

Describes the relation between variables

Regression models:- Relationship between one dependent variable & explanatory variables.

- 1. Use equations to set up relationship
Numerical dependent (Response) variable.
- 2. 1 or more Numerical or categorical Independent (Explanatory) variables
- 3. used mainly for prediction & Estimation

Regression modeling steps:

- 1. Hypothesize deterministic component
Estimate unknown parameters
- 2. Specify probability distribution of
Random Error terms.
Estimate sd of error
- 3. Evaluate the fitted model.
- 4. Use model for prediction & Estimation.

specifying the deterministic component

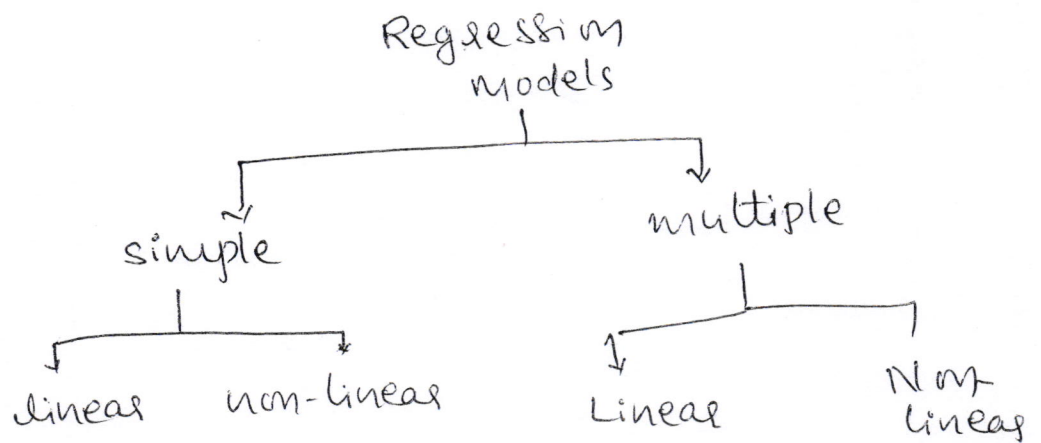
1. Define the dependent variable & independent variable.

2. Hypothesize nature of relationship
Expected effects (i.e. coefficients' signs)

Functional form (linear or non-linear)

Interactions

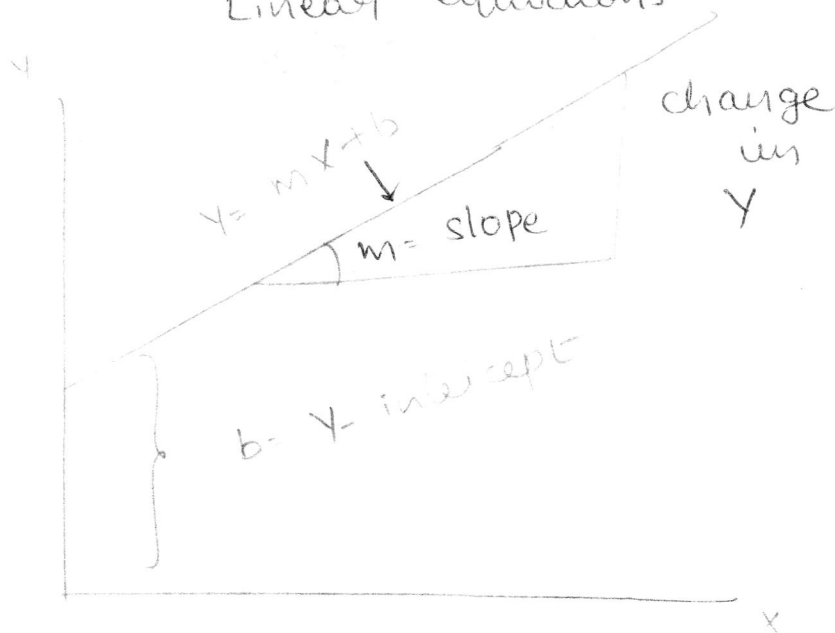
Types of regression models:



1. Explanatory variable

Two + explanatory variables.

Linear equations



Relationship between variables is a linear function

→ population slope

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

↳ Random error

Dependent
(Response variable)

e.g. CD + c.

Independent
(Explanatory variable)

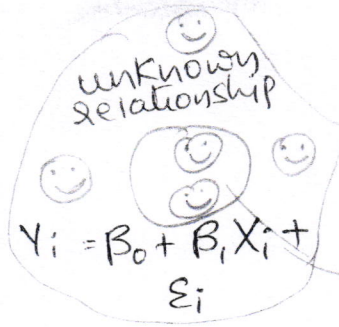
e.g. Years s. season)

Population

Y-intercept.

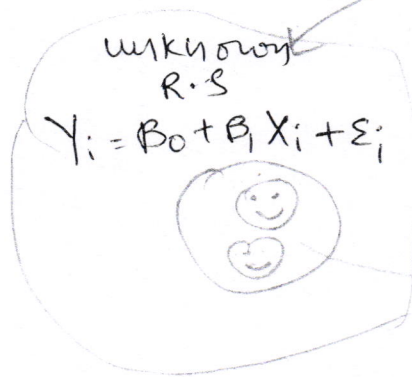
Population & Sample Regression Models

Population



Random sample

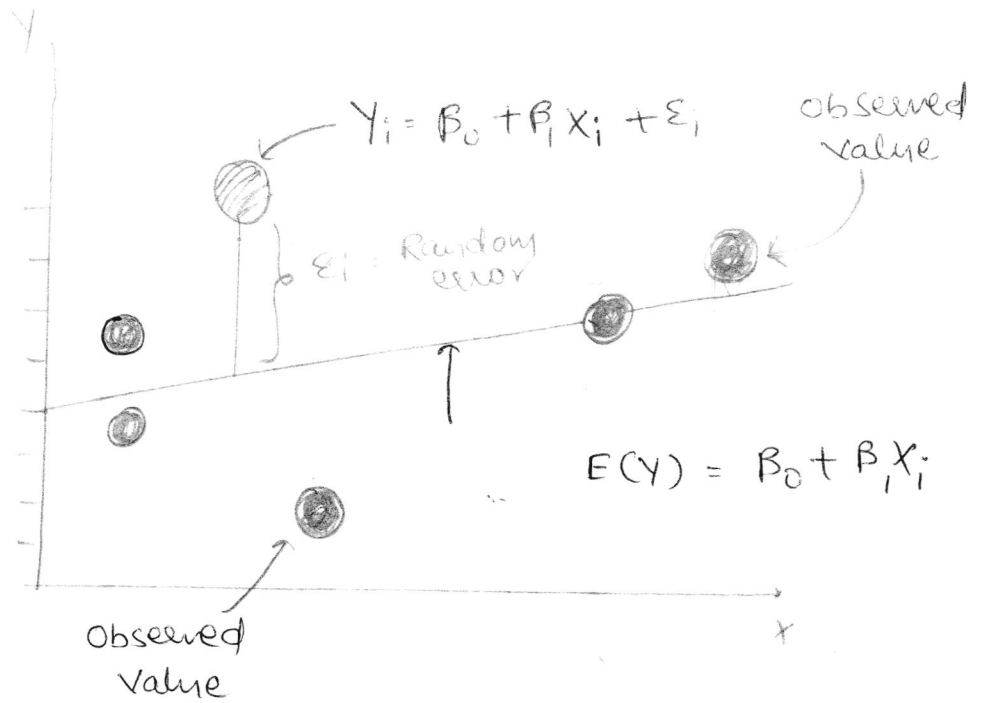
$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{\epsilon}_i$$



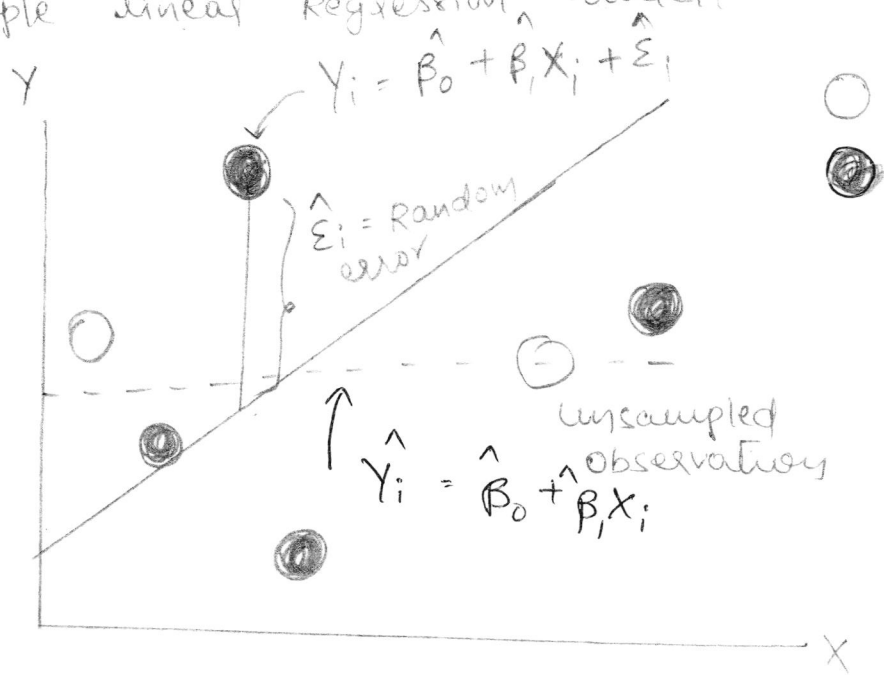
Random sample

Population

Population Linear Regression Model

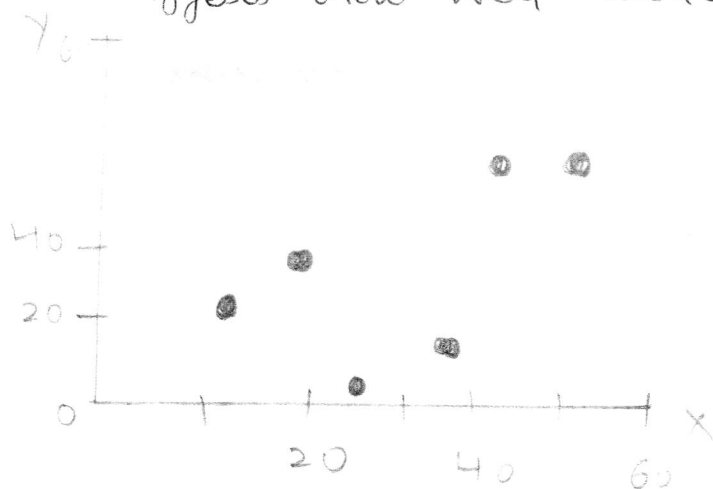


sample linear regression model:-



Scatter plot:

1. plot of All (x_i, y_i) pairs
2. Suggests how well model will fit



Thinking challenge:

How would you draw a line through the points? How do you determine which line fits best?



co-relationship coefficient

measures the strength of association b/w two variables. The most common correlation coefficient, called the Pearson Product Moment correlation coefficient

measures the strength of the linear association between variables.

correlation coefficient of a sample	refers to Pearson method.
r	a population
	P or R

How to interpret a correlation coefficient

The sign & the absolute value of a correlation coefficient describe the direction & the magnitude of the relationship between two variables.

1. The value of a correlation coefficient ranges between -1 & 1.
2. The greater the absolute value of a correlation coefficient, the stronger the linear relationship.

curve fitting with linear regression

We often think of a relationship between two variables as a straight line. That is, if you increase the predictor by 1 unit, the response always increases by x units. However, not all data have a linear relationship, and your model must fit the curves present in the data.

How do you fit a curve to your data?

→ curve-fitting methods in both linear regression & non-linear regression.

Goal: to accurately predict the o/p given the input.

difference b/w linear & non-linear equations in R.A?

Linear regression requires a linear model

A model is linear when each term is either a constant or the product of a parameter and a predictor variable.

A linear equation is constructed by adding the results for each term. This constrains the equation to just one basic form:

Response = constant + parameter \times predictor
+ parameter \times predictor

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_k X_k.$$

In statistics, a regression equation (or function) is linear when it is linear in the parameters. While the equation must be linear in the parameters, you can have the predictor variables in ways that produce curvature. For instance, you can include a squared variable to produce a U-shaped curve.

$$Y = b_0 + b_1 X_1 + b_2 X_1^2 + \dots$$

This model is still linear in the parameters even though the predictor variable is squared. You can also use log & inverse functional forms that are linear in the parameters to produce different types of curves.

Here is an example of a linear regression model that uses a square term to fit the curved relationship between BMI and body fat percentage.

Simple linear regression example

Problem starts last year, 5 randomly selected students took a math aptitude test before they began their statistics course. The statistics department has 3 equations.

- ① What linear regression equation best predicts statistics performance, based on math aptitude scores?
- 2) If a student made an 80 on the aptitude test, what grade would we expect her to make in statistics?
- 3) How well does the regression equation fit the data

How to find the regression equation.

In the table below, the x_i column shows scores on the aptitude test & y_i column shows statistic grades. The last ^{two} rows shows sums & mean scores that we will use to conduct the regression analysis.

Student	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$
1	95	85	17	8	289	64
2	85	95	7	18	49	324
3	80	70	2	-7	4	49
4	70	65	-8	-12	64	144
5	60	70	-18	-7	324	49
					<u>730</u>	630

$$(x_i - \bar{x})(y_i - \bar{y})$$

$$136$$

$$126$$

$$-14$$

$$96$$

$$126$$

$$470$$

The regression equation is a linear equation of the form $\hat{y} = b_0 + b_1x$.

To conduct a regression analysis, we need to solve for b_0 & b_1 . Computations are shown below

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

$$b_1 = \frac{470}{730} = 0.644$$

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

$$b_0 = 77 - (0.644)(78) = \underline{26.768}$$

Therefore, the regression equation is

$$\hat{y} = \underline{26.768 + 0.644x}$$

r

$$\text{slope} = r \cdot \frac{S_y}{S_x}$$

↓
Pearson correlation coefficient

How to use the regression equation?
Once you have the reg. eqn, choose a value for the independent variable (x), perform the computation, & you have an estimated value (\hat{y}) for the dependent variable.

In our example, the independent variable is the student's score on the aptitude test.

The dependent variable is the student's statistics grade. If a student made an 80 on the aptitude test, the estimated statistics grade would be

$$\begin{aligned}\hat{y} &= 26.768 + 0.644x = \\ &= 26.768 + 0.644 \times 80 \\ &= 26.768 + 51.52 \\ &= 78.288\end{aligned}$$

Extrapolation - When you use a regression equation, do not use values for the independent variable that are outside the range of values used to create the equation. That is called extrapolation and it can produce unreasonable estimates.

In this example, the aptitude test score used to create the regression equation ranged from 60 to 95. Therefore, ^{only} values inside that range to estimate statistics grade.

Using values outside that range (less than 60 or greater than 95) is problematic.

How to find the coefficient of determination

Whenever you use a regression equation, you should ask how well the equation fits the data. One way to assess fit is to check the coefficient of determination which can be computed from the following formula:

$$R^2 = \left\{ \frac{(1/N) * \sum [(x_i - \bar{x}) * (y_i - \bar{y})]}{(\sigma_x * \sigma_y)} \right\}^2$$

where N is the no. of observations used to fit the model, \sum is the summation symbol, x_i is the x value for observation i , \bar{x} is the mean value, y_i is the value for observation i , \bar{y} is the mean y value, σ_x is the s.d of x , σ_y is the s.d of y .

$$\sigma_x = \sqrt{\sum (x_i - \bar{x})^2 / N} \quad \sigma_y = \sqrt{\sum (y_i - \bar{y})^2 / N}$$

$$\begin{aligned} \sigma_x &= \sqrt{(730/5)} = \sqrt{(146)} = 12.083 \\ \sigma_y &= \sqrt{(630/5)} = \sqrt{(126)} = 11.225 \end{aligned}$$

$$\begin{aligned} R^2 &= \frac{1}{5} * 470 / (12.083 * 11.225)^2 \\ &= (94 / 135.632)^2 = (0.693)^2 = 0.48 \end{aligned}$$

A coefficient of determination equal to 0.48 indicates that about 48% of the variation in statistics grades (the dependent variable) can be explained by the relationship to math aptitude scores (the independent variable). This would be considered a good fit to the data, in the sense that it would substantially improve an educator's ability to predict student performance in statistics class.

Residual Analysis in Regression

Because a linear regression model is not always appropriate for the data. You should assess the appropriateness of the model by defining residuals & examining residual plots.

Residuals:

The difference between the observed value of the dependent variable (y) and the predicted value (\hat{y}) is called the residual (e). Each datapoint has one residual.

$$\text{Residual} = \text{observed value} - \text{predicted value}$$

$$e = y - \hat{y}$$

Both the sum & the mean of the residuals are equal to zero.

$$\text{that is } \sum e = 0 \text{ \& } \bar{e} = 0.$$

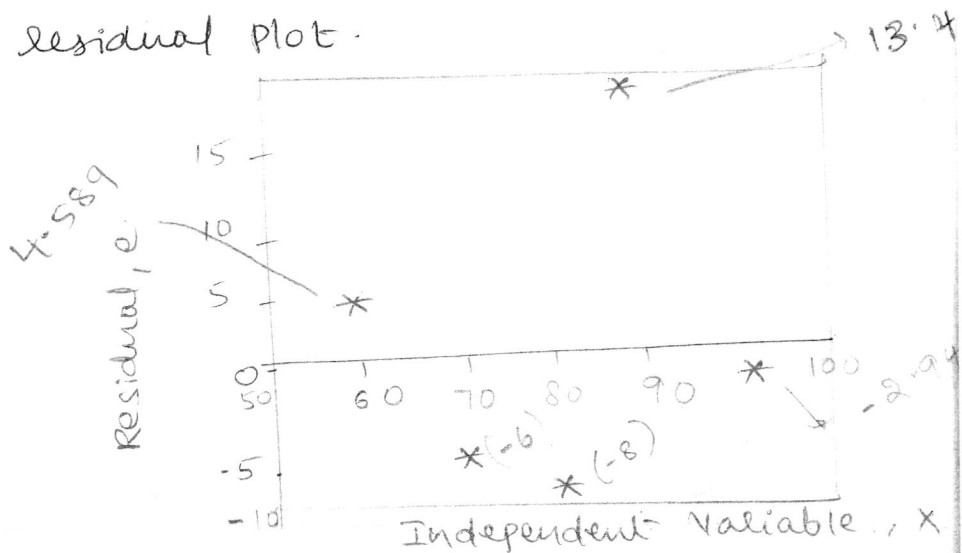
Residual plots

A residual plot is a graph that shows the residuals on the vertical axis & the independent variable on the horizontal axis. If the points in a residual plot are randomly dispersed around the horizontal axis, a linear regression is appropriate for the data, otherwise a non-linear model is more appropriate.

Below, the table shows inputs & o/p's from a simple linear regression analysis.

X	60	70	80	85	95
Y	70	65	70	95	85
\hat{y}	65.411	71.849	78.288	81.507	87.194
e	4.589	-6.849	-8.288	13.493	-2.194

And the chart below displays the residuals and the independent variable (x) as a residual plot.



The residual plot shows a fairly random pattern.

The first residual is +ve: 4.589

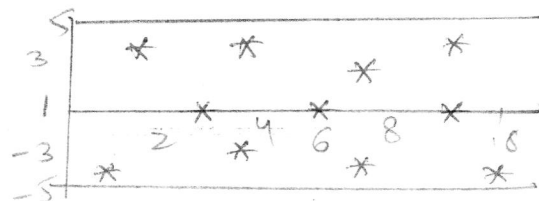
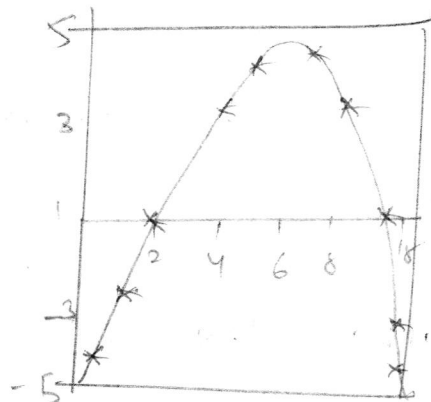
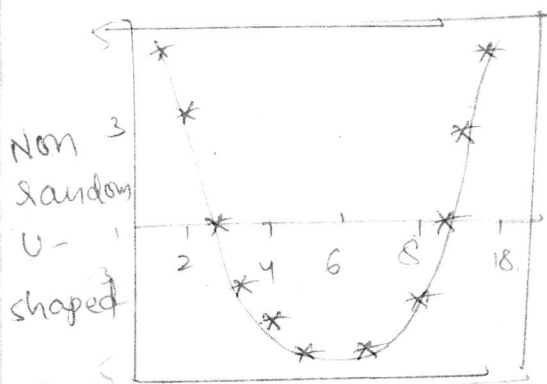
Next two are -ve: -6.549, -8.288

The last residual is -ve: -2.945

This random pattern indicates that a linear model provides a decent fit to the data.

Below, the residual plots show three typical patterns. The first plot shows a random pattern, indicating a good fit for a linear model.

Graph patterns - the other plots patterns are non-random (U-shaped & inverted U,) suggesting a better fit for a non-linear model.



Question:

In the context of R.A, which of the following statements are true?

1. When the sum of the residuals is greater than zero, the dataset is non-linear.
2. A random pattern of residuals supports a linear model.
3. A random pattern of residuals supports a non-linear model.

A. I only

B. II only

C. III only

D. I & II

E. I & III

Solution (B).

A random pattern of residuals supports a linear model.

A non-random pattern supports a non-linear model.

The sum of residuals is always zero
whether the data set is linear or
non-linear.

MULTI-COLLINEARITY

Collinearity: Two independent variables X_1 & X_2 are collinear, when they are correlated with each other.

In a multiple regression study, we assume that the X variables are independent of each other.

We also assume that each X -variable contains a unique piece of info about Y .

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

In the multiple regression model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

we believe that

β_1 = the change in Y for a 1-unit change in X_1 , while X_2 is held constant

Also

β_2 = while X_1 is held constant

when collinearity exists:

in that X_1 & X_2 are correlated

$\beta_1 \neq$ the change in y for a 1-
change in x_1 , with x_2 held
constant

$\beta_2 \neq$ change x_2 held constant

Effects of multicollinearity:

Variances (& standard errors) of
regression coefficient estimators
(i.e. the b_i) are inflated. This means
that $\text{Var}(b_i)$ is too large.

The magnitude of the b_i may be
different from what we expect.

The signs of b_i may be opposite of
what we expect.

Adding or removing any of the x -vars
produce large changes in the
values of remaining b_i or
their signs.

Sometimes removing a datapoint
causes large changes in the
value of b_i or their signs.

In some cases, F is significant,
but the t -values (for the b_i)
may not be significant.

t statistic for β_1

$$= \frac{b_1}{\text{SE of } b_1}$$

Tests for multicollinearity

calculate the correlation coefficient (r) for each pair of the X-variables. If any of the r -values is significantly different from zero, then the independent variables involved may be colinear.

$$r = \frac{\text{covariance}(X_i, X_j)}{\sigma_i \sigma_j}$$

Although the 'r' for any two X variables may be too small, three independent variables X_1, X_2 & X_3 may be highly correlated as a group.

2. VIF (Variance Inflation Factor)

check whether the VIF is too high.

Rule of thumb: collinearity exists if $VIF > 5$. A VIF of 10, for eg. means $\text{Var}(b_i)$ is 10 times what it should be if no collinearity existed (if no collinearity VIF should be 1).

VIF is more rigorous check for coll than correlation coefficient.

$$VIF = \frac{1}{(1 - R_i^2)}$$

In the R.M

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + e$$

R_i^2 is obtained from Regressing X_1 on X_2 & X_3 as follows

$$X_1 = \alpha_0 + \alpha_1 X_2 + \alpha_2 X_3 + e$$

similarly

$$X_2 = \alpha_0 + \alpha_1 X_1 + \alpha_2 X_3 + e \Rightarrow +$$

$$X_3 = \alpha_0 + \alpha_1 X_1 + \alpha_2 X_2 + e \Rightarrow -$$

R_3^2

Solutions for M.C

1. Drop the variables causing the problem.

If using a large no. of X -variables stepwise regression could be used to determine which of the variables to drop.

2.