

① Introduction To Regression

$$Z \sim N(0, 1)$$

↓ ↓ ↓
 fixed mean st. deviation

$$\frac{X - \bar{X}}{\sigma_{\text{standard}}}$$

↓ ↓
 Z score value mean

26/4/13

$$\sum_{i=1}^n (x_i - \bar{x})^2$$

$$\Rightarrow \sum_{i=1}^n (x_i^2 + \bar{x}^2 - 2x_i\bar{x})$$

$$\Rightarrow \sum_{i=1}^n (x_i^2 + n\bar{x}^2 - 2 \sum_{i=1}^n x_i\bar{x})$$

$$\Rightarrow \sum_{i=1}^n x_i^2 + n\bar{x}^2 - 2\bar{x} \sum_{i=1}^n x_i$$

$$\Rightarrow \sum_{i=1}^n x_i^2 + n\bar{x}^2 - 2\bar{x}(n\bar{x}) \quad \leftarrow \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

Variance always ADDED

If $x \sim N(0, 2^2)$ and $y \sim N(3, 1^2)$

$$1. \frac{1}{4} \text{Var}(x+y) = \frac{1}{4}(\text{Var}(x) + \text{Var}(y))$$

$$2. \text{Var}(4x+3y) = 16\text{Var}(x) + 9\text{Var}(y) \rightarrow \text{Constant gets squared}$$

$$3. \text{Var}(x-y) = \text{Var}(x) + \text{Var}(y) \rightarrow \text{Variance always ADDED}$$

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

$$E\left(\sum_{i=1}^n \frac{1}{n} \mu_i\right)$$

$$E(\bar{x}) = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \quad \text{Var}(\bar{x}) = \frac{1}{n} \text{Var}(x)$$

$$= \frac{1}{n} E\left(\sum_{i=1}^N x_i\right)$$
$$= \frac{1}{n} \sum_{i=1}^N E(x_i)$$
$$= \frac{1}{n} \sum_{i=1}^N \mu$$

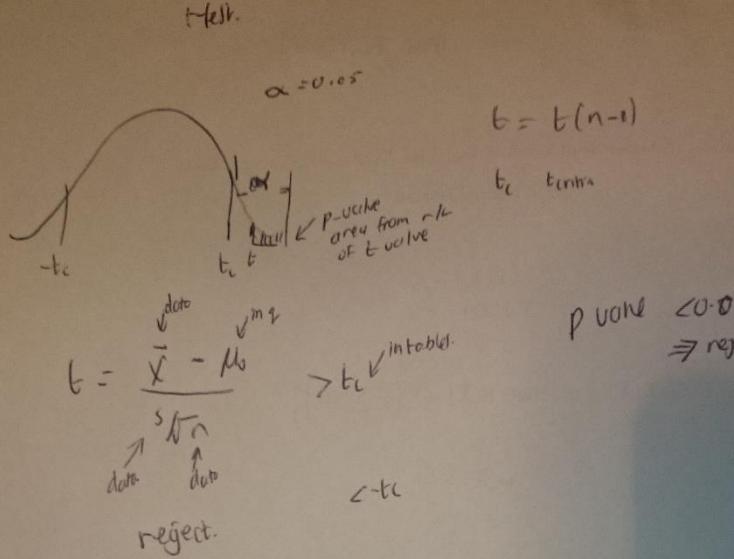
$$= \frac{1}{n} (\mu n)$$

$$= \mu$$

$$\text{Var}(\bar{x}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right)$$

$$= \left(\frac{1}{n}\right)^2 \text{Var}\left(\sum_{i=1}^n x_i\right)$$

$$\frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n x_i\right)$$
$$= \frac{1}{n^2} \sum_{i=1}^N \text{Var}(x_i)$$
$$= \frac{1}{n^2} (n \sigma^2)$$
$$= \frac{\sigma^2}{n}$$



Accept

$$-t_c < t < t_c$$

$$\Rightarrow -t_c < \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}} < t_c$$

$$\Rightarrow -t_c \frac{s}{\sqrt{n}} < \bar{x} - \mu_0 < t_c \frac{s}{\sqrt{n}}$$

$$\Rightarrow \bar{x} - t_c \frac{s}{\sqrt{n}} < \mu_0 < \bar{x} + t_c \frac{s}{\sqrt{n}}$$

$C1 =$

$$(\bar{x} - t_c \frac{s}{\sqrt{n}}, \bar{x} + t_c \frac{s}{\sqrt{n}})$$

Regression

$$R = \frac{cov(x, y)}{\sqrt{var(x)} \sqrt{var(y)}}$$

$$= \frac{cov(x, y)}{sd(x) \ sd(y)}$$

$$= \frac{\frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} \ \sqrt{\frac{1}{n-1} \sum (y_i - \bar{y})^2}}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \ \sqrt{\sum (y_i - \bar{y})^2}}$$

$$-1 \leq r \leq 1$$

$r > 0$ = positive correlation

$r < 0$ = negative

$r = 0$ = no correlation

No Linear Relationship

3

$$\text{Covariance: } \text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$$

$$\text{Sample: } \text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

If values of X lie above \bar{x} indicates values of Y lie above \bar{y}
then we have positive products. Similar for negative ones.

9/10/13

Regression

Correlation

$$R = \frac{\text{Cov}(x,y)}{\text{sd}(x)\text{sd}(y)}$$

1. Change in x cause change in y
2. " " " x
3. Change in some 3rd variable z , could change in x and y
4. The relationship is a coincidence

Correlation does not imply causality

Simple Linear Regression

Regression is a statistical tool that utilizes the relationship between two or more variables

Relationship

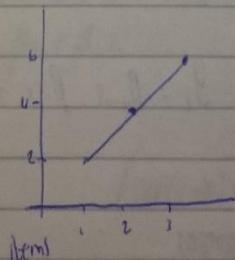
1. Functional
2. Statistical

Functional relationship is deterministic and can be exactly expressed by a formula

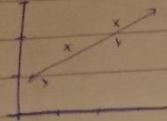
$$y = f(x)$$

e.g. $y = 2x$

Price of 1 item = £2.



2
~9/10/12
statistical relationship is not perfect like a function or



FR
 $y = F(x)$

SR
 $E(y) = f(x)$

Regression model is formal, mean of expressing a statistical relationship

A statistical relation can be summed up as:

1. The dependent variable y varies with respect to independent variables in some systematic way.

2. Observations (y_i) are scattered around the curve of relationship.

These 2 are represented in a regression model:

1. In the population of observation associated with the sampled ones, there is a probability distribution of y for every value of x .

2. The mean of y varies in some systematic relationship with respect to x .

$(x_1, y_1), \dots, (x_n, y_n)$

Simple linear equation: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$

$y_i \rightarrow i^{\text{th}}$ data point of y

$x_i \rightarrow$ "fixed" value for x

$\beta_0, \beta_1 \rightarrow$ unknown parameters

$\epsilon_i \rightarrow$ random error term.

11/01/13 Regression

$$1. Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

- i. $y_i \rightarrow i^{\text{th}}$ datapoint
- ii. $\beta_0, \beta_1 \rightarrow$ unknown parameters (need to be estimated)
- iii. X_i are constants
- iv. ε_i is a random error term $E(\varepsilon_i) = 0$
 $\text{Var}(\varepsilon_i) = \sigma^2$ for all i
 $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$ for all $i \neq j$.

Feature of the model

- 1. The observed value of Y at the i^{th} datapoint is the sum of
- 2. components

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

Systematic Chance (random term)

Error in prediction $U_i = \beta_0 + \frac{\beta_1}{\log P_i} + \varepsilon_i$

$$2. \text{ Since } E(\varepsilon_i) = 0 \\ E(Y_i) = E(\beta_0 + \beta_1 X_i + \varepsilon_i) = E(\beta_0 + \beta_1 X_i) + E(\varepsilon_i) \xrightarrow{\varepsilon_i=0} \\ = \beta_0 + \beta_1 X_i$$

$$3. \text{Var}(\varepsilon_i) = \sigma^2 \\ \text{Var}(Y_i) = \text{Var}(\beta_0 + \beta_1 X_i) \quad (\text{Var}(c+X) = \text{Var}(X) \text{ if } c \text{ is a constant}) \\ = \text{Var}(\varepsilon_i) = \sigma^2$$

$$4. \text{Important point for regression is } E(Y) = \beta_0 + \beta_1 X$$

$$5. Y_i \text{ and } Y_j \text{ are uncorrelated for } i \neq j. \text{ because } \text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \\ \Rightarrow \text{Cov}(Y_i, Y_j) = 0$$

$\text{Cov}(x, y) = 0 \Rightarrow x \text{ and } y \text{ are independent}$
 y_i and y_j are assumed to be independent

One more assumption which is not necessary $\Rightarrow y \sim N(\mu, \sigma^2)$

Estimate β_0 and β_1 only required.

i. $E(y_i) = 0 \quad \Rightarrow \quad y_i \sim N(0, \sigma^2)$

ii. $\text{Var}(x_i) = r^2$

iii. $\text{Cov}(y_i, y_j) = 0 \quad \Rightarrow \quad y_i \text{ and } y_j \text{ are independent}$

EXAM

ASSUMPTIONS LEADING TO LINEAR REGRESSION

1. y_i 's are assumed to be independent of each other.

2. y_i 's are normally distributed with mean 0.

3. All y_i 's have the same variance

4. It is assumed that $E(y_i)$ can be joined by a straight line. $E(y_i) = \beta_0 + \beta_1 x_i$

Example: Price Of Diamond Statistic

$$Y_i = -259.63 + 3721.02x_i + \epsilon_i$$

Price = $\beta_0 + \beta_1$ carat

$x_i = 0.23$ carats. Actual value \$591

$$E(y_i) = -259.63 + 3721.02(0.23)$$

$$= 596.20$$

$$\epsilon_i = -1.2 \quad (\text{difference - error - residual})$$

Moving on β_0 and β_1

β_0 = y intercept

β_1 = slope of line change in y with respect to unit change in x.

6/10/23 Introduction To Regression

Regression Model:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \quad i=1, \dots, n$$

Assumptions

1. X_i is the i^{th} value of the predictor variable and is constant.

2. ϵ_i is the error value and is a random variable.

(i) $E(\epsilon_i) = 0$ for all i

(ii) $\text{Var}(\epsilon_i) = \sigma^2$ for all i

(iii) $\text{Cov}(\epsilon_i, \epsilon_j) = 0$ for $i \neq j$

3. The mean of Y_i can be fitted by a straight line

$$E(Y_i) = \beta_0 + \beta_1 X_i \rightarrow \beta_0, \beta_1 \text{ are parameters}$$

(i) β_1 is the slope of the regression line. It indicates the change in mean of Y for unit change in X .

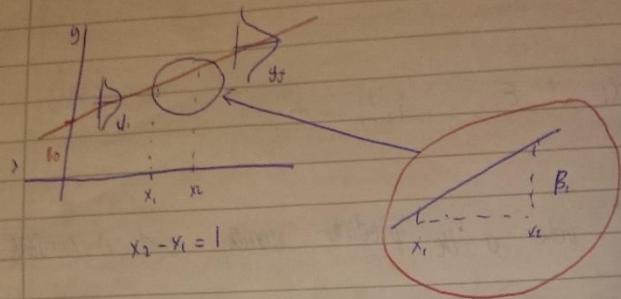
(ii) β_0 is the intercept.

Physically speaking, β_0 is the mean value of Y when $X=0$ and may not always have any meaning.

Simplified Assumptions

1. X_i is the ...
2. The observations y_i (or equivalently ϵ_i) are independent of each other.
3. For any i , y_i is normally distributed with μ
4. For any i , y_i 's have same standard deviation σ .
5. The mean of y_i can be ...

2 Regression



Diamonds Dataset

$$y_i = -259.63 + 3721.02x_i + \epsilon_i \quad y_i = 595 \quad \text{Price} \quad \text{const}$$

$$595 = -259.63 + 3721.02(0.23) + \epsilon_i \quad \leftarrow \text{error}$$

$$\epsilon_i = -1.20$$

$$E(\epsilon_i) = 0$$

Estimation of parameters β_0 and β_1

For the price of diamonds example, it was provided. How did one get:

$$\beta_0 = -259.63, \quad \beta_1 = 3721.01$$

If x followed a Normal distribution with unknown mean μ and σ^2 , μ was estimated by \bar{x} .

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

Procedure

One of the outcomes of estimating β_0 and β_1 is that they should result in a "line of best fit".

\Rightarrow one which reflects the data well.

Best in what sense?

16/10/13

3 Regression

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

$$\Rightarrow y_i - (\beta_0 + \beta_1 x_i) \text{ for all } i$$

We consider that line to be best for which the sum of squared error is minimal ($\sum_{i=1}^n \varepsilon_i^2$)

$$+ \quad +$$

Method of Least Squares

The method of LS gives estimates $\hat{\beta}_0$ (estimated for β_0) and $\hat{\beta}_1$ (for β_1) parameters

$$\text{In principle } \Theta(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2$$

$$= \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Taking partial derivatives

$$\frac{\partial \Theta}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial \Theta}{\partial \beta_1} = -2 \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i)$$

then solve the Normal Equations (equate partial derivative to 0)

Solve for β_0, β_1

We use b_0 and b_1 to denote the student

$$b_1 = \frac{\sum (x_i y_i) - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}$$

23/10/13

Regression

Estimation of β_0 and β_1

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

All the Y_i 's and ϵ_i 's have the same variance σ^2

What if σ^2 is unknown?

Estimate σ^2

When we had a single population

$$X_i \sim N(\mu, \sigma^2)$$

We estimate μ by \bar{x}

$$\text{Estimate } \sigma^2 \text{ by } \hat{\sigma}^2 = S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

$(n-1)$ is known as degrees of freedom

Estimating σ in SLR

$$Y \sim N(),$$

$$E[Y_i] = \beta_0 + \beta_1 X_i$$

$$\text{Var}[Y_i] = \sigma^2$$

$$\Rightarrow Y_i \sim N((\beta_0 + \beta_1 X_i), \sigma^2)$$

Firstly estimate β_0 and β_1 . We call them b_0 and b_1 .

We also know that a fitted \hat{Y}_i is defined as

$$\hat{Y}_i = b_0 + b_1 X_i$$

Define $SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ SSE = error sum of squares
(residual sum of squares)

$$= \sum_{i=1}^n (Y_i - (b_0 + b_1 X_i))^2$$

$$= \sum_{i=1}^n e_i^2 \text{ where } e_i \text{ is residual } (Y_i - \hat{Y}_i)$$

4.

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\begin{aligned} b_0 &= \frac{1}{n} (\sum y_i - b_1 \sum x_i) \\ &= \bar{y} - b_1 \bar{x} \end{aligned}$$

Some notation

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

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

\hat{y}_i → fitted value for the i^{th} observation x_i
 y_i → observed value

$$x_i = 0.17 \quad \hat{y}_i = -259.63 + 3721.02 \cdot 0.17 \\ = 372.97$$

Interpretation

The mean if many diamonds of carat $x_i = 0.17$ were sold under similar conditions as this data set, the average price of them would be 372.95 \$.

Of course the price of a single diamond with carat = 0.17 would not be 372.95 \$.

This is obvious, since there is chance variation in the pricing strategy.

Residuals

The i^{th} residual is the difference between the fitted value and the observed value of y_i .

$$e_i = \hat{y}_i - y_i$$

The residuals are the vertical deviations of y_i from the fitted regression line.

Difference between e_i and ϵ_i :

$$e_i = y_i - \hat{y}_i \quad \epsilon_i = y_i - (\beta_0 + \beta_1 x_i)$$

- (1) Known
(2) Difference between y_i and \hat{y}_i .
(1) Unknown quantity (theory \rightarrow hypothesis)
(2) Vertical deviation between y_i and "unknown" regression line

17/10/13 Regression

Estimating σ^2 by SLR
 $\text{Var}(\epsilon_i) = \sigma^2$ } unknown parameter, need to estimate
 $\text{Var}(y_i) = \sigma^2$

$y \sim (\mu, \sigma^2)$ both μ and σ^2 unknown
Mu estimated by \bar{y}

$$s^2 = \frac{\sum (y_i - \bar{y})^2}{n-1}$$

← Degrees of Freedom
- In estimating μ , we have lost 1 degree of freedom
- number of parameters that you have estimated

SLR simple linear rego

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

Since we have to estimate β_0 and β_1 before we can go on to estimate σ^2 , we lose 2 degrees of freedom

Define $y_i - \bar{y} = \epsilon_i$

The appropriate sum of squares is: SSE
$$\text{SSE} = \sum_{i=1}^n (\epsilon_i)^2$$

$$= \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

The estimate of σ^2 also known as MSE is

$$\hat{\sigma}^2 = \text{MSE} = \frac{\text{SSE}}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}$$

It can be shown that

$$E(\text{MSE}) = \sigma^2$$

For the diamonds example:

$$\hat{y}_i = -259.36 + 3721.03x_i \quad \text{regression equation}$$

$$\begin{array}{c} \text{SSE} = \\ y_i & x_i \\ 351 & 0.17 \\ 316 & 0.15 \end{array}$$

$$\text{SSE} = [351 - (-259.36 + 3721.03 \times 0.17)]^2 + \dots + [316 - (-259.36 + 3721.03 \times 0.15)]^2 = 46630$$

$$\text{MSE} = \frac{\text{SSE}}{n-2} = \frac{46630}{48-2} = 1013.82$$

Inference about a and b .
Fundamentally interested in distributions, t-test and confidence intervals.

Refer : If $x_i \sim N(\mu, \sigma^2)$
Then $\bar{x} \rightarrow \text{estimate of } \mu$
Remember $\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

So we look at sampling distribution of b_0 and b_1 .

Inference of b_1

Point estimate for b_1 is $b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$

23/10/13

Regression

ST2002

It can be shown $E(b) = \beta$

$$\text{Var}[b_i] = \frac{\sigma^2}{\sum(x_i - \bar{x})^2} = \frac{\sigma^2}{S_{xx}}$$

If σ^2 is known, $b_i \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$

However, if σ^2 is unknown, we need to estimate $\text{Var}[b]$ and
that is $\text{Var}[b_i] = S_{b_i}^2 = \text{MSE}$

standard error s.e. is the square root of $\text{Var}[b_i]$

$$\text{s.e.}(b_i) = \sqrt{\frac{\text{MSE}}{S_{xx}}} \quad \text{packages give out this}$$

For the diamond example:

$$\text{s.e.}(b_1) = \sqrt{\frac{101383}{0.156}} = 81.78$$

Hypothesis testing for β_1 / b_1 :

$$y_i = \beta_0 + \beta_1 x_i \quad \beta_1 \text{ has a physical meaning.}$$

Refers to change in y for unit change in x .

Experts think change in y unit change in x should be 7, whereas when we estimate β_1 from data it is 5

Test: $H_0: \beta_1 = 7$ vs $H_1: \beta_1 \neq 7$ at 5% DF

Golden rule for t-test

$$t_{\text{calc}} = \frac{(\text{estimator}) - (\text{value from } H_0)}{\text{s.e.}(\text{estimator})}$$

(Confidence Interval) $(\text{estimator}) \pm t_{\text{critical}} \text{s.e.}(\text{estimator})$

↑
are you get from table?

2. Regression

Diamonds data set.

$$H_0: \beta_1 = 0 \quad H_1: \beta_1 \neq 0$$

$$t_{\text{calc}} = \frac{b_1}{\text{se.}(b_1)} = \frac{372.02}{81.79}$$

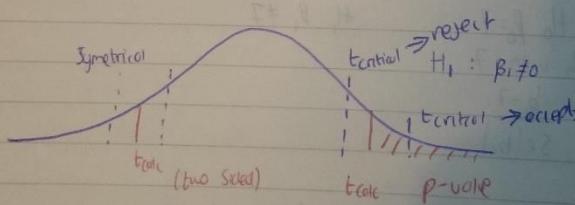
$$\text{From the } t\text{-table, } t_{0.975, 46} = 2.011$$

Since $|t_{\text{calc}}| > t_{\text{critical}}$ we reject null hypothesis

Interpretation: We can reject the null hypothesis that the rate of change of price of diamonds with respect to carats is equal to zero.

The test also provides p-values

P-value: Given the null hypothesis is true, the p-value measures the probability of observing a more extreme value than t_{calc} (in either direction)



If $p\text{ value} < \alpha \Rightarrow t_{\text{critical}} > t_{\text{calc}} \Rightarrow \text{reject}$

24/10/13 Regression

Inference on β_1

$$se(b_1) = \sqrt{\frac{MSE}{\sum \epsilon_i^2 (x_i - \bar{x})^2}} \text{ s.e.} = \sqrt{\frac{MSE}{S_{xx}}}$$

$\frac{b_1 - \beta_1}{se(b_1)} \sim t \text{ distribution with } (n-2) \text{ d.f.}$

Test. $H_0: \beta_1 = 0$ vs. $H_1: \beta_1 \neq 0$

$$t_{\text{calc}} = \frac{b_1 - \beta_1}{s.e.(b_1)} = \frac{b_1}{s.e.(b_1)} \quad (\beta_1 = 0 \text{ under } H_0)$$

~~$b_1 = \beta_1$~~

Example. $H_0: \beta_1 = 7$ $H_1: \beta_1 \neq 7$

Replace β_1 with 7.

$$t_{\text{calc}} = \frac{b_1 - 7}{s.e.(b_1)}$$

For a significance level α the rule is:

- If $|t_{\text{calc}}| \leq t_{\text{critical}}$, then accept H_0 .
- Reject otherwise
- $t_{\text{critical}} \Rightarrow$ value you get from t-tables at $\alpha=0.05$ df = n-2.

- Most often 5%.

$$t_{\text{critical}} = t_{[\alpha/2, n-2]} \text{ area df.}$$

4 REGRESSION

We are very interested in
 $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

$$[Y_i = \beta_0 + \beta_1 x_i + \epsilon_i]$$

$\beta_0 = 0$ no relationship

Note that $\beta_0 = 0$
⇒ no association between X and Y

So, under the null hypothesis, the sampling distribution of
 $b_1 - \beta_1 \sim t$ distribution with $(n-2)$ d.f.
 $SE(b_1)$

The $(n-2)$ d.f. is borrowed from $SE(b_1)$

For our example

$$p\text{-value} < 2 \times 10^{-16}$$

$$\alpha = 0.05$$

$$p < \alpha \Rightarrow \text{reject } H_0$$

Confidence Intervals for β_1 :

$$\text{estimator} \pm t_{\text{critical}} \times \text{s.e.}(\text{estimator})$$

The 95% CI for β_1 is

$$\hat{\beta}_1 \pm t_{0.975, n-2} \times \text{s.e.}(\hat{\beta}_1)$$

$$t_{0.975, 46} \text{ s.e.}(\hat{\beta}_1)$$

$$\text{Example: } 3.721.03 \pm 2.011 \times 81.79$$

$$= 3.721.03 \pm 164.48$$

$$95\% \text{ CI} = (3556.54, 3885.5)$$

This is an example of Interval estimation

Interpretation:

Firstly a 95% CI does not mean that there is a 95% chance that the true β_1 lies in the interval

(It just means that if we were to compute intervals in this way from similar data sets, then 95% of those intervals will contain the true slope β_1)

Correspondence between t-test and CI:

Both should give the same inference.

\Rightarrow t-test rejects H_0 then CI should not contain 0.
(for $H_0: \beta_1 = 0$) if $[H_0: \beta_1 = 0]$ should not contain 0]

6/13

Regression

Testing of hypothesis and C.I.

$$H_0: \mu = \mu^* \quad \text{vs} \quad H_1: \mu \neq \mu^*$$

One approach $\rightarrow t\text{-table or } p\text{-value}$ Other is to compute C.I. for μ Check if μ^* fall within the interval or notInference for β_1

$$\text{C.I.: } b_1 \pm t_{\text{crit}} s.e.(b_1)$$

$$(b_1 - t_{\text{crit}} s.e.(b_1); b_1 + t_{\text{crit}} s.e.(b_1))$$

$$\text{If } H_0: \beta_1 = 7 \text{ (say)} \quad \text{vs} \quad H_1: \beta_1 \neq 7$$

If 7 falls within C.I we accept H_0 Reject O/W.Inference about the intercept β_0 Mostly we are interested when $X=0$, makes sense

$$\text{point estimate of } \beta_0 \quad b_0 = \bar{y} - b_1 \bar{x}$$

$$\text{It can be shown } E(b_0) = \beta_0$$

$$\text{Var}(b_0) = s^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right]$$

$$\text{where } s_{xx} = \sum (x_i - \bar{x})^2$$

The sampling distribution of b_0 is

$$b_0 \sim N(\beta_0, s^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right])$$

$$\text{s.e.}[b_0] \quad (\text{square root of estimated variance}) = \sqrt{\text{Mse} \left[\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right]}$$

\downarrow
Estimated value

Regression

From the diamond data set

$$\text{s.e.}[b_0] = \sqrt{MSE \left[\frac{1}{n} + \frac{\sum x_i^2}{S_{xx}} \right]} \\ = \sqrt{101382 \left[\frac{1}{48} + \frac{0.20122}{0.1516} \right]} \\ = 17.32$$

The Sampling distribution of $\frac{b_0 - \beta_0}{\text{s.e.}(b_0)} \sim t$ distribution with $(n-2)$ d.f

Hypothesis testing:

$$H_0: \beta_0 = 0 \quad \text{vs} \quad H_1: \beta_0 \neq 0$$

Test Statistic $t_{\text{calc}} = \frac{b_0}{\text{s.e.}(b_0)}$

For α level of significance

- * If $|t_{\text{calc}}| \leq t_{\text{critical}}$ accept H_0 .
- Reject otherwise.

$$t_{\text{calc}} = \frac{-25.43}{17.32} = -14.94 \quad [\text{Intercept is negative}] \text{ fail.}$$

$$\text{For } \alpha = 0.05, t_{\text{critical}} = t_{0.95, 46} = 2.011$$

Since $|t_{\text{calc}}| \gg t_{\text{critical}}$ we reject the null hypothesis.

The true intercept is not equal to zero, in light of the data at hand

R gave us p-value of 2×10^{-16} which is < 0.05 .

C. for β_0

100(1- α)% CI for β_0 :

$$\beta_0 \pm t_{\text{critical}} \text{s.e.}[b_0]$$

If you consider the example

$$-25.43 \pm 2.011 \times (17.32) \\ \Rightarrow -25.63 \pm 34.38 \Rightarrow (-294.46, 2248)$$

13. Regression

(3)

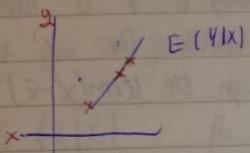
$$(-294.46, -224.8)$$

Since this does not contain 0, we can get the same conclusion as before: Reject H_0 .

Interpretation of C.I.: very important.

95% of time when calculating C.I., it will contain true parameters.

Inference on $E[Y]$ for $X = X$ some value,



One of the major goals in regression is to make inference about the mean or the distribution of Y at X .

We want to construct point and interval estimates for $E(Y|X = \text{some value})$

Let x' be some value of X for which we want to estimate the mean response

$$\text{known } E(Y|x) = \beta_0 + \beta_1 x'$$

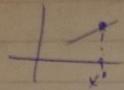
The point estimate, say \hat{y}' , of $E(Y|x')$

$$\hat{y}' = \hat{\beta}_0 + \hat{\beta}_1 x'$$

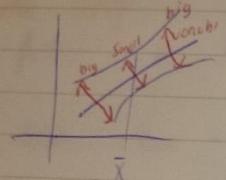
$$\begin{aligned} E(\hat{y}') &= E[\hat{\beta}_0 + \hat{\beta}_1 x'] \\ &= E[\hat{\beta}_0] + E[\hat{\beta}_1 x'] \\ &= E[\hat{\beta}_0] + E[\hat{\beta}_1] x' \quad (x' \text{ is a constant}) \\ &= \beta_0 + \beta_1 x' \quad \rightarrow \text{true value} \end{aligned}$$

4

$$\text{Var}(\hat{y}) = r^2 \left[\frac{1}{n} + \frac{(x' - \bar{x})^2}{s_{xx}} \right]$$



In term of a dataset, $x' = 0.33$



Note that variability of \hat{y}' is largely affected by how much far x' is from \bar{x} , through the term $(x' - \bar{x})^2$

Hence further away causes greater variability
Then values closer to \bar{x}
(Talking about variability of a point on the regression line)

$$\hat{y}' \sim N \left((\beta_0 + \beta_1 x'), \sigma^2 \left[\frac{1}{n} + \frac{(x' - \bar{x})^2}{s_{xx}} \right] \right)$$

$$s.e(\hat{y}') = \sqrt{\text{MSE} \left[\frac{1}{n} + \frac{(x' - \bar{x})^2}{s_{xx}} \right]}$$

Finally, sampling distribution of $\frac{\hat{y}' - (\beta_0 + \beta_1 x')}{s.e(\hat{y}')} \sim t \text{ distribution with } n-2 \text{ degrees freedom}$

HYPOTHESIS TESTING

$$H_0: \hat{y}' = y \quad H_1: \hat{y}' \neq y$$
$$t_{\text{calc}} = \frac{\hat{y}' - y}{s.e(\hat{y}')}$$

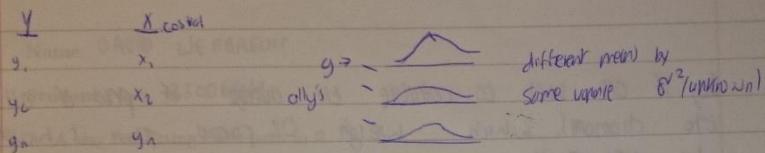
$$\text{CI: } \hat{y}' \pm \text{critical } s.e(\hat{y}')$$

10/B. Regression. ①

Assignment Q3c write down the values of SSE, MSE, and RSE

$$RSE = \sqrt{MSE}$$

residual standard error



If you can join these means by a straight line, then do linear regression

$$E(Y|X_i) = \beta_0 + \beta_1 x_i$$

$$\Rightarrow Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$= E(Y|X_i) + \epsilon_i$$

Estimate this line $\hat{Y} = b_0 + b_1 X$
estimated by b_0 and b_1

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \text{Sampling variation of } b_1 \text{ and } b_0$$

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

$$\text{Estimate of a point on a line } \hat{Y} = b_0 + b_1 x_i$$

$$E(\hat{Y}|X_i)$$

For some value $x = x'$, estimate of $E(Y|X) = \hat{Y}'$

$$\hat{Y}' = b_0 + b_1 x'$$

$$\text{s.e.}[\hat{Y}'] = \sqrt{MSE \left(\frac{1}{n} + \frac{(x' - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right)}$$

Diamond dataset

$$\bar{x} = 0.2012$$

$$1. Y' = 0.2 \cdot \text{plug in?}$$

$$\hat{Y}' = -2593 + 3721.43 \times 0.2 = 484.57$$

2.

fitted value here refers to the average price of diamonds with weight 0.2 carats.
If you collected data on every 0.2 carat diamond. This value is the average of all 0.2 carat diamonds.

So, if one was to calculate the average of price of all the diamonds which weigh 0.2 carat. Then it would be 484.57

$$S.E.(\hat{\gamma}') = \sqrt{MSE \frac{1}{n} + \frac{(x_{mean})^2}{\sum(x_i - \bar{x})^2}}$$

$$= \sqrt{1013.83 \left[\frac{1}{48} + \frac{(0.2 - 0.2002)^2}{0.1516} \right]}$$

$$= 4.609 \text{ closer to mean}$$

$$CI: \hat{\gamma}' \pm t_{0.975, 46} S.E(\hat{\gamma}') = (475.307, 493.85)$$

$$\bar{x} = 0.33$$

$$\hat{\gamma}' = 968.31$$

$$S.E. \hat{\gamma}' = \sqrt{1013.83 \left[\frac{1}{48} \frac{(0.23 - 0.2002)^2}{0.1516} \right]} \\ = 11.267 \text{ away from mean}$$

$$CI: \hat{\gamma}' \pm t_{0.975, 46} S.E(\hat{\gamma}') = (945.646, 990.9739)$$

diamond.lm \$ coefficients \rightarrow gives intercept b_0 and B_1

$b_0 = \text{diamond.lm } \$ \text{ coefficients}[1]$ First position.

RSE = summs.diamond.lm $\$/\sigma$

getting mean = $\frac{\text{mean}(\text{carots})}{\text{mean}(\text{price})}$

working view . name = c(1, 2, 3, 4)

ST2002: Introduction to regression

Computer laboratory 5

Name: DAVID WEITBRECHT

Student No.: 12300644

A way for setting one's path in R.

Regression

13/11/13 C.I. / t-test / Sampling Distribution

$$\text{of } E(Y|X=x) = y'$$

Prediction of a new observation

We consider prediction of a new observation y corresponding to some x .

Let us denote a new x as x' (x takes the value of x'), and we call the new observation as y_{new} .

y_{new} is a single observation corresponding to x' .

$$\begin{array}{c} \text{1. } y'_{\text{new}} = E(Y|X=x') = y_{\text{new}} \\ \hline x' \\ 2. z \sim N(\mu, \sigma^2) \\ E[z] = \bar{z} \sim N(\mu, \frac{\sigma^2}{n}) \end{array}$$

Big distinction between $E(Y|X=x) = y'$ in the previous section and our new response y_{new} .

First case refers to mean of distribution of Y for some $X=x'$.

In the later case, we predict an individual outcome drawn from the distribution of Y for $X=x'$.

\Rightarrow Greater variability for the 2nd case.

Point estimate: We predict y_{new} by its estimate \hat{y}_{new} , given by

$$\hat{y}_{\text{new}} = b_0 + b_1 x'$$
 - point on regression line.

Note that the point estimate is same as before

Properties

$$E(\hat{y}_{\text{new}}) = \hat{y}'_{\text{new}}$$
$$\text{Var}(\hat{y}_{\text{new}}) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x' - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

$$\text{For } E(y|x=x') = y'$$

$$\text{Var}(y') = \sigma^2 \left[\frac{1}{n} + \frac{(x' - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

Note that, $\text{var}(\hat{y}_{\text{new}})$ has 2 components,

1. The variance of the sampling distribution of fitted value y'

2. The variance of distribution of y at some x .

$$\begin{aligned}\text{Var}(\hat{y}_{\text{new}}) &= \sigma^2 + \sigma^2 \left[\frac{1}{n} + \frac{(x' - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] \\ &= \text{var}(y') + \text{var}(y)\end{aligned}$$
$$\text{Var}(y_{\text{new}}) \geq \text{var}(y')$$

Estimate of variance is:

$$\begin{aligned}\text{Var}(\hat{y}_{\text{new}}) &= s^2 / \hat{y}'_{\text{new}} \\ &= \text{MSE} \left[1 + \frac{1}{n} + \frac{(x' - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]\end{aligned}$$

$$\text{s.e.}(\hat{y}_{\text{new}}) = \sqrt{\text{MSE} \left[1 + \frac{1}{n} + \frac{(x' - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]}$$

$$x' = \bar{x} \Rightarrow \text{Var}(\hat{y}_{\text{new}}) = \text{MSE} \left[1 + \frac{1}{n} \right]$$

Intervals are now called prediction intervals!

$$\hat{y}_{\text{new}} \pm \text{t critical} \text{s.e.}(y_{\text{new}})$$

$$\hat{y}_{\text{new}} \pm \text{t critical} \sqrt{\text{MSE} \left[1 + \frac{1}{n} + \frac{(x' - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]}$$

$\text{lm}(\text{price} \sim \text{carats})$

$\text{lm}(\text{formula}, \text{dataframe})$

diamond.dataframe \rightarrow 2 columns price and carat

price = diamond.dataframe \$ price;

carat = diamond.dataframe \$ carat.

$\text{lm}(\text{price} \sim \text{carats})$

could also use $\text{lm}(\text{price} \sim \text{carats}, \text{diamond.dataframe})$

coef = diamond.lm\$coefficients

b0 = coef[1]

b1 = coef[2]

Significance code

* 0.05

** 0.01

*** 0.001

• 0.1

(Confidence Intervall for β_0 and $\beta_1 \Rightarrow \text{confint}(\text{lm object})$)

confint(lm object, level=0.95)

= confint(price, 0.95)

new.x = diamond \$ x = r(x1, x2, x3)

20/11/12

Regd's

2

The plot suggests a strong linear relationship between size of company and delay time. A ~~strong~~ high correlation coefficient of -0.84 supports this.

There is a causal relationship between the 2 variables, linear regression can be used.

b) Since Size = x and Delay = y
The "true" equation is given by

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i=1, \dots, n$$

where β_0 and β_1 are parameters and ϵ_i are the errors

The fitted regression line is:

$$\hat{y}_i = 25.47 - 0.18x_i$$

Assumption:

- i. The predictor variable X is assumed to be constant.
- ii. The observations y_i are independent.
- iii. For a given x_i , the corresponding y_i is normally distributed.
- iv. All y_i 's have the same variance σ^2 .
- v. The means of y_i given x_i can be joined by a straight line.

$$E(y_i | x_i) = \beta_0 + \beta_1 x_i$$

β_1 = Slope parameter. Signifies the change in Y for unit change in X .

β_0 = Intercept. Value of Y for $X=0$

c) For β_0

The hypothesis tested is: $H_0 : \beta_0 = 0 \quad H_1 : \beta_0 \neq 0$

The calculated t value is $t_{cal} = 11.687$.

From table $t_{0.975, 18} = 2.306$

$$\rightarrow t_{(1-\frac{\alpha}{2}, n-2)}$$

Since $|t| > t_{0.975, 18}$ we reject null hypothesis

3. Regression.

Example

Carat = 0.2

Instead of estimating interval for the true mean price of diamond that weigh 0.2 carat, we are interested in estimating the interval within which the true price of a single diamond of weigh 0.2 (r)

$$\text{For } x^1 = 0.2 \quad \hat{Y}_{\text{new}} \pm t_{0.975, 40} \text{ s.e.}(\hat{Y}_{\text{new}}) \\ \Rightarrow (4114.92, 549.34)$$

$$\text{For } x^1 = 0.33 \quad = \hat{Y}_{\text{new}} \pm t_{0.975, 41} \text{ s.e.}(\hat{Y}_{\text{new}}) \\ \Rightarrow (40032, 1036.30)$$

$$(1 \text{ for } x^1 = 0.2 \quad (475.31, 4485))$$

$$(1 \text{ for } x^1 = 0.33 \quad = (945.65, 990.98))$$

Prediction interval $\rightarrow P.I.$

Regression R

Vector

{2, 5, 7}
((2, 5, 7)

Matrix

1 4 7
2 5 8
3 6 9

Data frame

A B how to have name

1 4
2 5
3 6

A = c(1, 2, 3)

A = data.frame(A = c(1, 2, 3))

function to read data from file.

read.table (.txt, csv etc)

read.table (filename, header, sep)

header = true

excel

Price Cost

row 1 : :
row 2 : :

data.read = read.table (filename, header = TRUE, sep = " ")

If you say false it will assign names on its own

read.csv (filename, header, sep, col.names)
(col.names) = c ("A", "B")

CSV = comma separated values

Regression

26/11/13 ASSymetry | Answers

Q. ii) correlation coefficient $\frac{\frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} \sqrt{\frac{1}{n-1} \sum (y_i - \bar{y})^2}}$

 $= \frac{A}{\sqrt{B} \sqrt{C}}$

$$\begin{aligned} A &= \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n-1} \sum (x_i y_i - \bar{x} \bar{y} - x_i \bar{y} + \bar{x} y_i) \\ &= \frac{1}{n-1} \sum (x_i y_i - \sum x_i \bar{y} - \sum y_i \bar{x} + \sum \bar{x} \bar{y}) \\ &= \frac{1}{n-1} \sum x_i y_i - n \bar{x} \bar{y} - n \bar{y} \bar{x} + n \bar{x} \bar{y} \\ &= \frac{1}{n-1} \sum x_i y_i - n \bar{x} \bar{y} \\ &= \frac{1}{n-1} \sum x_i y_i - n (\frac{1}{n} \sum x_i) (\frac{1}{n} \sum y_i) \\ &= \frac{1}{n-1} \sum x_i y_i - \frac{\sum x_i \sum y_i}{n} \end{aligned}$$

$$\begin{aligned} B &= \frac{1}{n-1} \sum (x_i - \bar{x})^2 \\ &= \frac{1}{n-1} \sum (x_i^2 + \bar{x}^2 - 2x_i \bar{x}) \\ &= \frac{1}{n-1} \sum x_i^2 + n \bar{x}^2 - 2 \bar{x} \sum x_i \\ &= \frac{1}{n-1} \sum x_i^2 + n \bar{x}^2 - 2 \bar{x} (\bar{x}) \\ &= \frac{1}{n-1} \sum x_i^2 - n \bar{x}^2 \\ &= \frac{1}{n-1} \left(\sum x_i^2 - n \left(\frac{1}{n} \sum x_i \right)^2 \right) \\ &= \frac{1}{n-1} \left(\sum x_i^2 - \frac{(\sum x_i)^2}{n} \right) \\ &= \frac{1}{n-1} \left(\sum y_i^2 - \frac{(\sum y_i)^2}{n} \right) \end{aligned}$$

From data: $n=10 \quad \sum x_i = 435 \quad \sum x_i^2 = 2777$
 $\sum y_i = 175 \quad \sum y_i^2 = 3481 \quad \sum x_i y_i = 5990$

Plug into formula = -0.8429

2.

E. a) for mean respond size = 40

$$x'40 \quad y' = b_0 + b_1 x'$$

$$= 18.14149$$

$$\text{s.e.}(y') = \sqrt{MSE \left[\frac{1}{n} \frac{(x_i - \bar{x})^2}{\varepsilon(x_i - \bar{x})^2} \right]}$$

$$= 1.23896$$

$$(I) \quad y' \pm t_{0.975, 8} \text{s.e.}(y') = (15.2844, 20.9985)$$

$$+ 8.741024 \pm$$

2A. $\sum x_i = 27 \quad \sum x_i y_i = 160$
 $\sum x_i^2 = 330 \quad \sum y_i = 210$
 $n = 16 \quad SSE = 1750$

Regression model 1) error: $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

$$E(y_i | x_i) = \beta_0 + \beta_1 x_i$$

$$\hat{\beta}_1 = b_1 = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = 4.6605$$

$$\hat{\beta}_0 = b_0 = \bar{y} - b_1 \bar{x}$$

$$= \frac{1}{n} (\sum y_i - b_1 \sum x_i) = 5.26$$

b) Hypothesis $H_0: \beta_1 = 6$ vs $H_1: \beta_1 \neq 6$

$$t = \frac{b_1 - 6}{\text{s.e.}(b_1)}$$

$$\text{s.e.}(b_1) = \sqrt{MSE \left(\frac{1}{n} + \frac{\bar{x}^2}{\varepsilon(x_i - \bar{x})^2} \right)} = \sqrt{\frac{MSE}{\varepsilon(x_i - \bar{x})^2}}$$

III/D 3 Regressions

$$S.e(b_1) = 0.6624$$

$$t = \frac{b_1 - b}{S.e(b_1)} = \frac{4.66 - 6}{0.6624} = -2.52$$

$$t_{0.975, 14} = 2.145 \quad |t| < t_{0.975, 14} \text{ accept null hypothesis}$$

accept H₀ that slope or parameter 1 is equal to 0

SOLUTION:

```
box.office.dataframe = read.csv("boxoffice.csv", header=TRUE)
```

```
box.office.lm = lm(currentweek ~ LastWeek, boxoffice.dataframe)
```

```
coef(box.office.lm)
```

```
summary.box.office = summary.lm(box.office.lm)
```

```
summary.box.office $ coefficients
```

RSE = summary.box.office \$ sigma

MSE = RSE ^ 2

```
confint(box.office.lm)
```

```
fitted(box.office.lm)[1]
```

```
fitted(box.office.lm)[length(box.office.dataframe) $ currentweek]
```

```
new.boxoffice = data.frame(10week = 400000)
```

```
predict.lm(box.office.lm, newdata = new.boxoffice, se.fit=T, interval=c("confidence"))
```

```
predict.lm(box.office.lm, newdata = new.boxoffice, se.fit=T, interval=c("prediction"))
```

7/11/13 Regression.

Problem in Assignment.

Q1E $H_0: \beta_0 = 0$

$H_1: \beta_0 \neq 0$ β_0 and β_1 cannot be in same hypotheses!

Should be $H_0: \beta_0 = 0$ $H_0: \beta_1 = 0$

$H_1: \beta_0 \neq 0$ $H_1: \beta_1 \neq 0$

β_0, β_1 are theoretical values of parameters

b_0, b_1 are estimates of the parameters

$s.e.(b_1)$ wrong

s.e.(b_0) or $s.e.(b_1)$ correct.

2C. If $x_i = 0$ Yield = β_0

According to the question the yield will be 0 tonnes TEST.

$H_0: \beta_0 = 0$ vs $H_1: \beta_0 \neq 0$

$$s.e.(b_0) = \sqrt{MSE \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right]}$$

$$= \sqrt{\frac{SSE}{n-2} \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right]}$$

$$= 3.01$$

$$t = \frac{b_0}{s.e.(b_0)} = \frac{5.26}{3.01} = 1.74$$

$$t_{0.975, 14} = 2.145$$

$$\Rightarrow |t| < t_{0.975, 14} \text{ Do not reject the null hypothesis}$$

The analyst was right in pointing out that allocating 0 hectare could be the yield to be 0 tonnes

20/11/13 3 Regression

t σ ∞

Thus, we can conclude that the true intercept is not equal to 0, given the data.

In other words the mean delay time for a company of size 0 is not equal to 0; as seen from the test.

For B_1

$$t = -4.432$$

$$t_{0.975, 18} = 2.306$$

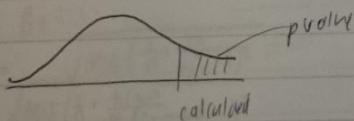
$\Rightarrow |t| > t_{0.975, 18}$ we reject the null hypothesis that the slope is 0.

In other words, the change in ~~mean~~ delay time for unit change in size is not equal to 0.

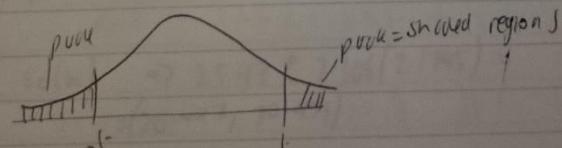
P-value:

Given that the null hypothesis is true, p-value measures the probability of observing a more extreme value than that of the calculated one.

One Sided test:



Two Sided test:



It is the area in tails, more extreme than the calculated test statistic. If p-value is less than α , we reject the null hypothesis.

$$|t| > t_{\text{critical}} \equiv p < \alpha \quad \text{In this case, p-value for both } B_1 \text{ and } B_0 \text{ is } \alpha$$

| t_1 | > $t_{0.975, 18}$ hence we reject null hypothesis

1/13.

Regression Assignment Solution

$|t| > t_{\text{critical}} \Leftrightarrow$ related to p-value (< 0.05)

If $t = \text{reject}$ p will also reject

$$1.8. \text{ Estimate of } \sigma^2 = \text{MSE} = \frac{\text{SSE}}{n-2}$$

$$\text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2$$

$$\begin{aligned} \frac{y_i}{29} &= \frac{\hat{y}_i}{29} = 25.47 \pm 0.08x_i & y_i - \hat{y}_i \\ \vdots & \rightarrow \times & \text{BSE} = 0.08 \quad F.S. = 1.15 \quad A.S. \\ \vdots & \rightarrow \times & 0.08 = 0.08 \\ & \vdots & 0.08 = 0.08 \\ & = 121.14071 & \text{S.E.} = 0.08 \\ \text{MSE} &= \frac{121.14071}{(29-2)} \approx 15.14 & \end{aligned}$$

R output:

$$\text{RSE} = 3.891 \quad \text{RSE}^2 = \text{MSE} = (3.891)^2 = 15.139 = 15.14$$

C.I. for β_0 :

$$\begin{aligned} \text{s.e.}(\beta_0) &= \sqrt{\text{MSE} \left(\frac{1}{n} + \frac{x_i^2}{\sum(x_i - \bar{x})^2} \right)} \\ &= \sqrt{\text{MSE} \left(\frac{1}{29} + \frac{1/(29)^2}{\sum(x_i - 25.47)^2} \right)} \\ &= 2.1795 \end{aligned}$$

$$\begin{aligned} \beta_0 &\pm t_{\text{critical}} \text{s.e.}(\beta_0) \Rightarrow 25.47 \pm 2.306(2.1795) \\ &= (20.447, 30.499) \end{aligned}$$

C.I. for β_1 :

$$\text{s.e.}(\beta_1) = \sqrt{\frac{\text{MSE}}{\sum(x_i - \bar{x})^2}} = s_{x_1}$$

$$\begin{aligned} \text{s.e.}(\beta_1) &= 0.04136 \\ \beta_1 &\pm t_{\text{critical}} \text{s.e.}(\beta_1) = (-0.279, -0.0879) \end{aligned}$$

Anova - Analysis of Variance

Total Variation = Systematic Variation + Chance Variation.
Aim to max systematic and min chance variation.

Given any random variable y_1, \dots, y_n the variability is denoted by variance.

$$\text{Var}(y) = \frac{1}{n-1} \sum (y_i - \bar{y})^2$$

In other words, a measure of variability is $\sum (y_i - \bar{y})^2$
→ refers to total variation.

In regression context, this is known as Total Sum of Squared or SSTO

$$= \sum (y_i - \bar{y})^2 = \text{Unadjusted sum of squares}$$

Note: Unadjusted sum of squares = $\sum y_i^2$

$$\text{SSTO} = 0$$

⇒ All y_i 's are equal to \bar{y} (all same number e.g. 5)

The greater the value of SSTO, the greater the variability of the data.

Aim:

Basic aim is to partition total variability (SSTO) into 2 components:
① Variability due to model
② Variability caused due to Error or chance

① Variability explained by regression equation
This is used to find out if our model was good enough in explaining variation in the data

② We have done this already using Error Sum of Squares SSE

3.
7/11/13 Regression

The value quantified the variability of data around the fitted line

↳ Refer to chance variation

The rest of the variability is defined as Regression

sum of Squared or SSR

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

The greater the SSR, the better is the regression eqn in accounting for variability in the data

$$SSTD = SSR + SSE$$

$$\text{Total Variation} = \text{Systematic variation} + \text{chance variation}$$

* If SSE = 0, means all the data lie perfectly on the regression line

$$SSE = \sum (y_i - \hat{y}_i)^2 \text{ equals } 0 \text{ when } y_i = \hat{y}_i$$

Partitioning the total sum of squares

$$y_i - \bar{y} = \underbrace{\hat{y}_i - \bar{y}}_{\text{due to regn}} + \underbrace{y_i - \hat{y}_i}_{\text{chance}}$$

Squaring and summing them on both sides gives us

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

The intersection/cross-product then is equal to 0.

To prove this we the following properties:

① The regression line can be written as $\hat{y}_i = \bar{y} + b_1(x_i - \bar{x})$

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

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

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

4

2 The sum of observed values y_i is equal to sum of fitted values.

$$\sum y_i = \sum \hat{y}_i$$

3 ~~inx~~ If $e_i = y_i - \hat{y}_i$

$$\sum e_i = 0$$

4 $\sum \hat{y}_i e_i = 0$

Regression

Newdat = data.frame (r nsize = c(100, 20, 30))

$$SST_0 = \text{sum}(y_i - \bar{y})^2$$

term = due to regression + error sum of square

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

better model is SSE is smaller

The right hand side is used to judge new model

Also used to compare with other model

Degrees of Freedom

$\sum (g_i - \bar{y})^2$ has $(n-1)$ df SSE only 1 less estimating \bar{y} .

$\sum (y_i - g_i)^2$ has $(n-2)$ df SSE

$$df[SST_0] = df[SSR] + df[SSE]$$

$$\Rightarrow df[SSR] = 1.$$

Partitioning of degrees of freedom

Mean Sum of Squares

This is defined as the ratio of sum of square and the corresponding d.f.

$$MSE = \frac{SSR}{n-2}$$

Mean squared error.

$$MSR = \frac{SSR}{n-2}$$

Mean squared regression

$$\begin{array}{lll}
 \text{Source of Variation} & \text{SS} & \text{DF} \\
 \text{Due to regression} & \text{SSR} = \sum (y_i - \bar{y})^2 & 1 \\
 \text{Error} & \text{SSE} = \sum (y_i - \hat{y}_i)^2 & n-2 \\
 \text{Total} & \text{SSTO} = \sum (y_i - \bar{y})^2 & n-1
 \end{array}$$

Mean Square MS / F

$$MSR = \frac{SSR}{1}$$

$$MSE = \frac{SSE}{n-2}$$

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

The F value is used to test the following hypothesis:

$$H_0: \beta = 0 \quad \text{vs} \quad H_1: \beta \neq 0$$

Tells whether it is meaningful to do a regression model

If $F_{\text{calc}} \leq F_{1-\alpha, df_1, df_2}$ then accept H_0

df_2	1	2	3	$F(0.95, 1, 3)$
1	"	"	"	
2	"	"	"	
3	"	"	"	

Diamond	Dataset	Above		
	SS	DF	MS	F
Regression	2144218	1	2144218	
Error	46636	46	1014	2114.99
Total	2145232	47	-	

$$F(0.95, 1, 46) = 7.95.$$

Since $F(2114.99) \geq F(0.95, 1, 46)$ we reject H_0

R command: lm (dataset, lm)

$H_0: \beta = 0$ vs $H_1: \beta \neq 0$
 Relationship between total and f-test has been used

2 Residual Analysis

Basis idea is to check if assumption regarding SLR is followed or not

Assumption:

1. Mean of y_i given x_i can be joined by a straight line
2. Mean of y_i given x_i = $\beta_0 + \beta_1 x_i$
3. y_i 's are independent.
4. y_i 's have a constant variance
5. y_i 's come from a normal distribution

Equivalent:

1. Mean of ϵ_i equal to zero
2. ϵ_i 's are independent
3. Variance of ϵ_i 's is constant
4. ϵ_i follows a normal distribution

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \Rightarrow \text{imply } \epsilon_i = y_i - E[y_i]$$

$$E(y_i) = \beta_0 + \beta_1 x_i$$

$$\text{or } g_i = \beta_0 + \beta_1 x_i \quad \checkmark \text{ residual}$$

Residual Analysis

New with residuals e_i

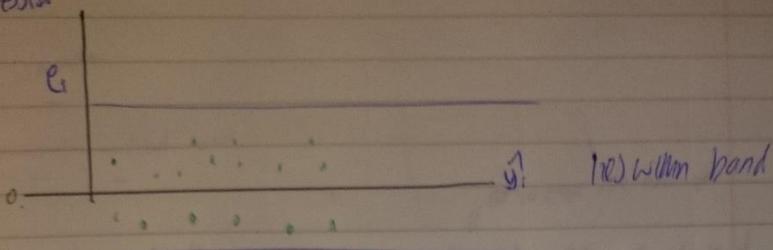
1. Check for mean error term e_i equal to 0 and error term having constant variance.

- Plot residuals against fitted values e_i vs y_i .
Look for patterns.

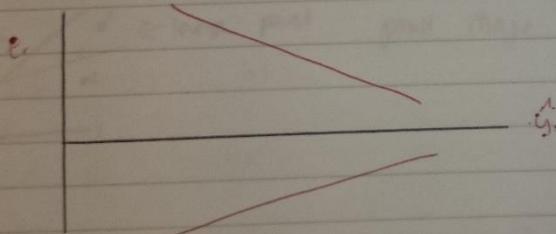
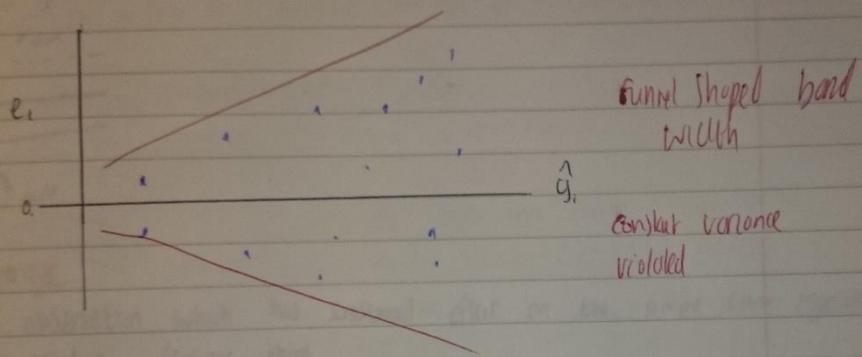
- No pattern implies \Rightarrow Randomly distributed around 0 means assumption hold

3 Regression

Good



Bad



i. U-shaped

ii. constant varm

3. Regressions

Relationship between t-test and F-test.

Can be shown that $F_{\text{calc}} = (t_{\text{calc}})^2$

For diamond example: $t_{\text{calc}} = t_{\text{obs}} = (45.44)^2 = 2110.99 = F$.

So for same test we can use either the t-test or the F-test.
for $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

t-test is more versatile, since it allows us to test for one sided alternative

(read) $H_0: \beta_1 = 0$ $H_1: \beta_1 > 0$ using t-test.
or < 0

In hypothesis β_0 or β_1

In t-test use b_0 or b_1

13 Regression Analysis

Sources of variation	SS	df	Mean square	F
Regression	SSR	1	MSR	$\frac{MSR}{MSE}$
Error	SSE	n-2	MSE	
Total	SSTo	n-1	-	-

F is used to test $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

Compare F with $F_{\text{critical}} = F_{1-\alpha/2, 1, n-2}$

R^2 statistic

$$R^2 = \frac{\text{Sum of squares due to Regression}}{\text{Total sum of squares.}} = \frac{SSR}{SSTo} = 1 - \frac{SSE}{SSTo}$$

R^2 measures the proportion of total variation around the mean as defined by the regression equation.

R^2 lies between 0 and 1

R^2 is the correlation coefficient between y and \hat{y} .

In SLR, $R^2 = r^2$, where r is the correlation of x and y .

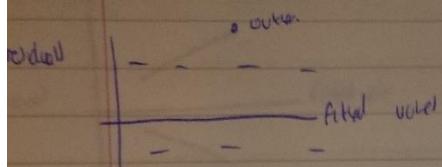
Summary.lm() → in output multiply R-sq is R^2 .

For a good fit, R^2 is close to 1. If

For diamond dataset, $R^2 = 0.9783$

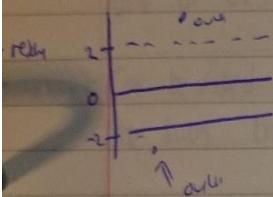
Bad for bad fit, $R^2 < 0.7$ no relation

2/11 Regression



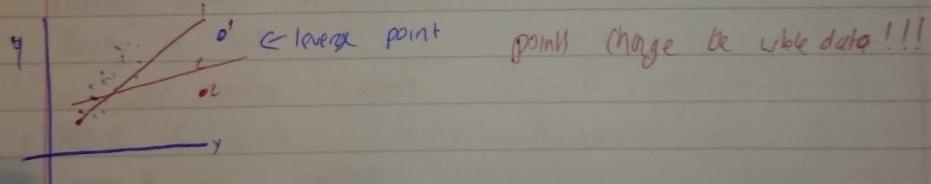
not a proper rule more than

Usually points outside the range of -2, +2 are considered outliers



Leverage

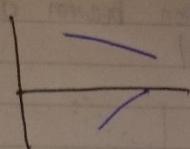
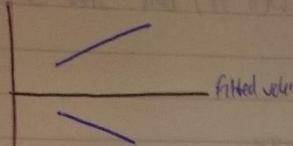
An observation which has substantial effect on the simple linear regression fit is called a leverage point.



2/13

Residuals

Residuals



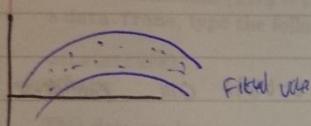
This problem means variance is not constant.

None of the problem is "heteroscedasticity" $\Rightarrow \text{Var}(\epsilon_i) \neq \sigma^2$ for all i .

Sometimes can be fixed by transforming y and x

One such transformation is square root transform

Residuals



A non-linear model needs to be fit

$$\text{Eg try } Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i.$$

$$\text{or try } " = " + \beta_3 x_i^3 + \epsilon_i$$

Plotting residuals vs fitted is similar to plotting residuals vs x
strictly in SLR.

We were checking:

i. mean of ϵ_i is 0

ii. Variance of ϵ_i is constant

Assumptions

- All ϵ_i 's are independent of each other

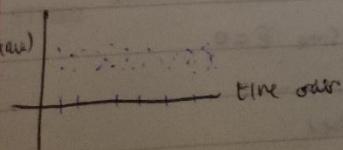
- It is only possible that to tell this when there is some sort of time order in the data

\Rightarrow plot residuals against time order

No pattern \Rightarrow OK

Pattern \Rightarrow not independent

Residuals

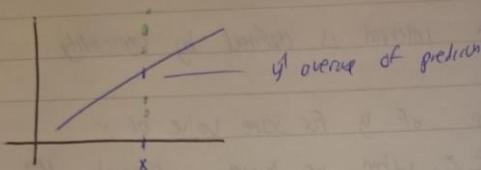


2

Alternative

Average r.v. $y_i = \text{mean} + \text{chance variation}$

Hence, variability for a single y_i is greater than variability of mean



(CI of mean in some sense refers to how much the point on line regression line can vary)

PI of a single point refers to in some sense variability of any of those green dots

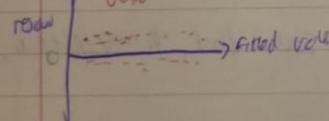
2010/2011 QL

a 1b

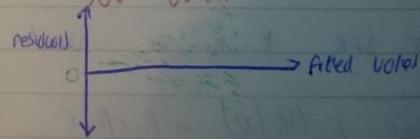
i. constant.

Constant variance for E_i for all i .

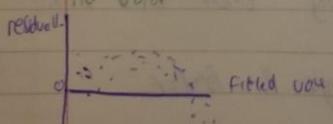
Valid



NOT valid

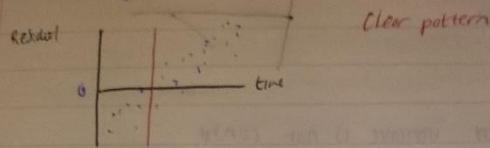


Not valid



2

Fit a regression between diameter of a well (x) and shear strength of the well (y)



Clear pattern

The observation are not independent.

There is correlation between successive observations (points) (auto-correlation)

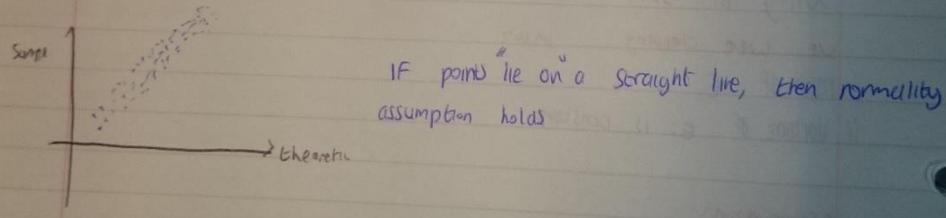
ASSUMPTION:

The error e_i are normally distributed.

-Normal distribution Necessary for computing t-tests and confidence intervals

-Plot Q-Q plot or normal probability plot. Is used to check for normality.

Here, the ordered residual are plotted against ordered expected normal scores



Outliers

Outliers are extreme observations

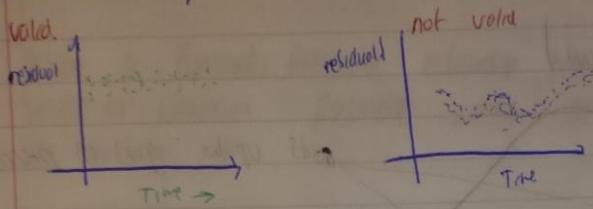
-Detection of outliers is better done with Standardized residuals,

-Standardized residuals $\rightarrow \frac{e_i - \bar{e}}{\sqrt{MSE}} = \frac{e_i}{\sqrt{MSE}}$ since $\bar{e} = 0$

-Do a residual plot residual vs fitted values

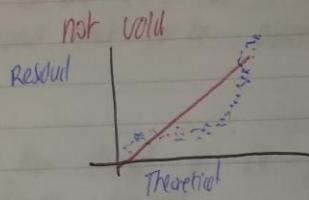
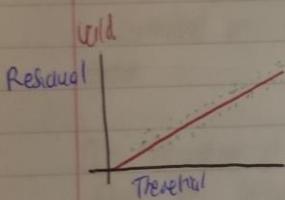
3 Regression

(ii) Independence of ϵ_i
Time order plot



(iii) Normality of ϵ_i

A-Q plot



C i. $H_0: \beta_0 = 0$ vs $H_1: \beta_0 \neq 0$

From Computer $t_{\text{calculated}}$ with t_{critical}

$$t_{\text{critical}} \rightarrow t_{0.475, 8} = 2.306$$

$|t| = 11.69 \gg t_{0.475, 8}$ reject H_0 , ~~thus~~ ~~the~~ we
reject ~~that~~ that time taken to adopt the innovation
of a company of size = 0.

ii. $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$

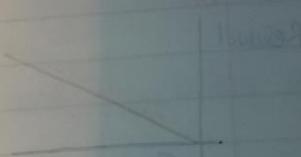
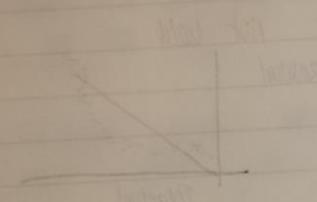
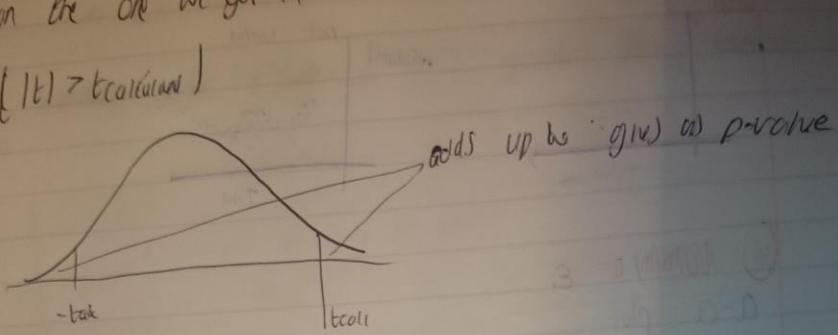
$$t_{0.475, 8} = 2.306$$

$$|t| = 4.43 \gg t_{\text{critical}} \text{ reject } H_0.$$

reject H_0 but β_1 is zero. Change in time taken to
adopt innovation for 1 unit increase in size

4
c. P-value is the probability of getting a t value more extreme than the one we got in the test.

$$P(|t| > t_{\text{calculated}})$$



Regression

Confidence and Prediction Intervals for predictions

$$\text{Interval} = \text{point estimate} \pm t_{\text{critical}} * \text{s.e. (estimate)}$$

\ variation or variability

The width or range of an interval is defined by variability.

(I i) for mean prediction of y for some value of x

For a certain value of x , when we have several y 's then we construct the mean (the average) of those y 's and finally get the standard error of average of y .

This way we get I or Average of y for a certain value of x .

y_1, \dots, y_n are normal r.v. such that $y_i \sim N(\mu, \sigma^2)$

$$\text{Then } S \sim N \left(\mu, \frac{\sigma^2}{n} \right)$$

\Rightarrow variability of \bar{y} < variability of y_i for all i

For SLR, y_1, y_n are normal r.v. such that $y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$

$$\text{Now } E(Y_i | X_i) = \beta_0 + \beta_1 X_i$$

Average of y for $X = x_i$

\Rightarrow average of y for $x=x_i$, has less variability than y_i for $x=x_i$.

In SLR, the P.I. is constructed for a single Y_i for $X=x$.

The PI is wider than CI for prediction