

Module Code ST3451

Module Name APPLIED LINEAR STATISTICAL METHODS I

Module Short Title STATISTICAL METHODS I

**ECTS 5
weighting**

Semester/term taught Michaelmas

Contact Hours 3 hours per week, some of which will be tutorials

Module Personnel Jason Wyse

wyse@tcd.ie

Learning Outcomes When students have successfully completed this module they should be able to:

1. Derive and apply estimators, tests and confidence intervals for the parameters for a range of linear regression models.
2. Derive and construct ANOVA tables.
3. Examine the fit of a regression model through regression diagnostics and test the assumptions of the model.
4. Build an appropriate linear regression model for a given data set.

Module Learning Aims The student will learn about the simple linear regression (SLR) model in detail. This will include derivation of least squares estimators and their properties, sampling distributions of the estimators in the case of Gaussian errors, and tests of significance. The student will also learn about ANOVA- decomposition of the error sum of squares. The matrix approach to linear regression will follow where multiple regression will be discussed. Various diagnostics of fit will be explored, with illustration of how these can be used in practice. Some modifications of the usual regression model will be discussed as well as model building through

variable selection.

Module Content

1. Simple linear regression
2. Multiple regression
3. Regression diagnostics
4. Variable selection

Recommended

Reading List

- A second course in statistics regression analysis- Mendenhall & Sincich
- Classical & modern regression with applications- R. H. Meyers
- Introduction to linear regression analysis- Montgomery, Peck & Vining
- Applied regression analysis- Draper & Smith

Assessment Final exam (90%) and assignments (10%)
Details

Academic Year 2015/16
of Data

28/04/15

ALSM 1

1.0 SIMPLE LINEAR REGRESSION

Basic Model

- Let y be a random variable which follows some mean μ and variance σ^2
- We can write y as $y = \mu + \epsilon$

where μ is a fixed parameter and ϵ is a rv

$$\mathbb{E}[\epsilon] = 0 \quad \text{Var}[\epsilon] = \sigma^2$$

$$y = \mu + \epsilon \quad \mathbb{E}[y] = \mathbb{E}[\mu + \epsilon] = \mu$$

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

- If $y \sim N(\mu, \sigma^2)$ then $\epsilon \sim N(0, \sigma^2)$, then the term ϵ is referred to as the error

- We could write $\epsilon = y - \mu$ so that ϵ is the random deviation of y about its mean.

- Suppose there is an iid sample y_1, y_2, \dots, y_n

- We could use this basic model by writing $y_i = \mu + \epsilon_i \quad i=1, \dots, n$

- When $y_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$, $\epsilon_i \sim N(0, \sigma^2)$ (by independence)

- Then using this model is equivalent to situation where a one sample t-test is appropriate for testing the hypothesis: $H_0: \mu = \mu_0$ vs $H_1: \mu \neq \mu_0$

- Gives a $100(1-\alpha)\%$ CI for the mean:

$$\bar{y} \pm (t_{n-1, \alpha/2}) (s/\sqrt{n})$$

Estimators: parameter \rightarrow Estimator

$$\mu \rightarrow \bar{y}$$

$$\sigma^2 \rightarrow s^2$$

$$\text{With } y_1, y_2, \dots, y_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2) \quad \bar{y} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

Simple model $y_i = \mu + \epsilon_i \quad i=1, \dots, n$

$$\mathbb{E}[\epsilon_i] = 0 \quad \text{Var}[\epsilon_i] = \sigma^2$$

$$\text{Estimator for } \mu: \hat{\mu} = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\text{Estimator for } \sigma^2: s^2 = \frac{1}{n-1} \sum_{i=1}^{n-1} (y_i - \bar{y})^2$$

$$\hat{y}_i = \hat{\mu} + \hat{\epsilon}_i \quad \Leftrightarrow \quad \hat{\epsilon}_i = y_i - \bar{y}$$

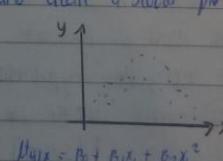
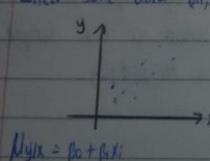
$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n \hat{\epsilon}_i^2$$

Simple linear Regression Model

- When we have one independent variable, we call the model simple linear regression SLR
- Dependent variable y_i (continuous)
- Independent variable X (continuous)
- X is incorporated into the model as a predictor or explanatory variable
 - It explains "in some way" the value of obtained y
 - The values of the independent variable are assumed known X is not a r.v.
- For a given value of x , we consider y to be distributed with mean $\mathbb{E}[y|x] = \mu_{yx}$ i.e. $y = \mu_{yx} + \epsilon$

- Now, how does μ_{yx} depend on x ?

- Collect some data $(x_1, y_1), (x_2, y_2)$ and create a scatter plot of x only



- Simplest model is SLR model; $\mu_{yx} = \mathbb{E}[y|x] = \beta_0 + \beta_1 x$

β_0 and β_1 are fixed constants (parameters) ⇒ call these regression coefficients

→ β_0 : intercept parameter - mean of y when x is zero

→ β_1 : slope parameter - change in mean value of y for a one unit increase in x

$$\mu_{y|x+1} - \mu_{yx} = \beta_0 + \beta_1(x+1) - (\beta_0 + \beta_1 x) = \beta_1$$

Note: When we say linear models, we mean linear in the parameters

e.g. $y_i = e^x + \epsilon_i$ is not linear

28/09/15

ALSM 1

Assumptions:

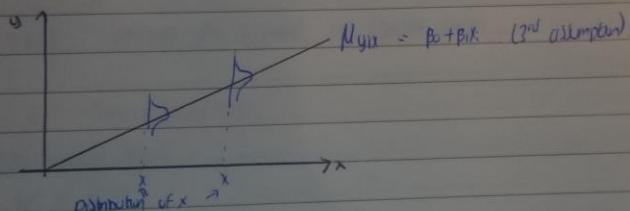
1. For a fixed value of x , y is a random variable with a finite mean and variance:
 $|\mathbb{E}[y|x]| < \infty$ $\text{Var}[y|x] < \infty$ (i.e. it exists)

2. The values of y are uncorrelated

3. The model is linear i.e. $\mathbb{E}[y|x] = \beta_0 + \beta_1 x$

4. Homoscedasticity - Variance of y is not dependent on x , i.e. variance is constant.

$$\text{Var}[y|x] = \text{Var}[y] = \sigma^2$$



Statement of the model:

- Data $(x_1, y_1), \dots, (x_n, y_n)$
- $y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i=1, \dots, n$
- $\mathbb{E}[\epsilon_i] = 0 \quad \forall i$
- $\text{Var}[\epsilon_i] = \sigma^2 \quad \forall i$ (Homoscedastic)
- $\text{Cov}[\epsilon_i, \epsilon_j] = 0 \quad \text{if } i \neq j \quad (\epsilon_i \text{ are uncorrelated})$

Models as data descriptors

- In data analysis, we use a model to describe the data generating process
- This will allow us to say something about the patterns in the data generation
- We make inference about the model parameters
- "All models are wrong, but some are useful"

Basic Model: $y = \mu + \varepsilon_i$ $i=1, \dots, n$
 $E[\varepsilon_i] = 0$ $\text{Var}[\varepsilon_i] = \sigma^2$
 $\text{Cov}[\varepsilon_i, \varepsilon_j] = 0$ $i \neq j$

SLR: $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ (instead of μ in BM)
Data (x_i, y_i) , (x_n, y_n)

1.5 ALSM I

1.5 PARAMETER ESTIMATION

- We will use the collected data $(x_i, y_i), (x_n, y_n)$ to estimate β_0, β_1 .
- The objective is to find the "best" straight line through the data.

The Method of Least Squares

- Suppose that $\hat{\beta}_0, \hat{\beta}_1$ are the estimates of β_0, β_1 (in general $\hat{\beta} \neq \beta$, $\hat{\beta} \in \mathbb{R}$)
- When $x = x_i$, the model predicts $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$. Prediction error is $y_i - \hat{y}_i$.
- Want to minimize this error for all points $i=1, \dots, n$.

- The Least Squares (LS) criterion is to choose $\hat{\beta}_0, \hat{\beta}_1$ to minimize:

$$SS[\text{Error}] = SSE = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = Q(\hat{\beta}_0, \hat{\beta}_1)$$

- In the optimization literature, this is called an objective function.

- SSE known as:
• Sum of squared errors

• Residual sum of squares

• Error sum of squares

• Sum of squared errors about the regression line

- $SSE = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$ SSE a function of $\hat{\beta}_0, \hat{\beta}_1$

- Minimize SSE w.r.t. $\hat{\beta}_0, \hat{\beta}_1$.

$$\frac{d(SSE)}{d(\hat{\beta}_0)} = d \left[\frac{(y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2}{d(\hat{\beta}_1)} \right]$$

$$= -2(y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1) + \dots + -2(y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)$$

$$= -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \quad (1)$$

$$\frac{d(SSE)}{d(\hat{\beta}_1)} \rightarrow \sum \frac{d(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{d(\hat{\beta}_1)}$$

$$= \sum (-2x_i)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)$$

(2)

2.

If we set (1) and (2) = 0, we get normal equations

$$(1) \sum y_i - n\bar{y} - \beta_0 \sum x_i = 0$$

$$(2) \sum x_i y_i - \beta_0 \sum x_i - \beta_1 \sum x_i^2 = 0$$

$$(1) \beta_0 = \frac{\sum y_i}{n} - \beta_1 \frac{\sum x_i}{n} = \bar{y} - \beta_1 \bar{x}$$

$$\text{Using (1) in (2)} \quad \sum x_i (y_i - \bar{y} + \beta_1 \bar{x} - \beta_1 x_i) = 0$$

$$\sum x_i (y_i - \bar{y}) - \beta_1 \sum x_i (\bar{x} - x_i) = 0$$

$$\beta_1 = \frac{\sum (x_i)(y_i - \bar{y})}{\sum (x_i)(\bar{x} - x_i)} = \frac{S_{xy}}{S_{xx}}$$

$$\sum x_i (\bar{x} - \bar{x})$$

$$= \sum (x_i)(\bar{x}) - \sum (x_i)(\bar{x})$$

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

$$\sum (x_i)(\bar{x}) = \sum x_i - n\bar{x} = n\bar{x} - n\bar{x} = 0$$

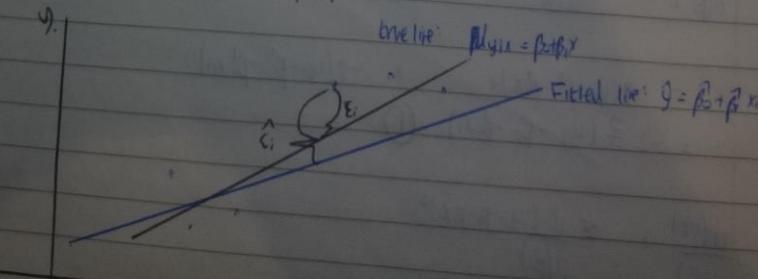
$$\text{and similarly: } \sum x_i (y_i - \bar{y}) = \sum (y_i)(\bar{x} - \bar{x}) = \sum (x_i - \bar{x})(y_i - \bar{y})$$

The least squares estimators are:

$$\hat{\beta}_1 = S_{xy}/S_{xx} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Fitted regression line is: $y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

4.



But fitted values \hat{y}_i random variable (statistic) with one formula of all data

05/10/15 ALSM1

Last Week:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad i=1, \dots, N$$

Dependent Variable Parameter Independent Variable error

$$\mathbb{E}[\varepsilon_i] = 0 \quad \text{Var}[\varepsilon_i] = \sigma^2$$

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

Want to predict mean value of y at any given x . $\hat{y} = \beta_0 + \beta_1 x$

Example

Basic Model: $y_i = \beta_0 + \varepsilon_i \quad (\beta_1 = 0)$

What is the LS estimator of β_0 ?

$$\sum \varepsilon_i^2 = \sum (y_i - \bar{y})^2$$
$$\frac{\partial L}{\partial \beta_0} = -2 \sum (y_i - \bar{y}) = 0$$

$$\sum y_i - n\bar{y} = 0$$

$$\beta_0 = \bar{y}$$

Example:

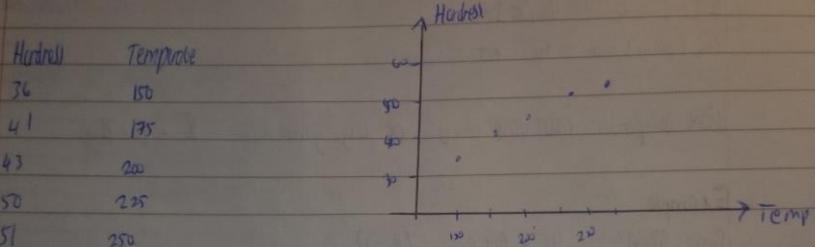
The corrected sum of cross product and square can be computed in an efficient manner. E.g. $S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$ [corrected \rightarrow mean has been subtracted]

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

$$\text{Similarly, } S_{xx} = \sum (x_i - \bar{x})(x_i - \bar{x}) = \sum x_i^2 - \frac{(\sum x_i)^2}{n}$$

Example:

A smartphone has an outer shield protective case which is made of hybrid material. The hardness of the material depends on temperature (which materials are merged). Durability is indicated by hardness. The materials were merged at 5 different temperatures & measured hardness.



Fit the model $y = \text{hardness}$ $x = \text{temperature}$

$$\begin{aligned} n &= 5 & \hat{\beta}_1 &= \frac{\sum y_i}{\sum x_i} & S_{xy} &= \sum x_i y_i - \frac{\sum x_i \sum y_i}{n} \\ \sum x_i &= 1000 & & & &= 4575 - \frac{(225)(100)}{5} = 975 \\ \sum y_i &= 221 & S_{xx} &= \sum x_i^2 - \frac{(\sum x_i)^2}{n} & &= 20620 - \frac{100^2}{5} = 6220 \\ \sum x_i^2 &= 20620 & \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} & &= 4575 - \frac{975}{6220} = 0.156 \\ \sum x_i y_i &= 4575 & & & & \end{aligned}$$

$$\hat{y} = \bar{y} - \hat{\beta}_1 \bar{x} = \frac{221}{5} - \hat{\beta}_1 \frac{\sum x_i}{n} = \frac{221}{5} - 0.156 \left(\frac{1000}{5} \right) = 13.0$$

Model Fitted line = $y = 13 + 0.156x$

- 0.156 is the estimated increase in mean hardness for every extra degree Celsius

- 13 is the estimated mean hardness at temperature zero

- In this case, intercept does not have meaningful interpretation

- However, by having the intercept in the model, we are allowing the line more flexibility in its fit of the data

10/15 ALSM1

$$\begin{aligned}
 & \text{Try to find the SSE for } \hat{\beta}_0, \hat{\beta}_1 \quad \text{SSE} = \sum (y_i - \bar{y})^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\
 &= \sum (y_i - \bar{y} + \hat{\beta}_1 (\bar{x} - \hat{\beta}_1 x_i))^2 \\
 &\leq (y_i - \bar{y})^2 - \hat{\beta}_1^2 (\bar{x} - \hat{\beta}_1 x_i)^2 \\
 &= \sum (y_i - \bar{y})^2 - 2\hat{\beta}_1^2 \sum (\bar{x} - \hat{\beta}_1 x_i)(y_i - \bar{y}) + \hat{\beta}_1^2 \sum (\bar{x} - \hat{\beta}_1 x_i)^2 \\
 &= S_{yy} - 2\hat{\beta}_1 S_{xy} + \hat{\beta}_1^2 S_{xx} \\
 &= S_{yy} - \frac{2S_{xy} S_{yy}}{S_{xx}} + \frac{S_{xy}^2}{S_{xx}}
 \end{aligned}$$

For our example: $S_{yy} = 9927 - \frac{221^2}{5} = 1588$

$\text{SSE} = 67$ which is the sum of squared residuals about the fitted line

$$\text{SSE} = S_{yy} - \frac{S_{xy}^2}{S_{xx}} \rightarrow S_{yy} = \text{SSE} + \frac{S_{xy}^2}{S_{xx}}$$

/

variability in the data variability around fit by fitting the line

05/10/15

1 (2)

ALSM 1

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

↑
SSE for basic model ↓
SSE for SLE

1.7 ANALYSIS OF VARIANCE (ANOVA) IN THE DEPENDENT VARIABLE

The variation in y can be quantified by either:

$$\sum (y_i - \bar{y})^2 = \sum y_i^2 \rightarrow \text{Total Uncorrected Sum of Squared SS (Uncorrected)}$$

or more usually

$$\sum (y_i - \bar{y})^2 = S_{yy}^2(n-1) \rightarrow \text{Total corrected sum of squared}$$

$$\begin{aligned} \text{Now } \sum (y_i - \bar{y})^2 &= \sum (y_i^2 - 2y_i\bar{y} + \bar{y}^2) \\ &= \sum y_i^2 - 2n\bar{y}^2 + n\bar{y}^2 \\ &= \sum y_i^2 - n\bar{y}^2 \\ &= SS(\text{uncorrected}) - \text{correction} \end{aligned}$$

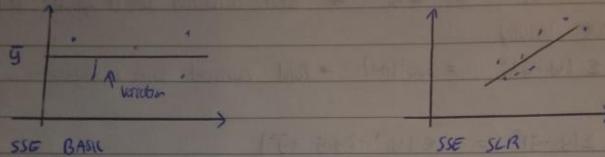
Total corrected SS = Total Uncorrected SS - Correction

Consider writing SS(uncorrected) in terms of fitted values and residuals ($y_i = \hat{y}_i + e_i$)

$$\begin{aligned} SS(\text{uncorrected}) &= \sum y_i^2 = \sum (\hat{y}_i + e_i)^2 \\ &= \sum (\hat{y}_i^2 + 2\hat{y}_i e_i + e_i^2) \\ &= \sum \hat{y}_i^2 + \sum e_i^2 + 2 \sum \hat{y}_i e_i \\ \\ &\sum \hat{y}_i e_i = \sum \hat{y}_i (y_i - \bar{y}) \\ &= \sum \hat{y}_i (y_i - \beta_0 - \beta_1 x_i) \\ &= \sum \hat{y}_i (y_i - \bar{y} + \beta_1 \bar{x} - \beta_1 x_i) \\ &= \sum \hat{y}_i [(y_i - \bar{y}) + \beta_1 (\bar{x} - x_i)] \\ &= [\beta_0 + \beta_1 (\bar{x} - \bar{y})] \sum [(y_i - \bar{y}) + \beta_1 (\bar{x} - x_i)] \\ &= \sum [\beta_0 + \beta_1 (\bar{x} - \bar{y})] [(y_i - \bar{y}) + \beta_1 (\bar{x} - x_i)] \\ &= \sum \bar{y} (y_i - \bar{y}) + \sum \bar{y} \beta_1 (\bar{x} - \bar{y}) + \beta_1 \sum (\bar{x} - \bar{y})(y_i - \bar{y}) \quad \text{← } -\hat{\beta}_1^2 \sum (\bar{x} - \bar{y})(\bar{x} - x_i) \\ &= \bar{y} \left[\sum (y_i - \bar{y}) - \bar{y} \beta_1 \sum (\bar{x} - \bar{y}) + \beta_1 \sum (\bar{x} - \bar{y})(y_i - \bar{y}) - \hat{\beta}_1^2 \sum (\bar{x} - \bar{y})^2 \right] \\ &= \bar{y} \left[\sum (y_i - \bar{y}) - \bar{y} \bar{\beta}_1 \sum (\bar{x} - \bar{y}) + \beta_1 \sum (\bar{x} - \bar{y})(y_i - \bar{y}) - \hat{\beta}_1^2 \sum (\bar{x} - \bar{y})^2 \right] \end{aligned}$$

$$= \frac{SSA}{SSR} - \frac{SSR^2}{SSA} = \frac{\frac{SSA}{n}}{\frac{SSR}{n}} - \frac{\frac{SSR}{n}^2}{\frac{SSA}{n}} = 0$$

$$\begin{aligned} SS(\text{Unadjusted}) &= \sum y_i^2 + \sum \hat{y}_i^2 + 2 \sum y_i \hat{y}_i \\ &= \sum y_i^2 + \sum \hat{y}_i^2 \\ &= SS(\text{Model}) + SSE \end{aligned}$$



$$\begin{aligned} SS(\text{Model}) &= \sum \hat{y}_i^2 = \sum [y_i - \bar{y} + \hat{y}]^2 \\ &= \sum [(y_i - \bar{y})^2 + 2(y_i - \bar{y})(\hat{y} - \bar{y}) + \hat{y}^2] \\ &= \sum (y_i - \bar{y})^2 + 2 \sum (y_i - \bar{y})(\hat{y} - \bar{y}) + n\bar{y}^2 \end{aligned}$$

$$\begin{aligned} \text{Aside: } \bar{y} \sum (y_i - \bar{y}) &= \bar{y} \sum y_i - \bar{y} \sum \bar{y} \\ &= \bar{y} \sum y_i - n\bar{y}\bar{y} \\ &= \frac{\sum y_i}{n} - \frac{n \sum y_i}{n} = 0 \end{aligned}$$

$$\begin{aligned} \text{Aside: } \bar{y} \sum (\hat{y} - \bar{y}) &= \bar{y} \sum (\beta_0 + \beta_1 x_i - \bar{y}) \\ &= \bar{y} \sum (\bar{y} - \beta_1 \bar{x} + \beta_1 x_i - \bar{y}) \\ &= \bar{y} \sum (x_i - \bar{x}) \beta_1 \\ &= \beta_1 \bar{y} \sum (x_i - \bar{x}) \\ &\stackrel{\bar{y} = 0}{=} 0 \end{aligned}$$

$$\begin{aligned} SS(\text{Model}) &= \sum (y_i - \bar{y})^2 + n\bar{y}^2 \\ &= SSE + \text{correction} \end{aligned}$$

The regression sum of squares represents the reduction in variation (around the predicted mean) by adding the $\beta_1 x_i$ term to a model containing β_0

05/10/15

ALSM

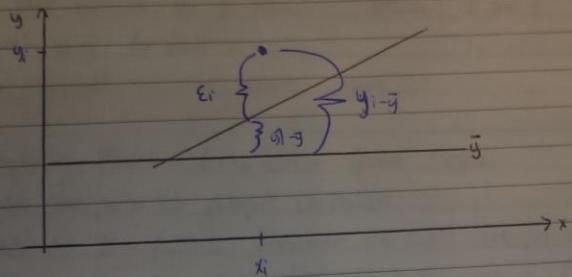
2 (v)

$$SS(\text{Uncorrected}) = SS(\text{Reg}) + n\bar{y}^2 + SSE$$

IF we subtract LR correction from each side:

$$SS(\text{Uncorrected}) - n\bar{y}^2 = SS(\text{Reg}) + SSE$$

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



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

07/05/15

ALM 1

$$\begin{aligned}\sum (y_i - \bar{y})^2 &= \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 \\ \text{SS(Corr)} &= \text{SS(Reg)} + \text{SSE}\end{aligned}$$

- If there is little or no relationship between X and Y . In this case \hat{y}_i will be "close" to \bar{y} for each x_i .
- Then SS(Reg) will be smaller than SSE .

Comparing Models - Extra Sum of Squares Methods

Suppose we wish to compare the models:

$$\text{Model 1: } y_i = \beta_0 + \epsilon_i \quad (\text{horizontal line model}) \quad i=1, \dots, n \quad (\text{reduced model})$$

$$\text{Model 2: } y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (\text{sloping line model}) \quad i=1, \dots, n \quad (\text{full model})$$

Model 1 is a special case of Model 2 ($\beta_1=0$) and can be viewed as having less structure as it only describes deviations about a mean β_0 .

$$\text{Model 1: } \hat{\beta}_0 = \bar{y} \quad \text{SSE} = \sum (y_i - \bar{y})^2$$

$$\text{Model 2: } \hat{\beta}_1 = \frac{\sum x_i y_i}{\sum x_i^2} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \text{SSE} = \sum (y_i - \hat{y}_i)^2 \quad \text{where } \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Note: β_0 is different in Model 1 and Model 2.

Model 2 is less restricted than Model 1 $\text{SSE}_2 \leq \text{SSE}_1$. Always can get closer to the data with a complex model, and $\text{SSE}_1 - \text{SSE}_2$ is a measure of how much better model 2 is than model 1. \rightarrow better in the sense of accounting for the variability in the data.

$$\begin{aligned}\text{SSE}_1 - \text{SSE}_2 &= \sum (y_i - \bar{y})^2 - \sum (y_i - \hat{y}_i)^2 \\ &= \text{SS(Corrected)} - \text{SSE}_2 \\ &= \text{SS(Regression)} \rightarrow \text{Extra Sum of Squares}\end{aligned}$$

- Hence SS(Reg) indicates the importance of the $\beta_1 x_i$ term in the model.
- We can use the "R" notation, $\text{SS(Reg)} = R(\beta_1 | \beta_0)$ adds β_1 already have β_0

Now a second comparison:

$$\text{Model 0: } y_i = \epsilon_i \quad i=1, \dots, N \quad (\text{Reduced Model})$$

$$\text{Model 1: } y_i = \beta_0 + \epsilon_i \quad i=1, \dots, N \quad (\text{Full Model})$$

$$\text{Model 0: } SSE_0 = \sum (y_i - \bar{y})^2 = \sum (y_i)^2$$

$$\text{Model 1: } SSE_1 = \sum (y_i - \bar{y})^2$$

Thus $SSE_1 \leq SSE_0$ and $SSE_0 - SSE_1$ is a measure of how much better model 1 is.

$$SSE_0 - SSE_1 = SS(\text{Uncorrected}) - SS(\text{corrected}) = n\bar{y}^2 \quad (\text{the correction})$$

Hence, $n\bar{y}^2$ represents the importance of β_0 in the model.

$$R(\beta_0) = n\bar{y}^2 \quad (\text{the reduction})$$

Finally, let's compare:

$$\text{Model 0: } y_i = \epsilon_i \quad i=1, \dots, n$$

$$\text{Model 2: } y_i = \beta_0 + \beta_1 x_i \quad i=1, \dots, n$$

Where Model 0 is a special case of Model 2 ($\beta_0=0, \beta_1=0$)

$$\text{Model 0: } SSE_0 = \sum y_i^2$$

$$\text{Model 2: } SSE_2 = \sum (y_i - \bar{y})^2$$

$$\begin{aligned} \text{We have } SSE_2 &\leq SSE_0 \quad \text{and } SSE_0 - SSE_2 = SS(\text{uncorrected}) - SS_2 \\ &= SS(\text{Resid}) + n\bar{y}^2 \end{aligned}$$

= SS(Model) - difference from adding for β_1 to model with no structure

Here $SS(\text{Model})$ represents the importance of β_0 and $\beta_1 x_i$ in the model and

$$R(\beta_0, \beta_1) = SS(\text{Model}) = R(\beta_0) + R(\beta_1 | \beta_0)$$

$$SS(\text{uncorrected}) = \sum y_i^2 \quad n$$

$$SS(\text{corrected}) = \sum (y_i - \bar{y})^2 \quad n-1$$

$$\text{Correlation} = n\bar{y}^2 \quad 1$$

$$SS(\text{Model}) = R(\beta_0, \beta_1) \quad 2$$

$$SS(\text{Resid}) = \sum (y_i - \bar{y})^2 \quad 1$$

$$SS_b = \sum (y_i - \bar{y})^2 \quad n-2$$

12/10/15

1 (1)

ALSM 1

Deviation between predicted mean from regression line and sample mean
 $SS(\text{Reg}) = \sum (y_i - \bar{y})^2 \quad (n-2)$

$$\begin{aligned} SSE &\rightarrow \text{residuals (not explained by model)} = \sum (y_i - \hat{y}_i)^2 \quad (n-2) \\ SS(\text{Corrected}) &= \sum (y_i - \bar{y})^2 \quad (n-1) \end{aligned}$$

$$SS(\text{Reg}) = \sum (y_i - \bar{y})^2 = \beta_1^2 S_{xx}$$

$$\hat{y}_i = \beta_0 + \beta_1 x_i$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\beta_1^2 = \frac{S_{xy}}{S_{xx}}$$

$$= \sum (\beta_0 + \beta_1 x_i - \bar{y})^2$$

$$= \sum (\bar{y} - \beta_1 \bar{x} + \beta_1 x_i - \bar{y})^2$$

$$= \sum (-\beta_1 (\bar{x} - x_i))^2$$

$$= \beta_1^2 \sum (x_i - \bar{x})^2$$

$$= \beta_1^2 S_{xx}$$

$$= \frac{S_{xy}^2}{S_{xx}^2} S_{xx} = \beta_1^2 S_{xy}$$

Anova Table

An analysis of variance (ANOVA) table gives a summary of the sources and contributions to $SS(\text{Corrected})$ or $ss(\text{uncorrected})$ by various parts of the model.

Source	DF	SS	Mean Square = SS/DF
Regression on x	1	$\sum (y_i - \bar{y})^2$	$MS(\text{Reg})$
Residual Error	$n-2$	$\sum (y_i - \hat{y}_i)^2$	MSE
Total Corrected	$n-1$	$\sum (y_i - \bar{y})^2$	

Source	DF	SS	MS = SS/DF
Correction R (β_0)	1	$n\bar{y}^2$	
Regression on x R ($\beta_1 x$)	1	$\sum (y_i - \bar{y})^2$	$MS(\text{Reg})$
Residual Error	$n-2$	$\sum (y_i - \hat{y}_i)^2$	MSE
Total Uncorrected	n	$\sum y_i^2$	

The ANOVA Table gives $R(\beta_0, \beta_1)$ and $R(\beta_1)$ without having to fit extra models.

Coefficient of Determination

- The contribution of the $\beta_1 x_i$ term in the model, measured by $SS(\text{reg})$ and this has an upper limit of $SS(\text{corrected})$

- Hence we can quantify the contribution of x by the market means of the coefficient of determination R^2

$$\text{Where } R^2 = \frac{SS(\text{Reg})}{SS(\text{corr})} \quad SS(\text{corr}) = SS(\text{reg}) + SSE$$

$$\Rightarrow SS(\text{Reg}) = SS(\text{corrected}) - SSE$$

$$R^2 = \frac{SS(\text{corr}) - SSE}{SS(\text{corr})} = 1 - \frac{SSE}{SS(\text{corr})} \quad 0 \leq R^2 \leq 1$$

If $R^2 = 0.76$, then 76% of the variation in y is "explained" by its linear relationship with x .

NOTE: We can get large R^2 for poor model

Example

$$\sum X_i = 100 \quad \sum Y_i = 221 \quad \sum X_i^2 = 20620 \quad \sum Y_i^2 = 9427 \quad \sum X_i Y_i = 48775$$

$$SS(\text{corr}) = SS(\text{uncorrected}) - n\bar{y}^2$$

$$= 9972 - \frac{221^2}{5} = 158.8$$

$$\beta_1 = 0.156 \quad S_{xy} = 97.5$$

$$SS(\text{Reg}) = \beta_1 S_{xy} = 0.156(97.5) = 15.21$$

$$SSE = SS(\text{corr}) - SS(\text{Reg}) = 158.8 - 15.21 = 6.7$$

ANOVA Table

Source	DF	SS	MS
Regression	1	15.21	15.21
Residual	3	6.7	2.23
Corrected	4	158.8	

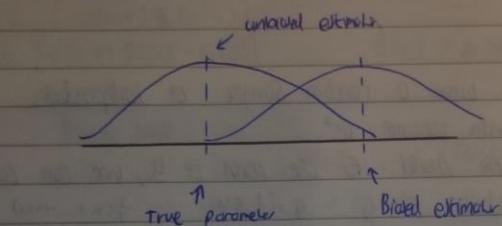
$$R^2 = \frac{SS(\text{Reg})}{SS(\text{corr})} = \frac{15.21}{158.8} = 0.0957$$

12/10/15 ALSM 1

36.1

1.8 PROPERTIES OF ESTIMATORS

- An estimator is a sample statistic $\hat{\beta}_1 = \bar{y} - \hat{\beta}_0 \bar{x}$, $\hat{\beta}_0 = \bar{y}_{\text{avg}}$ and as a result has a sampling distribution
- If the mean of the sampling distribution is NOT equal to the true parameter value, being estimated, then the estimator is biased.
- Otherwise unbiased.

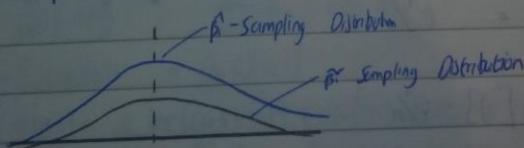


- Apart from being unbiased, it is also desirable for an estimator to have the smallest possible variance.

(could do: $\hat{\beta}_1 = \frac{\sum y_i - \bar{y}}{\sum x_i - \bar{x}}$, and $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ (suggested estimator))

$$\mathbb{E}[\hat{\beta}_1] = \mathbb{E}\left[\frac{\sum y_i - \bar{y}}{\sum x_i - \bar{x}}\right] = \frac{1}{\sum x_i - \bar{x}} \mathbb{E}[\sum y_i - \bar{y}]$$
$$= \frac{1}{\sum x_i - \bar{x}} [\beta_0 + \beta_1 x_1 - \beta_0 - \beta_1 \bar{x}]$$

$= \beta_1 \rightarrow$ unbiased estimator



Revision of Random Variable Properties

- Let μ and v be 2 r.v.

$\mathbb{E}[\mu]$: mean of μ

$$\text{Var}[\mu] = \mathbb{E}[\mu^2] - [\mathbb{E}[\mu]]^2$$

$$\text{Corr}[\mu, v] = \frac{\text{Cov}[\mu, v]}{\sqrt{\text{Var}[\mu] \text{Var}[v]}}$$

- If μ and ν are uncorrelated or independent, then $\text{Cov}[\mu, \nu] = 0$

$$\text{Cov}[\mu, \nu] = \text{Var}[\mu]$$

$$\mathbb{E}[a\mu + b\nu] = a\mathbb{E}[\mu] + b\mathbb{E}[\nu]$$

$$\text{Var}[a\mu + b\nu] = a^2\text{Var}[\mu] + b^2\text{Var}[\nu] + 2ab\text{Cov}[\mu, \nu]$$

- Consider another random variable z and constant c

$$\text{Cov}[a\mu + b\nu, cz] = ac\text{Cov}[\mu, z] + bc\text{Cov}[\nu, z]$$

For example:

- Consider a situation where a random sample of independent observations y_1, \dots, y_n each with variance σ^2

- Then, irrespective of the model for the mean of y_i , we can consider the variance of the sampling distribution of \bar{y} $\bar{y} = \frac{1}{n}\sum y_i = \frac{1}{n}(y_1 + \dots + y_n)$

$$= \frac{1}{n}y_1 + \dots + \frac{1}{n}y_n$$

$$\mathbb{E}[\bar{y}] = \frac{1}{n}\mathbb{E}[y_1] + \dots + \frac{1}{n}\mathbb{E}[y_n]$$

$$= \frac{1}{n}\mu + \dots + \frac{1}{n}\mu$$

$$= \mu$$

$$\begin{aligned}\text{Var}[\bar{y}] &= \text{Var}\left[\frac{1}{n}y_1 + \dots + \frac{1}{n}y_n\right] \\ &= \text{Var}\left[\frac{1}{n}y_1\right] + \dots + \text{Var}\left[\frac{1}{n}y_n\right] + \text{Cov}\left[\frac{1}{n}y_1, \frac{1}{n}y_2\right] + \dots + \text{Cov}\left[\frac{1}{n}y_m, \frac{1}{n}y_n\right] \\ &= \frac{1}{n^2}\text{Var}[y_1] + \dots + \frac{1}{n^2}\text{Var}[y_n] + 2 \sum_{i=1, j>i}^{n-1} \frac{1}{n^2} \text{Cov}[y_i, y_j] \\ &= \frac{1}{n^2} \sigma^2 + \dots + \frac{1}{n^2} \sigma^2 + 0 \quad (\text{Cov}=0 \text{ by way of independence}) \\ &= n\sigma^2/n^2 = \sigma^2/n\end{aligned}$$

$$\mathbb{E}[\bar{y}] = \mu, \text{Var}[\bar{y}] = \sigma^2/n$$

$$y_i \sim N(\mu, \sigma^2) \Rightarrow \bar{y} \sim N(\mu, \sigma^2/n)$$

$$\sigma^2 \text{ unknown: } \frac{\bar{y} - \mu}{\sigma/\sqrt{n}} \sim t_{n-1}$$

12/10/15

ALSM1

SC)

1.9 PROPERTIES OF THE LEAST SQUARES ESTIMATOR

Consider the Sampling distribution of $\hat{\beta}_1$.

$$\begin{aligned} E[\hat{\beta}_1] &= E\left[\frac{\sum y_i}{\sum x_i}\right] = E\left[\frac{\sum \varepsilon_i(x_i - \bar{x})(y_i - \bar{y})}{\sum x_i}\right] \\ &= \frac{1}{\sum x_i} E\left[\sum \varepsilon_i(x_i - \bar{x})(y_i - \bar{y})\right] \quad x_i's \text{ known} \Rightarrow \text{assumed to be fixed} \\ &= \frac{1}{\sum x_i} E\left[(x_1 - \bar{x})(y_1 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y})\right] \quad x_i's \text{ assumed non-random} \\ &= \frac{1}{\sum x_i} \left[\sum (x_i - \bar{x}) E[y_i - \bar{y}] \right] \\ &= \frac{1}{\sum x_i} \left[\sum (x_i - \bar{x}) E[y_i] - E[y_i] \right] \quad E[y_i] = \beta_0 + \beta_1 x_i \quad (\text{from Model}) \end{aligned}$$

Assuming Model is correct:

$$\begin{aligned} E[y_i] &= E\left[\frac{\sum y_i}{n}\right] = \frac{1}{n} E[y_i] = \frac{1}{n} \sum \varepsilon_i(\beta_0 + \beta_1 x_i) \\ &= \frac{1}{n} [n\beta_0 + \beta_1 \sum x_i] = \beta_0 + \beta_1 \bar{x} \\ &= \frac{1}{\sum x_i} \left[\sum (x_i - \bar{x}) \left[\frac{E[y_i]}{E[y_i]} - \beta_0 - \beta_1 \bar{x} \right] \right] \\ &= \frac{1}{\sum x_i} \sum (x_i - \bar{x}) \beta_1 \\ &= \beta_1 \frac{\sum (x_i - \bar{x})}{\sum x_i} \\ &= \beta_1 \end{aligned}$$

 $\hat{\beta}_1$ is an unbiased estimator of β_1 .

Variance

$$\begin{aligned} \text{Var}[\hat{\beta}_1] &= \text{Var}\left[\frac{\sum y_i}{\sum x_i}\right] = \frac{1}{\sum x_i^2} \text{Var}\left[\sum \varepsilon_i(x_i - \bar{x})y_i\right] \\ &= \frac{1}{\sum x_i^2} \text{Var}\left[\sum (x_i - \bar{x})(y_i) + \dots + (x_n - \bar{x})(y_n)\right] \\ &= \frac{1}{\sum x_i^2} \left[\text{Var}[(x_1 - \bar{x})(y_1)] + \dots + \text{Var}[(x_n - \bar{x})(y_n)] + \text{Cov}[(x_1 - \bar{x})(y_1), (x_2 - \bar{x})(y_2)] + \dots + \text{Cov}[(x_{n-1} - \bar{x})(y_{n-1}), (x_n - \bar{x})(y_n)] \right] \\ &= \frac{1}{\sum x_i^2} \left[\sum \text{Var}[(x_i - \bar{x})(y_i)] + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n \text{Cov}[(x_i - \bar{x})(y_i), (x_j - \bar{x})(y_j)] \right] \\ &= \frac{1}{\sum x_i^2} \left[\sum (x_i - \bar{x})^2 \text{Var}[y_i] + 2 \sum (x_i - \bar{x})(x_j - \bar{x}) \text{Cov}[y_i, y_j] \right] \\ &= \frac{1}{\sum x_i^2} \left[\sum (x_i - \bar{x}) \sigma^2 + 0 \right] \quad (\text{by independence}) \\ &= \sigma^2 \frac{\sum x_i^2}{\sum x_i^2} = \sigma^2 \bar{x} \end{aligned}$$

19/10/15 ALSM 1

$$\mathbb{E}[\hat{\beta}_0] = \beta_0 \quad \mathbb{E}[\hat{\beta}_1] = \beta_1 \\ \text{Var}[\hat{\beta}_0] = ? \quad \text{Var}[\hat{\beta}_1] = \frac{\sigma^2}{S_{xx}}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\begin{aligned} \text{Var}[\hat{\beta}_0] &= \text{Var}[\bar{y} - \hat{\beta}_1 \bar{x}] \\ &= \text{Var}[\bar{y}] + \text{Var}[\hat{\beta}_1 \bar{x}] - 2\text{Cov}[\bar{y}, \hat{\beta}_1 \bar{x}] \\ &= \underline{\text{Var}[\bar{y}]} + \underline{\text{Var}[\hat{\beta}_1 \bar{x}]} - 2\bar{x}^2 \text{Cov}[\bar{y}, \hat{\beta}_1 \bar{x}] \quad \text{Assume } \bar{x} \text{ known} \Rightarrow \text{not random} \end{aligned}$$

$$\begin{aligned} \text{Cov}[\bar{y}, \hat{\beta}_1] &= \text{Cov}\left[\frac{1}{n}\sum y_i, \frac{S_{xy}}{S_{xx}}\right] \quad (x \text{ assumed not random}) \\ &= \frac{1}{n} S_{xx} \text{Cov}[y_i, S_{xy}] \\ &= " \text{Cov}[y_1 + \dots + y_n, S_{xy}] \\ &= " \text{Cov}[y_1, S_{xy}] + \dots + \text{Cov}[y_n, S_{xy}] \\ &= " \left[\sum \text{Cov}[y_i, S_{xy}] \right] \\ &= " \sum \text{Cov}[y_i, \frac{1}{n}(x_i - \bar{x})y_i] \\ &= " \sum \text{Cov}[y_i, (x_i - \bar{x})y_i] + \dots + \text{Cov}[y_i, (x_n - \bar{x})y_n] \\ &= " \sum \text{Cov}[y_i, (x_i - \bar{x})y_i] \\ &= \frac{1}{n} S_{xx} \sum (x_i - \bar{x}) \text{Cov}[y_i, y_i] \\ &= \frac{n}{n} \frac{S_{xx}}{S_{xx}} = 0 \quad \Rightarrow (\text{Cov}[y_i, y_i] = 0) \end{aligned}$$

$$\begin{aligned} \text{Var}[\hat{\beta}_1] &= \text{Var}[\bar{y}] + \text{Var}[\hat{\beta}_1 \bar{x}] \\ &= \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{S_{xx}} + 0 \\ &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right) \end{aligned}$$

$$\begin{aligned} \text{Note: } \text{Cov}[\hat{\beta}_0, \hat{\beta}_1] &= \text{Cov}[\bar{y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1] \\ &= \text{Cov}[\bar{y}, \hat{\beta}_1] - \bar{x} \text{Cov}[\hat{\beta}_1, \hat{\beta}_1] \\ &= 0 - \bar{x} \text{Var}[\hat{\beta}_1] \\ &= -\frac{\bar{x} \sigma^2}{S_{xx}} \end{aligned}$$

1.10 GAUSS-MARKOV THEOREM

Assuming $y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i=1, \dots, n$
 $\mathbb{E}[\epsilon_i] = 0 \quad \text{Var}[\epsilon_i] = \sigma^2 \quad \text{Cov}[\epsilon_i, \epsilon_j] = 0 \quad i \neq j$

- Then the LS estimators have minimum variance amongst all unbiased linear estimators (i.e. LS estimator)

- BLUE : Best Linear Unbiased Estimator

III ESTIMATION OF σ^2

When the model is $y_i = \mu + \epsilon_i$ we use $s^2 = \frac{1}{n-2} \sum (y_i - \bar{y})^2$ as an estimator of σ^2 .

When the model is $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ we use $\frac{1}{n-2} \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$
 $MSE = SSE/(n-2)$

Is this estimator unbiased?

$$\mathbb{E}[MSE] = \frac{1}{n-2} \sum \mathbb{E}[(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2]$$

$$\text{Remember } \mathbb{E}[u^2] = \text{Var}[u] + (\mathbb{E}[u])^2$$

$$\begin{aligned} \mathbb{E}[MSE] &= \frac{1}{n-2} \sum \left[\text{Var}[y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i] + (\mathbb{E}[y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i])^2 \right] \\ &= \frac{1}{n-2} \leq \left[\text{Var}[y_i] + \text{Var}[\hat{\beta}_0] + x_i^2 \text{Var}[\hat{\beta}_1] - 2 \text{Cov}[y_i, \hat{\beta}_0] - 2x_i \text{Cov}[y_i, \hat{\beta}_1] + 2x_i \text{Cov}[\hat{\beta}_0, \hat{\beta}_1] \right. \\ &\quad \left. + [(\hat{\beta}_0 + \hat{\beta}_1 x_i) - \beta_0 - \beta_1 x_i]^2 \right] \\ &= \frac{1}{n-2} \leq \left[\sigma^2 + \sigma^2 \left(\frac{1}{n-2} \sum x_i^2 \right) + \frac{x_i^2 \sigma^2}{S_{xx}} - 2 \text{Cov}[y_i, \hat{\beta}_0]^{(2)} - 2x_i \text{Cov}[y_i, \hat{\beta}_1] + 2x_i \text{Cov}[\hat{\beta}_0, \hat{\beta}_1] \right] \end{aligned}$$

$$\textcircled{1}: \text{Cov}[y_i, \hat{\beta}_0] = \text{Cov}[y_i, \frac{S_{yy}}{S_{xx}}]$$

$$= \frac{1}{S_{xx}} \text{Cov}[y_i, \sum (x_i - \bar{x}) y_i] \quad \text{will only pick up } i^{\text{th}} \text{ term (linearity)}$$

$$= \frac{1}{S_{xx}} \text{Cov}[y_i, (x_i - \bar{x}) y_i]$$

$$= \frac{(x_i - \bar{x}) \text{Cov}[y_i, y_i]}{S_{xx}} = \frac{\sigma^2(x_i - \bar{x})}{S_{xx}}$$

$$\textcircled{2}: \text{Cov}[y_i, \hat{\beta}_1] = \text{Cov}[y_i, g - \beta_1 \bar{x}]$$

$$= \text{Cov}[y_i, \frac{S_{yy}}{S_{xx}} - \beta_1 \bar{x}]$$

$$= \text{Cov}[y_i, g] - \text{Cov}[y_i, \beta_1 \bar{x}]$$

$$\hookrightarrow = \frac{1}{n} \text{Cov}[y_i, g y_i]$$

$$\downarrow = \frac{1}{n} \text{Cov}[y_i, y_i] = \sigma^2 n$$

$$\frac{\sigma^2}{n} - \sigma^2 (x_i - \bar{x})$$

1/10/15 ALSM 1

3.

$$\begin{aligned}\mathbb{E}[\text{MSE}] &= \mathbb{E} \left[\sigma^2 + \sigma^2 \left(\frac{1}{n} + \frac{S_{xx}}{S_{xx}} \right) + X_1^2 + \frac{X_1^2 \sigma^2}{S_{xx}} - 2 \left(\frac{\sigma^2}{n} - \frac{\mathbb{E}(X_1 - \bar{x})\sigma^2}{S_{xx}} \right) \right. \\ &\quad \left. - \frac{2\sigma^2 (X_1 - \bar{x})\sigma^2}{S_{xx}} + 2\sigma^2 \left(\frac{-\bar{x}\sigma^2}{S_{xx}} \right) \right] \\ &= \frac{1}{n-2} \mathbb{E} \left[\sigma^2 + \sigma^2 - \frac{2\sigma^2}{n} - \frac{X_1^2 \sigma^2}{S_{xx}} - \frac{\bar{x}^2 \sigma^2}{S_{xx}} - \frac{2\bar{x}X_1 \sigma^2}{S_{xx}} \right] \\ &= \dots \mathbb{E} \left[\sigma^2 + \sigma^2/n - \sigma^2 S_{xx} (X_1^2 - 2\bar{x}X_1 + \bar{x}^2) \right] \\ &= \dots \mathbb{E} \left[\sigma^2 - \sigma^2/n - \sigma^2 S_{xx} (\mathbb{E}(X_1 - \bar{x})^2) \right] \\ &= \frac{1}{n-2} \left[n\sigma^2 - \sigma^2 - \frac{\sigma^2}{S_{xx}} S_{xx} \right] \\ &= \frac{1}{n-2} [n\sigma^2 - \sigma^2 - \sigma^2] \\ &= \frac{n-2}{n-2} \sigma^2 = \sigma^2 \Rightarrow \text{Unbiased estimator of } \sigma^2 \text{ is the MSE.}\end{aligned}$$

Aside:

$$y_i \sim N(\mu, \sigma^2) \quad y_1, \dots, y_n$$
$$z_i = \frac{y_i - \mu}{\sigma} \sim N(0, 1) \quad z^2 \sim \chi^2_1 \quad \text{chi-square, 1 d.f.}$$
$$z_1^2 + \dots + z_n^2 \sim \chi^2_n \quad n \text{ d.f.}$$

$$g \sim N(\mu, \sigma^2)$$
$$\mu = \frac{g \mu}{\sigma} = \sqrt{n(g \mu)} / \sigma \sim N(0, 1)$$

$$\mu^2 \sim \chi^2_1$$

$$R = \frac{1}{\sigma^2} \sum (y_i - \mu)^2 = \sum \left(\frac{y_i - \mu}{\sigma} \right)^2 \sim \chi^2_n \quad n \text{ d.f.}$$

$$\begin{aligned}R &= \frac{1}{\sigma^2} \sum (y_i - \mu)^2 = \frac{1}{\sigma^2} \sum [(y_i - g) + (g - \mu)]^2 \\ &= \frac{1}{\sigma^2} \left[\sum (y_i - g)^2 + 2\sum (y_i - g)(g - \mu) + n(g - \mu)^2 \right] \\ &= \frac{1}{\sigma^2} \sum (y_i - g)^2 + \frac{n}{\sigma^2} (g - \mu)^2 \\ &= \frac{1}{\sigma^2} [(n-1)s^2] + \frac{\sum (y_i - g)^2}{\sigma^2} \uparrow N(0, 1)\end{aligned}$$

$$R = \frac{(n-1)s^2}{\sigma^2} + n$$

$$R \sim \chi^2_{n-1} \quad n \sim \chi^2_1$$

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1}$$

Show $\chi^2_n \leftrightarrow$ sum of n squared $N(0,1)$

$$F_{n_1, n_2} = \frac{\chi^2_{n_1}/\mu_1}{\chi^2_{n_2}/\mu_2} \quad F\text{-test for comparison of variances}$$

$$t_d = \frac{\sqrt{\chi^2_d/d}}{\sqrt{2/\delta}} \quad t\text{-distribution with } d \text{ degrees of freedom}$$

10/15 ALSM

From last week:

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

$$z_i = \frac{y_i - \mu}{\sigma} \sim N(0, 1)$$

$$z_i^2 \sim \chi^2_1$$

$$\sum z_i^2 \sim \chi^2_n$$

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1}$$

1.2 STATISTICAL INFERENCE FROM THE LINEAR MODEL

- In order to make inferences we must make an assumption about the distribution of y_i : are:
- We will assume the distribution is normal

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

$$\varepsilon_i \sim N(0, \sigma^2)$$

$$1 \quad \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \sum d_i y_i$$

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^2 S_{xx})$$

$$2 \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \sum d_i y_i$$

$$\hat{\beta}_0 \sim N(\beta_0, \sigma^2 (1 + \hat{\beta}_1^2 S_{xx}))$$

$$3 \quad \text{It can be shown that } \frac{SSE}{\sigma^2} \sim \chi^2_{n-2} \quad (\text{two parameters } \hat{\beta}_0, \hat{\beta}_1)$$

$$\frac{(n-2)SSE}{\sigma^2} = \frac{(n-2)MSE}{\sigma^2} \sim \chi^2_{n-2}$$

4. It can be shown that $\hat{\beta}_0$ and SSE are independent and $\hat{\beta}_1$ and SSE are independent.

1.3 HYPOTHESIS TESTING

$$H_0: \beta_1 = \text{some value} \quad H_1: \beta_1 \neq \text{value} \quad \text{or} \quad \beta_1 > \text{value} \quad \beta_1 < \text{value}$$

Usually we would like to test $H_0: \beta_1 = 0$

$$\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$$

Assume that H_0 is true i.e. $\beta_1 = m$

$$\frac{\hat{\beta}_1 - m}{\text{S.E.}} \sim N(0,1) \text{ distribution (standardizing } \hat{\beta}_1)$$

But we don't know what or is

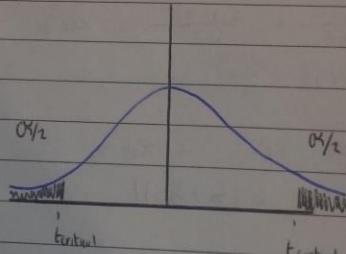
Student's t distribution is defined as follows:

$$t_{n-2} = \frac{z}{\sqrt{\frac{2/n}{z^2}}} \quad z \sim N(0,1) \quad \text{independent}$$

$$\text{So } t_{n-2} = \frac{\frac{\hat{\beta}_1 - m}{\text{S.E.}}}{\sqrt{\frac{(n-2)\text{MSE}}{\sigma^2/(n-2)}}}$$

$$t_{n-2} = \frac{\hat{\beta}_1 - m}{\sqrt{\text{MSE}/S_{xx}}}$$

Look at distribution:



- If we said "I'll reject H_0 if I was less than 5% likely to have seen this t-value if H_0 was true"

- If $|t_{n-2}| > t_{\text{critical}}$ then the probability of observing this data when the null hypothesis is true is less than 0.05 or α .

- We call the region where we would reject, the rejection region

- We can find critical values using tables, statistical packages, computer p-value

Using a similar argument:

$$H_0: \beta_0 = c \quad H_1: \beta_0 \neq c \quad \text{or} \left(\frac{\beta_0 - c}{\text{S.E.}} \right)$$

Using the test statistic:

$$t = \frac{\hat{\beta}_0 - c}{\sqrt{\text{MSE} \left(\frac{1}{s_{xx}} + \frac{x^2}{s_{xx}} \right)}}$$

NOTE: The estimated standard deviation of an estimator is called its standard error.

$$\text{SE}(\hat{\beta}_0) = \sqrt{\frac{\text{MSE}}{s_{xx}}}$$

$$\text{SE}(\hat{\beta}_1) = \sqrt{\text{MSE} \left(\frac{1}{s_{xx}} + \frac{x^2}{s_{xx}} \right)}$$

28/10/15 ALSM I

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{S_{xx}}\right) \quad \sigma^2 = \text{Var}[e]$$
$$\frac{(n-2)MSE}{\sigma^2} \sim \chi^2_{n-2}$$

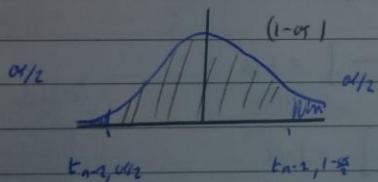
Construct a t -distribution from basis of normal distribution

$$t = \frac{z}{\sqrt{\frac{\sigma^2}{S_{xx}}}} \quad z \sim N(0,1)$$

$$= \frac{(\hat{\beta}_1 - \beta_1) / \sqrt{\frac{\sigma^2}{S_{xx}}}}{\sqrt{\frac{(n-2)MSE}{(n-2)\sigma^2}}} = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{MSE}{S_{xx}}}} \sim t_{n-2}$$

If we know that $\frac{\hat{\beta}_1 - \beta_1}{\sqrt{MSE/S_{xx}}} \sim t_{n-2}$ can use it to construct a CI

1.14 CONFIDENCE INTERVALS



$$P[-t_{n-2, \alpha/2} \leq t_{n-2} \leq t_{n-2, 1-\alpha}] = 1-\alpha$$

$$\text{Replace } t \rightarrow P\left[t_{n-2, \alpha/2} \leq \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{MSE}{S_{xx}}}} \leq t_{n-2, 1-\alpha}\right]$$

$$-t_{n-2, \alpha/2} \sqrt{\frac{MSE}{S_{xx}}} \leq \hat{\beta}_1 - \beta_1 \leq t_{n-2, 1-\alpha} \sqrt{\frac{MSE}{S_{xx}}}$$

$$\hat{\beta}_1 - t_{n-2, \alpha/2} \sqrt{\frac{MSE}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{n-2, 1-\alpha} \sqrt{\frac{MSE}{S_{xx}}}$$

$$\rightarrow \hat{\beta}_1 - t_{n-2, \alpha/2} \sqrt{\frac{MSE}{S_{xx}}}$$

$$\hat{\beta}_1 \pm t_{n-2, 1-\alpha} \sqrt{\frac{MSE}{S_{xx}}}$$

Confidence Interval for β_0 ?

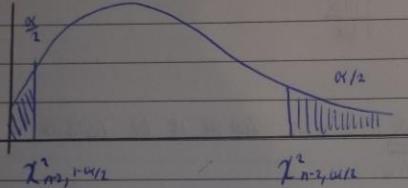
$$\frac{\hat{\beta}_0 - \beta_0}{\text{MSE}(\frac{1}{n} \mathbf{x}' \mathbf{x})} \sim t_{n-2}$$

$$P[-t_{n-2, \alpha/2} \leq \frac{\hat{\beta}_0 - \beta_0}{\text{MSE}(\frac{1}{n} \mathbf{x}' \mathbf{x})} \leq t_{n-2, \alpha/2}] = 1 - \alpha \text{ area}$$

$$\hat{\beta}_0 \pm \sqrt{\text{MSE}(\frac{1}{n} \mathbf{x}' \mathbf{x})} t_{n-2, \alpha/2}$$

(Confidence Interval for σ^2)?

$$\frac{(n-2) \text{MSE}}{\sigma^2} \sim \chi^2_{n-2} \text{ use this fact}$$



Not dealing with a
symmetric distribution

$$P[\chi^2_{n-2, 1-\alpha/2} \leq \frac{(n-2) \text{MSE}}{\sigma^2} \leq \chi^2_{n-2, \alpha/2}] = 1 - \alpha$$

$$P\left[\frac{1}{\chi^2_{n-2, 1-\alpha/2}} \leq \frac{\sigma^2}{(n-2) \text{MSE}} \leq \frac{1}{\chi^2_{n-2, \alpha/2}}\right] = 1 - \alpha$$

$$P\left[\frac{(n-2) \text{MSE}}{\chi^2_{n-2, 1-\alpha/2}} \leq \sigma^2 \leq \frac{(n-2) \text{MSE}}{\chi^2_{n-2, \alpha/2}}\right] = 1 - \alpha$$

Which is a $100(1-\alpha)\%$ confidence interval for σ^2 .

IS AN ALTERNATIVE METHOD FOR TESTING THE HYPOTHESIS $H_0: \beta_1 = 0 \vee H_1: \beta_1 \neq 0$

- Have possibility of directional alternatives in t-test like $\beta_1 > 0$ etc

- Cannot do this in the F-test

Recall the model comparison of section 1.7 using the "extra sum of squares" method.

28/10/15 ALSM 1

3

Model 1: $y_i = \beta_0 + \epsilon_i$ Intercept only model

Model 2: $y_i = \beta_0 + \beta_1 x_{i1} + \epsilon_i$

$R(\beta_1 | \beta_0)$ = Reduction in error sum of squares by adding $\beta_1 x_{i1}$ term to the model
= $SSE(\text{Reg})$ "Extra Sum of Squares"

$SSE(\text{Reg})$ and SSE are independent

When H_0 is true ($\beta_1 = 0$) then $SSE(\text{Reg}) / \sigma^2 \sim \chi^2$

- But by definition the F distribution is the ratio of two independent χ^2 distributions divided by their df.

- An F-dist has 2 df.

$$F_{v_1, v_2} = \frac{\chi^2_{v_1} / v_1}{\chi^2_{v_2} / v_2} \quad \chi^2_{v_1}, \chi^2_{v_2} \text{ are independent}$$

So when H_0 is true $\frac{SSE(\text{Reg})}{\sigma^2} \sim \chi^2$

$MSE(n-2) / \sigma^2 \sim \chi^2_{n-2}$

$$F = \frac{\frac{SSE(\text{Reg})}{\sigma^2} / 1}{\frac{MSE(n-2)}{\sigma^2} / n-2} = \frac{SSE(\text{Reg})}{MSE} \sim F_{1, n-2}$$

- Using this F statistic, we can compare our observation with what would be true if H_0 was true

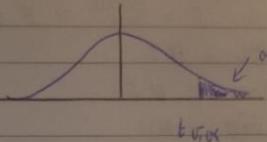
- look up critical F value

- If $F > F_{1, n-2, \alpha}$ then the probability of this happening under H_0 would be or or less, so in that case we'd reject $H_0 = 0$.

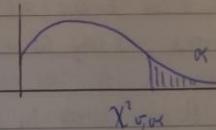
- Evidence for the alternative H_A .

02/11/15 ALSM I

- t distribution with v d.f.
- Want to get point such that area in the tail to the right is α
- Or as a percentage 100 $(1 - \alpha)$ of the area
- This percentage point is denoted $t_{v,\alpha}$



- χ^2 with v d.f.
- Area to right of $\chi^2_{v,\alpha}$ is α



Left point: $t_{n-2, 1-\alpha/2} = -t_{n-2, \alpha/2}$ (symmetrical)

Right point: $t_{n-2, \alpha/2}$

Example - Material Durability

$$n=5 \quad \beta_0 = 0.156 \quad \beta_1 = 13 \quad S_{xx} = 6290$$
$$\bar{x} = 200 \quad SS(R_{xy}) = 152.1 \quad SSE = 6.7 \quad MS(R_{xy}) = 152.1 \quad MSE = 2.23$$

Non-Directional Hypothesis

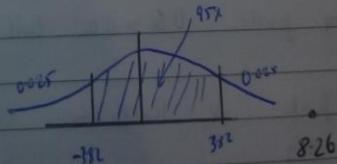
$H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$ Testing for relationship between 2 variables

$$t = \frac{\hat{\beta}_1 - 0}{\sqrt{\frac{MSE}{S_{xx}}}} = \frac{0.156}{\sqrt{\frac{2.23}{6290}}} = 8.26 \quad t_{\text{crit}}$$

- Go to t-distribution with 3 d.f.

-Significance level of α

-Critical value $\pm t_{3, 0.025} = \pm 3.182$



-Test Statistic 8.261 lands in the rejection region

-Evidence to support the alternative $H_1: \beta_1 \neq 0$

-Conclude: Durability of material depends on temperature

F-statistic

- Under H_0 : $F = \frac{MSE_{\text{real}}}{MSE_{\text{null}}} \sim F$ distribution with 6, 7 d.f.

$$F = \frac{52.3}{2.2} = 68.2$$

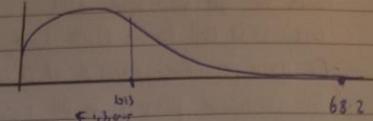
- Significance level of 5%

$$F_{1, 6, 0.05} = 10.13$$

→ Reject H_0 , evidence against $\beta_1 = 0$

NOTE: $t = 8.26$ $F = 68.2 \approx t^2$

The F-test is only for a non directional hypothesis



Does Durability increase by 0.15 units for extra degree in temperature or is it less?

$$H_0: \beta_1 = 0.15$$

$$H_1: \beta_1 < 0.15$$

$$t_{\text{calc}} = \sqrt{\frac{0.15}{MSE_{\text{null}}}} = \sqrt{\frac{0.15 - 0.15}{2.2}} = 0.32$$

- T-table: $t_{3, 0.05} = 2.35$

- Rejection region is any $t_{\text{calc}} < -2.35$

- Don't reject H_0 : no evidence $\beta_1 < 0.15$

P-Value

- P-value: For a test Pr [Obtaining test statistic as extreme as that observed if H_0 is true]

$$\text{Hence p-value} = \Pr [t < 0.32] = 0.61 \quad (\text{t-table})$$

Using p-value, $0.61 > 0.05$; Fail to reject H_0

$$H_0: \beta_1 = b \quad v \quad H_1: \beta_1 \neq b$$

$$t = \frac{b - \bar{b}}{\sqrt{MSE_{\text{null}}}}$$

$$\text{Assume } t \text{ is positive} \quad \Pr [|t_{n-2}| \geq t] = \Pr [t_{n-2} \leq -t] + \Pr [t_{n-2} \geq t]$$

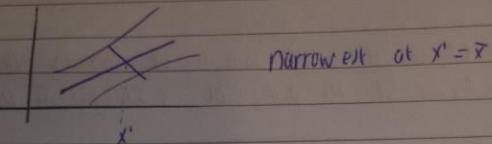
$$\text{p-value} = 2 \Pr [t_{n-2} \leq -t]$$

$$\text{p-value} = 2 \Pr [t_{n-2} \geq t]$$

100(1- α)% confidence interval for mean of y at x'

$$P\left[-t_{n-2, \alpha/2} \leq \frac{\hat{y}' - E[y|x']}{\sqrt{MSE\left(\frac{1}{n} + \frac{(x'-\bar{x})^2}{S_{xx}}\right)}} \leq t_{n-2, \alpha/2}\right] = 1-\alpha$$

$$\hat{y}' \pm t_{n-2, \alpha/2} \sqrt{MSE\left(\frac{1}{n} + \frac{(x'-\bar{x})^2}{S_{xx}}\right)}$$



1.17 PREDICTION

Suppose we wish to predict the value of a new observation when $x=x'$

Predicted value \rightarrow not mean value $\hat{y}' = \beta_0 + \beta_1 x'$ more uncertainty

- Point prediction of y' would be $\hat{y}' = \hat{\beta}_0 + \hat{\beta}_1 x'$

- Consider the prediction error:

$y' - \hat{y}'$ (Deviation between actual and predicted)

$$E[y' - \hat{y}'] = \beta_0 + \beta_1 x' - (\hat{\beta}_0 + \hat{\beta}_1 x')$$

= 0 expected value of 0

$$Var[y' - \hat{y}'] = Var[y'] + Var[\hat{y}'] + 0 \quad \text{uncorrelated} \Rightarrow \text{independence}$$

$$= Var[\varepsilon'] + Var[\hat{y}']$$

$$= \sigma^2 + \sigma^2 \left(\frac{1}{n} + \frac{(x'-\bar{x})^2}{S_{xx}}\right)$$

$$= \sigma^2 \left(1 + \frac{1}{n} + \frac{(x'-\bar{x})^2}{S_{xx}}\right)$$

Since y' and \hat{y}' are independent normals, then $y' - \hat{y}'$ is normal

$$\Rightarrow (y' - \hat{y}') \sim 0$$

$$\sqrt{MSE\left(1 + \frac{1}{n} + \frac{(x'-\bar{x})^2}{S_{xx}}\right)} \quad \text{follows a t-distribution } n-2 \text{ df}$$

02/10/15 AISM1

How would you do a test of $H_0: \beta_1 = 0$ v $H_A: \beta_1 < 0$ or ($\beta_1 > 0$) if you had the p-value above?

Example Continued

$$95\% \text{ CI for } \beta_1 = \hat{\beta}_1 \pm t_{2,0.025} \sqrt{\text{MSE}} \\ = 0.156 \pm 3.182 \sqrt{2.23/6290} \\ 0.96 \leq \beta_1 \leq 0.216$$

(Comparison between CI and test): 0 isn't in the 95% CI then we would reject $H_0: \beta_1 = 0$, in favour of $H_A: \beta_1 \neq 0$ at a 5% significance level.

$$95\% \text{ CI for } \sigma^2: \hat{\sigma}^2 \pm t_{3,0.025} \sqrt{\text{MSE} \left(\frac{1}{n} + \frac{x^2}{Sxx} \right)} \\ 0.0743 \leq \sigma^2 \leq 3.836$$

$$95\% \text{ CI for } \sigma^2: \frac{(n-2)\text{MSE}}{\chi^2_{3,0.025}} \leq \sigma^2 \leq \frac{(n-2)\text{MSE}}{\chi^2_{3,0.975}}$$

$$\chi^2_{3,0.025} = 9.348 \quad \chi^2_{3,0.975} = 2.158 \\ 0.716 \leq \sigma^2 \leq 3.1 \quad \text{Wide interval but small } n$$

1.1b INFERENCE CONCERNING $E[y|x]$

From SLR $E[y|x] = \mu_{y|x} = \beta_0 + \beta_1 x'$

Point estimate $\hat{g} = \hat{\beta}_0 + \hat{\beta}_1 x'$

$$E[\hat{g}] = E[\hat{\beta}_0] + x' E[\hat{\beta}_1] = \beta_0 + \beta_1 x'$$

$$\text{Var}[\hat{g}] = \text{Var}[\hat{\beta}_0] + x'^2 \text{Var}[\hat{\beta}_1] + 2x' (\text{cov}[\hat{\beta}_0, \hat{\beta}_1]) \\ = \sigma^2 \left(\frac{1}{n} + \frac{x^2}{Sxx} \right) + x'^2 \left(\frac{\sigma^2}{Sxx} \right) - 2x' \frac{2\sigma^2}{Sxx} \\ = \frac{\sigma^2}{n} + \frac{\sigma^2(x^2 - \bar{x}^2)}{Sxx}$$

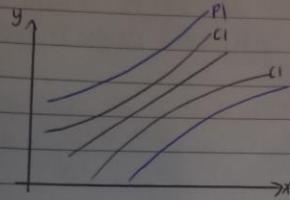
$$SE: \sigma \sqrt{\frac{1}{n} + \frac{(x^2 - \bar{x}^2)}{Sxx}}$$

$$\Rightarrow \frac{\hat{y}' - E[y|x]}{\sqrt{\text{MSE} \left(\frac{1}{n} + \frac{(x^2 - \bar{x}^2)}{Sxx} \right)}} \text{ follows a t-distribution with } n-2 \text{ df}$$

4 02/10/15 ALSM 1 5

$$P\left[\hat{y}' - t_{n-2, \alpha/2} \sqrt{MSE\left(1 + \frac{(x-x')^2}{S_{xx}}\right)} \leq \hat{y}' \leq \hat{y}' + t_{n-2, \alpha/2} \sqrt{MSE\left(1 + \frac{(x-x')^2}{S_{xx}}\right)}\right] = 1 - \alpha$$

and this gives a $100(1-\alpha)\%$ PRECISION interval for y at x'



Since we are estimating the population mean of x' rather than a single value of y , the $C1$ is always narrower than the prediction interval

Both are narrowest at $x = \bar{x}$

Example:

$x = x'$

95% CI for mean durability when temperature is 200°C

$$\bar{x} = 200 \quad \hat{y} \pm t_{2, 0.05} \sqrt{MSE/n}$$

$$42.07 \text{ to } 46.32$$

95% PI: $\hat{y} \pm t_{2, 0.05} \sqrt{MSE(1 + \frac{1}{S_{xx}})}$

$$38.945 \text{ to } 49.405$$