

## Module-1: Predictive Modelling

### MODULE: PREDICTIVE MODELLING

## Outlines

- Correlation
  - Pearson's correlation
  - Kendall rank correlation
  - Spearman rank correlation
- Regression
  - Types of regression
  - Fitting a function – Criterion for best fit
  - Least squares
- Simple regression
- Multiple regression
- Model assessment and validation

### CORRELATION

## Preliminaries

- $n$  observations for  $x$  and  $y$  variables ( $x_i, y_i$ )
- Sample means  $\bar{x}$  and  $\bar{y}$

$$\bar{x} = \frac{\sum x_i}{n} \quad \bar{y} = \frac{\sum y_i}{n}$$

- Sample variances  $S_{xx}$  and  $S_{yy}$

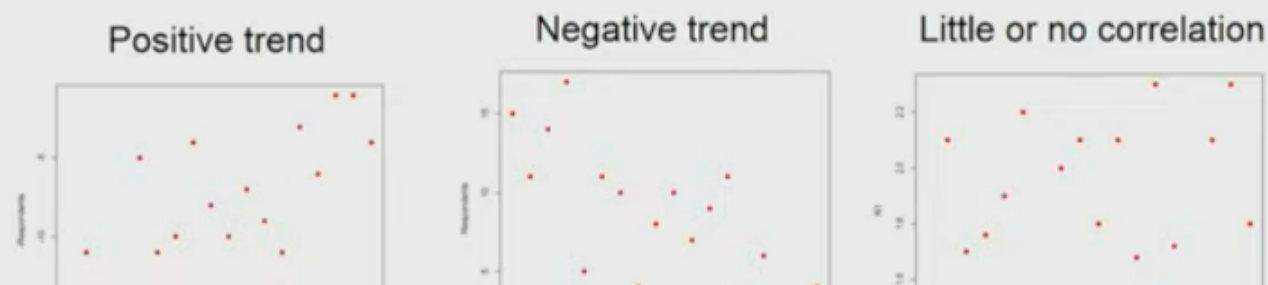
$$S_{xx} = \frac{1}{n} \sum (x_i - \bar{x})^2 \quad S_{yy} = \frac{1}{n} \sum (y_i - \bar{y})^2$$

- Sample covariance  $S_{xy}$

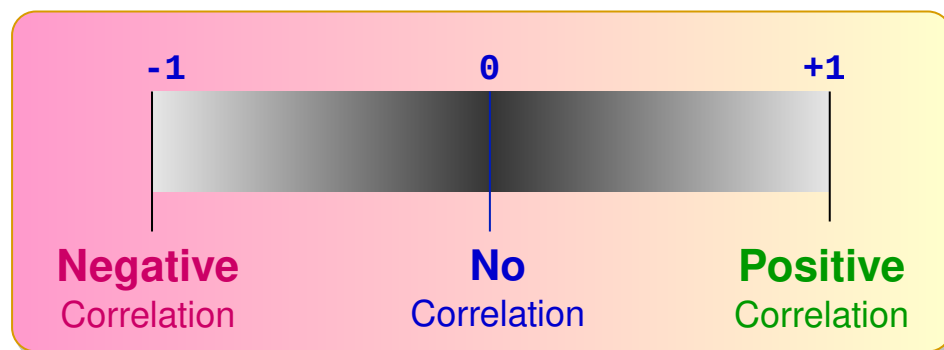
$$S_{xy} = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})$$

# Correlation

- Correlation: the strength of association between two variables
- Correlation does not imply causation
- Visual representation of correlation: Scatter grams



## Pearson's Correlation



## Pearson's Correlation

- $n$  observations for  $x$  and  $y$  variables ( $x_i, y_i$ )
- Pearson's product-moment correlation coefficient ( $r_{xy}$ )

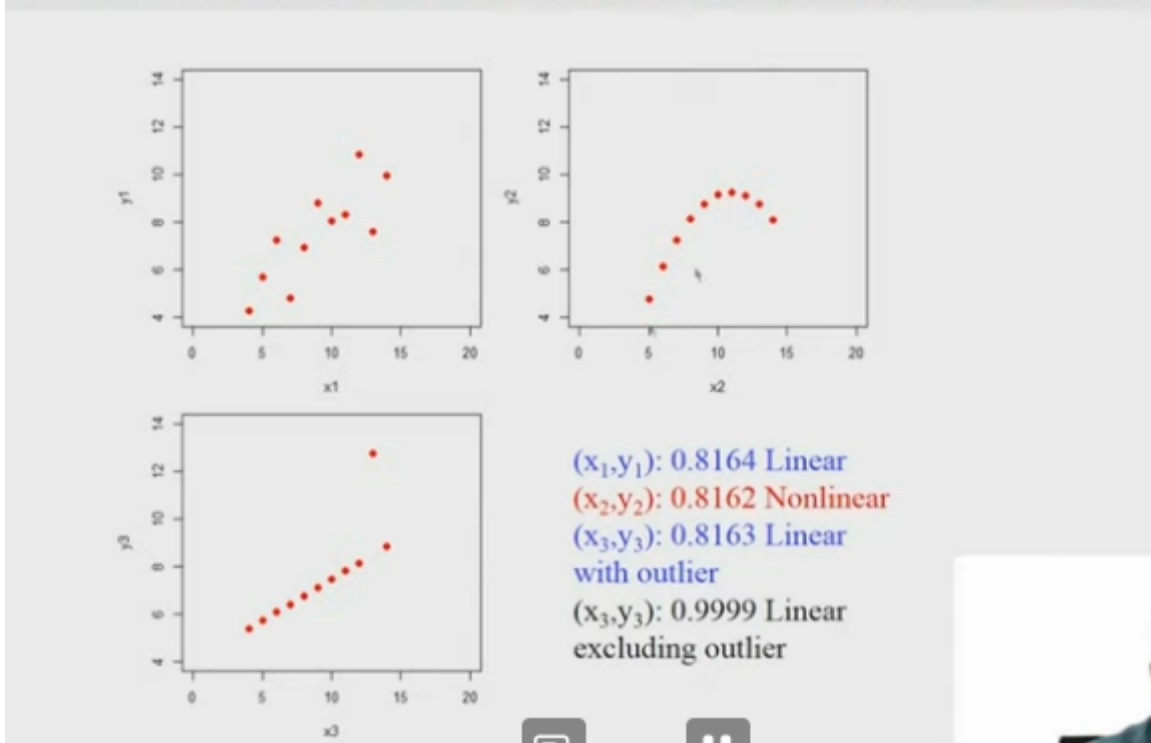
$$r_{xy} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sqrt{(\sum x_i^2 - n \bar{x}^2)} \sqrt{(\sum y_i^2 - n \bar{y}^2)}} = \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}}$$

- $r_{xy}$  takes a value between -1 (negative correlation) and 1 (positive correlation)
- $r_{xy} = 0$  means no correlation

## Pearson's Correlation (Cont.)

- A measure for the degree of linear dependence between  $x$  and  $y$  ,  
Means.... the variables which hold the **Order**. Like indices of a data, rank..
- Cannot be applied to ordinal variables
- Sample size: Moderate (20-30) for good estimate
- Robustness: Outliers can lead to misleading values

# Pearson's Correlation: Anscombe's data



-- notice that, it gives almost same value of *linear* and

*non-linear*.

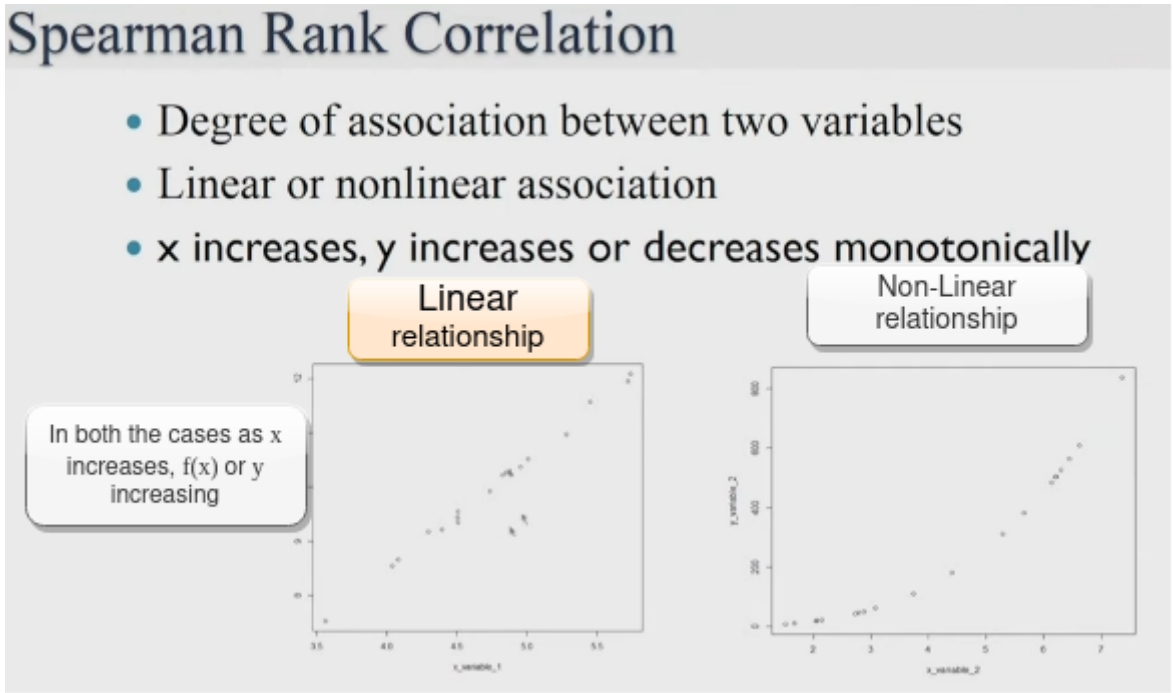
It is a dataset, which contains only 8-13 points.

Example:

## Pearson's Correlation (Cont.)

- Example: Nonlinear

## Spearman Rank Correlation



### Spearman Rank Correlation

- Spearman rank correlation computation for n observations:

$$r_s = 1 - \frac{6\sum d_i^2}{n(n^2-1)}$$

$d_i$  is the difference in the ranks given to the two variables values for each item of the data

- Example:

**Special Case...**  
 When multiple occurrences of same value, their allotted ranks are averaged.  
 Here... for first 7 -- rank 6, and for next 7 -- rank 10, as both are repeated (twice) its  $\frac{6+7}{2} = 6.5$  in all it occurrences

Number	1	2	3	4	5	6	7	8	9	10
X <sub>1</sub>	7	6	4	5	8	7	10	3	9	2
Y <sub>1</sub>	5	4	5	6	10	7	9	2	8	1
Rank X <sub>1</sub>	6.5	5	3	4	8	6.5	10	2	9	1
Rank Y <sub>1</sub>	4.5	3	4.5	6	10	7	9	2	8	1
d <sup>2</sup>	4	4	2.25	4	4	0.25	1	0	1	0

$r_s = 0.88$

So on...

Next lowest order --> Next lowest rank

Lowest Order --> Lowest rank

And the same for SY, 155 and 55 text(Rank) X, 155

Check...

```
In [14]: import numpy as np
lst = np.array([4, 4, 2.25, 4, 4, 0.25, 1, 0, 1, 0])
n=10
print(1-sum(6*lst)/(n*(n**2-1)))
```

0.8757575757575757

```
In [17]: lst = c(4, 4, 2.25, 4, 4, 0.25, 1, 0, 1, 0)
n=10
1-sum(lst*6)/(n*(n**2-1))
```

0.8757575757575756

## Spearman Rank Correlation

- $r_s$  takes a value between -1 (negative association) and 1 (positive association)
- $r_s = 0$  means no association
- Monotonically increasing  $r_s = 1$
- Monotonically decreasing  $r_s = -1$
- Can be used when association is nonlinear
- Can be applied for ordinal variables

## Kendall Rank Correlation

```
In [1]: from IPython.display import IFrame
IFrame('resources/KendallRankCoeff.html', width=950, height=800)    ## The graphs in the image, are vi.
```

Out[1]:

### Kendall rank correlation coefficient

- Correlation coefficient to measure association between two ordinal variables
- Concordant Pair: A pair of observations  $(x_1, y_1)$  and  $(x_2, y_2)$  that follows the property  $x_1 > x_2$  and  $y_1 > y_2$  or  $x_1 < x_2$  and  $y_1 < y_2$
- Discordant Pair: A pair of observations  $(x_1, y_1)$  and  $(x_2, y_2)$  that follows the property  $x_1 > x_2$  and  $y_1 < y_2$  or  $x_1 < x_2$  and  $y_1 > y_2$

### Kendall rank correlation coefficient

- Kendall rank correlation coefficient

$$\tau_k = \frac{\text{Number of concordant pairs} - \text{Number of discordant pairs}}{n(n-1)/2}$$

- The pair for which  $x_1 = x_2$  and  $y_1 = y_2$  are not classified as concordant or discordant and are ignored.

Ignorable Cases

### Concordant Pairs

### Discordant Pairs

An example...

**How the matrix was filled...?**\_(Looks like, for simple understanding.. consider either Expert1 column or Expert2 col)

>> -- here consider items as  $x_i$  and Expert\_1 as  $y_1$ 's ad Expert\_2 as  $y_2$  -- try making pairs in this way -- and compare like.. take two rows(item's instances) and compare the values (of both experts) only if both concordant -> Concordant, else discordant pairs.. .. its respectively filled in the matrix table beside.

**NOTE** that.. for 1,1  $x$  pairs(i.e, rowS), there is no change .. so as per previous slide -- this is ignored, -- this can be seen as empty value.

- High +ve value indicates -- A strong agreement between associates.
- High -ve value indicates -- A strong disagreement..



# Kendall rank Correlation: Anscombe's data

## Kendall rank correlation coefficient

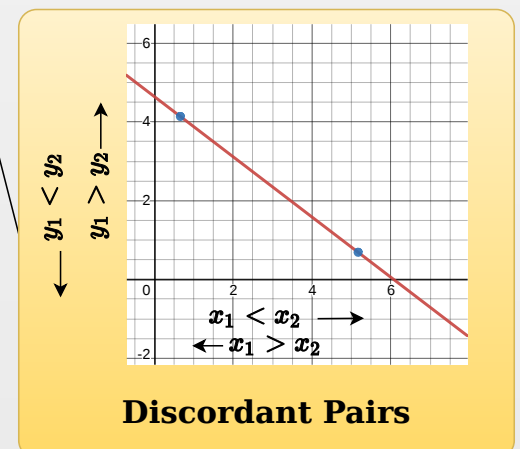
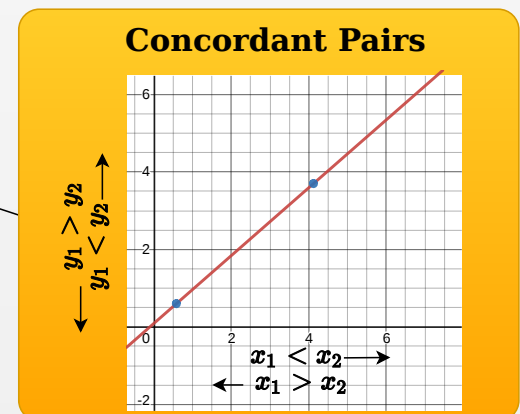
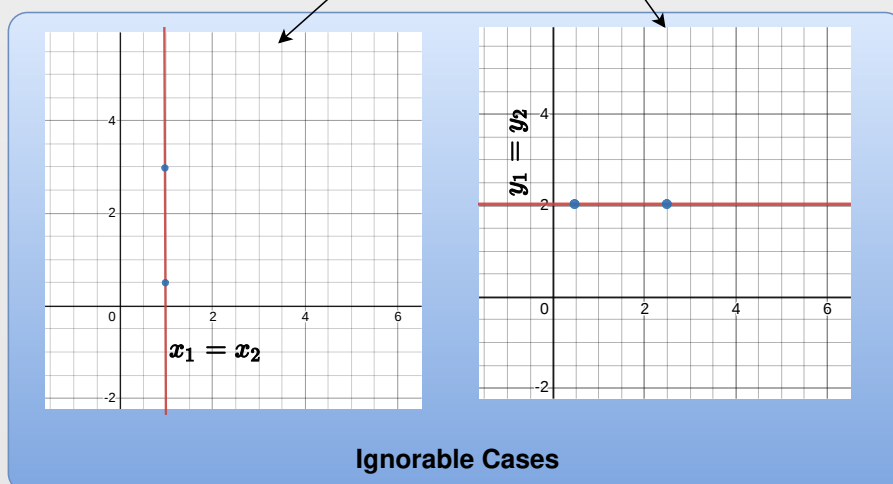
- Correlation coefficient to measure association between two ordinal variables
- Concordant Pair: A pair of observations  $(x_1, y_1)$  and  $(x_2, y_2)$  that follows the property  $x_1 > x_2$  and  $y_1 > y_2$  or  $x_1 < x_2$  and  $y_1 < y_2$
- Discordant Pair: A pair of observations  $(x_1, y_1)$  and  $(x_2, y_2)$  that follows the property  $x_1 > x_2$  and  $y_1 < y_2$  or  $x_1 < x_2$  and  $y_1 > y_2$

## Kendall rank correlation coefficient

- Kendall rank correlation coefficient

$$\tau_k = \frac{\text{Number of concordant pairs} - \text{Number of discordant pairs}}{n(n-1)/2}$$

- The pair for which  $x_1 = x_2$  and  $y_1 = y_2$  are not classified as concordant or discordant and are ignored.



summary..

## Module-2: Linear Regression

### LINEAR REGRESSION

Let's start with the motivation..

#### Motivation

- Purpose is to build a functional relationship (model) between *dependent variable(s)* and *independent variable(s)*
- Example
  - Business : What is the effect of price on sales? (Can be used to fix the selling price of an item)
  - Engineering : Can we infer difficult to measure properties of a product from other easily measured variables? (mechanical strength of a polymer from temperature, viscosity or other process variables) – also known as a soft sensor

In Engineering... (*for the above examples*) -- finding their data via measurement instruments is difficult, but it can be inferred with some given parameters -- and its needed continuously, to make inferences. In these cases, the usage of model is absolutely needed.

Brushing up basic terms...

#### Regression - Basics

- One of the widely used statistical techniques
- Dependent variables also known as *Response variable*, *Regressand*, *Predicted variable*, *output variable* - denoted as variable/s  $y$
- Independent variable also known as *Predictor variable*, *Regressor*, *Exploratory variable*, *input variable* - denoted as variable/s  $x$

#### Regression types

- Classification of Regression Analysis
  - Univariate vs Multivariate
    - *Univariate*: One dependent and one independent variable
    - *Multivariate*: Multiple independent and multiple dependent variables
  - Linear vs Nonlinear
    - *Linear*: Relationship is linear between dependent and independent variables
    - *Nonlinear*: Relationship is nonlinear between dependent and independent variables
  - Simple vs Multiple
    - Simple: One dependent and one independent variable (SISO)
    - Multiple: One dependent and many independent variables (MISO)

Types of Regression:

How can we go with choosing Regression technique?

## Regression analysis

- Is there a relationship between these variables?
- Is the relationship linear and how strong is the relationship?
- How accurately can we estimate the relationship?
- How good is the model for prediction purposes?

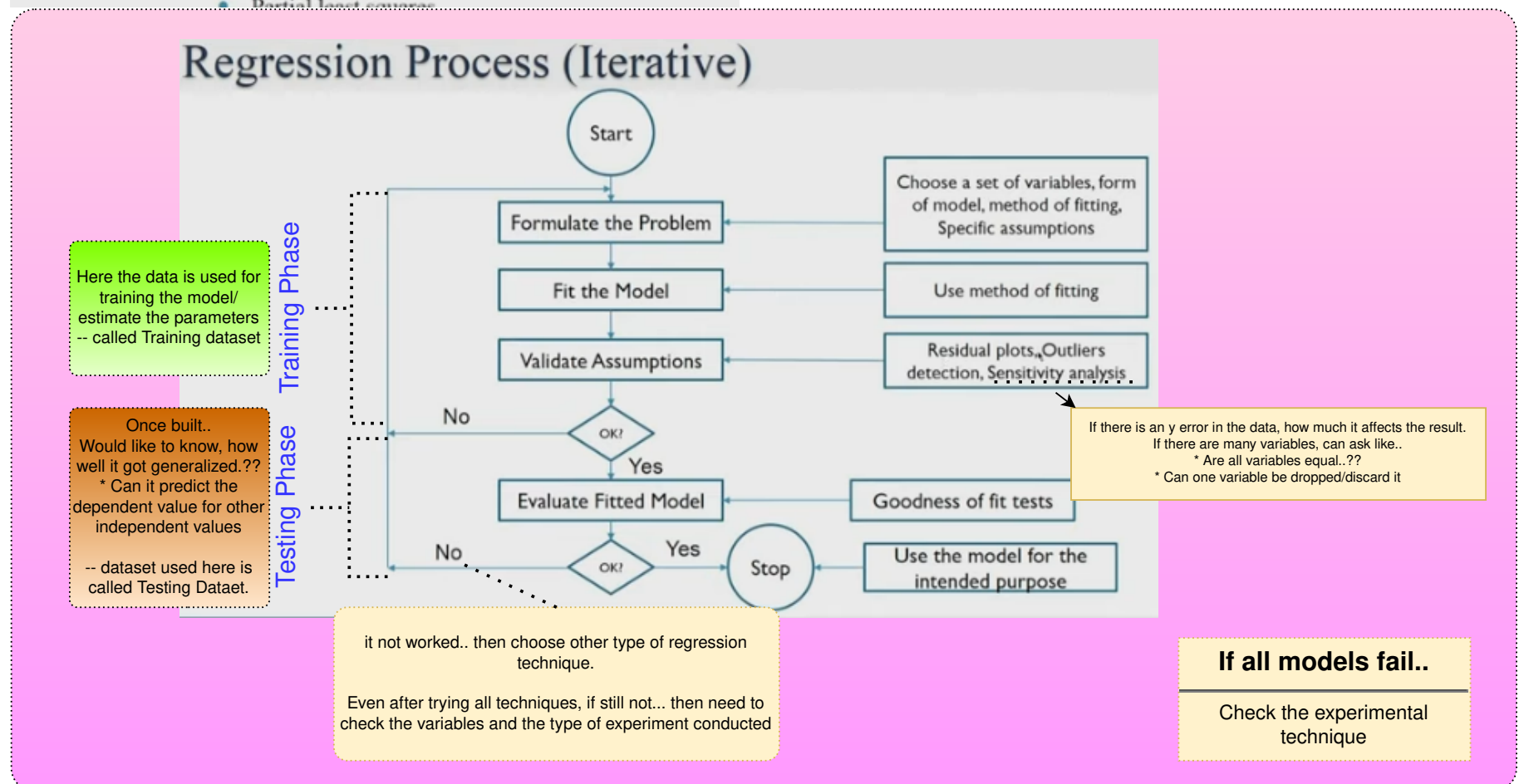
-- needed at choosing the model

and even at developing the model..<br/? Various techniques to go with regression..

## Regression methods

- Linear regression methods

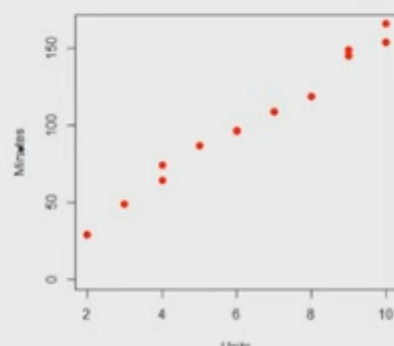
- Simple linear regression
- Multiple linear regression
- Ridge regression
- Principal component regression
- Lasso
- Partial least squares



Let's see some of the techniques..

## Ordinary Least Squares (OLS)

- Fourteen observations obtained on time taken in minutes for service calls and number of units repaired
- Objective is to find relationship between these variables (useful for judging service agent performance)



-- In this (say) a service agent

reports to his boss -- how many he repaired in the given time. Such way here 14 points are taken. -- This can be used for the purposes like,



- Measuring the efficiency of the service-man or decision on increasing his salary.. or
- Productivity level of the comany.. like some of those..

Why error..( $\epsilon_i$ )??

- May be due to the model we've chosen is inadequate for this data.
- May be some errors in the measurement of  $x$  and  $y$ 's..'

Here we assume that,  $x$  has no error -- it has perfectly measured.  $y$  can contain.  
This gives the decision of choosing the independent and dependent variable..

The one which had no errors --> Independent variable  
Which had some errors -> Dependent variable

Here in the given example...

**Units** is chosen as *Independent variable* -- with an idea that, he bills all the customers and the same copy-slips are returned in the office, so this can't be an error unless someone transcribes it.  
**Minutes** as *Dependent variable* as his measurement may include the timings like traffic, location(if far, takes much time than the service-time).. and many such.

But, its also a arguable point that, some even tells to use **The one to be predicted as *Dependent variable* and with which we can do that as *Independent variable***, but, after building the model with the obtained equation, one can go in either way right..

If both contain error, then we need to go for some other methods like *Sensitive Linear regression*..

Ordinary Least Squares (OLS)

Linear model between  $y_i$  and  $x_i$ ,  $i = 1, \dots, n$

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

$\epsilon_i$  captures...  
The vertical-difference between  
the predicted\_line and  
observed\_points(i)

Error in only dependent variable and  
no error in independent variable:

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

The sum of squares of errors (SSE)

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

Squaring and totalling all such those..

The minimization of SSE gives estimates of  $\beta_0$  and  $\beta_1$

By some calculus techniques..  
we approach at..

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

Cross co-variance between x and y  
Like the Pearson's correlation

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

Co-variance of x

If the intention is to pass the line  
through the origin, then put, it to 0.  
But should not force always....

Minutes

Units

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

tells the vertical-difference

$$\epsilon$$

$$\beta_0$$

tells the intercept  
(where it touches y-axis when  $x=0$ )

OLS: Testing Goodness of Fit

If 'x' had influence on y, then it should able to reduce the variability.  
i.e., Should able to do a better prediction.and difference  
between

$$y_i \& \hat{y}$$

should be low.

Prediction using the regression equation:  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

Coefficient of determination -  $R^2$  is a measure of variability in output variable explained by input variable

Observed Value

Predicted value

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

Variance

Explains the variability which is explained by the individual variable  $X_i$

Total variability in y

Tells how much variability is present in the given data.  
~~ This much variability exists in data.

$R^2$  values: Between 0 and 1

➤ Values close to 0 indicates poor fit

➤ Values close to 1 indicates a good fit (However, should not be used as sole criterion to judge that a linear model is adequate)

Adjusted  $\bar{R}^2$

To balance the no. of variables used in the above equation....

$$\bar{R}^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2 / (n - p - 1)}{\sum (y_i - \bar{y})^2 / (n - 1)}$$

Equation Interpretation

If the  $R^2$  is approx. == to the Dr. then we get near to 1, upon SUB with 1, yields to 0.

-----This happens when 'x' has very little impact on explaining 'y' (i.e., probably no relationship)

If  $R^2$  is approx == 0, then it yields a value close to 0.  
Hence, when this subtracted with 1, it yields a value close to 1.

--- This explains that  $x_i$  can explain the variation in  $y_i$ . -- i.e., there is a strong relationship.

complicates

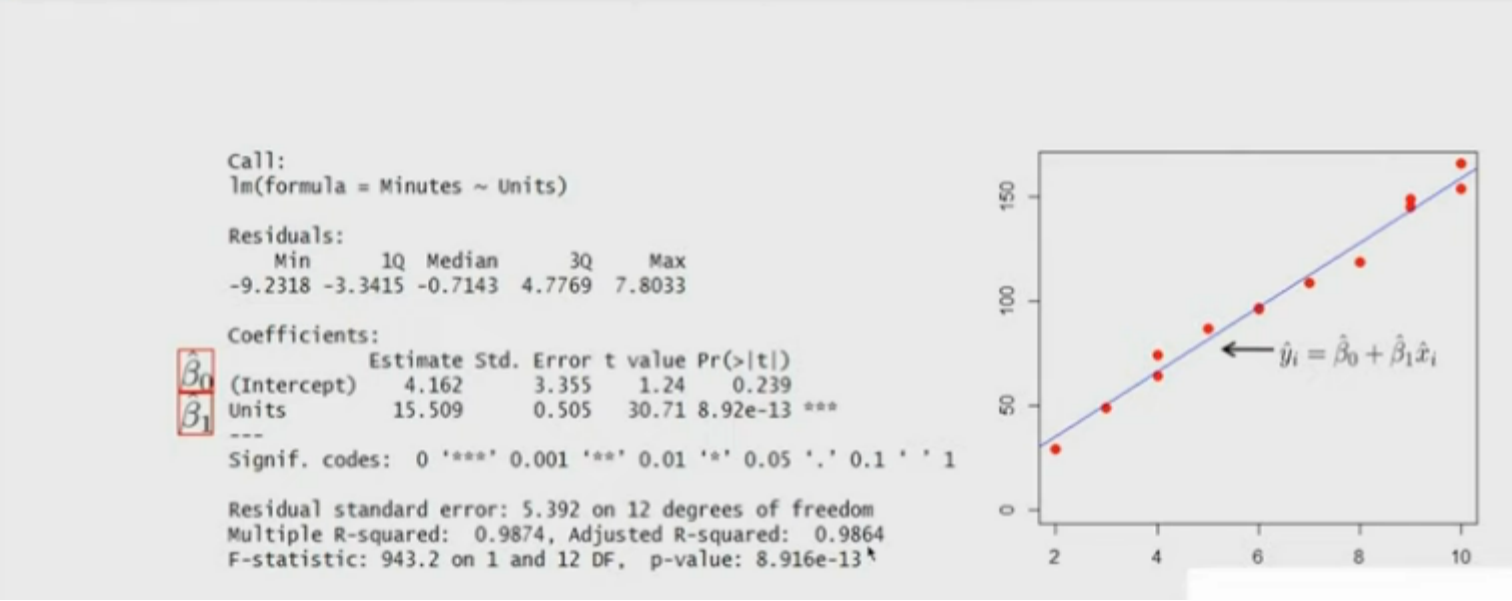
Poor Fit

complicates

Good Fit

Getting the value closer to 1, doesn't say the model is adequate to data, its an assurance to go with and check with other too.

OLS: Example using R



Module-3: Model Assessment

Now, we are done with the (basic version) of model building. Now lets assess it.. and there are various ways to achieve that.. So, what are the questions that one can take..

OLS Model Assessment and Improvement

These are focused in this lecture..

- How good is a linear model?
- Which coefficients of the linear model are significant (Identify important variables)
- Can we improve quality of linear model?
  - Are assumptions made about errors reasonable?
    - Normality: Errors are normality distributed
    - Homoscedasticity: Errors in different samples have same variance
- Are there bad measurements in the data (outliers)

If not, try fitting with other model

Basically, this is mostly important in **Multi-Variable Regression**. --where we have multiple independent variables, and need to decide, whether all are important..?? Can we drop any..

Before starting, need to estimate some of the parameters..

**OLS: Properties of Estimates**

Hat symbols denote that: **These are estimates**

i.e.,  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are the estimates of true intercept and slope

Intercept

Slope

Both  $\hat{\beta}_0$  and  $\hat{\beta}_1$  estimates are unbiased

"E"  $E[\hat{\beta}_0] = \beta_0, E[\hat{\beta}_1] = \beta_1$

means **Estimation or Mean**

Variance of the estimates

Variance of independent variable  $var[\hat{\beta}_1] = \frac{\sigma^2}{S_{xx}}, var[\hat{\beta}_0] = \sigma^2 \frac{\sum x_i^2}{n S_{xx}}$

Estimate of  $\sigma^2$

Variance of independent variable (x) i.e.,  $\sum (\bar{x} - x_i)^2 / n$

Distribution of slope estimate  $\hat{\beta}_1 \sim \mathcal{N}(\beta_1, \frac{\sigma^2}{S_{xx}})$

**These stmts tell that...**

If could repeat the experiment for some n times, and average(mean) all those, then it reaches the true values.

-- NOTE: These values are unknown to us intially, and the obtained ones are called **UnBiased estimates**.

We get different different estimates, depending on sample So, what's the strength of these.. We can show from the assumptions that .. (as present as expression..)

*This depends on the instrument measured. Like we measured for errors and accuracy in for Vernier calipers, Screw Gauge tools usage in Physics Lab (Masabtank-clg)*

Already two of the estimates are used (beta\_0-cap and beta\_1- capt) for beta\_0 and beta\_1. Therefore only remaining n-2 samples are available for estimating it.

-If had only 2 samples, then surely result is 0.

Why n"-2"??

Mostly, this won't be given.. But, it can be estimated as...

**SSE: Sum Squared Error**

**Normal Distribution**

Once we derived the distribution of paramters, we can go for hypothesis testing to decide, whether are these significantly equal to 0. And can also derive the confidence intervals.

OLS: Confidence Intervals on regression coefficients

Recollect that.. an interval contains the lower and upper boundaries

And also the concepts and worked-out problems in COSM subject(JNTUH)

□ 95% two-sided confidence intervals (CI) for  $\hat{\beta}_0$  and  $\hat{\beta}_1$

Lower critical value

Upper critical value

Degrees of Freedom(DoF). As two were already used, left-over are 14-2=12 (Anscombe's DataSet)

$\beta_1 \in [\hat{\beta}_1 - 2.18 s_{\hat{\beta}_1}, \hat{\beta}_1 + 2.18 s_{\hat{\beta}_1}]$

Standard Deviation of  $\hat{\beta}_1$

Here its for 2.5% critical value.

$s_{\hat{\beta}_1} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{(n-2)S_{xx}}}$

o longer al values stribution. na^2 is ata, which priorly.

Once after constructing 95% interval, we can test whether unknown  $\beta_0$  and unknown  $\beta_1$  are equal to 0 or not.. Let's do that.. But.. **Why needed this Hypothesis test..??**

We've fitted the linear model, assuming that,

- 1. linear dependence between x and y and obtained an estimate of  $\beta_1$ , and also...
- 2. fitted the intercept term.

Now, we may want to ask,

- 1. **Is the intercept term ( $\beta_0$ ) Significant:** \* - May be line should pass through origin or not
- 2. **Slope term ( $\beta_1$ ):** May be variable y does not depend on x -- i.e., not depending in significant manner(i.e.,  $\beta_1$  is close to 0)

\* Significant

Considering the equation's view.... evaluating with 0 (for intercept term, as its the const) is as same as not writing the intercept right..?? -- so, does it had significance..?? if 0 not, and !=0 , yes...

For graphical view..

- 1. **For intercept:** Does the line doesn't pass through origin or not.
- 2.

Now, by NULL Hypothesis, we are testing for  $\beta_1 = 0$  vs  $\beta_1 \neq 0$ .

If  $\beta_1 = 0$  ( $H_0$  is accepted) true, indicates that...**Independent variable ( x ) has no effect on dependent variable( y )**.  
If rejected NULL Hypothesis (means accepting Alternative Hypothesis).... concludes that.. **Independent variable ( x ) has some effect on dependent variable( y )**.