06/08/2020    ( Pre-Processing Data)

{ Missing Values}

missing value occur when no data value is stored for a variable (feature) in an observations.

Represented as →  ?, N/A, "0".

Dealing with missing Values →

Check with data collection Source :→

* Drop the missing Values :-
   • drop the variable
   • drop the data entry

* Replace the missing value -
   • replace it with an average ( of similar datapoint)
   • replace it by frequency (if data is categorical)

* Leave it as missing data.

How to drop missing Values in Python →

use →'
       dataframes.dropna( ) :
       axis=0  drop the entire row
       axis = 1   drop the entire cloumn

## Data Formatting »

Non-formatted:
- Confusing
- hard to aggregate
- hard to compare

formatted:-
- more clear
- easy to aggregate
- easy to compare

## Correcting Data types »

* To identify data types:-
  (1) dataframe.dtype s ( )
* To convert data types:-
  dataframe.astype ( )

for e.g   df ["Price"] = df[price].astype ('int')

## Data Normalization »

| age | income |
|-----|--------|
| 20  | 100000 |
| 30  | 90000  |
| 40  | 500000 |

$\Rightarrow$

| age | income |
|-----|--------|
| 0.2 | 0.2    |
| 0.3 | 0.09   |
| 0.7 | 1      |

## method of normalization data »

(1.) Simple feature Scaling

$$X_{new} = \frac{X_{old}}{X_{max}}$$

(iv) Min-max →

$$x_{new} = \frac{x_{old} - x_{min}}{x_{max} - x_{min}}$$

(v) Z-Score →

$$x_{new} = \frac{x_{old} - \mu}{\sigma}$$

## Binning In Python →

- Binning → Grouping of values into "bins"
- Convert numeric into categorical variables
- Group a set of numerical values into a set of "bins".

Turning Categorical variable into Quantative Variable in Python →

One-hot-encoding →

using Pandas library →

(I) Pd. get dummies ()  :
convert categorical variable to dummy variable (0 or 1).

e.g - Pd. get dummies ( df ('fuel').