1. **Applied Regression Analysis 3rd edition**

Norman R. Draper, Harry Smith

Response variable = Model Function + Random error

Errors assumed to be independent

Linear first order model:
$$Y = \beta_0 + \beta_1 X_i + \varepsilon$$

For a given X, a corresponding Y consists of the value $\beta_0 + \beta_1 X_i$ plus an amount $\varepsilon$, the increment by which an individual Y may fall off the regression line

## Meaning of linear Model

− When we say a model is linear or nonlinear we are referring to linearity or nonlinearity in the parameters

− Value of highest power of a predictor variable in the model is called the order of the model.
$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \varepsilon$$
is a second order (x) linear ($\beta s$) regression model.

## Least Squares Estimation:

− $\beta_0, \beta_1$ and $\varepsilon$ are unknown & difficult to discover as it changes in each observation Y.

− We calculate estimate of $\beta_0$ and $\beta_1$, $b_0, b_1$:
$$\hat{Y} = b_0 + b_1 X_i$$

$\hat{Y}$ predicted value of Y for a given X.

Suppose we have n sets of observations $(X_1, Y_1) \dots (X_n, Y_n)$ then we can write
$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \qquad \text{for } i = 1, 2 \dots n$$

So that the sum of squares of deviation from the true line is:
$$S = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_i)^2$$

2

- S is called the Sum of Squares function
- We shall choose our value of $b_0$ and $b_1$ that, when substituted for $\beta_0$ and $\beta_1$, produce the least possible value of S

- We determine $b_0$ and $b_1$ by differentiating first with respect to $\beta_0$ and then with respect to $\beta_1$ and setting results equal to zero:

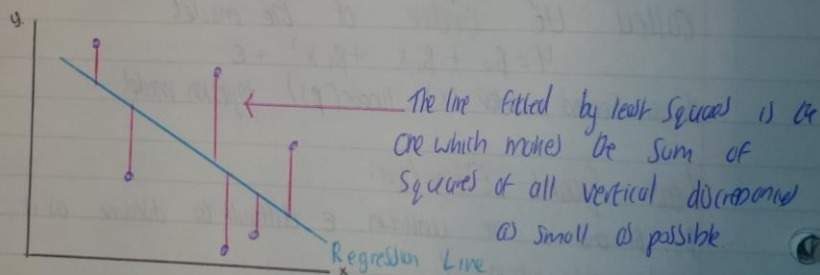$$\frac{dS}{d\beta_0} = -2 \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{dS}{d\beta_1} = -2 \sum_{i=1}^{n} x_i (Y_i - \beta_0 - \beta_1 x_i)$$

So that the estimated $b_0$ and $b_1$ are solutions to 2 eq'ns:

$$\sum (Y_i - b_0 - b_1 x_i) = 0$$
$$\sum x_i (Y_i - b_0 - b_1 x_i) = 0$$

We sub $(b_0, b_1)$ for $(\beta_0, \beta_1)$ when we equate to zero



The line fitted by least squares is the one which makes the sum of squares of all vertical discrepancies as small as possible

Regression Line

The solution to these eq'ns yield:

$$b_1 = \frac{\sum (x_i - \bar{x})(Y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \qquad \text{where } \bar{x} = \frac{\sum x_i}{n}$$

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

Also: $\sum (x_i - \bar{x})(Y_i - \bar{y}) = \sum x_i Y_i - \bar{x} \sum Y_i - \bar{y} \sum x_i + n\bar{x}\bar{y}$

$$= \sum x_i Y_i - n\bar{x}\bar{y}$$

$$= \sum x_i Y_i - \left( \frac{\sum x_i \sum Y_i}{n} \right)$$

$\sum x^2$ is known as uncorrected sum of squares of x's

$\frac{(\sum x)^2}{n}$ is the correction for the mean of x's

$\sum x^2 - \frac{(\sum x)^2}{n}$

The difference is called corrected sum of squares of x's $S_{xx}$

$\sum x_i y_i$ is called the uncorrected sum of products

$\frac{\sum x_i \sum y_i}{n}$ is called correction for the means

$\sum (x - \bar{x})(y - \bar{y}) = S_{xy}$

The difference is called corrected sum of product of X and Y

## Notation.

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

$$S_{xx} = \sum (x_i - \bar{x})^2$$
$$= \sum (x_i - \bar{x}) x_i$$
$$= \sum x_i^2 - (\sum x)^2 / n$$
$$= \sum x_i^2 - n \bar{x}^2$$

$$S_{yy} = \sum (y_i - \bar{y})^2$$
$$= \sum (y_i - \bar{y}) y_i$$
$$= \sum y_i^2 - (\sum y)^2 / n$$
$$= \sum y_i^2 - n \bar{y}^2$$

The easily remembered formula for $b_1$ is then: $\frac{S_{xy}}{S_{xx}}$

$b_0 = \bar{y} - b_1 \bar{x}$   Sub in $b_0$   $\hat{y} = b_0 + b_1 x_i$

$$\implies \hat{y} = \bar{y} + b_1 (x - \bar{x})$$

If we set $x = \bar{x}$, we see $\hat{y} = \bar{y}$. The mean point $(\bar{x}, \bar{y})$ lies on the fitted line.

In other words, the least squares line contains the centre of gravity of the data

4

$(y_i - \hat{y}_i)$ equals the residual $e_i$

Since $\hat{y}_i = \bar{y} + b_1(x_i - \bar{x})$

$$y_i - \hat{y}_i = (y_i - \bar{y}) - b_1(x_i - \bar{x})$$

Which we can sum to give:

$$\Sigma(y_i - \hat{y}_i) = \Sigma(y_i - \bar{y}) - b_1\Sigma(x_i - \bar{x}) = 0$$
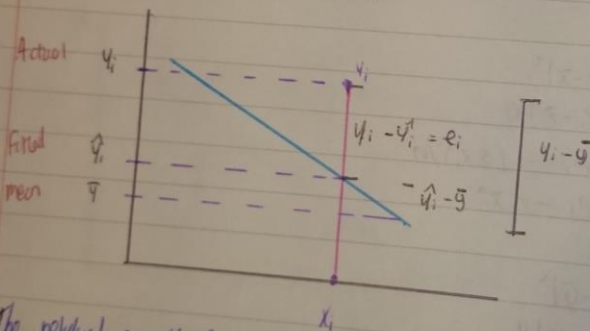
Residuals Sum to zero in these

## The Analysis Of Variance

We tackle question of how much of the between variance in the data has been explained by the regression line.

Consider $y_i - \hat{y}_i = y_i - \bar{y} - (\hat{y}_i - \bar{y})$

What this means geometrically is



The residual $e_i = y_i - \hat{y}_i$ is difference between 2 quantities (1) the deviation of the observed $y_i$ from the overall mean $\bar{y}$ and (2) the deviation of the fitted $\hat{y}_i$ from the overall mean $\bar{y}$.

Note that average of the $\hat{y}_i$ is the same as average of $y_i$:

$$\frac{\Sigma \hat{y}_i}{n} = \Sigma(b_0 + b_1 x_i)/n$$

$$= (nb_0 + b_1 n\bar{x})/n$$

$$= b_0 + b_1 \bar{x}$$

$$= \bar{y}$$

5.

This fact also reconfirms that $\sum e_i = \sum(y_i - \hat{y}_i) = n\bar{y} - n\bar{y} = 0$

We can rewrite as $(y_i - \bar{y}) = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)$

If we square both sides of this and sum from $i = 1, 2, \dots n$ we obtain
$$\sum(y_i - \bar{y})^2 = \sum(\hat{y}_i - \bar{y})^2 + (y_i - \hat{y}_i)^2$$

The cross product term $(PT = 2\sum(\hat{y}_i - \bar{y})(y_i - \hat{y}_i)$ can be shown to vanish by applying $\hat{y} = \bar{y} + b_1(x - \bar{x})$
$$\hat{y}_i - \bar{y} = b_1(x_i - \bar{x})$$
$$y_i - \hat{y}_i = y_i - \bar{y} - b_1(x_i - \bar{x})$$

It follows that the cross product term is:
$$PT = 2\sum b_1(x_i - \bar{x})\{(y_i - \bar{y}) - b_1(x_i - \bar{x})\}$$
$$= 2b_1(S_{xy} - b_1 S_{xx})$$
$$= 0$$

It is also clear: $\sum(\hat{y}_i - \bar{y})^2 = \sum b_1^2 S_{xx}$
$$= b_1 S_{xy}$$
$$= S^2_{xy}/S_{xx}$$

## Sum of Squares

The quantity $(y_i - \bar{y})$ is the deviation of the $i$th observation from the overall mean and so, the LHS $\sum(y_i - \bar{y})^2$ is the sum of squares of deviations of the observations from the mean.

This is shorted to SS about the mean and is also the corrected sum of squares of the $y$'s namely $S_{yy}$

Since $\hat{y}_i - \bar{y}$ is the deviation of predicted value of the $i$th observation from the mean and $y_i - \hat{y}_i$ is the deviation of the $i$th observation from its predicted/fitted value we can express:

$$6 \text{ SSTO} = \text{SSR} \quad \text{SSE}$$

$$\begin{pmatrix} \text{Sum of Square} \\ \text{about the mean} \end{pmatrix} = \begin{pmatrix} \text{Sum of Squar} \\ \text{due to regression} \end{pmatrix} + \begin{pmatrix} \text{Sum of Squar} \\ \text{about regression} \end{pmatrix}$$

This shows that of the variation in the Y's about their mean, some of the variation can be ascribed to the regression line and some $\sum(Y_i - \hat{Y}_i)^2$ to the fact that the actual observations do not lie on the regression line, if they did at all, the sum of square would be zero.

We see that a sensible way to assess how useful the regression line will be as a predictor is to see how much of the SS about the mean has fallen into the SS due to regression and how much about the regression.

We will be happy if SS due to regression is much greater than the SS about regression, or what amounts to the same thing, if the ratio $R^2 = \dfrac{\text{SS due to regression}}{\text{SS about mean}}$

is not far from unity.

## Degrees Of Freedom

This number indicates how many independent piece of information including the $n$ independent number $Y_1 \ldots Y_n$ are needed to compute the sum of squares.

For example, SS about mean needs $(n-1)$ independent piece [ of the numbers $Y_1 - \bar{Y}, Y_2 - \bar{Y} \ldots Y_n - \bar{Y}$, only $(n-)$ are independent since all $n$ numbers sum to zero by definition of the mean ]

We can compute SS due to regression from a single function of $Y_1 \ldots Y_n$ namely $b_1$ [ Since $\sum(\hat{Y}_i - \bar{Y})^2 = b_1^2 \sum(x_i - \bar{x})^2$ as $\boxed{}$ and so has one degree of freedom.

7.

By Subtraction, the SS about regression, which we shall in future call the residual Sum of Squares, & has $n-2$ degrees of freedom (df).

Two parameters are estimated

In general, residual Sum of Squares is based on (number of observations - number of parameters estimated) degrees of freedom

$$n - 1 = 1 + (n-2)$$

Analysis of Variance Table

We can construct an analysis of variance table. The "mean Square" column is calculated by dividing each Sum of Squares entry by its corresponding degrees of freedom.

| 1·3 | Source of Variation | Degrees of freedom | Sum of Squares SS | Mean Square MS |
|---|---|---|---|---|
| | Due to regression | 1 | $\sum(\hat{y}_i - \bar{y})^2$ SSR | $MS_{Reg} = \frac{SSR}{1}$ |
| | About regression (residual) | $n-2$ | $\sum(y_i - \hat{y}_i)^2$ SSE | $s^2 = \frac{SS}{(n-1)}$ · $\frac{SSE}{n-2}$ |
| | Total corrected for mean $\bar{y}$ | $n-1$ | $\sum(y_i - \bar{y})^2$ SSTo | $f = \frac{MSR}{MSE}$ |

A more general form of the analysis of variance table is obtained by incorporating the correction factor for the mean of the Y's into the table where it is called SS $(b_0)$.

Incorporating SS $(b_0)$

| 1·4 | Source | df | SS | MS = SS/df |
|---|---|---|---|---|
| | Due to $b_1|b_0$ | 1 | $SS(b_1|b_0) = \sum(\hat{y}_i - \bar{y})^2$ | $MS_{reg}$ |
| | Residual | $n-2$ | $\sum(y_i - \hat{y}_i)^2$ | $s^2$ |
| | Total corrected. | $n-1$ | $\sum(y_i - \bar{y})^2$ | |
| | | | | |
| | correction factor due to $b_0$ | 1 | $SS(b_0) = \frac{(\sum y)^2}{n} = n\bar{y}^2$ | |
| | Total | $n$ | $\sum y_i^2$ | |

8.

Alternative way to display 14 is to drop the line "total corrected".
Total line is sum of remaining 3 entries.

| Source | df | SS | $MS = \frac{SS}{df}$ |
|--------|-----|------------------|-------|
| $b_0$ | 1 | $n\bar{y}^2$ | — |
| $b_1\vert b_0$ | 1 | $S_{xy} b_{xx}$ | $MS_{reg}$ |
| Residual | $n-2$ | By Subtraction | $s^2$ |
| TOTAL | $n$ | $\Sigma y^2$ | |

First SS entry due to $b_0$ is the amount of variation in $n\bar{y}^2$ explained by a horizontal straight line $\hat{y}' = \bar{y}$.

If the model $Y = \beta_0 + \varepsilon$ via least squares, fitted model is $\hat{y} = \bar{y}$.

If we subsequently fit the "with slope $b_1$" model $Y = \beta_0 + \beta_1 x_i + \varepsilon$, the "due to $b_1\vert b_0$" SS entry $\frac{S_{xy}}{S_{xx}}$ is the extra variation picked up by the slope term over and above that picked up by the intercept alone.

$$SS(b_1\vert b_0) = \Sigma(\hat{y}_i - \bar{y})^2 = b_1\{\Sigma(x_i - \bar{x})(y_i - \bar{y})\} = b_1 S_{xy}$$

$$= \frac{(\Sigma(x_i - \bar{x})(y_i - \bar{y}))^2}{\Sigma(x_i - \bar{x})^2} = \frac{S^2_{xy}}{S_{xx}}$$

$$= \frac{(\Sigma x_i y_i - (\Sigma x_i \Sigma y_i)/n)^2}{\Sigma x_i^2 - (\Sigma x_i)^2/n} = \frac{S^2_{xy}}{S_{xx}}$$

$$= \frac{(\Sigma(x_i - \bar{x}) y_i)^2}{(\Sigma x_i - \bar{x})^2}$$

Note, total corrected SS, $\Sigma(y_i - \bar{y})^2$ can be written as:

$$S_{yy} = \Sigma y_i^2 - (\Sigma y_i)^2/n$$

or $S_{yy} = \Sigma y_i^2 - n\bar{y}^2$

The mean square about regression, $s^2$ will provide an estimate based on the degrees of freedom of the variance about the regression, which we will call $\sigma^2_{y \cdot x}$

9

If regression equation were estimated from on indefinitely large number of observations the variance about the regression would represent a measure of the error with which any observed value of Y could be predicted from a given value of X using the determined eq⁻

## Skeleton Analysis of Variance Table.
A Skeleton analysis of variance table consists of 4 "Source" and "df" columns only.

## $R^2$ Statistic
A useful statistic to check is the $R^2$ value of a regression fit:

$$R^2 = \frac{(SS \text{ due to regression given } b_0)}{(\text{Total SS corrected for the mean } \nabla)} \quad \frac{SSR}{SSTO}$$

$$= \frac{\Sigma(\hat{y}_i - \bar{y})^2}{\Sigma(y_i - \bar{y})^2} \quad 1.3.15. \quad i = 1, 2 \dots n$$

$R^2$ measures 'proportion of total variance about the mean Y explained by the regression".

In fact R is the correlation between Y and $\hat{Y}$ and is usually called the multiple correlation coefficient

$R^2$ is then "the square of the multiple correlation coefficient".
For a straight line fit: $R^2 = SS(b_1|b_0)/S_{yy}$
$$= S_{xy}^2/(S_{xx}S_{yy})$$

Thus $r^2$ explains G% of the total variation in the data about the overall $\bar{Y}$. This is quite a large proportion

$R^2$ can take values as high as $1 \rightarrow 100$ when all the X values are different

When repeat runs exist in the data, the value of $R^2$

11

is $\quad V(a) = a_1^2 V(Y_1) + a_2^2 V(Y_2) + \ldots + a_n^2 V(Y_n)$

if the $Y_i$ are pairwise uncorrelated and the $a_i$ are constant; furthermore,
if $V(Y_i) = \sigma^2$

$$V(a) = (a_1^2 + a_2^2 + \ldots + a_n^2)\sigma^2$$
$$= (\Sigma a_i^2)(\sigma^2)$$

In the expression for $b_1$, $\quad a_i = (x_i - \bar{x}) / \Sigma (x_i - \bar{x})^2$ since $x_i$ can be regarded as constant. Hence after reduction:

$$V(b_1) = \frac{\sigma^2}{\Sigma(x_i - \bar{x})^2} = \frac{\sigma^2}{S_{xx}}$$

The Standard deviation of $b_1$ is the square root of the variance:

$$sd(b_1) = \frac{\sigma}{\sqrt{\Sigma(x_i - \bar{x})^2}} = \frac{\sigma}{\sqrt{S_{xx}}}$$

or if $\sigma$ is unknown, and we use the estimate $s$ in its place, assuming the model is correct, the estimated standard deviation of $b_1$:

$$est \ sd(b_1) = \frac{s}{\sqrt{\Sigma(x_i - \bar{x})^2}} = \frac{s}{\sqrt{S_{xx}}}$$

Alternative terminology for estimated standard deviation is the standard error

Confidence Interval for $\beta_1$

If we assume that the variation of the observation about the line are normal - that is the errors $\varepsilon_i$ are all from the same normal distribution, $N(0, \sigma^2)$ - it can be shown that we can assign $100(1-\alpha)\%$ confidence limits for $\beta_1$ $\Rightarrow$

$$b_1 \pm \frac{t_{(n-2, 1-\frac{1}{2}\alpha)} s}{\sqrt{\Sigma(x_i - \bar{x})^2}}$$

where $t_{(n-2, 1-\frac{1}{2}\alpha)}$ is the $100(1-\frac{1}{2}\alpha)$ percentage point of a $t$ dist, with $n-2$ df.

12

Test for $H_0$ $\beta_1 = \beta_{10}$ Vs $H_1$ $\beta_1 \neq \beta_{10}$

Test null hypothesis that $\beta_1$ is equal to $\beta_{10}$, where $\beta_{10}$ is a specified value that could be zero, against alternative hypothesis that $\beta_1$ is different to $\beta_{10}$ by calculating

$$t = \frac{b_1 - \beta_{10}}{se(b_1)} = \frac{(b_1 - \beta_{10})(\sqrt{\Sigma(x_i - \bar{x})^2})}{s}$$

and comparing $|t|$ with $t_{n-2, 1-\frac{1}{2}\alpha}$ from a t-table with $n-2$ d.f. The test will be a 2 sided test conducted at $100 \alpha \%$ level in this case.

The value lies with $(\sim, \sim)$ and this statement is made u 95% confidence in

If $|t|$ exceeds central $H_0$ is rejected.

We could also examine CI to see if it includes the value or not.

Reject or do not reject.
- If $|t|$ value had been smaller than the critical value we could not reject the hypothesis.
- The most we can say is that on basis of certain observed data we cannot reject it.
- It may be that in another set of data we can find evidence that is contrary to our hypothesis and so reject it.

Confidence Interval Repeats a Set of Test.
Once we have CI for $\beta_1$ we do not actually have to compute the $|t|$ value for a particular 2 sided-t-test at the same or less

It is easier to examine interval for $\beta_1$ and see if it contains $\beta_{10}$
If it does then hypothesis $\beta_1 = \beta_{10}$ cannot be rejected

13

If $|b_1 - \beta_{10}| > t(n-2, 1-\tfrac{1}{2}\alpha) \dfrac{s}{\sqrt{\Sigma(x_i-\bar{x})^2}}$

that is, $\beta_{10}$ lie outside the limit of eq 1.4d

## Standard Deviation of the Interval for $\beta_0$

A CI for $\beta_0$ and a test of whether or not $\beta_0$ is equal to some specified value can be constructed in a way similar to that just described for $\beta_1$

$$Sd(b_0) = \sigma \sqrt{\dfrac{\Sigma x_i^2}{n \, \Sigma(x_i-\bar{x})^2}}$$

Replacement of $\sigma$ by $s$ provides the estimated $sd(b_0)$, that is $s.e(b_0)$. Thus $100(1-\alpha)\%$ confidence limits are given by:

$$b_0 \pm t(n-2, 1-\tfrac{1}{2}\alpha) \; s \sqrt{\dfrac{\Sigma x_i^2}{n \, \Sigma(x_i-\bar{x})^2}}$$

A test for $H_0: \beta_0 = \beta_{00}$ against $H_1: \beta_0 \neq \beta_{00}$ will be rejected at the $\alpha$ level if $\beta_{00}$ falls outside the CI, or will not be rejected if it lies inside or may be conducted separately by finding quantity:

$$t = \dfrac{(b_0 - b_{00})}{s \cdot \sqrt{\dfrac{\Sigma x_i^2}{n \, \Sigma(x_i-\bar{x})^2}}}$$

and comparing it with $t(n-2, 1-\alpha\tfrac{1}{2})$. $s$

# Regression    F-Test

## F-Test for Significance of Regression

Since $y_i$ are random variables, any function of them is also a random variable; two particular functions are $MS_{reg}$, the mean square due to regression, and $s^2$ the mean square due to residual variation, which are in the analysis of variance table,

These functions then have their own distribution, mean, variance and moments.

It can be shown that the mean values are:

$$E[MS_{reg}] = \sigma^2 + \beta_1^2 \Sigma(x_i - \bar{x})^2$$
$$E(s^2) = \sigma^2$$

Where if $Z$ is a random variable, $E[Z]$ denotes its mean or expected value.

Suppose that the errors $\varepsilon_i$ are independent $N(0, \sigma^2)$ variables.

It can be shown that if $\beta_1 = 0$, the variable $MS_{reg}$ multiplied by its degrees of freedom (here one) and divided by $\sigma^2$ follow a $x^2$ dist with the same (one) df

In addition $(n-2)s^2/\sigma^2$ follow $x^2$ dist with $(n-2)$ df.

Since 2 variables are independent, a statistical theorem tell us that the ratio     $F = \dfrac{MS_{reg}}{s^2}$

follow an F distribution with 1 and $n-2$ df provided $\beta_1 = 0$

This fact can thus be used as a test $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

We compare the ratio $F = MS_{reg}/s^2$ with the $100(1-\alpha)\%$ point of the tabulated $F_{(1, n-2)}$ dist in order to determine whether $\beta_1$ can be considered nonzero on the basis of the data we have see

$$F = t^2$$

We have 2 tests for the test of $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$ a t-test and F-td.

In fact, the two test are equivalent and mathematically related here, due to the theoretical fact that $F_{(1, v)} = \{t_{(v)}\}^2$ that is, the square of a t-variable with $v$ d.f is an F-variable with $1$ and $v$ df.

NOTE: This only happens when 1st df of F is 1.

$$F = \frac{MS_{reg}}{s^2} = \frac{b_1 \Sigma(x_i - \bar{x})S(y_i - \bar{y})}{s^2}$$

$$= \frac{b_1^2 \Sigma(x_i - \bar{x})^2}{s^2} \quad \text{(by def of } b_1\text{)}$$

$$= \left[\frac{b_1 \sqrt{\Sigma(x_i - \bar{x})^2}}{s}\right]^2$$

$$= t^2$$

Since variable $F_{(1, n-2)}$ is the square of the $t_{(n-2)}$ variable, and this carried over to the percentage points (upper or tail of F and 2-tailed t tailed of $\alpha$). Exact same result

When there are more regression coef's the overall F-test for reg, which is the extension of the one given here, does not correspond to the t-test of a coef. This is why we need to know both t and F-test.

However tests for an individual coef can be made in either t or F by a similar argument.

10

cannot attain 1 no matter how well the model fit.

This is because no model however good, can explain the variation in the data - due to pure error.

We now make the basic assumptions of Model

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \qquad i = 1, 2, \ldots n$$

1  $\varepsilon_i$ is a random variable with mean zero and variance $\sigma^2$ (unknown) that is: $E(\varepsilon_i) = 0 \qquad V(\varepsilon_i) = \sigma^2$

2  $\varepsilon_i$ and $\varepsilon_j$ are uncorrelated, $i \neq j$, so that $Cov(\varepsilon_i, \varepsilon_j) = 0$ thus:
$$E(Y_i) = \beta_0 + \beta_1 X_i \qquad V(Y_i) = \sigma^2$$

A further assumption, which is not immediately necessary and will be recalled when used is that:

3  $\varepsilon_i$ is normally distributed random variable, with mean zero and variance $\sigma^2$ by assumption (D); that is:
$$\varepsilon_i \sim N(0, \sigma^2)$$
Under this additional assumption, $\varepsilon_i, \varepsilon_j$ are not only uncorrelated but necessarily independent.

Standard Deviation of Slope $b_1$; Confidence Interval for $\beta_1$.
We know $b_1 = \sum(X_i - \bar{X})(Y_i - \bar{Y}) / \sum(X_i - \bar{X})^2$
$$= \sum(X_i - \bar{X})Y_i / \sum(X_i - \bar{X})^2$$
[Since the other term removed from numerator is $\sum(X_i - \bar{X})\bar{Y} = \bar{Y}\sum(X_i - \bar{X}) = 0$

$$b_1 = \{(X_1 - \bar{X})Y_1 + \ldots + (X_n - \bar{X})Y_n\} / \sum(X_i - \bar{X})^2$$

Now the variance of a function
$$a = a_1 Y_1 + a_2 Y_2 + \ldots + a_n Y_n$$