| Preprocessing Steps | R | Python |
|---|---|---|
| **Reading_file** | read.csv(file = "file_name" , header=TRUE , sep="," , stringsAsFactors = F ,  na.strings=" " ) | import pandas as pd<br>dataframe = pd.read_csv("file_name ") |
| **Statistics for columns** | summary(dataframe) | dataframe.describe() |
| **Summary of dataframe** | str(dataframe) | dataframe.info() |
| **Column names** | colnames(dataframe) | dataframe.columns |
| **Changing Index** | rownames(dataframe) = dataframe$column<br>dataframe$column=NULL | dataframe.set_index('column_name', inplace=True) |
| **Checking columns  datatypes** | is.numeric(dataframe$column) | dataframe.column.dtype |
| **Changing the datatypes** | as.numeric(dataframe$column) | dataframe.column.astype('int') |
| **Finding missing values** | colSums(is.na(dataframe))<br>*# Returns no. of missing values per column*<br>sum(is.na(dataframes))<br>*# Returns total number of missing values per dataframe* | dataframe.isnull().sum()<br>*# Returns no. of missing values per column* |
| **Imputing Missing values  (Numeric)** | library(DMwR)<br>df_imputed  = centralImputation(dataframe) | from sklearn.preprocessing import **Imputer**<br>imputer = Imputer(missing_values='NaN', strategy='mean',axis =0)<br> imputer.fit_transform(dataframe.column)<br>*# Returns numpy array* |
| **Imputing Missing values (Categorical)** | NA | dataframe.column.apply(lambda x : x.fillna(x.value_counts().index[0])) |
| **Categorical  to Numeric** | library(dummies)<br>dummy(dataframe$column) | pd.get_dummies(dataframe.column,  prefix= 'column_name') |
| **Scaling the Numerical Attribute (Range)** | library(vegan)<br>dataframe = decostand(x =dataframe[ , c(all_numeric_columns)], method = "range",MARGIN = 2) | from sklearn.preprocessing import **MinMaxScaler**<br>minmax_scaler = MinMaxScaler()<br>minmax_scaler.fit_transform(dataframe[all_numeric_columns])<br>*# Returns Numpy array* |
| **Scaling the Numerical Attribute (Standardize)  (Mean = 0, SD = 1)** | dataframe = decostand(x =dataframe[ , c(all_numeric_columns)], method ="standardize",MARGIN = 2) | from sklearn.preprocessing import **StandardScaler**<br>std_scalar = StandardScaler()<br>std_scalar.fit_transform(dataframe[all_numeric_columns])<br> *# Returns Numpy array* |