{22} & & X_{p2} \\ \vdots & & & & \\ 1 & X_{1n} & X_{2n} & \cdots & X_{pn} \end{bmatrix}}_{n \times (p+1)}
\quad
\underbrace{\begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}}_{(p+1) \times 1}
+
\underbrace{\begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}}_{n \times 1}
$$

$\beta$ is a vector of unknown parameters to be estimated from the data

$\beta_j$ is the change in the mean value of $y$ per unit change in $X_J$ assuming all other independent variables are held constant

$$\mathbb{E}[y] = X\beta \qquad \mathbb{E}[\varepsilon] = 0 \qquad Var[y] = \sigma^2 I = Var[\varepsilon] \quad \text{(independent errors)}$$
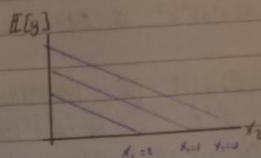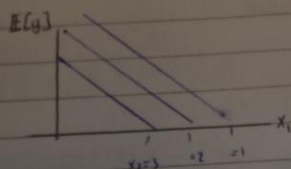
When $p = 2$ the model is 2 dimensional plane in a 3d space



Slope $P_2$   Slope $P_1$   $\mathbb{R}^3$

$X_1$

$X_2$   Fixed $X_2$ held constant

23/11/15    ALM 1

Previous Diagram, would have $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$

Could write model as $y = x\beta + \varepsilon$



When $Y_2$ held constant, all lines have

Slope $\beta_1$

- The model   $y = x\beta + \varepsilon$   is known as the general linear model - not to be confused with GLM's

- Includes:   SLR, multiple regression, analysis of variance (one way classification), others.

## 2.4  LEAST SQUARES ESTIMATORS

The sum of squared errors can be written in matrix form

$$SSE = \sum \hat{\varepsilon}_i^2 = [\hat{\varepsilon}_1, \dots, \hat{\varepsilon}_n] \begin{bmatrix} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_n \end{bmatrix} = \hat{\varepsilon}^T \hat{\varepsilon}$$

$$(y^T - \hat{\beta}x)$$

$$= (y - x\hat{\beta})(y - x\hat{\beta})$$

$$= y^T y - y^T x \hat{\beta} - \hat{\beta}^T x^T y + \hat{\beta}^T x^T x \hat{\beta}$$

We want to minimise SSE WRT $\hat{\beta}$: ie. least squares Estimates

$$\frac{d\,SSE}{d\hat{\beta}} = 0 - x^T y - x^T y + 2x^T x \hat{\beta} = 0$$

$$x^T x \hat{\beta} = x^T y \quad \text{(normal equations)}$$

Assuming $x^T x$ is non-singular (ie. it is invertible) then $(x^T x)^{-1} x^T x \hat{\beta} = (x^T x)^{-1} x^T y$

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

## 2.5 LEAST SQUARES PLANE

Once we have the least squares estimates of $\beta$, then the predicted mean of $y$

is $\hat{y} = x\hat{\beta}$ $= x(x^Tx)^{-1}x^Ty = H y$

$H = x(x^Tx)^{-1}x^T$ is called the Hat Matrix

## 2.6 ANALYSIS OF VARIATION IN $y$

$$SSE = \hat{\varepsilon}^T\hat{\varepsilon} = (y-x\hat{\beta})^T(y-x\hat{\beta})$$

$$= y^Ty - y^Tx\hat{\beta} - \hat{\beta}^Ty + \hat{\beta}^Tx^Tx\hat{\beta}$$
$$\underset{\text{scalar}}{\underbrace{\quad\quad}}$$

$$= y^Ty - \hat{\beta}^Tx^Ty - \hat{\beta}^Ty + \hat{\beta}^Tx^Tx\hat{\beta} \quad \text{from above}$$

$$= y^Ty - 2\hat{\beta}^Tx^Ty + \hat{\beta}^Tx^Tx [ (x^Tx)^{-1}x^Ty ]$$
$$\underset{\text{cancel each other}}{\underbrace{\qquad\qquad}}$$

$$= y^Ty - 2\hat{\beta}^Tx^Ty + \hat{\beta}^Tx^Ty$$

$$= \underline{y^Ty - \hat{\beta}^Tx^Ty}$$

$$y^Ty = [y_1 \dots y_N]\begin{bmatrix}y_1 \\ \vdots \\ y_N\end{bmatrix} = y_1^2 + \dots + y_N^2 = \Sigma y_i^2 = SS(\text{uncorrected})$$

$SSE = SS(\text{Uncorrected}) - SS(\text{Model})$

$SS(\text{uncorrected}) = SS(\text{Model}) + SSE$

$n\bar{y}^2 = \text{correction}$

$SS(\text{Uncorrected}) - \text{Correction} = SS(\text{Reg}) + SSE$

$SS(\text{Corrected}) = SS(\text{Reg}) + SSE$

## 2.7 PROPERTIES OF THE ESTIMATORS

SLR    $E[\hat{\beta_0}] = \beta_0$          $Var[\hat{\beta_0}] = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)$

$\quad\quad\ E[\hat{\beta_1}] = \beta_1$          $Var[\hat{\beta_1}] = \frac{\sigma^2}{S_{xx}}$

$$\hat{\beta} = (X^TX)^{-1}X^Ty$$

$$E[\hat{\beta}] = E[(X^TX)^{-1}X^Ty]$$

$$= (X^TX)^{-1}X^T E[y] \quad x\text{'s are fixed quantities} \Rightarrow \text{not random}$$

$$= (X^TX)^{-1}X^T [X\beta]$$

$$= \underbrace{(X^TX)^{-1}X^T(X^TX)}_{I}\ \beta$$

$$= \beta \quad\quad LS \text{ is unbiased estimator}$$

For SLR model:

$$Var[\hat{\beta}] = \begin{bmatrix} Var[\hat{\beta_0}] & Cov[\hat{\beta_0}, \hat{\beta_1}] \\ Cov[\hat{\beta_0}, \hat{\beta_1}] & Var[\hat{\beta_1}] \end{bmatrix}$$

In general, $Var[\hat{\beta}] = Var[\overbrace{(X^TX)^{-1}X^T}^{A}y]$

$$= A\, Var[y]\, A^T$$

$$= A\, [\sigma^2]\, A^T$$

$$= \sigma^2 AA^T$$

$$= \sigma^2 \left[(X^TX)^{-1}X^T\right]\left[X(X^TX)^{-1}\right]$$

$$= \sigma^2 (X^TX^{-1})$$

In multiple regression, Gauss-Markov theorem holds: LS estimates are Best linear unbiased estimators (BLUE)

## Example: Cereals Dataset

- Nutritional info and shelf location for 77 breakfast cereals
- A rating was calculated from consumer reports
- "Middle Shelf" cereals tended to have lowest rating
- R Script: investigate relationship b/w rating, sugar content (per 100g) and fat content (per 100g)

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$$

Rating    Intercept    Sugar/100g    Fat/100g

$$\hat{\beta} = \begin{bmatrix} \hat{\beta_0} \\ \hat{\beta_1} \\ \hat{\beta_2} \end{bmatrix} = \begin{bmatrix} 61.084 \\ -2.213 \\ -3.066 \end{bmatrix}$$

- Mean rating decreasing by 2.2 for every extra gram of sugar per 100g, keeping the fat content fixed.

- Similarly, the mean rating decreases by ~3 for every extra gram of fat, keeping sugar fixed.

ANOVA Table    n = 77

| Source | DF | SS | ms |
|---|---|---|---|
| $X_1$ | 1 | 8654.7 | 8654.7 |
| $X_2$ | 1 | 670.5 | 670.5 |
| Residual | 74 | 5671.5 | 76.6 |

Add together for pre-table [ ... ]

MSE = 76.6

$$Var[\hat{\beta}] = MSE (X'X)^{-1} = \begin{bmatrix} 3.813 & -0.315 & -0.632 \\ -0.315 & 0.055 & -0.066 \\ -0.632 & -0.066 & 1.074 \end{bmatrix}$$

S.E. $[\hat{\beta_0}]$ = 3.813

SE $[\hat{\beta_1}]$ = 0.055

SE $[\hat{\beta_2}]$ = 1.074

ANOVA Table

| Source | DF | SS | MS |
|--------|-----|-----|-----|
| Regression | $p$ | $\hat{\beta}^T x^T y - n\bar{y}^2$ | $SS(Reg)/p$ |
| Residual | $n-p-1$ | $(y-x\hat{\beta})^T(y-x\hat{\beta})$ | $SSE/(n-p-1)$ |
| Total (corrected) | $n-1$ | $y^T y - n\bar{y}^2$ | |

SLR: $SSE = \Sigma(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \rightarrow n-2$ df

Here we have $p$ predictors (independent variables). We estimate $p+1$ parameters. So the DF associated with SSE is $n-(p+1)$

$y^T y n\bar{y}^2 = \Sigma y_i^2 - n\left(\frac{\Sigma y_i}{n^2}\right) = \Sigma y_i^2 - \frac{(\Sigma y_i)^2}{n} = \Sigma(y_i - \bar{y})^2 \quad df = n-1$

$R^2 = \frac{SS(Reg)}{SS(corrected)}$   $0 < R^2 < 1$   The square of the multiple correlation between $y$ and the predictors $(x_1, \ldots, x_p)$

- Gives proportion of ~~variance~~ variation in $y$ explained by the predictors

$\hat{\underline{\varepsilon}} = \underline{y} - \hat{\underline{y}}$   vector of residuals
$= \underline{y} - H\underline{y} = I_{n \times n} \underline{y} - H\underline{y}$
$= (I - H)\underline{y}$

$Cov[\hat{\beta}_0, \hat{\beta}_1] = -0.315$
$Cov[\hat{\beta}_1, \hat{\beta}_2] = -0.066$

$H_0 : \beta_2 - \beta_1 = 0$   v   $H_1 : \beta_2 - \beta_1 \neq 0$   need estimate for S.E. for $[\hat{\beta}_2 - \hat{\beta}_1]$
$Var[\hat{\beta}_2 - \hat{\beta}_1] = Var[\hat{\beta}_2] + Var[\hat{\beta}_1] - 2Cov[\hat{\beta}_1, \hat{\beta}_2]$   ←
Used in one-way classification models

2.8 INFERENCE FOR MULTIPLE REGRESSION

In order to make inferences we again have to make an assumption about the distribution of $y$. We usually assume normality

ASIDE: $y$ has a multivariate normal distribution with parameters $\mu$ (mean vector) and $\Sigma$ (covariance matrix)

$$\underset{n \times 1}{y} \sim N_r \; (\underset{n \times 1}{\mu}, \underset{n \times m}{\Sigma})$$

Density function of $y$

$$f(y) = \frac{1}{(2\pi)^{n/2}} |\Sigma|^{-1/2} \exp\left[-\tfrac{1}{2}(y-\mu)^T \Sigma^{-1}(y-\mu)\right]$$

↳ multivariate distribution

$\mu = x\beta$

$\Sigma = \sigma^2 I$

$$f(y) = \frac{1}{(2\pi)^{n/2}} |\sigma^2 I|^{-1/2} \exp\left[-\tfrac{1}{2}(y-x\beta)^T (\sigma^2 I)^{-1}(y-x\beta)\right]$$

$$\frac{1}{(2\pi)^{n/2}(\sigma^2)^{n/2}} \exp\left[\frac{-1}{2\sigma^2}(y-x\beta)^T(y-x\beta)\right]$$

When $n=1$, $y \sim N(\mu, \sigma^2)$

$$f(y) = \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left[\frac{1}{2\sigma^2}(y-\mu)^2\right]$$

## CONFIDENCE INTERVALS

Confidence Intervals on the regression coefficients: $\hat{\beta}$ is a linear estimator hence

$$\hat{\beta} \sim N_{p+1}(\beta, (x^Tx)^{-1}\sigma^2)$$

$$\hat{\beta}_j \sim N(\beta_j, C_{jj}\sigma^2)$$

Where $C_{jj}$ is the $j+1$st diagonal entry of $(x^Tx)^{-1}$ for $j = 0, 1, \ldots, p$

Consequently $\dfrac{\hat{\beta}_j - \beta_j}{\sqrt{MSE \; C_{jj}}} \sim t_{n-p-1}$

The degrees of freedom of the $t$ dist is that associated with the MSE using a similar argument to Section 1-2

$$P\left[\hat{\beta}_j - t_{n-p-1,\,\alpha/2}\sqrt{MSE\,C_{jj}} \leq \beta_j < \hat{\beta}_j + t_{n-p-1,\,\alpha/2}\sqrt{MSE\,C_{jj}}\right] = 1-\alpha$$

This is a $100(1-\alpha)\%$ confidence interval for $\beta_j$

$\hat{\beta} \sim N_{p+1} \left( \beta, \sigma^2 (x^T x)^{-1} \right)$

$\uparrow_{c_{jj}}$ (j+1)th diagonal element

## Confidence Interval For the mean response

Define $x_0 = \begin{bmatrix} 1 \\ x_{10} \\ x_{20} \\ \vdots \\ x_{p0} \end{bmatrix}$   (p+1) x 1   col vector

The mean value of $y$ at this point $x_0$ is $x_0^T \beta$ which is estimated

by $x_0^T \hat{\beta} = \hat{y}_0$

$E[\hat{y}_0] = E[x_0^T \hat{\beta}] = x_0^T \beta = \mu_0$

$Var[\hat{y}_0] = x_0^T \, Var[\hat{\beta}] \, x_0$

$\qquad = \sigma^2 \, x_0^T \, (x^T x)^{-1} \, x_0$

$\hat{y}_0 \sim N \left( \mu_0, \; \sigma^2 x_0^T (x^T x)^{-1} x_0 \right)$

Replace $\sigma^2$ by MSE

CI for $\mu_0$:   $\dfrac{\hat{y}_0 - \mu_0}{\sqrt{\sigma^2 \, x_0^T (x^T x)^{-1} x_0}} \sim N(0,1)$

$\dfrac{\hat{y}_0 - \mu_0}{\sqrt{MSE \, x_0^T (x^T x)^{-1} x_0}} \sim t_{n-p-1}$   distribution

$100(1-\alpha)\%$ CI:   $\hat{y}_0 \pm t_{n-p-1, \alpha/2} \left( \sqrt{MSE \, x_0^T (x^T x)^{-1} x_0} \right)$

## Prediction

Prediction Interval for a new observation $y_0$ is the value of a future observation

of $x_0$, it is estimated by $\hat{y}_0 = x_0^T \hat{\beta}$ and the prediction interval is

$\hat{y}_0 \pm t_{n-p-1, \alpha/2} \left( \sqrt{MSE \left(1 + x_0^T (x^T x)^{-1} x_0 \right)} \right) = 1-\alpha$

# Hypothesis Testing

An overall test of significance of regression is given by:

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_p = 0 \quad vs \quad H_A: \beta_j \neq 0 \text{ for at least one } j$$
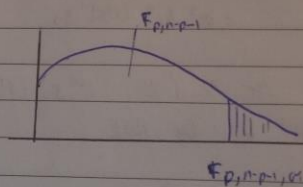
$$E[y_i] = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{pi} \quad \text{if all } \beta's, \, j > 0 \text{ are zero, } E(y_i) = \beta_0$$

## Test Statistic

$$F = \frac{MS(Reg)}{MSE} = \frac{R(\beta_1, \cdots, \beta_p \mid \beta_0)/p}{MSE}$$

If $H_0$ is true, then (ie. $\beta_1 = \cdots = \beta_p = 0$), then F follows an F-Dist with $p$ and $n-p-1$ D.F.

If $F > F_{p,n-p-1,\alpha}$ then
reject $H_0$ at $100\alpha\%$ Significance level
with a controlling probability of
type I error.



## Testing On Individual Regression Coefficient

$$H_0: \beta_j = 0 \quad v \quad H_A: \beta_j \neq 0 \quad \text{(can be directional)}$$

$$t\text{-test} = \frac{\hat{\beta}_j}{\sqrt{MSE \, C_{jj}}} \qquad C_{jj}: j+1^{th} \text{ diagonal element of } (x^Tx)^{-1}$$

$n-p-1$ D.F.

This is the test of the contribution of $x_j$ given all the other independent variables are the model. Compare test statistic with critical value and conclude

## "Extra Sum of Squares"

Reduced model: $E[y_i] = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_{j-1} x_{j-1,i} + \beta_{j+1,i} + \cdots + \beta_p x_{pi}$

   Just removed the $j^{th}$ term

Full Model: $E[y_i] = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{pi}$

$R(\beta_J \mid \beta_{J-1}, \beta_{J-1}, \beta_{J+1}, \beta_P)$ = Extra Sum of Squares

Here, the extra sum of squared is the partial SS for $x_J$ and represents the contribution of $x_J$, adjusted for all other independent variables in the model

$$F = \frac{R(\beta_J \mid \beta_0, \beta_P)/1}{MSE} \rightarrow D.F \ 1, \ n-p-1$$

is the test statistic for $H_0: \beta_J = 0$ vs $H_A: \beta_J \neq 0$ (can't do directional testing).

This is called the partial F-test for $\beta_J$. Equivalent to the two tailed t-test $[F = t^2]$

How to do a joint test for regression coefficients to measure overall usefullness of the regression model

$H_0: \beta_1 = \beta_2 = \ldots = \beta_P = 0$     (Model Utility Test)
$H_A:$ At least one of the $\beta_J \neq 0$

Test For a Subset OF Regression Coefficients
  Reduced Model: $E[y_i] = \beta_0 + \beta_1 x_{1i} + \ldots + \beta_q x_{qi}$     $q < p$
  Full Model: $E[y_i] = \beta_0 + \beta_1 x_{1i} + \ldots + \beta_p x_{pi}$

$H_0: \quad \beta_{q+1} = \beta_{q+2} = \ldots = \beta_p = 0$     (testing a subset)
$H_A: \quad$ At least one of these $\beta$'s is $\neq 0$

$R(\beta_{q+1}, \beta_p \mid \beta_0, \beta_q) =$ Extra Sum of Squares (how much more info you explain by) including them in model

Use F test: $F = \dfrac{R(\beta_{q+1}, \beta_p \mid \beta_0, \beta_q)/(p-q)}{MSE}$ — Difference

If $H_0$ is true, then $F$ follows an $F_{p-q, \ n-p-1}$ distribution
Compute $F$ and compare it to a critical value

## 2.9 SEQUENTIAL SUM OF SQUARES

The extra Sum of Squares partition $SS(Reg)$ in the ANOVA Table

$$SS(Reg) = R(\beta_1, \beta_2, \dots, \beta_p \mid \beta_0)$$
$$= R(\beta_1, \dots, \beta_q \mid \beta_0) + R(\beta_{q+1}, \dots, \beta_p \mid \beta_0, \dots, \beta_q)$$

| Source | DF |
|---|---|
| $R(\beta_1, \dots, \beta_q \mid \beta_0)$ | $q$ |
| $R(\beta_{q+1}, \dots, \beta_p \mid \beta_0, \dots, \beta_q)$ | $p-q$ |
| $R(\beta_1, \dots, \beta_p \mid \beta_0)$ | $p$ |
| Residual | $n-p-1$ |
| Total | $n-1$ |

ADD

We can extend this partitioning as follows:

| Source | DF |
|---|---|
| $R(\beta_1 \mid \beta_0)$ | 1 |
| $R(\beta_2 \mid \beta_0, \beta_1)$ | 1 |
| $\vdots$ | $\vdots$ |
| $R(\beta_p \mid \beta_0, \dots, \beta_{p-1})$ | 1 |
| $SS(Reg) = R(\beta_1, \dots, \beta_p \mid \beta_0)$ | $p$ |
| Residual | $n-p-1$ |
| Total | $n-1$ |

$R(\beta_J \mid \beta_0, \dots, \beta_{J-1})$ is known as the sequential sum of squares for $X_J$ → The amount by which you will reduce the residual Sum of Squares by adding $X_J$ to the model given $X_1, \dots, X_{J-1}$ → and represent the contribution of $X_J$ in the model adjusted for $X_1, \dots, X_{J-1}$ but NOT $X_{J+1}, \dots, X_p$. It also depends on the order.

We can see $\sum_{J=1}^{p} R(\beta_J \mid \beta_0, \dots, \beta_{J-1}) = R(\beta_1, \dots, \beta_p \mid \beta_0)$

Example:

$E[y] = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i}$

ALSM 1

| Var | Partial SS | Sequential SS |
|-----|-----------|---------------|
| $X_1$ | $R(\beta_1 \mid \beta_0, \beta_2, \beta_3)$ | $R(\beta_1 \mid \beta_0)$ |
| $X_2$ | $R(\beta_2 \mid \beta_0, \beta_1, \beta_3)$ | $R(\beta_2 \mid \beta_0, \beta_1)$ |
| $X_3$ | $R(\beta_3 \mid \beta_0, \beta_1, \beta_2)$  ← equal → | $R(\beta_3 \mid \beta_0, \beta_1, \beta_2)$ |

- Partial SS - Find the info contained in $X_5$ which is not contained in $X_k$ for $k \neq 5$

- Sequential SS - Info in $X_5$ not contained in $X_1, \ldots, X_{5-1}$

$R(\beta_1 \mid \beta_0) + R(\beta_2 \mid \beta_0, \beta_1) + R(\beta_3 \mid \beta_0, \beta_1, \beta_2) = R(\beta_1, \beta_2, \beta_3 \mid \beta_0) = SS(Reg)$

In general, Sum of Partial SS $\neq$ SS(Reg)

Example: Cereals

1. 95% CI: $\beta_0$: $\hat{\beta}_0 \pm t_{74, 0.025} \; S.E.[\hat{\beta}_0] = 61.084 \pm 1.993 \sqrt{3.813}$

$\beta_1$: $\hat{\beta}_1 \pm t_{74, 0.025} \; S.E.[\hat{\beta}_1] = -2.213 \pm 1.993 \sqrt{0.055}$

$\beta_2$: $\hat{\beta}_2 \pm t_{74, 0.025} \; S.E.[\hat{\beta}_2] = -3.066 \pm 1.993 \sqrt{1.074}$

2. $H_0$: $\beta_1 = \beta_2 = 0$    v   $H_A$: Either $\beta_1$ or $\beta_2 \neq 0$

$F = \dfrac{9325.2/2}{76.6} = 60.864$    df = 2, 74

    5% critical value: 3.12 → test is highly significant, reject $H_0$

3. $H_0$: $\beta_1 = 0$    (given $\beta_2$ is in the model)

$H_A$: $\beta_1 \neq 0$

Partial SS → $R(\beta_1 \mid \beta_0, \beta_2)$

$t = \dfrac{\hat{\beta}_1 - 0}{SE[\hat{\beta}_1]} = \dfrac{-2.213}{\sqrt{0.055}} = -9.4362$

Critical value 5% → $\pm 1.993$, highly significant, reject $H_0$

2.10 THE GENERAL LINEAR HYPOTHESIS

$H_0$: $L\beta = c$    $H_A$: $L\beta \neq c$

$L$ : A $k \times (p+1)$ matrix of coefficients

$\beta$ : $(p+1) \times 1$ vector of parameters

$s$ : $k \times 1$ vector of constants

All the hypotheses we discussed are special cases of this

Example: $E[y_i] = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i}$

$H_0: \beta_2 = 0$     v $H_1: \beta_2 \neq 0$

$$[0 \ 0 \ 1 \ 0] \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = [0]$$

$\underset{L}{} \quad \underset{\beta}{} \quad \underset{c}{}$

$H_0: \beta_1 = \beta_2 = 0$

$$\begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$\underset{L}{} \quad\quad \underset{\beta}{} \quad\quad \underset{c}{}$

- The extra sum of squares method can be used to test $H_0: L\beta = c$
- Reduced model ($H_0$) is constrained so that $H_0$ is true. Full model is not constrained

$F = \dfrac{\text{extra SS} / k}{MSE}$    with $k$, $n-p-1$ df

## 2.11 ORTHOGONALITY

Consider for example : $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$

$$y = X\beta + \varepsilon$$

$X \rightarrow$ Design matrix $n \times 3$

$$X = \begin{bmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & & \\ 1 & x_{1n} & x_{2n} \end{bmatrix} = [X_0 \vdots X_1 \vdots X_2]$$

Consider the situation where the columns of X are orthogonal

$X_0^T X_1 = 0$

$X_1^T X_2 = 0$       Orthogonality

$X_0^T X_2 = 0$

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

$$X^T = \begin{bmatrix} X_0^T \\ X_1^T \\ X_2^T \end{bmatrix} \qquad X^T X = \begin{bmatrix} X_0^T \\ X_1^T \\ X_2^T \end{bmatrix} [X_0 \ X_1 \ X_2]$$

$$= \begin{bmatrix} X_0^T X_0 & X_0^T X_1 & X_0^T X_2 \\ X_1^T X_0 & X_1^T X_1 & X_1^T X_2 \\ X_2^T X_0 & X_2^T X_1 & X_2^T X_2 \end{bmatrix}$$

$$= \begin{bmatrix} X_0^T X_0 & 0 & 0 \\ 0 & X_1^T X_1 & 0 \\ 0 & 0 & X_2^T X_2 \end{bmatrix} \qquad \text{because of orthogonality}$$

$$(X^T X)^{-1} = \begin{bmatrix} (X_0^T X_0)^{-1} & 0 & 0 \\ 0 & (X_1^T X_1)^{-1} & 0 \\ 0 & 0 & (X_2^T X_2)^{-1} \end{bmatrix}$$

$$X^T y = \begin{bmatrix} X_0^T \\ X_1^T \\ X_2^T \end{bmatrix} [\underline{y}] = \begin{bmatrix} X_0^T y \\ X_1^T y \\ X_2^T y \end{bmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T y = \begin{bmatrix} (X_0^T X_0)^{-1} X_0^T y \\ (X_1^T X_1)^{-1} X_1^T y \\ (X_2^T X_2)^{-1} X_2^T y \end{bmatrix}$$

$$SS(\text{model}) \quad \hat{\beta}^T X^T y$$

$$= [\hat{\beta}_0 \ \hat{\beta}_1 \ \hat{\beta}_2] \begin{bmatrix} X_0^T \\ X_1^T \\ X_2^T \end{bmatrix} \underline{y}$$

$$= [\hat{\beta}_0 \ \hat{\beta}_1 \ \hat{\beta}_2] \begin{bmatrix} X_0^T y \\ X_1^T y \\ X_2^T y \end{bmatrix}$$

$$= \hat{\beta}_0 X_0^T y + \hat{\beta}_1 X_1^T y + \hat{\beta}_2 X_2^T y \qquad \text{Partitions according to each column of design matrix}$$

So there is a contribution (clean) for each column of the design matrix

Consider a reduced model: $y_i = \beta_0 + \beta_1 x_{1i} + \varepsilon_i$

Can be shown, as above:

$$\hat{\beta}_0 = (X_0^T X_0)^{-1} X_0^T y$$

$$\hat{\beta}_1 = (X_1^T X_1)^{-1} X_1^T y$$

So the full and reduced models have the same parameter estimates for $\beta_0$ and $\beta_1$

The reduced model has:

$$SS_{Red}(\text{model}) = \hat{\beta}_0 X_0^T y + \hat{\beta}_1 X_1^T y$$

Orthogonally gives us separation we don't normally get in design matrix

Therefore: $SS(\text{model}) - SS_{Red}(\text{model}) = \hat{\beta}_2 X_2^T y$

$$R(\beta_2 | \beta_0, \beta_1) = \hat{\beta}_2 X_2^T y$$

Which is independent of $\beta_0$, $\beta_1$ and it's regression sum of squares we'd get if we regress $y$ on $x_2$ Alone

$$R(\beta_2 | \beta_0 \beta_1) = R(\beta_2 | \beta_0)$$

In Orthogonal model: Partial SS = Sequential SS

Example: Regression Through The Origin
Consider problem 3 on problem Sheet 2

M1 · Full  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$
    Reduced:  $y_i = \beta_1 x_i + \varepsilon_i$

M2: Full:  $y_i = a_0 + a_1(x_i - \bar{x}) + \varepsilon_i$
    Reduced:  $y_i = a_1(x_i - \bar{x}) + \varepsilon_i$

M1  Full:  $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$     $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$
    Reduced:  $Q = \sum(y_i - \beta_1 x_i)^2$
       $\frac{dQ}{d\beta} = 2\sum(y_i - \beta_1 x_i)(x_i) = 0$
          $\hat{\beta}_1 = \sum x_i y_i / \sum x_i^2$   (Different from full model)
    $\frac{S_{xy}}{S_{xx}} = \frac{\sum(x_i - \bar{x})\bar{x}}{\sum(x_i - \bar{x})^2}$

M2:  Full  $a_1 = \frac{S_{xy}}{S_{xx}} = \beta_1 \,(\text{Full})$     $\hat{a}_0 = \bar{y}$
    Reduced:  $y_i = a_1(x_i - \bar{x}) + \varepsilon_i$
       $Q = \sum(y_i - a_1(x_i - \bar{x}))^2$
       $\frac{dQ}{da_1} = \sum(x_i - \bar{x})(y_i - a_1 x_i - \bar{x}))$
          $= \sum(x_i - \bar{x})y_i - a_1 \sum(x_i - \bar{x})^2$
       $\hat{a}_1 = \frac{\sum(x_i - \bar{x})y_i}{\sum(x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$   Orthogonal

So the values of $a_1$ are the same as the full and reduced model
Making the transformation $z_i = x_i - \bar{x}$ gives an orthogonality transform

Write down design matrix for full M2

$$\begin{pmatrix} 1 & x_1 - \bar{x} \\ 1 & x_2 - \bar{x} \\ \vdots & \vdots \\ 1 & x_n - \bar{x} \end{pmatrix}$$

$$x_0^T x = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \begin{pmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_n - \bar{x} \end{pmatrix} = \Sigma(x_i - \bar{x}) = 0$$

$x_0 \ x_1$

## 2.12 MULTICOLLINEARITY

Recall the normal equations in matrix form $(X^T X)\hat{\beta} = X^T y$

We made the assumption that $X^T X$ is invertible

For $(X^T X)^{-1}$ to exist, the columns of $X^T X$ have to be linearly independent

If the columns are not linearly independent it can suggest that there is some redundancy of information in the predictors (model) → i.e. the predictors are giving the same information about the mean of $y$

If $x_1$ and $x_2$ are the proportion of water and solids in beer, then $x_1 + x_2 = 1$ and we have a linear dependency
- The correct approach would be to fit an intercept and either $x_1$ or $x_2$

- If there is linear dependence in the columns of $X$, we can't invert $X^T X$
- This is called multicollinearity

- One could have approximate or close to linear dependence between columns. Leads to large values in $(X^T X)^{-1}$ which in turn leads to large standard errors for $\hat{\beta}$
- Unreasonably large CI for $\beta$

$$X = \begin{bmatrix} 1 & x_{11} & 1 - x_{11} \\ 1 & x_{21} & 1 - x_{21} \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} & 1 - x_{n1} \end{bmatrix}$$