# CS 593
## Khasha Dehnad

## Simple + Multiple Linear Regression

## Class restarts
## at 7:45

# Simple Regression

# Introduction to Regression Analysis

- Regression analysis is used to:
  - Predict the value of a dependent variable based on the value of at least one independent variable
  - Explain the impact of changes in an independent variable on the dependent variable

Dependent variable: the variable we wish to explain
(also called the endogenous variable)

Independent variable: the variable used to explain the dependent variable
(also called the exogenous variable)

# Aims

- Describe the relationship between an independent variable X, and a continuous dependent variable Y as a straight line in $R^2$
  - Two Cases:
    - Fixed X: values of X are preselected by investigator
    - Variable X: a random sample of (X,Y) pairs
- Draw inferences regarding the relationship
- Predict the value of Y for a given X

# Simple Linear Regression Model

The population regression model:

Dependent Variable

Population Y intercept

Population Slope Coefficient

Independent Variable

Random Error term

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

Linear component

Random Error component

# Linear Regression Assumptions

- The true relationship form is linear (Y is a linear function of X, plus random error)

- The error terms, $\varepsilon_i$, are independent of the X values

- The error terms are random variables with mean 0 and constant variance, $\sigma^2$

  (the uniform variance property is called homoscedasticity)

$$E[\varepsilon_i] = 0 \quad \text{and} \quad E[\varepsilon_i^2] = \sigma^2 \quad \text{for } (i = 1, \ldots, n)$$

- The random error terms, $\varepsilon_i$, are not correlated with one another, so that

$$E[\varepsilon_i \varepsilon_j] = 0 \quad \text{for all } i \neq j$$

# Graphically (p 85)

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Figure 6.2: *Simple Linear Regression Model for Fixed X's*

# Simple Linear Regression Model

*(continued)*

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Y

Assignment Project Exam Help

https://powcoder.com

Observed Value
of Y for $x_i$

Add WeChat powcoder

$\varepsilon_i$

Slope = $\beta_1$

Predicted Value
of Y for $x_i$

Random Error for this
value

Intercept = $\beta_0$

$x_i$

X

# α and β (p 86)



Figure 6.3: Theoretical Regression Line Illustrating $\alpha$ and $\beta$

# Simple Linear Regression Equation

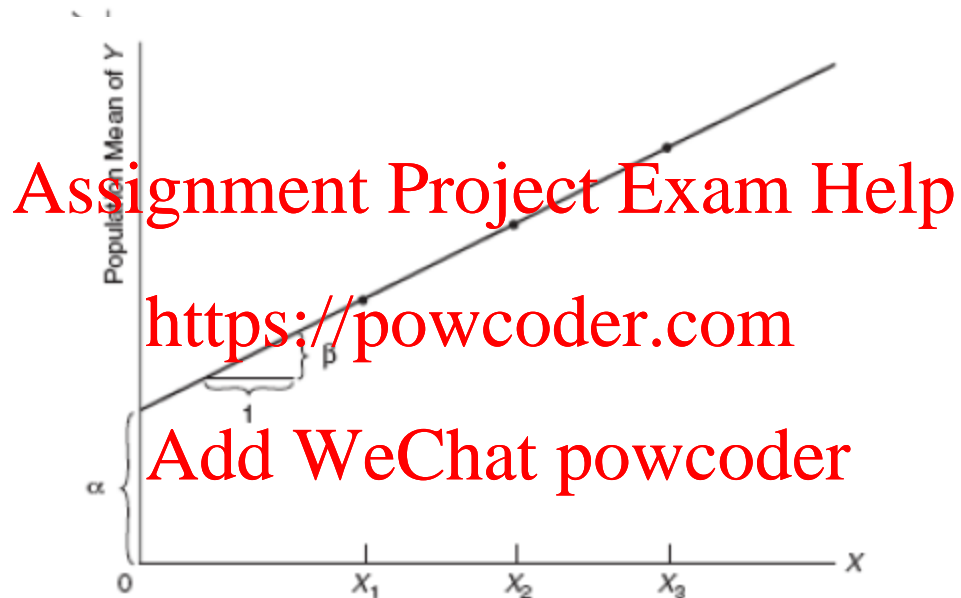The simple linear regression equation provides an estimate of the population regression line

Assignment Project Exam Help

Estimated (or predicted) y value for observation i

Estimate of the regression intercept

Estimate of the regression slope

https://powcoder.com

Add WeChat powcoder

Value of x for observation i

$$\hat{y}_i = b_0 + b_1 x_i$$

The individual random error terms $e_i$ have a mean of zero

$$e_i = (y_i - \hat{y}_i) = y_i - (b_0 + b_1 x_i)$$

# Least Squares Coefficient Estimators

- $b_0$ and $b_1$ are obtained by finding the values of $b_0$ and $b_1$ that minimize the sum of the squared residuals (errors), SSE:

$$\min \; SSE = \min \sum_{i=1}^{n} e_i^2$$

$$= \min \sum (y_i - \hat{y}_i)^2$$

$$= \min \sum [y_i - (b_0 + b_1 x_i)]^2$$

Differential calculus is used to obtain the coefficient estimators $b_0$ and $b_1$ that minimize SSE

# Prediction

- The regression equation can be used to predict a value for y, given a particular x

- For a specified value, $x_{n+1}$, the predicted value is

$$\hat{y}_{n+1} = b_0 + b_1 x_{n+1}$$

# Least Squares Coefficient Estimators

_(continued)_

- The slope coefficient estimator is

$$b_1 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} = \frac{Cov(x, y)}{s_x^2} = r\frac{s_y}{s_x}$$

Assignment Project Exam Help
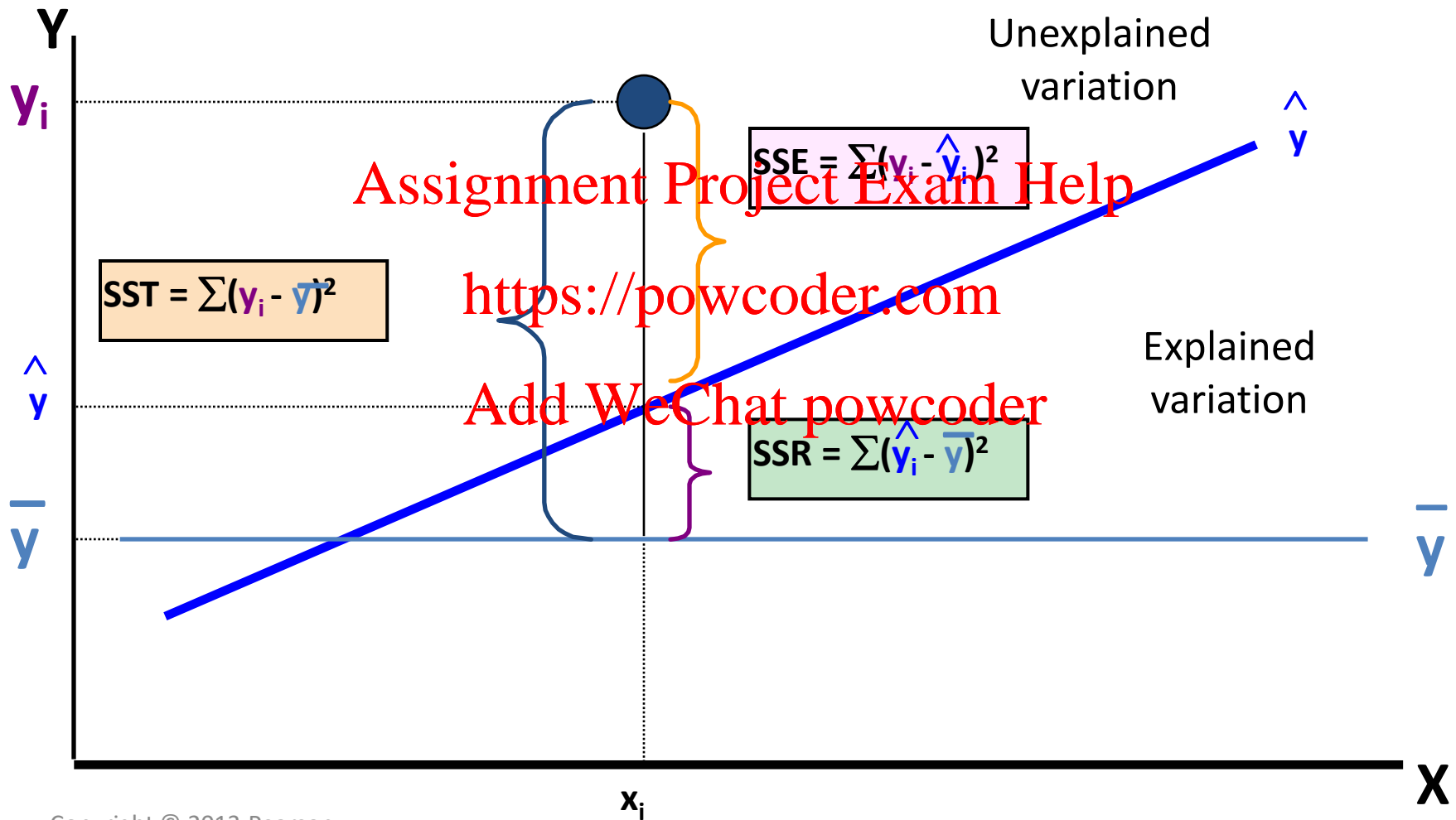
https://powcoder.com

Add WeChat powcoder

- And the constant or y-intercept is

$$b_0 = \bar{y} - b_1\bar{x}$$

- The regression line always goes through the mean $\bar{x}, \bar{y}$

# Analysis of Variance

$$SST = \sum(y_i - \bar{y})^2$$

$$SSE = \sum(y_i - \hat{y}_i)^2$$

$$SSR = \sum(\hat{y}_i - \bar{y})^2$$

Unexplained variation

Explained variation

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Explanatory Power of a Linear Regression Equation

11.4

- Total variation is made up of two parts:

$$SST = SSR + SSE$$

| **Total** Sum of Squares | **Regression** Sum of Squares | **Error** (residual) Sum of Squares |
|---|---|---|

$$SST = \sum (y_i - \bar{y})^2$$

$$SSR = \sum (\hat{y}_i - \bar{y})^2$$

$$SSE = \sum (y_i - \hat{y}_i)^2$$

where:

$\bar{y}$ = Average value of the dependent variable

$y_i$ = Observed values of the dependent variable

$\hat{y}_i$ = Predicted value of y for the given $x_i$ value

# Proof

$$\sum_{i=1}^{n}(Y_i - \bar{Y})^2 = \sum_{i=1}^{n}(\hat{Y}_i - \bar{Y} + Y_i - \hat{Y}_i)^2$$

$$= \sum_{i=1}^{n}\{(\hat{Y}_i - \bar{Y})^2 + (Y_i - \hat{Y}_i)^2 + 2(\hat{Y}_i - \bar{Y})(Y_i - \hat{Y}_i)\}$$

$$= SSR + SSE + 2\sum_{i=1}^{n}(\hat{Y}_i - \bar{Y})(Y_i - \hat{Y}_i)$$

$$= SSR + SSE + 2\sum_{i=1}^{n}(\hat{Y}_i - \bar{Y})e_i$$

$$= SSR + SSE + 2\sum_{i=1}^{n}(b_0 + b_1 X_i - \bar{Y})e_i$$

$$= SSR + SSE + 2b_0\sum_{i=1}^{n}e_i + 2b_1\sum_{i=1}^{n}X_i e_i - 2\bar{Y}\sum_{i=1}^{n}e_i$$

$$= SSR + SSE$$

# Hypothesis Test for Population Slope Using the F Distribution

- F Test statistic:

$$F = \frac{MSR}{MSE}$$

where

$$MSR = \frac{SSR}{k}$$

$$MSE = \frac{SSE}{n - k - 1}$$

where F follows an F distribution with  k  numerator  and (n – k - 1) denominator degrees of freedom

(k = the number of independent variables in the regression model)

# Computer Analysis

Results:

- estimates of slope ($\beta_1$) and intercept ($\beta_0$), using least squares

- residual mean square = estimate of variance ( $S^2$ )

- test if $\beta = \beta_0$

  - Usually, test $\beta = 0$, i.e. X has no effect on Y

# Hypothesis Test for Population Slope Using the F Distribution

*(continued)*

- An alternate test for the hypothesis that the slope is zero:

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

- Use the F statistic

$$F = \frac{MSR}{MSE} = \frac{SSR}{s_e^2}$$

- The decision rule is

$$\text{reject } H_0 \text{ if } F \geq F_{1,n-2,\alpha}$$

Copyright © 2013 Pearson Education, Inc. Publishing as Prentice Hall
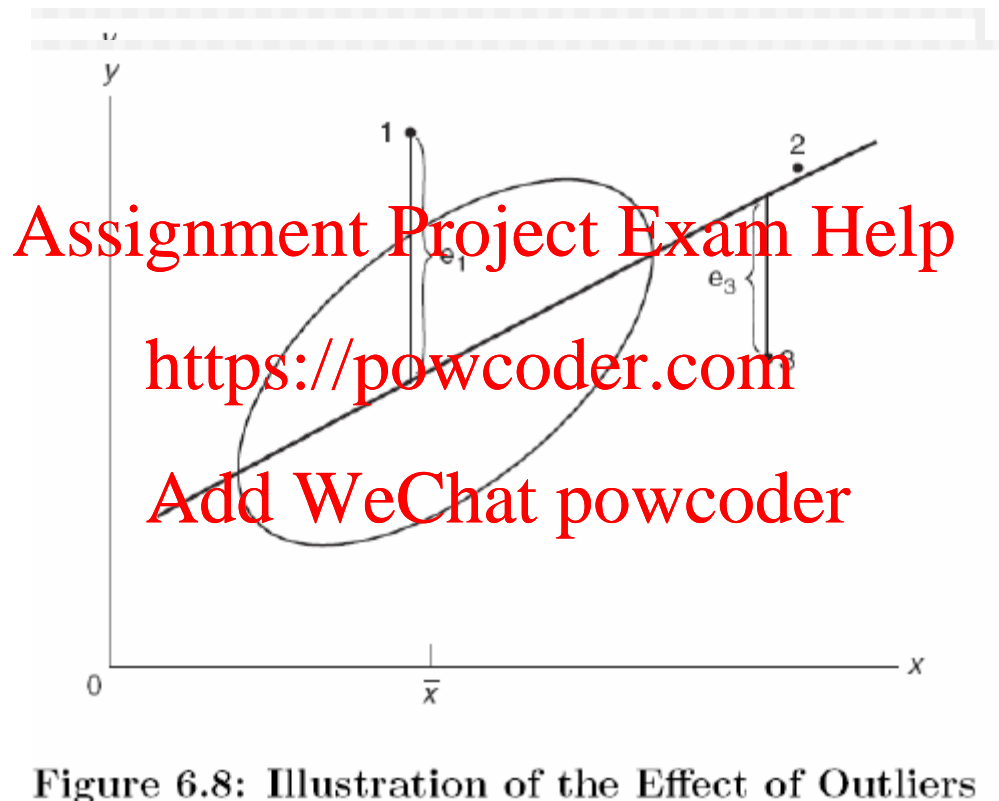
# Steps in Simple Regression

1.  State the research hypothesis.

2.  State the null hypothesis

3.  Gather the data

4.  Assess each variable separately first (obtain measures of central tendency and dispersion; frequency distributions; graphs); is the variable normally distributed?

5.  Calculate the regression equation from the data

6.  Calculate and examine appropriate measures of association and tests of statistical significance for each coefficient and for the equation as a whole

7.  Accept or reject the null hypothesis

8.  Reject or accept the research hypothesis

9.  Explain the practical implications of the findings

# Effect of Outliers (p 102)



Figure 6.8: Illustration of the Effect of Outliers

# Leverage

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2}$$

# Influence Measures

- Cook's distance:  "distance" between B with and without the $i^{th}$ observation

- DFFITS:  "distance" between Ŷ with and without the $i^{th}$ observation

# Cook's Distance

$$D_i = \frac{(y_i - \hat{y_i})^2}{(m+1)s^2} \frac{h_i}{(1 - h_i)^2}$$

# Influential observations

An observation is influential if:

– It is an outlier in X and Y

– Cook's distance > $F_{0.5}(P+1, N-P-1)$

– DFFITS > $\dfrac{2\sqrt{P+1}}{\sqrt{N-P-1}}$

Try analysis with and without influential observations and compare results.

# Explanatory Power of a Linear Regression Equation

- Total variation is made up of two parts:

$$SST = SSR + SSE$$

| Total Sum of Squares | Regression Sum of Squares | Error (residual) Sum of Squares |

$$SST = \sum (y_i - \bar{y})^2 \qquad SSR = \sum (\hat{y}_i - \bar{y})^2 \qquad SSE = \sum (y_i - \hat{y}_i)^2$$

where:

$\bar{y}$ = Average value of the dependent variable

$y_i$ = Observed values of the dependent variable

$\hat{y}_i$ = Predicted value of y for the given $x_i$ value

# Confidence & Prediction Intervals

- Confidence interval (CI) for mean of Y

- Prediction interval (PI) for individual Y

PI is wider than CI

# Confidence Interval for the Average Y, Given X

Confidence interval estimate for the
**expected value of y** given a particular $x_i$

Confidence interval for $E(Y_{n+1} \mid X_{n+1})$:

$$\hat{y}_{n+1} \pm t_{n-2,\alpha/2} s_e \sqrt{\left[ \frac{1}{n} + \frac{(x_{n+1} - \overline{x})^2}{\sum (x_i - \overline{x})^2} \right]}$$

Notice that the formula involves the term $(x_{n+1} - \overline{x})^2$

so the size of interval varies according to the distance $x_{n+1}$ is

from the mean, $\overline{x}$

# Prediction Interval for an Individual Y, Given X

Confidence interval estimate for an **actual observed value of y** given a particular $x_i$

Assignment Project Exam Help

Confidence interval for $\hat{y}_{n+1}$ :

https://powcoder.com

Add WeChat powcoder

$$\hat{y}_{n+1} \pm t_{n-2,\alpha/2} s_e \sqrt{\left[1 + \frac{1}{n} + \frac{(x_{n+1} - \overline{x})^2}{\sum (x_i - \overline{x})^2}\right]}$$

This extra term adds to the interval width to reflect the added uncertainty for an individual case

# Estimating Mean Values and Predicting Individual Values

Goal: Form intervals around y to express uncertainty about the value of y for a given $x_i$

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Confidence Interval for the expected value of y, given $x_i$

$\hat{y} = b_0 + b_1 x_i$

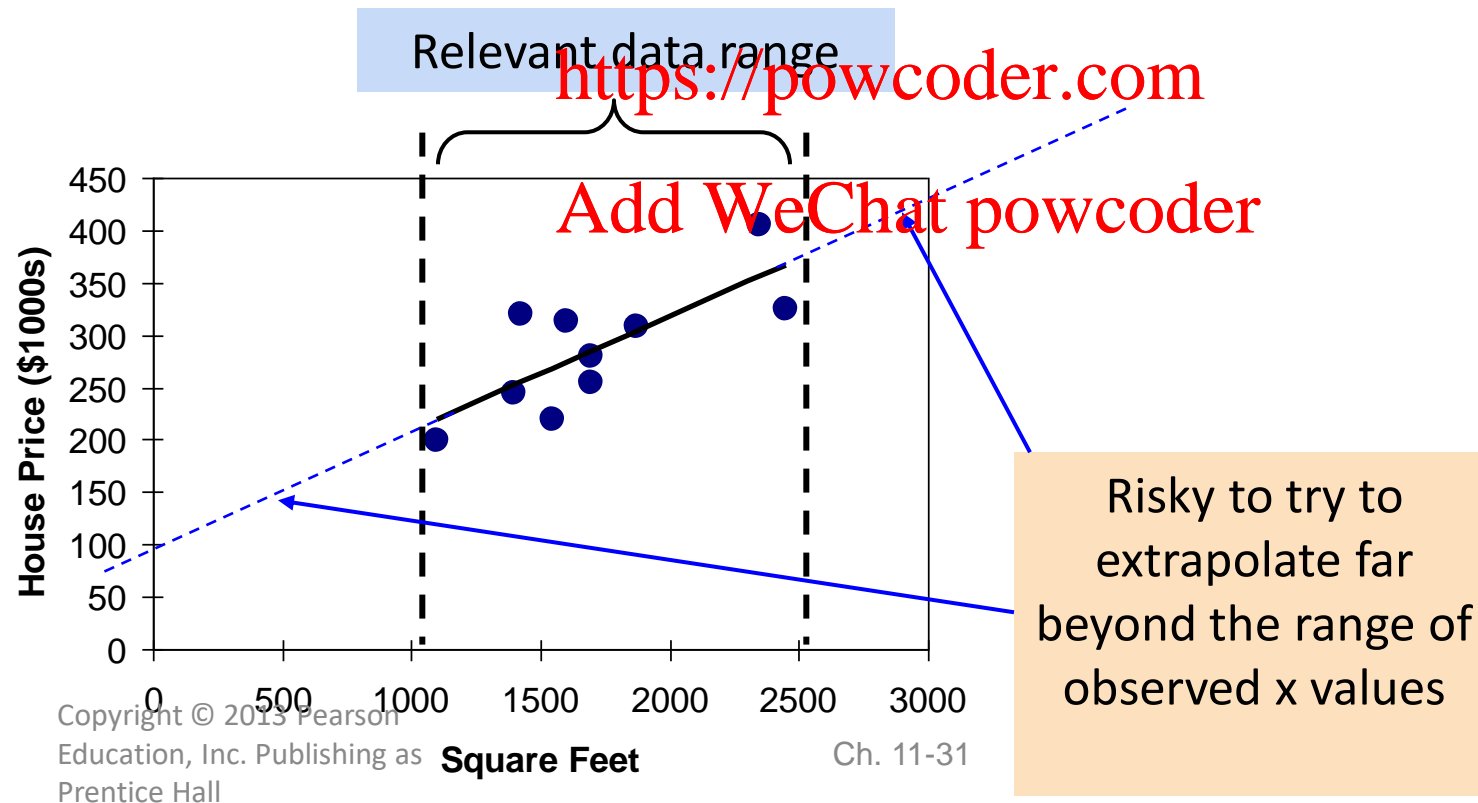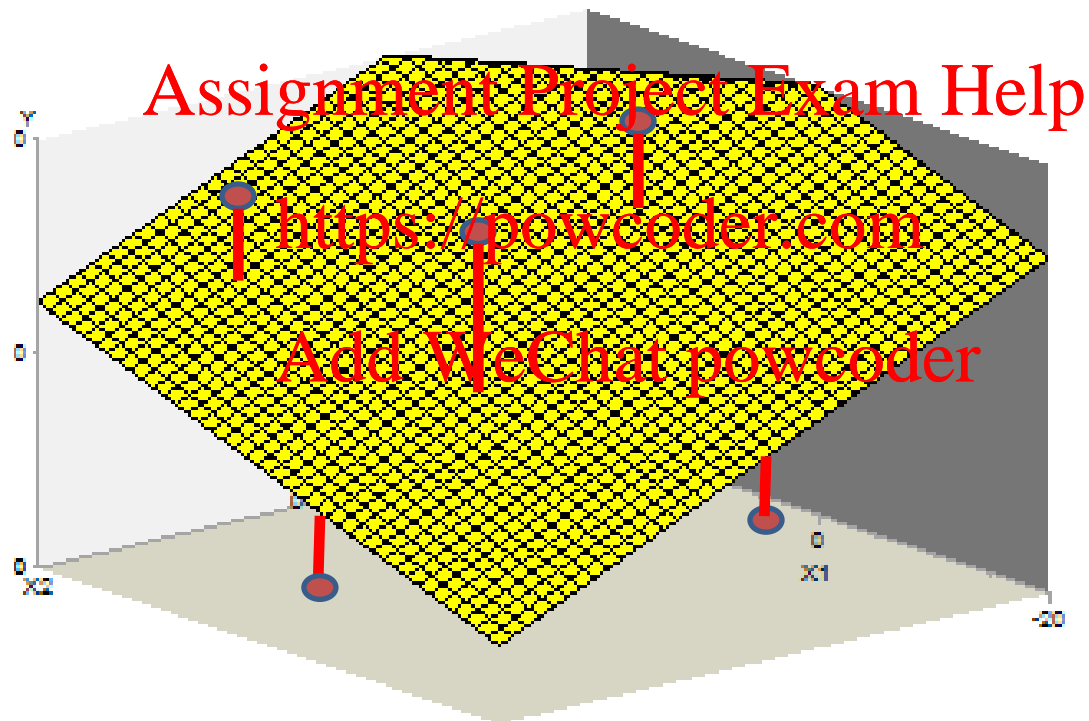Prediction Interval for an single observed y, given $x_i$

Y

$\hat{y}$

$x_i$

X

# Relevant Data Range

- When using a regression model for prediction, only predict within the relevant range of data

Assignment Project Exam Help
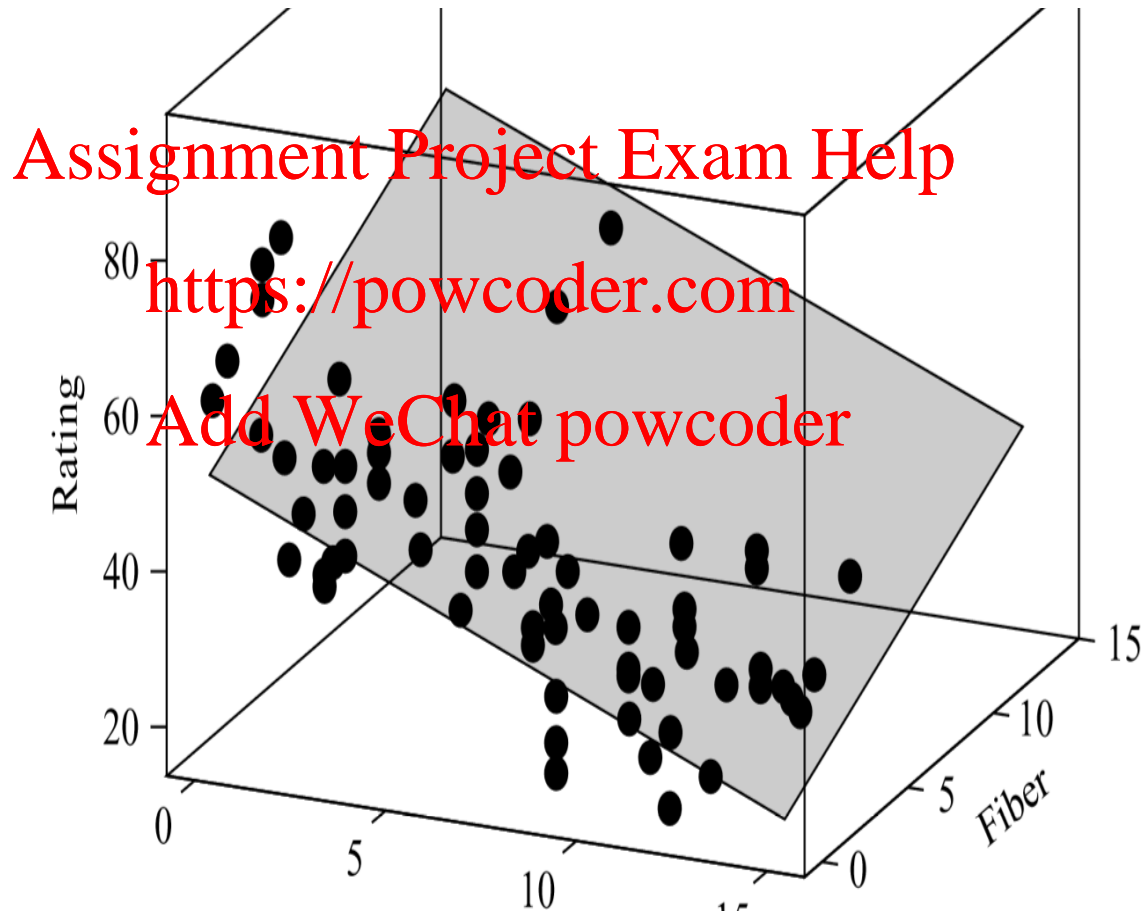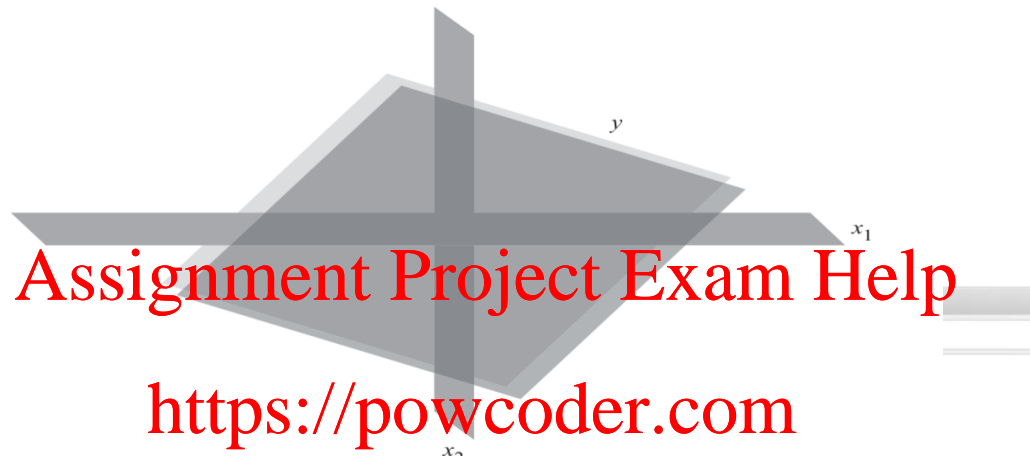
Relevant data range

https://powcoder.com

Add WeChat powcoder

Risky to try to extrapolate far beyond the range of observed x values

# Multiple Regression



Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Multiple Regression

# Adjusted R-Sqr

$$R^2_{\text{adj}} = 1 - (1 - R^2)\frac{n - 1}{n - m - 1}$$

# VIF

# VIF

$$VIF_i = \frac{1}{1 - R_i^2}$$

20

(indicates deviation, or residual

0    5    10    X

Data Example illustrating computation of output deviation...

The number $N - 2$, called the **residual degrees of freedom**, is the sample size minus the number of parameters in the line (in this case, $\alpha$ and $\beta$). Using $N - 2$ as a divisor in computing $S^2$ produces an **unbiased estimate** of $\sigma^2$. In this example,

$$RES\ MS = \frac{17...}{...}$$

The square root of the residual mean square is called the **standard error of the estimate** and is denoted by $S$.

Software regression programs will also produce the **standard errors** of $A$ and $B$. These statistics are computed as

$$SE(A) = S\left[\frac{1}{N} + \frac{\bar{X}^2}{\Sigma(X - \bar{X})^2}\right]^{1/2}$$

$$SE(B) = \frac{S}{[\Sigma(X - \bar{X})^2]^{1/2}}$$

# Correlation Coefficient - ρ

- Correlation coefficient measures the strength of linear association between X and Y in the population (ρ).

- it is estimated by sample ( r )

# Correlation Analysis

- Correlation analysis is used to measure strength of the association (linear relationship) between two variables

  - Correlation is only concerned with strength of the relationship

  - No causal effect is implied with correlation

  - Correlation was first presented in Chapter 4

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Correlation Analysis

- The population correlation coefficient is denoted $\rho$ (the Greek letter rho)

- The sample correlation coefficient is

$$r = \frac{s_{xy}}{s_x s_y}$$

where

$$s_{xy} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

# Calculating the value of ρ

- 100 $(1 - \rho^2)^{1/2}$ = % of Standard Deviation NOT "explained" by X

$$\sigma^2 = \sigma_y{}^2 (1 - \rho^2)$$

$$\Rightarrow \sigma = \sigma_y \sqrt{1 - \rho^2}$$

$$\Rightarrow \rho^2 = \frac{\sigma_y{}^2 - \sigma^2}{\sigma_y{}^2}$$

# Graphically (p 92)



Figure 6.5: Ellipses of Concentration for Various $\rho$ Values

# Calculating the value of ρ

- 100 $(1 - \rho^2)^{1/2}$ = % of Standard Deviation NOT "explained" by X
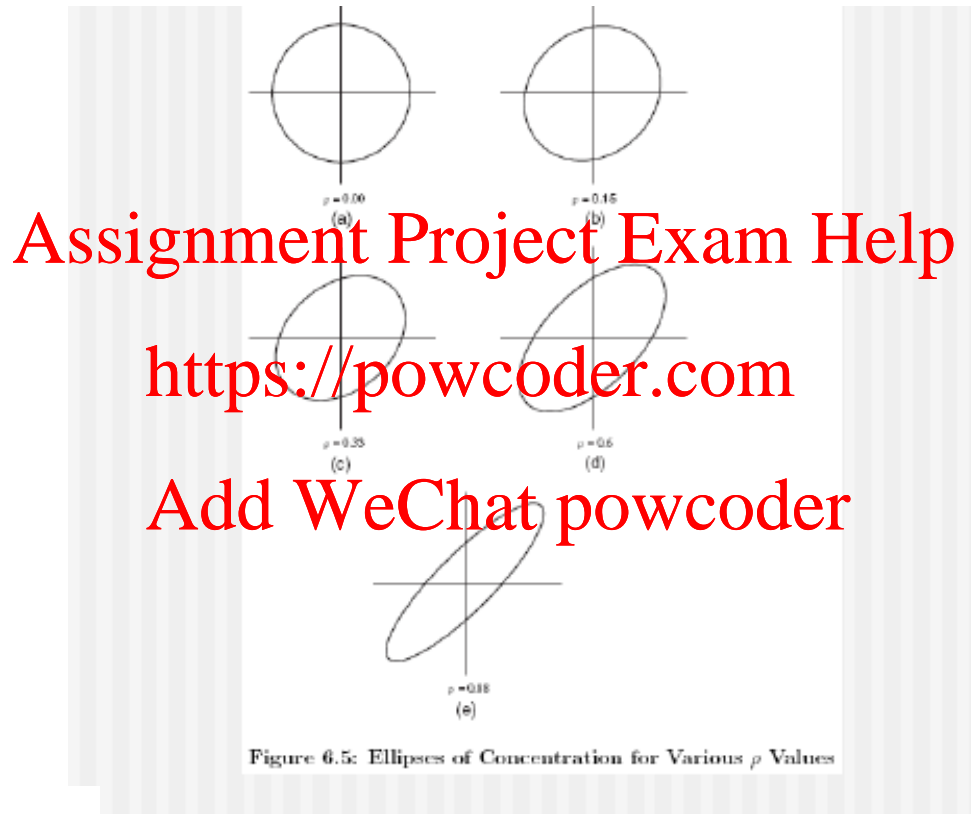
$$\sigma^2 = \sigma_y^2 (1 - \rho^2)$$

$$\Rightarrow \ \sigma = \sigma_y \sqrt{1 - \rho^2}$$

$$\Rightarrow \rho^2 = \frac{\sigma_y^2 - \sigma^2}{\sigma_y^2}$$

# Interpretation of ρ

- $\rho^2$ = reduction in variance of Y associated with knowledge of X/original variance of Y

- $100\rho^2$ = % of variance of Y "explained by X"

Caveat:  correlation vs causation

# Estimating the value of ρ
# (Pearson's Correlation Coefficient)

$$\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

$$r = \frac{S_{XY}}{S_X S_Y}$$

$$S_{XY} = \sum (X - m(X))(Y - m(Y))/(N-1)$$

# Interpretation of ρ

| ρ | % of variance "explained" | % of variance not "explained" | % of SD "explained" | % of SD not "explained" |
|---|---|---|---|---|
| ±0.3 | 9% | 91% | 5% | 95% |
| ±0.5 | 25% | 75% | 13% | 87% |
| ±0.71 | 50% | 50% | 29% | 71% |
| ±0.95 | 90% | 10% | 69% | 31% |

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Test for Zero Population Correlation

- To test the null hypothesis of no linear association,

$$H_0 : \rho = 0$$

the test statistic follows the Student's t distribution with $(n-2)$ degrees of freedom:

$$t = \frac{r\sqrt{(n-2)}}{\sqrt{(1-r^2)}}$$

# Example from Text:  Lung Function

- Data from an epidemiological study of households
    - living in four areas with different amounts and types of air pollution (Appendix A)
- Data only on non-smoking fathers
    - X = height in inches
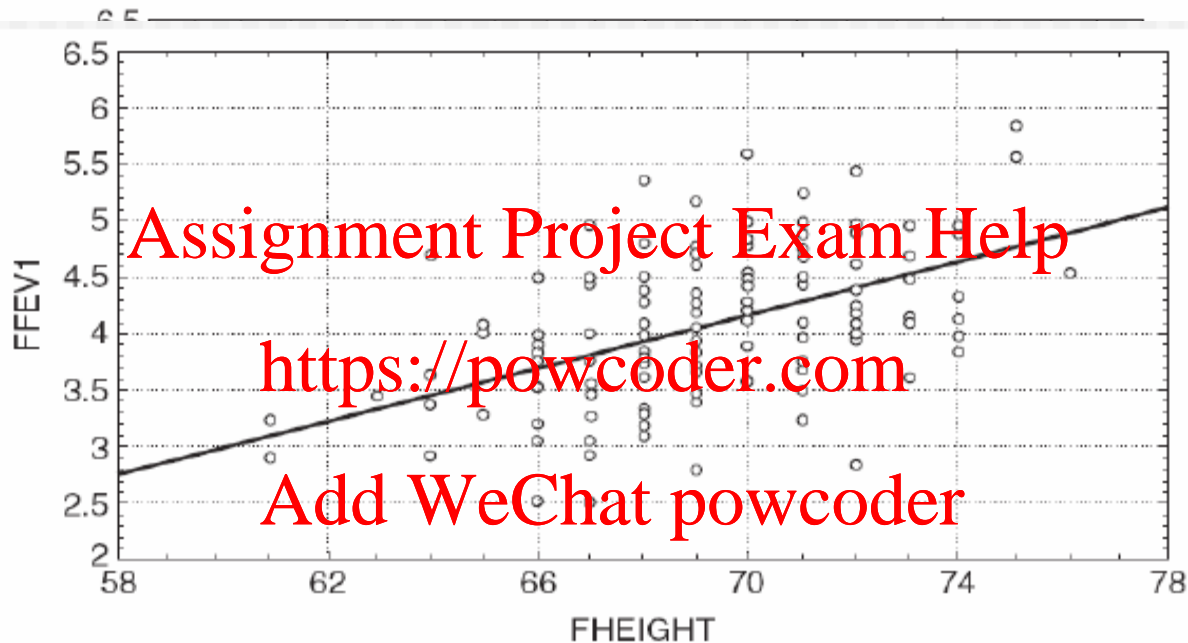    - Y = forced expiratory volume in 1 second (FEV1)

# Scatter Plot (p 83)



Figure 6.1: Scatter Diagram and Regression Line of FEV1 Versus Height for Fathers

# Example Results

- Least Squares Equation:  $Y = -4.087 + 0.118X$

- Correlation $r = 0.504$

- Test p = 0,

    - t = 7.1 (p 94),  $\rho < 0.0001$

    - t test can be one or two sided

# Analysis of Variance

- SST = total sum of squares
  - Measures the variation of the $y_i$ values around their mean, $\bar{y}$

- SSR = regression sum of squares
  - Explained variation attributable to the linear relationship between x and y

- SSE = error sum of squares
  - Variation attributable to factors other than the linear relationship between x and y

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Explanatory Power of a Linear Regression Equation

- Total variation is made up of two parts:

$$SST = SSR + SSE$$

| Total Sum of Squares | Regression Sum of Squares | Error (residual) Sum of Squares |
|---|---|---|

$$SST = \sum (y_i - \bar{y})^2 \qquad SSR = \sum (\hat{y}_i - \bar{y})^2 \qquad SSE = \sum (y_i - \hat{y}_i)^2$$
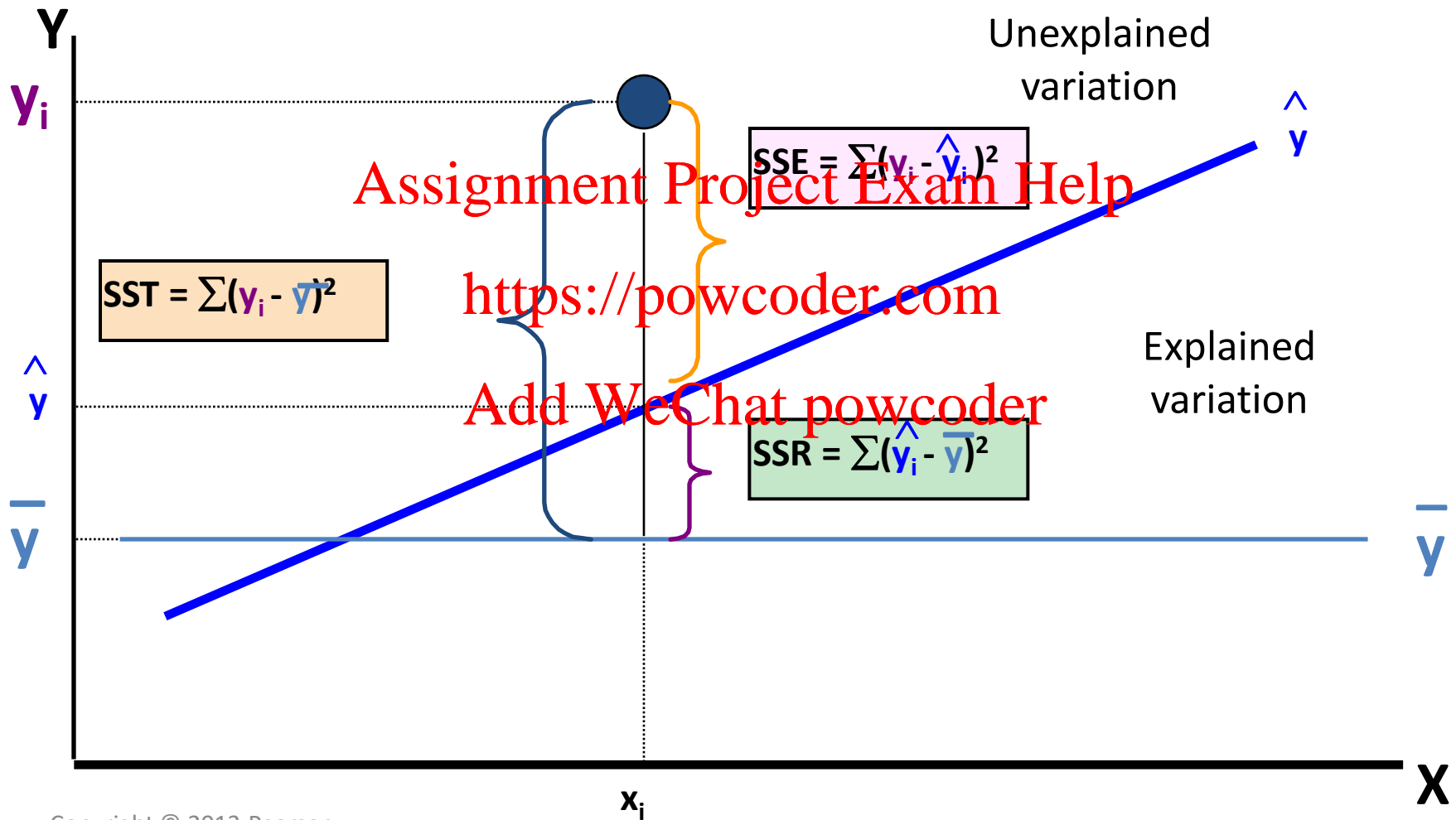
where:

$\bar{y}$ = Average value of the dependent variable

$y_i$ = Observed values of the dependent variable

$\hat{y}_i$ = Predicted value of y for the given $x_i$ value

# Analysis of Variance

*(continued)*



$$SST = \sum(y_i - \bar{y})^2$$

$$SSE = \sum(y_i - \hat{y}_i)^2$$

$$SSR = \sum(\hat{y}_i - \bar{y})^2$$

Unexplained variation

Explained variation

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Coefficient of Determination, $R^2$

- The coefficient of determination is the portion of the total variation in the dependent variable that is explained by variation in the independent variable

- The coefficient of determination is also called R-squared and is denoted as $R^2$

$$R^2 = \frac{SSR}{SST} = \frac{\text{regression sum of squares}}{\text{total sum of squares}}$$

note: $\boxed{0 \leq R^2 \leq 1}$
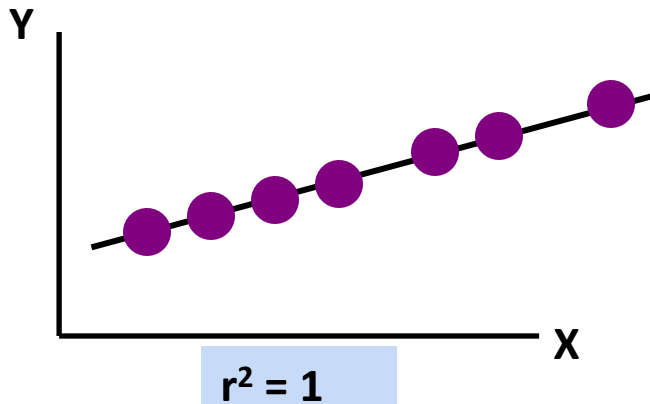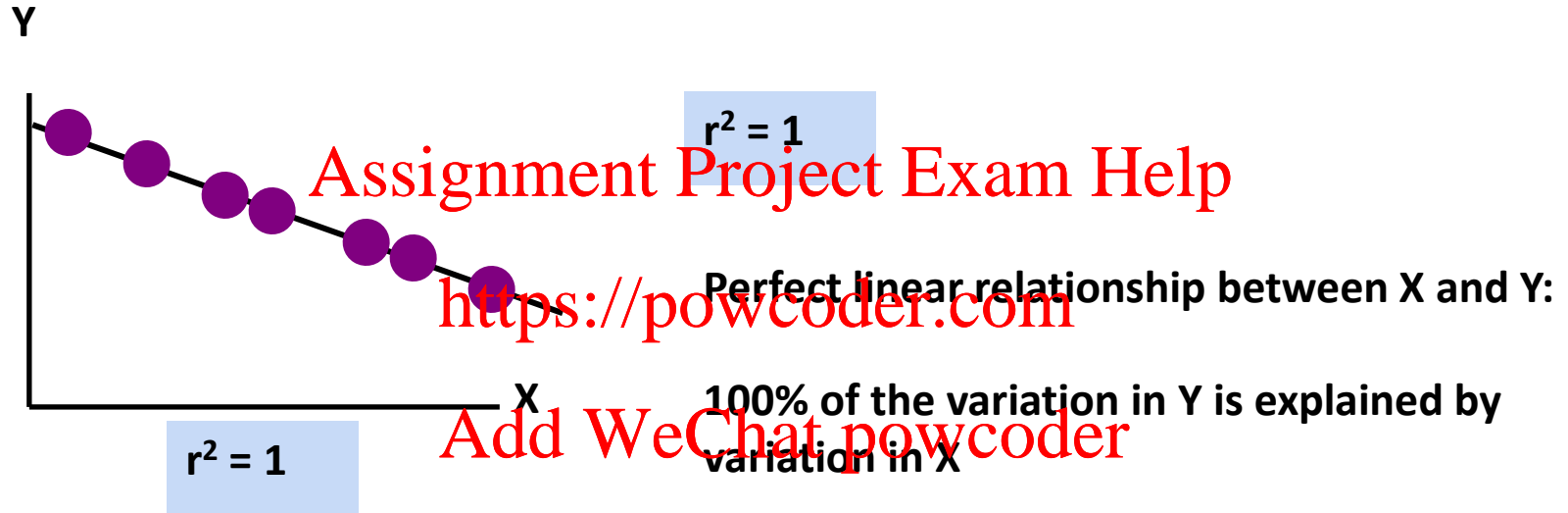
# Correlation and $R^2$

- The coefficient of determination, $R^2$, for a simple regression is equal to the simple correlation squared
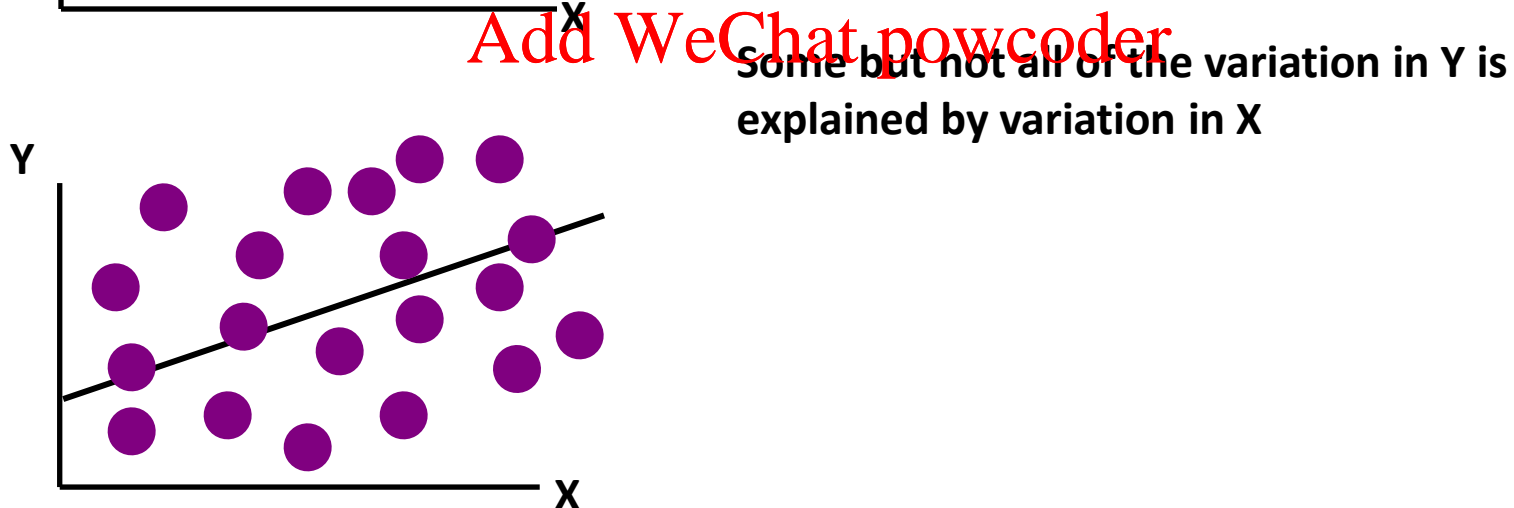
$$R^2 = r^2$$

# Examples of Approximate r² Values

Y

r² = 1

Assignment Project Exam Help

https://powcoder.com

Perfect linear relationship between X and Y:

X

Add WeChat powcoder

100% of the variation in Y is explained by variation in X

r² = 1

Y

X

r² = 1

# Examples of Approximate r² Values



**0 < r² < 1**

Weaker linear relationships between X and Y:

Some but not all of the variation in Y is explained by variation in X

# Examples of Approximate r² Values



**Y**

r² = 0

r² = 0

**X**

No linear relationship between X and Y:

The value of Y does not depend on X. (None of the variation in Y is explained by variation in X)

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Estimation of Model Error Variance

- An estimator for the variance of the population model error is

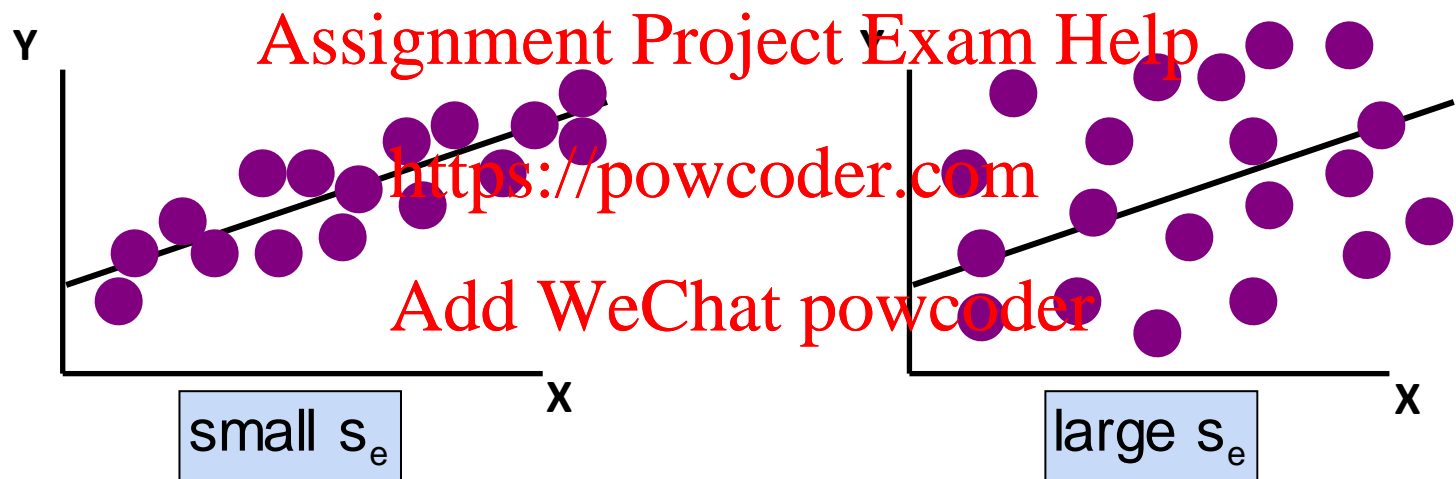$$\hat{\sigma}^2 = s_e^2 = \frac{\sum_{i=1}^{n} e_i^2}{n-2} = \frac{SSE}{n-2}$$

- Division by n – 2 instead of n – 1 is because the simple regression model uses two estimated parameters, $b_0$ and $b_1$, instead of one

is called the standard error of the estimate

$$s_e = \sqrt{s_e^2}$$

# Comparing Standard Errors

$s_e$ is a measure of the variation of observed y values from the regression line



small $s_e$

large $s_e$

The magnitude of $s_e$ should always be judged relative to the size of the y values in the sample data

# Statistical Inference: Hypothesis Tests and Confidence Intervals

- The variance of the regression slope coefficient ($b_1$) is estimated by

$$s^2_{b_1} = \frac{s^2_e}{\sum(x_i - \bar{x})^2} = \frac{s^2_e}{(n-1)s^2_x}$$

where:

$s_{b_1}$ = Estimate of the standard error of the least squares slope

$s_e = \sqrt{\dfrac{SSE}{n-2}}$ = Standard error of the estimate

# Example Results

- Least Squares Equation:  Y = -4.087 + 0.118X

- Correlation r = 0.504

- Test p = 0,

  – t = 7.1 (p 94),  ρ < 0.0001

  – t test can be one or two sided

# Hypothesis Test for Population Slope Using the F Distribution

- F Test statistic:

$$F = \frac{MSR}{MSE}$$

where

Assignment Project Exam Help

https://powcoder.com

$$MSR = \frac{SSR}{k}$$

Add WeChat powcoder

$$MSE = \frac{SSE}{n - k - 1}$$

where F follows an F distribution with  k  numerator  and (n – k - 1) denominator degrees of freedom

(k = the number of independent variables in the regression model)

# Hypothesis Test for Population Slope Using the F Distribution

- An alternate test for the hypothesis that the slope is zero:

Assignment Project Exam Help

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

https://powcoder.com

- Use the F statistic

Add WeChat powcoder

$$F = \frac{MSR}{MSE} = \frac{SSR}{s_e^2}$$

- The decision rule is

$$\text{reject } H_0 \text{ if } F \geq F_{1,n-2,\alpha}$$

# ANOVA Overview

**Table 6.1: ANOVA table for simple linear regression**

| Source of variation | Sums of squares | df | Mean square | F |
|---|---|---|---|---|
| Regression | $\sum(\hat{Y} - \bar{Y})^2$ | 1 | $SS_{reg}/1$ | $MS_{reg}/MS_{res}$ |
| Residual | $\sum(Y - \hat{Y})^2$ | $N - 2$ | $SS_{res}/(N-2)$ | |
| Total | $\sum(Y - \bar{Y})^2$ | $N - 1$ | | |

**Table 6.2: ANOVA example from Figure 6.1**

| Source of variation | Sums of squares | df | Mean square | F |
|---|---|---|---|---|
| Regression | 16.0532 | 1 | 16.0532 | 50.50 |
| Residual | 47.0451 | 148 | 0.3179 | |
| Total | 63.0983 | 149 | | |

# Test β = 0

- From ANOVA table:  F – 50.5

  – Gives 2-sided test, p-value < 0.0001

- One sided test is:  t = F$^{1/2}$ = 7.1

Same as test for ρ = 0

# Outliers

- Outlier in Y is studentized (or deleted studentized) residual >2

- Leverage = h $= \frac{1}{N} + \frac{(X - \bar{X})^2}{\sum (X - \bar{X})^2}$

  - X's far from the mean of X have large leverage (h)

  - Observations with large leverage have large effect on the slope of the line.

- Outlier in X if h > 4/N

# Residual Analysis

- Residual $= e = Y - \hat{Y}$

- Studentized residual $= e/S(1 - h)^{1/2}$

  - h called "leverage"

- Deleted studentized residual = studentized residual with observation for computing regression and S deleted.

# Influential observations

An observation is influential if:

- It is an outlier in X and Y
- Cook's distance > $F_{0.5}(2,N-2)$
- DFFITS > $\dfrac{2\sqrt{2}}{\sqrt{N-2}}$

Try analysis with and without influential observations and compare results.

# Observations

- Point 1 is an outlier in Y with low leverage
  - impacts estimate of intercept but not slope
  - Tends to increase the estimates of S & SE of B

- Point 2 has high leverage; not an outlier in Y
  - doesn't impact estimate of B or A

- Point 3 has high leverage and is an outlier in Y
  - impacts the values of B, A, and S

70

# Assumptions

- Homogeneity of variance (same $\sigma^2$)
  - Not extremely serious
  - Can be achieved through transformations if necessary
- Normal residuals
  - Slight departures ok
  - Can use transformations to achieve it
- Randomness
  - Serious
  - Can use hierarchical models for clustered samples

# Checking Assumptions

- Plot residuals vs X or vs the predicted Y to check linearity and homogeneity of variance

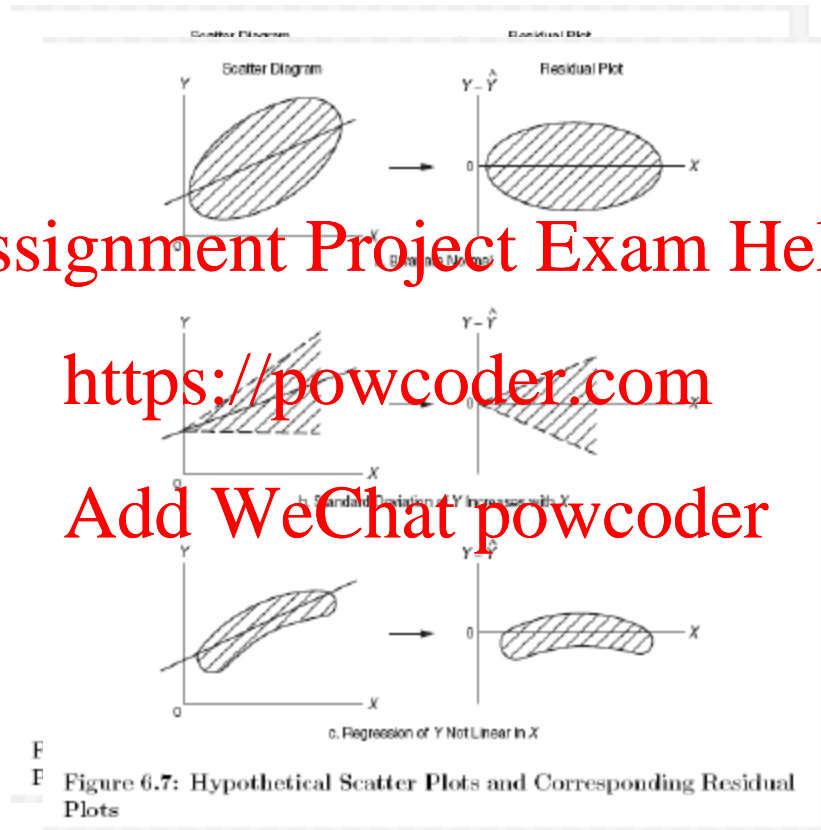- Create normal probability plots of residuals to check for normality

# Residual Plots (p 98)



Scatter Diagram

Residual Plot

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

c. Regression of Y Not Linear in X

Figure 6.7: Hypothetical Scatter Plots and Corresponding Residual Plots
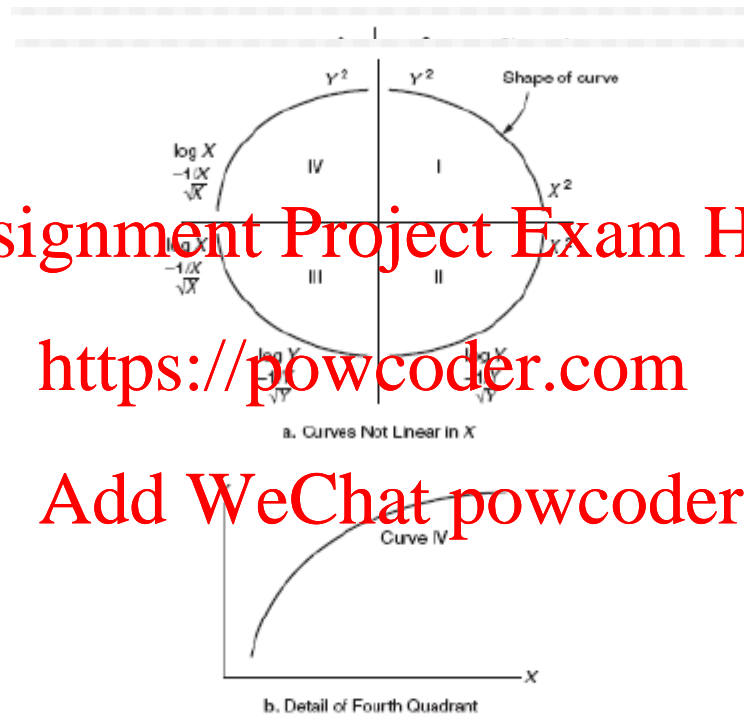
# Transformations (p 105 )



Figure 6.9:  Choice of Transformation: Typical Curves and Appropriate Transformation

# Weighted Regression

- If $\sigma^2$ are not equal, use weight for each residual in the sum of squares used in Least Squares process.

- Weight = $1/\sigma^2$

- Gives unbiased estimate with smaller variance

# Weighted Regression - Caveat

- Solution,, standardize weight (w) to add up to the sample size (N)

  - e.g. N = 5, w = 4,1,8,2,4, sum of w = 19

  - define standardized weight (sw) = w*5/19

  - sum of sw = 5

  - = 1.05 + .26 + 2.11 + .53 + 1.05 = 5

# What to watch for

- Need representative sample
- Range of prediction should match observed range in X in sample
- Use of nominal or ordinal, rather than interval or ration data
- Errors in variables
- Correlation does not imply causation
- Violation of assumptions
- Influential points
- Appropriate model

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Multiple Linear Regression

# Keywords for OUTPUT Statement

| Keyword | Description |
|---|---|
| COOKD=*names* | Cook's influence statistic |
| COVRATIO=*names* | standard influence of observation on covariance of betas |
| DFFITS=*names* | standard influence of observation on predicted value |
| H=*names* | leverage, |
| LCL=*names* | lower bound of a % confidence interval for an individual prediction. This includes the variance of the error, as well as the variance of the parameter estimates. |
| LCLM=*names* | lower bound of a % confidence interval for the expected value (mean) of the dependent variable |
| PREDICTED \| P=*names* | predicted values |
| PRESS=*names* | th residual divided by , where is the leverage, and where the model has been refit without the th observation. |
| RESIDUAL \| R=*names* | residuals, calculated as ACTUAL minus PREDICTED |
| RSTUDENT=*names* | a studentized residual with the current observation deleted |
| STDI=*names* | standard error of the individual predicted value |
| STDP=*names* | standard error of the mean predicted value |
| STDR=*names* | standard error of the residual |
| STUDENT=*names* | studentized residuals, which are the residuals divided by their standard errors |
| UCL=*names* | upper bound of a % confidence interval for an individual prediction |
| UCLM=*names* | upper bound of a % confidence interval for the expected value (mean) of the dependent variable |

# Aims

- Extend simple linear regression to multiple dependent variables.

- Describe a linear relationship between:
  - A single continuous Y variable, and
  - Several X variables

- Draw inferences regarding the relationship

- Predict the value of Y from $X_1$, $X_2$, …, $X_p$.

- Research Questions:  To what extent does some combination of the IVs predict the DV?

- E.g. To what extent does age, gender, type/amount of food consumption predict low density lipid level.

# Assumptions

- Level of Measurement:
  - IVs – two or more, Continuous or dichotomous
  - DV - continuous
- Sample Size – Enough cases per IV
- Linearity:  Are bivariate relationships linear
- Constant Variance (about line of best fit) – Homoscedasticity
- Multicollinearity:  Between the IVs
- Multivariate outliers
- Normality of residuals about predicted value

# Approaches

- Direct:  All IVs entered simultaneously
- Forward:  IVs entered one by one until there are no significant IVs to be entered.
- Backward: IVs removed one by one until there are no significant IVs to be removed.
- Stepwise:  Combination of Forward and Backward
- Hierarchical:  IVs entered in steps.

# Write ups

- Assumptions:  How tested, extent met

- Correlations:  What are they, what conclusions

- Regression coefficients: Report and interpret

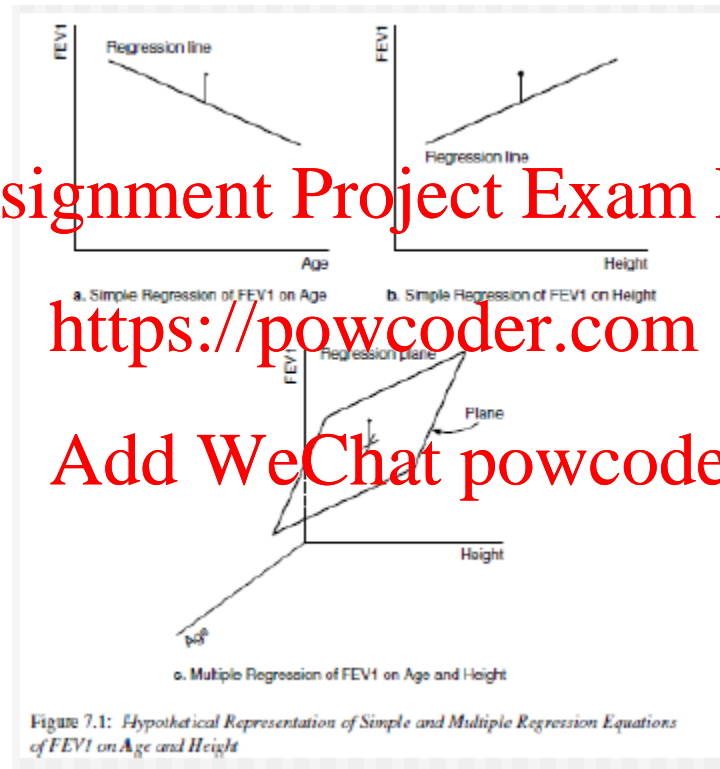- Conclusions and Caveats

# Steps in Multiple Regression

1. State the research hypothesis.
2. State the null hypothesis
3. Gather the data
4. Assess each variable separately first (obtain measures of central tendency and dispersion; frequency distributions; graphs); is the variable normally distributed?
5. Assess the relationship of each independent variable, one at a time, with the dependent variable (calculate the correlation coefficient; obtain a scatter plot); are the two variables linearly related?
6. Assess the relationships between all of the independent variables with each other (obtain a correlation coefficient matrix for all the independent variables); are the independent variables too highly correlated with one another?
7. Calculate the regression equation from the data
8. Calculate and examine appropriate measures of association and tests of statistical significance for each coefficient and for the equation as a whole
9. Accept or reject the null hypothesis
10. Reject or accept the research hypothesis
11. Explain the practical implications of the findings

# Example (p 121)



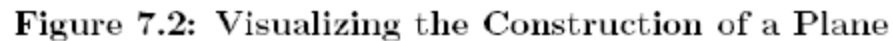Figure 7.1: *Hypothetical Representation of Simple and Multiple Regression Equations of FEV1 on Age and Height*

# Example (p 122)



Figure 7.2: Visualizing the Construction of a Plane

# Mathematical Model

- The mean of Y values at a given X is:

$$\alpha + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_p X_p$$

- Variance of Y values at any set of X's is $\sigma^2$

  (For all X)

- Y values are normally distributed at each X

  (needed for inference)

# Types of X (independent) variables

- Fixed: selected in advance

- Variable: as in most studies

- X's can be continuous or discrete (categorical)

- X's can be transformations of other X's, e.g., polynomial regression.

# Computer Analysis

- Estimates of: $\alpha$, $\beta_1$, $\beta_2$, …, $\beta_p$ using least-squares.

- Residual mean square ( $S^2$ ) is estimate of variance $\sigma^2$

- Confidence intervals for mean of Y

- Prediction intervals for individual Y

# Example of Bonferroni

- Test 3 hypotheses
- P-values are:  0.014, 0.036, 0.075
- Let nominal significance level = 0.15
  - $\therefore$ first 2 are significant
- Bonferroni Adjusted p-values   multiply by 3, giving: 0.042, 0.108, 0.225
  - Only first is significant
  - Probablility of at rejecting at least 1 out of m hypotheses

$$FWER = Pr\left\{\bigcup_{i_o}(p_i \leq \frac{\alpha}{m})\right\} \leq \sum_{i_o}\{Pr(p_i \leq \frac{\alpha}{m})\} \leq m_0\frac{\alpha}{m} \leq m\frac{\alpha}{m} = \alpha$$

# Analysis of variance (p 132)

- Does regression plane help in predicting values of Y?

- Test hypothesis that all $\beta_I$'s = 0

Table 7.1: ANOVA Table for multiple regression

| Source of variation | Sums of squares | df | Mean square | $F$ |
|---|---|---|---|---|
| Regression | $\sum(\hat{Y}-\bar{Y})^2$ | $P$ | $\text{SS}_{\text{reg}}/P$ | $\text{MS}_{\text{reg}}/\text{MS}_{\text{res}}$ |
| Residual | $\sum(Y-\hat{Y})^2$ | $N-P-1$ | $\text{SS}_{\text{reg}}/(N-P-1)$ | |
| Total | $\sum(Y-\bar{Y})^2$ | $N-1$ | | |

# Example: Reg of FEV1 on height and weight (p 132)

Table 7.2: ANOVA example from the lung function data (fathers)

| Source of variation | Sums of squares | df | Mean square | $F$ |
|---|---|---|---|---|
| Regression | 21.0570 | 2 | 10.5285 | 36.81 |
| Residual | 42.0413 | 147 | 0.2860 | |
| Total | 63.0983 | | | |

- $F = 36.81$; df = 2, 147; p-value <0.0001
- Use percentile link from web site: http://faculty.vassar.edu/lowry/tabs.html#f

# Venn Diagrams

- Multiple $R^2$
- Bivariate Correlation between IV1 and DV
- Bivariate Correlation between IV2 and DV
- Correlation between IV1 and IV2

- Target:  IV's that highly correlate with the DV, but don't highly correlate with each other

DV

IV2

IV1

# Correlation Coefficient

- The multiple correlation coefficient (R) measures the strength of association between Y, and the set of X's in the population.

- It is estimated as the simple correlation coefficient between the Y's and their predicted values ( Ŷ's )

# Coefficient of Determination

- $R^2$ = Coefficient of determination

  = SS due to regression/SS total

- $R^2$ = (reduction in variance of Y due to X's) / (original variance of Y).

- Therefore $100R^2$ = % of variance of Y "explained by X's".

- And $100(1 - \rho^2)^{1/2}$ = % of Standard Deviation NOT "explained" by X's

# Regression

# Standard Deviation of bet1

$$\sigma_{b_1} = \frac{\sigma}{\sqrt{\sum x^2 - (\sum x)^2 / n}}$$

# Confidence Interval Mean Value

**CONFIDENCE INTERVAL FOR THE MEAN VALUE OF *y* FOR A GIVEN VALUE OF *x***

$$\hat{y}_p \pm t_{n-2}(s)\sqrt{\frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum(x_i - \bar{x})^2}}$$

where $\hat{y}_p$ is the point estimate of *y* for a particular value of *x*, *t* a multiplier associated with the sample size and confidence level, *s* the standard error of the estimate, and $x_p$ the particular value of *x* for which the prediction is being made.

# Confidence Interval for prediction

**PREDICTION INTERVAL FOR A RANDOMLY CHOSEN VALUE OF *y* FOR A GIVEN VALUE OF *x***

$$\hat{y}_p \ \pm \ t_{n-2}(s)\sqrt{1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

# Adjusted R-square

$$R^2_{\text{adj}} = 1 - (1 - R^2)\frac{n-1}{n-m-1}$$

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder



$y$

$x_1$

$x_2$

**Figure 3.8**    When the predictors $x_1$ and $x_2$ are uncorrelated, the response surface $y$ rests on a solid basis, providing stable coefficient estimates.

101

$x_1$

$x_2$

$y$

**Figure 3.9**

# Sequential SS vs. Partial SS

**TABLE 3.14  Difference Between Sequential and Partial SS**

| Variable | Sequential SS | Partial SS |
|----------|---------------|------------|
| $x_1$ | $SS(x_1)$ | $SS(x_1 \mid x_2, x_3, x_4)$ |
| $x_2$ | $SS(x_2 \mid x_1)$ | $SS(x_2 \mid x_1, x_3, x_4)$ |
| $x_3$ | $SS(x_3 \mid x_1, x_2)$ | $SS(x_3 \mid x_1, x_2, x_4)$ |
| $x_4$ | $SS(x_4 \mid x_1, x_2, x_3)$ | $SS(x_4 \mid x_1, x_2, x_3)$ |

Assignment Project Exam Help

https://powcoder.com

**120**    **CHAPTER 3   MULTIPLE REGRESSION AND MODEL BUILDING**

Add WeChat powcoder

What effect do these changes in $\text{VIF}_i$ have on $s_{b_i}$, the variability of the $i$th coefficient? We have

$$s_{b_i} = s c_i = s \sqrt{\frac{1}{(n-1)\,s_i^2} \frac{1}{1 - R_i^2}} = s \sqrt{\frac{VIF_i}{(n-1)\,s_i^2}}$$

If $x_i$ is uncorrelated with the other predictors, $\text{VIF}_i = 1$, and the standard error of the coefficient $s$ will not be inflated. However, if $x_i$ is correlated with the other

# VIF

rge when $x_i$ is highly correlated with the other

the first factor, $1/((n-1)s_i^2)$, measures

f $i$th predictor, $x_i$. It is the second factor, $1/(1$

ween the remaining pr

cond factor is denoted as the *variance inflatio*

$$\text{VIF}_i = \frac{1}{1 - R_i^2}$$

navior of the VIF? Suppose that $x_i$ is com

g predictors, so that $R^2 = 0$. Then we wi

# Interpretation of R

| R | % of variance "explained" | % of variance not "explained" | % of SD "explained" | % of SD not "explained" |
|---|---|---|---|---|
| ±0.3 | 9% | 91% | 5% | 95% |
| ±0.5 | 25% | 75% | 13% | 87% |
| ±0.71 | 50% | 50% | 29% | 71% |
| ±0.95 | 90% | 10% | 69% | 31% |

# Partial Correlation

- The correlation coefficient measuring the degree of dependence between two variables

  - after adjusting for the linear effect of one or more of the other X variables
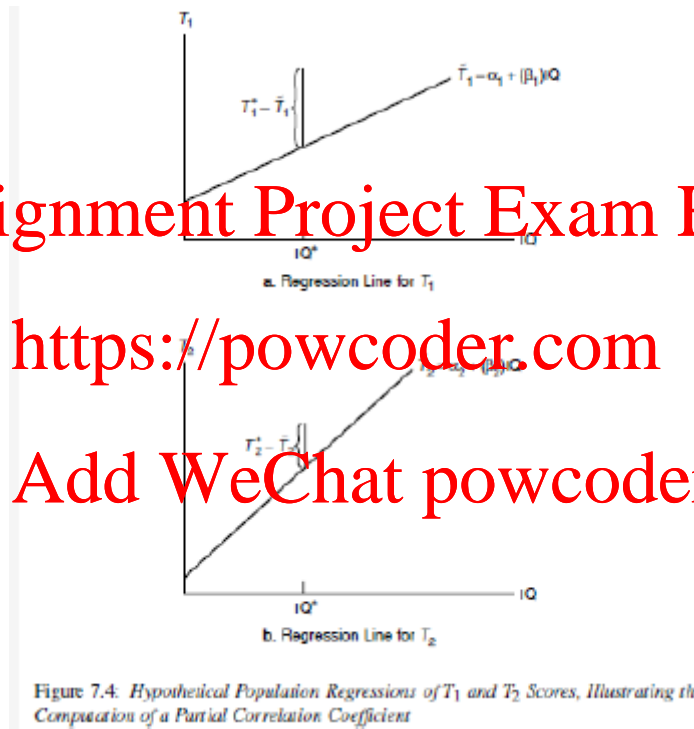
Example:  $T_1$  and  $T_2$ are test scores

- Find partial R between $T_1$  and  $T_2$ after adjusting for IQ

# Visually ( p 130)



Figure 7.4: *Hypothetical Population Regressions of $T_1$ and $T_2$ Scores, Illustrating the Computation of a Partial Correlation Coefficient*

- Partial R = simple R between the two residuals

# Interpretation of regression coefficients

- In the model: $\alpha + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_p X_p$ if $\rho$ is the partial correlation between Y and $X_1$, given $X_1, X_2, \ldots, X_p$, then

- Testing that $\beta_1 = 0$ is equivalent to testing that $\rho = 0$

Hence, $\beta_I$ is called the partial regression coefficient of Y on $X_1$, given $X_1, X_2, \ldots, X_p$

# Values of regression coefficients

- Problem: Values of $\beta_i$ 's are not directly comparable

- Hence:  Standardized coefficients:
  - Standardized $\beta_i$ = $\beta_i$ * (SD ($X_i$) / SD (Y))

- Standardized $\beta_i$ are directly comparable.

# Multicollinearity

$$[SE(B_i)]^2 = \frac{S^2}{(N-1)(S_i)^2} \times \frac{1}{1-(R_i)^2}$$

- The case where some of the X variables are highly correlated

- This will impact estimates, and their SE's (p 143)

- Consider Tolerance, and its inverse, Variance Inflation Factor.

- Target Tolerance < 0.01, or VIF > 100

- Remedy:  use variable selection to delete some X variables, or a dimension reduction techniques such as Principal Components.

# Misleading Correlations

- Example (Lung Function data, Appendix A): FEV1 vs height and age

- Depends on gender

# Total vs Stratified Correlation
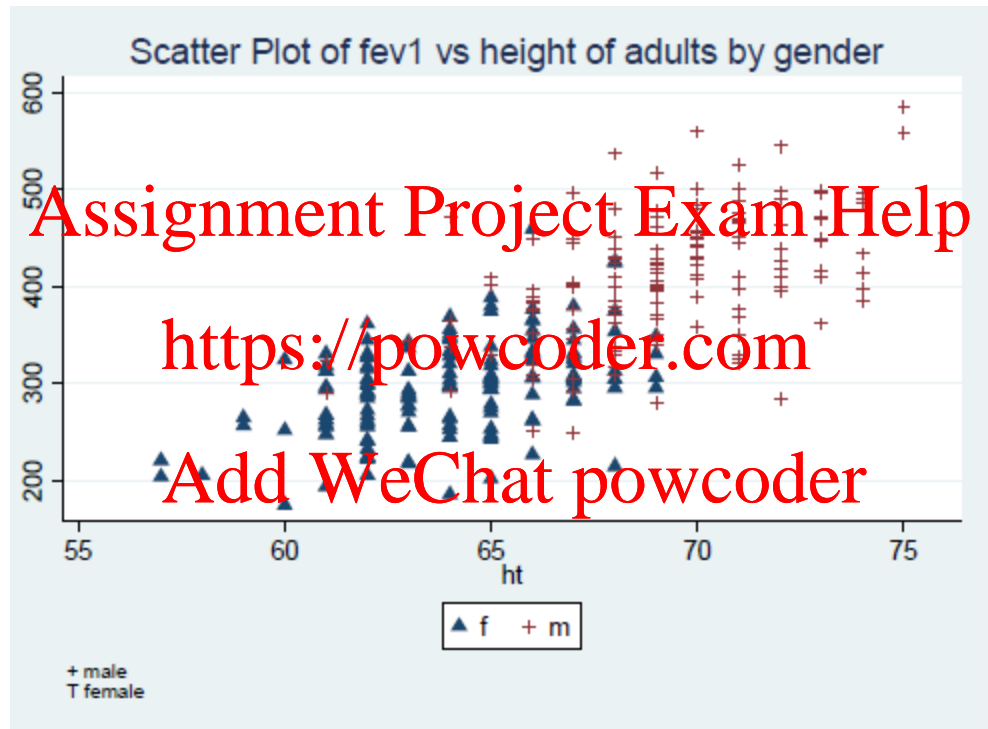
| Gender | Correlation between FEV1 and | |
| --- | --- | --- |
| | Height | Age |
| Total | 0.739 | -0.073 |
| Male | 0.504 | -0.310 |
| Female | 0.465 | -0.267 |
| | | |

# FEV1 vs height



Scatter Plot of fev1 vs height of adults by gender

# FEV1 vs height – Regression lines



Scatter Plot of fev1 vs height of adults by gender

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# FEV1 vs age



Scatter Plot of fev1 vs age of adults by gender

# FEV1 vs age– Regression lines



Scatter Plot of fev1 vs age of adults by gender

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Residual Analysis

- Residual = e = Y - Ŷ

- Studentized residual = $e/S(1 - h)^{1/2}$

  - h called "leverage"

- Deleted studentized residual = studentized residual with observation for computing regression and S deleted.

# Outliers

- Outlier in Y is studentized (or deleted studentized) residual >2 (same as simple case)

- Outlier in X if h > 2(p+1)/N

# Some Caveats

- See list for simple regression
- Need representative sample
- Violations of assumptions, outliers
- Multicollinearity: coefficient of any one variable can vary widely, depending on what others are included in the model
- Missing values
- Number of observations in the sample should be large enough relative to number of variables in the model.

# Outline

- Matrix Review: $(A - \lambda I) X = 0$; Eigenvalues
- Simple linear regression
- Visit
http://www.ats.ucla.edu/stat/sas/output/reg.htm
- Assign HW 6.1,2,5 for next week
- If we get to Chapter 7, assign HW 7.2, 7.4, 7.5, 7.6 (Hand in 7.2,4,5)  7.7 Will be assigned next week.
- Start Multiple Regression Lecture
- Go over Multiple Regression Example – 7.1

# Quick Matrix Review

$(A - \lambda\, I)\, X = 0$

$A = \begin{pmatrix} 3 & 1 \\ 2 & 2 \end{pmatrix}$
$\qquad\qquad (A - \lambda\, I) = \begin{pmatrix} 3 - \lambda & 1 \\ 2 & 2 - \lambda \end{pmatrix}$

$\lambda = 1, 4$

$\lambda = 1 \Rightarrow y = -2x$

$\lambda = 4 \Rightarrow y = x$

# Quick Matrix Review

(A − λ I) X = 0

A =

| 3 | 1 | 3 |
|---|---|---|
| 2 | 2 | 5 |
| 1 | 3 | 2 |

(A − λ I ) =

| 3-λ | 1 | 3 |
|---|---|---|
| 2 | 2-λ | 5 |
| 1 | 3 | 2-λ |

(3-λ)(2-λ)(2-λ)+1*5*1+3*3*2 − ((1*(2-λ)*3) + (2*1*(2-λ)) + ((3-λ)*5*3) = 0

(12 -16λ + 7λ^2 − λ^3 +5 +18) − ((6 - 3λ) + (4 - 2λ) + (45 - 15λ)) = 0

-20 + 4λ +7λ^2 − λ^3  = 0

λ = 7.17, -1.76, 1.59

# Analysis of Variance

Y

Observed Value
of Y for three
Groups

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder