Dependent & Independent
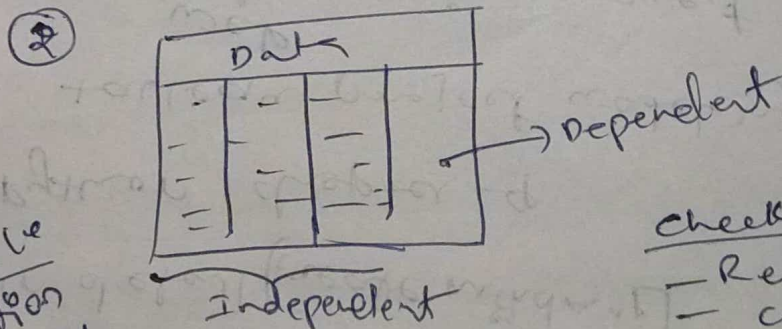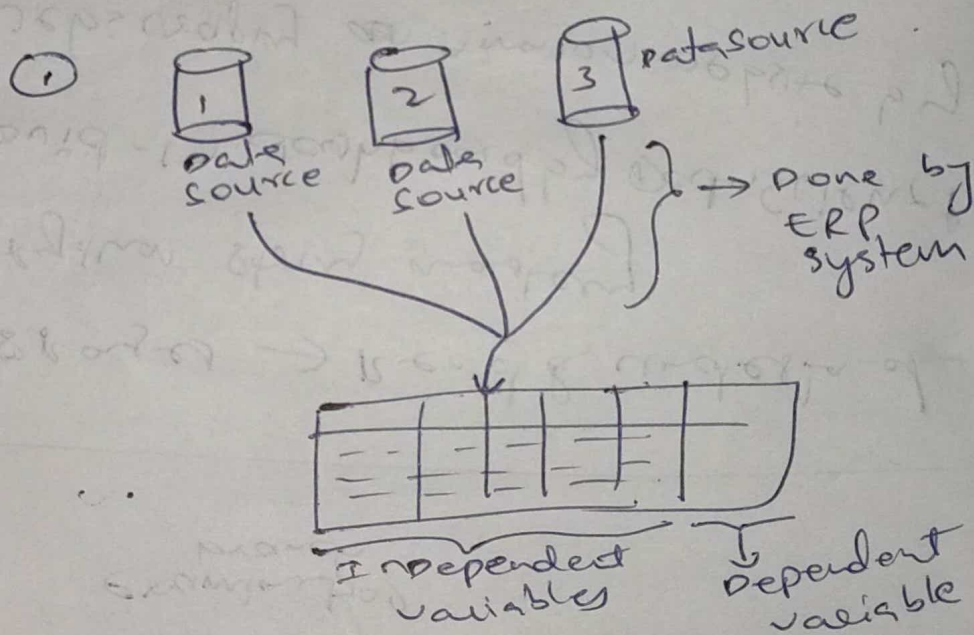variable { variable
    → Association
      relationship


Linear Regression [supervised]
              (predictive Analytic

8 steps in building CR-model
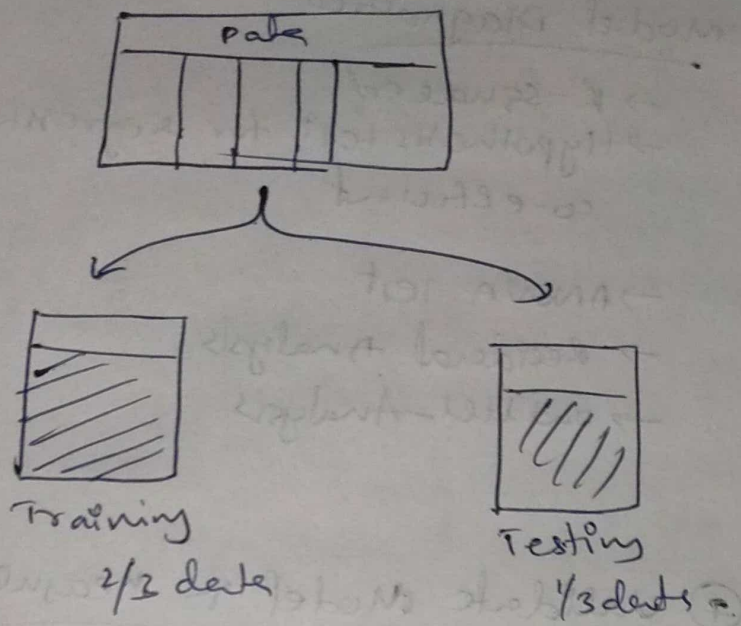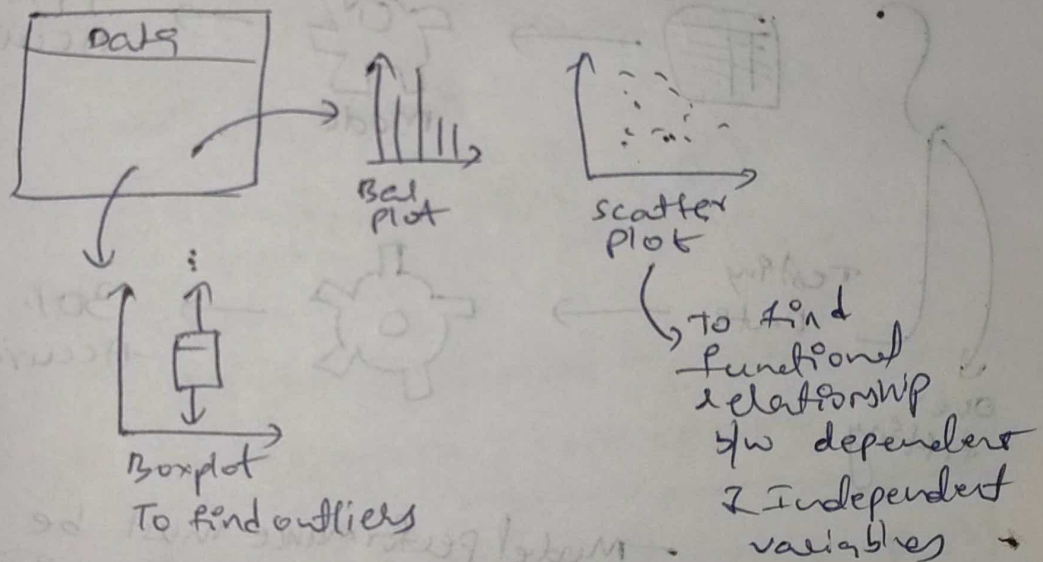

→ several iterations
  needed

① 
[1]        [2]        [3]  → Data source
Data      Data
Source    Source        } → Done by
                            ERP
                            system

Independent          Dependent
variables            variable

② 
| Data | | |
| --- | --- | --- |
| — | — | — |
| — | — | — |
| — | — | — |
              → Dependent

Independent

Check for
 — Reliability
 — completeness
 — Accuracy
 — missing data
 — outliers

Techniques
To handle
1) Data Imputation
2) New variables
   deriving
3) Handle
   categorical data by
   encoding

③



Data

Training
2/3 data

Testing
1/3 data

④



Data

Bar
plot

scatter
plot

Boxplot
To find outliers

→ To find
functional
relationship
b/w dependent
& Independent
variables

⑤ OLS → To find Regression
parameters

$$y = \beta_0 + \beta_1 x + \epsilon$$
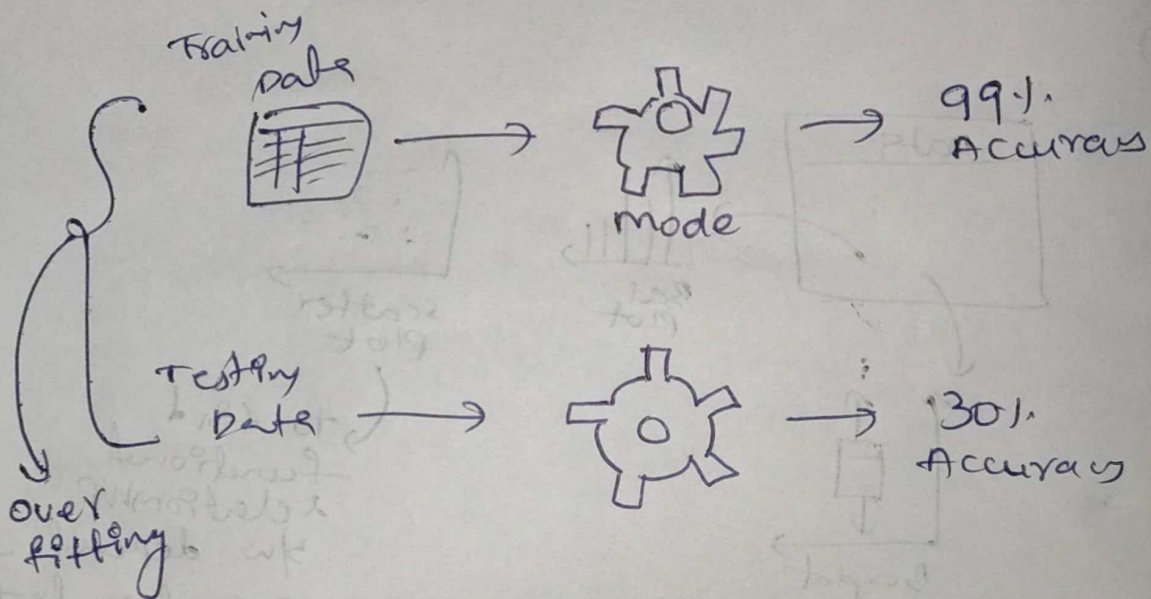
const        grade
             in 10th

you will get $\beta_0$ & $\beta_1$

(6) Model Diagnostics

→ R-squared
→ Hypothesis test for regression co-efficient
→ ANOVA Test
→ Residual Analysis
→ outlier Analysis

(7) Validate Model & measure Accuracy

Training Data



→ 99% Accuracy

mode

Testing Data

→ 30% Accuracy

over fitting

— Model performance must be consistent on both training & testing dataset

— cross validate the model using multiple training and test datasets

(8) Deploy the model

— According to the Business rules

$$\boxed{Y = \beta_0 + \beta_1 X + \varepsilon}$$

$Y$ = Dependent variable

$(\beta_0, \beta_1)$ = Regression parameters

$\varepsilon$ = Random Error

$X$ = Independent variable

For n observations,

$$\boxed{Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i}$$

where $i = 1, 2, 3 \cdots n$

$$\boxed{\varepsilon_i = Y_i - \beta_0 + \beta_1 X_i}$$

$\beta_0$ & $\beta_1$ can be estimated by minimizing sum of squared Errors (SSE)

$$SSE = \boxed{\sum_{i=1}^{n} \varepsilon_i^2 = \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X)}$$

values of $\beta_0$ & $\beta_1$ are taken by partial derivatives of SSE

$$\hat{\beta_1} = \sum_{i=1}^{n} \frac{(X_i - \bar{X})(Y_i - \bar{Y})}{(X_i - \bar{X})^2}$$

$$\hat{\beta_0} = \bar{Y} - \beta_1 \bar{X}$$

Above method is called as OLS method.
It yields BLUE (Best Linear Unbiased Estimates)

## Myths

→ Residuals follow Normal Distribution

→ $(\varepsilon_i)$ is constant for varous independent variables

→ $\varepsilon_i$ & $x$ are correlated

→ functional relationship b/w outcome variable & feature is correctly defined.

## Properties

① $Mean(Y_i) = \hat{\beta_0} + \hat{\beta_1} x$

② $Y_i$ follows Normal Distribution with mean $\hat{\beta_0} + \hat{\beta_1} x$ & variace $VAR\,\varepsilon_i$

$$Y = \textcircled{$\beta_0$} + \textcircled{$\beta_1$} X + \varepsilon$$

constant
(1)

$\hookrightarrow$ OLS API estimates
$(\beta_1)$ only to estimate
$\beta_0$ a constant term
need to be added
as new feature.

$$Y_i = \beta_0 + \beta_1 x + \varepsilon$$

OLS gives $(\beta_0 \& \beta_1)$

$$Y_i = 30587.285652 + \begin{pmatrix} 3560.587 \\ * 62.1 \end{pmatrix}$$

for every 1% increase in grade '10' salary
increases by 3560.587

Rahman