Milestone 2: Data Collection & Preparation

ML depends heavily on data. It is the most crucial aspect that makes algorithm training possible. So this section allows you to download the required dataset.

Activity 1: Collect the dataset

There are many popular open sources for collecting the data. Eg: kaggle.com, UCI repository, etc.

In this project we have used .csv data. This data is downloaded from kaggle.com. Please refer to the link given below to download the dataset.

Link: https://www.kaggle.com/datasets/mansoordaku/ckdisease

As the dataset is downloaded. Let us read and understand the data properly with the help of some visualisation techniques and some analysing techniques.

Note: There are several techniques for understanding the data. But here we have used some of it. In an additional way, you can use multiple techniques.

Activity 1.1: Importing the libraries

Import the necessary libraries as shown in the image. (optional) Here we have used visualisation style as fivethirtyeight.

```
Importing Libaries

pandas pd

murpy np

collections Counter c

matplotlib.pyplot plt

seaborn sms

missingno mano

skinaro.metrics accuracy_score, confusion_matrix

sklearo.model_selection train_test_split

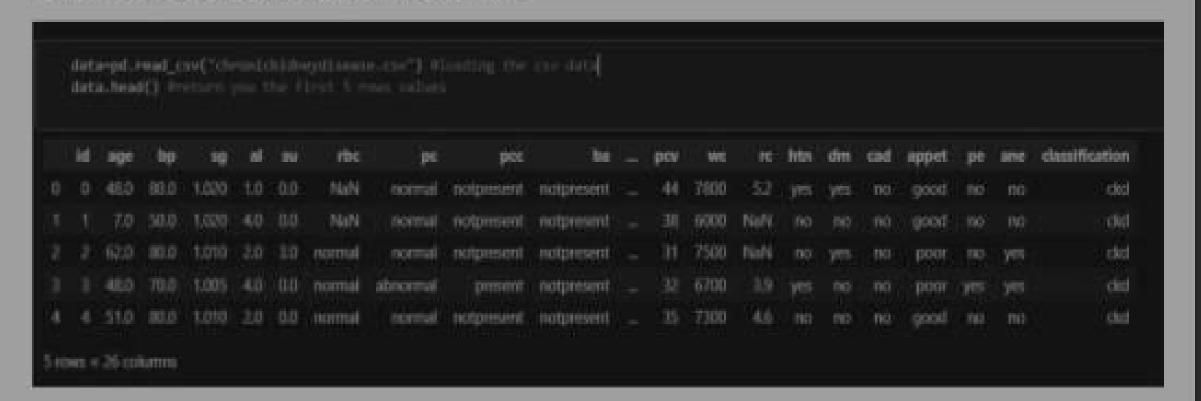
sklearo.preprocessing LabelEncoder

pickle
```

Activity 1.2: Read the Dataset

Our dataset format might be in .csv, excel files, .txt, .json, etc. We can read the dataset with the help of pandas.

In pandas we have a function called read_csv() to read the dataset. As a parameter we have to give the directory of the csv file.



Activity 2: Data Preparation

As we have understood how the data is, let's pre-process the collected data.

The download data set is not suitable for training the machine learning model as it might have so much randomness so we need to clean the dataset properly in order to fetch good results. This activity includes the following steps.

- Rename the columns
- Handling missing values
- Handling categorical data
- Handling Numerical data

Note: These are the general steps of pre-processing the data before using it for machine learning. Depending on the condition of your dataset, you may or may not have to go through all these steps.

Activity 2.1: Rename the columns

```
data.columns
index(['age", 'bs', 'ag", 'ai", 'su', 'rbc', 'ac', 'scc', 'be', 'bgr', 'su',
     "appet", 'pe', 'ane', 'classification'),
     dtypes"object")
   data.columns#['age', blood pressure', specific gravity', albumin',
                    'sugar', 'red_blood_cells', 'pus_cell', 'pus_cell_clumps', 'bacteria',
                    'blood glucose random', 'blood urea', 'serum creatinine', 'sodium', 'potassium',
                    'hemoglobin', 'packed_cell_volume', 'white_blood_cell_count', 'red_blood_cell_count',
                    'hypertension', 'diabetesmellitus', 'coronary artery disease', 'appetite',
                    'pedal_edema', 'anemia', 'class'] # minustry giving the name of the columns
   data_columns
Infex[] "sge", "blood pressure", "specific gravity", "altumin", "super",
      "red_blood_sells", "pus_cell", "pus_cell_clumps", "bacteris",
      "blood glucone rundom", "blood pres", "serum tresticine", "sodium",
      'potassium', "hemoglobin', "packed cell volume",
      "white bland cell count", "red blood cell count", "hypertension",
      "Winbetesmellition", "communey_artery_disease", "appetite",
      'pedal_edoma', 'aventa', 'class'),
```

Activity 2.2: Handling missing values

 Let's find the shape of our dataset first. To find the shape of our data, the df.shape method is used. To find the data type, df.info() function is used.

```
### State | Section | Sect
```

```
data.ismull().any() fit will return from if my release is furting out! wither
                          Trus
blood pressure
                          True
apecific gravity
                          TTUE
High
red himse sells
pun_intkl
                          TPM
                          True
hilood_seems
serum_creatistime
                          True
                          True
potasitum
                          True
howing Torbital
                          Trus
packed cell volume
white binet call court
                          True
red_blood_cell_comet
                          17340
hypartensise.
                          7798
diobetenecLibbo
                          True
coronary_artery_disease
                          True
                          77100
profet nimes
anest a
```

```
data['blood glucose random'].fillna(data['blood glucose random'].mean(),inplace="read")
data['blood_pressure'].fillna(data['blood_pressure'].mean(),inplace="read")
data['blood_urea'].fillna(data['blood_urea'].mean(),inplace="read")
data['packed_cell_volume'].fillna(data['packed_cell_volume'].mean(),inplace="read")
data['potassium'].fillna(data['potassium'].mean(),inplace="read")
data['red_blood_cell_count'].fillna(data['red_blood_cell_count'].mean(),inplace="read")
data['serum_creatinine'].fillna(data['serum_creatinine'].mean(),inplace="read")
data['sodium'].fillna(data['sodium'].mean(),inplace="read")
data['sodium'].fillna(data['sodium'].mean(),inplace="read"),inplace="read")
```

```
data['age'].filina(data['age'].mode()[0],inplace=[row]
data['hypertension'].filina(data['hypertension'].mode()[0],inplace=[row]
data['pus_cell_clumps'].filina(data['pus_cell_clumps'].mode()[0],inplace=[row]
data['appetite'].filina(data['appetite'].mode()[0],inplace=[row]
data['albumin'].filina(data['pus_cell'].mode()[0],inplace=[row]
data['pus_cell'].filina(data['pus_cell'].mode()[0],inplace=[row]
data['red_blood_cells'].filina(data['red_blood_cells'].mode()[0],inplace=[row]
data['appetite'].filina(data['bacteria'].mode()[0],inplace=[row]
data['appetite'].filina(data['appetite'].mode()[0],inplace=[row]
data['appetite'].filina(data['appetite'].mode()[0],inplace=[row]
data['diabetesmellitus'].filina(data['diabetesmellitus'].mode()[0],inplace=[row]
data['pedal_edema'].filina(data['pedal_edema'].mode()[0],inplace=[row]
data['appetite'].filina(data['pedal_edema'].mode()[0],inplace=[row]
data['appetite'].filina(data['pedal_edema'].mode()[0],inplace=[row]
data['appetite'].filina(data['pedal_edema'].mode()[0],inplace=[row]
```

Let's now check the count of null values after filling all null values using isnull.sum()

Activity 2.3: Handling Categorical columns

The below code is used for fetching all the object or categorical type of columns from our data and we are storing it as set in variable catcols.

As, you can observe that it gives us the same count of columns which we find previously.

```
i i catcols:
                                              print("Columns :",1)
                                               print(c(data[1])) Funding country for charaking the matter of citaness in the relu-
                                                print('*'*120+'\n')
    (Silvers | Ropertancian
    Constitution 1 201, 'yes': 147, Nat: 233
   Sulumn : packed cell volume
  Counter ([name 18, 1821; 21, 1851; 22, 1861; 29, 1851; 29, 1851; 38, 1851; 38, 1451; 33, 1421; 32, 1351; 32, 1861; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 32, 1351; 3
    "SWY 1 25, 1271 25, 1841 25, 1861 3, 1261 5, 1261 5, 1261 5, 1261 6, 1261 6, 1261 7, 1261 7, 1261 8, 1261 6, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1261 4, 1
    "$2" A, "$4" A, "$7" A, "$7" A, "$7" A, "$5" A, "$6" A
     "5548"1 $, "$"1 $33.
  Counter ( "cas" | 254, "acticle" | 154)
  Columns : consumary artery disease
  Charter [ "NO": NEX, "yet": DA, "(SRO": I, NEW: 2)]
   Counter([ 'be'( 328, 'yes': 48, test 1])
      influence of registrood cell count
Constitut Class: $18, "5.2" | 18, "4.5" | 18, "4.5" | 54, "4.2" | 53, "5.2" | 18, "5.2" | 18, "4.5" | 28, "4.5" | 2, "5.4" | 2, "5.4" | 2, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4" | 3, "5.4
5"1 No. "SCATE NO. "SCATE NO. "SCATE NO. "SCATE NO. "SCATE NO. "ACRES NO. "SCATE NO.
     $5.500 m. ($1.500 m.) ($1.500 m.)
```

```
slumme | red blood balls
            Courter(["mormalf: 281, nen: 152; 'sbrormal'; 47)}
             Desirant Control of the Control of t
      Deligent | people edites
          Disorter [ [ 167 | 323, "yes" | 74, con: 11]
       Enlament i apportite
       Counter (1 good 1 157, 'poor': $2, new 339.
     Collegery 1 pers nucl.
     Courter(("sermal"): 259, "absormal") 75, nam: 423)
         Colomno o distate medilitos
          Innotary ( "en't 258, "pen't 258, "ttel 1 2, "ttel 1 2, "ttel 1 2, "pen't 21)
          Column | puscelli-Compa
       (months (("motoressect": USA; "present": 43; spot 4)).
         Calbumny I while bland rell count
  Counter(Count Set, 1988): 22, Table 2 26, Table 2 26, Table 2 3, T
 BY: 7, "2000": 7, "ADMON B, "SAME": 8, "SAME": 8, "SAME": 8, "SAME": 5, "SAME
 S. TARRELL E. TRONG C. S. TORRELL E. TORRELL E. TRANSIC A. TRANSIC
     WINDS TO THROUGH S. TRANSFE S. TR
 MYS 2, "EXEMPTR 2, "DOME") 2, "ROME") 2, "EXEMPTR 2, "
 #FO E, "$4000"C.E, "25000"C.E, "EXEMPTO E, "COMMON E, "STANO"C.E, "EXEMPTO E, 
         L. CAMARRIO S. TERRORIUS, "EDNOT E. TERRORIUS, "APPRILIE E. "EDNOT E. TERRORIUS, "EDNOT E. TERRORIUS, "E. TERRO
Wining Tabband S. Tabb
```

In the above we are looping with each categorical column and printing the classes of each categorical columns using counter function so that we can detect which columns are categorical and which are not.

If you observe some columns have a few classes and some have many, those columns are having many classes can be considered as numerical column and we have to remove it and add it to the continuous columns.

```
catcols.remove('red_blood_cell_count') # ......... is used for removing a particular column
catcols.remove('packed_cell_volume')
catcols.remove('white_blood_cell_count')
print(catcols)

["spertension', 'class', 'coronary_ortery_disease', 'smemia', 'red_blood_cells', 'bacteria', 'pedal_come', 'specific', 'pu_cell', 'dish
etesmillion', 'pus_cell_clamp')
```

As we store our columns as set, we can make use of remove function which is used to remove the element in our case we can take it as columns.

Activity 2.3.1: Label Encoding for categorical columns

Typically, any structured dataset includes multiple columns with combination of numerical as well as categorical variables. A machine can only understand the numbers. It cannot understand the text. That's essentially the case with <u>Machine Learning algorithms</u> too. We need to convert each text category to numbers in order for the machine to process those using mathematical equations.

How should we handle categorical variables? There are Multiple way to handle, but will see one of it is LabelEncoding.

Label Encoding is a popular encoding technique for handling categorical variables. In this technique, each label is assigned a unique integer based on alphabetical ordering.

Let's see how to implement label encoding in Python using the scikit-learn library.

we have to convert only the text class category columns; we first select it then we will implement Label Encoding to it.

In the above code we are looping through all the selected text class categorical columns and performing label encoding.

```
LABEL INCODENG OF Lawrence
Enuclarity and I had, "yes" | 6825
LABOL DECCRISO OF 1 pedal screen
Doubline ( 1 'mm' : 324, "yes" : 763)
LABEL ENCORENG OF L Apportate
Countre((0) 33%, 2) 823)
LARGE FRETERING OF PRICESOLS
Counter(("metpresent": 378, 'present": 22))
Country (18: 178, 1: 121)
LABEL CHCCCCNE BE: Class
Country ( | 'ctd' | 250, 'wefiche' | 150)
Courter([8: 258, 1: 158]]
*********************************
LABEL ENCOMENS OF COMMUNICATIONS ACTIONS
Courter (Coor: Det, "yes": 3433
Country (CR: 200, T: 345)
LABEL ENCOREMS OF | disbetsomeilities
Country ( 'mo' : 30%, "pee": 52735
Counter ( (8) 383, 31 33773
LABEL ENCORONG OF Paper Conclusion
Counter(['me': 25%, 'pag': 347])
Counter((6: 253, 1: 367))
LARGE THETEODNE DET DUT THEE
Country (("mormal": Nid., "stoormal": 765)
Constant((%: 304, m: 70))
LABEL DECKENG OF PART PART CHAPTE
Counter(('metpresent': 35%, 'present': 43))
Counter((0: 35E, 2: 67])
LARFL ENCORENE OF: red blass calls
Country ( consult: 95%, "stetersel": 67))
```

As you can see here, after performing label encoding alphabetical classes is converted to numeric.

Activity 2.4: Handling Numerical columns

Same as we did with categorical columns, we are majing use of dtypes for finding the continuous columns

```
print("Continous Columns :",i)
print(c(data[i]))
print("""+128+"\n")
```

If we observe the output of the above code we can observe that some columns have few values or you can say classes which can be considered as categorical columns. So, let's remove it and add the columns which we observed into their respective variables.

```
contcols.remove('specific_gravity')
contcols.remove('albumin')
contcols.remove('sugar')
print(contcols)
```

With the help of add() function we can add an element.

```
contcols.add('red_blood_cell_count') // """
contcols.add('packed_cell_volume')
contcols.add('white_blood_cell_count')
print(contcols)

('blood_cell_count', 'packed_cell_volume', 'blood_pressure', 'blood_glucous random', 'sustom', 'becomplaint', 'red_blood_cell_count')
_count', 'age', 'polantiom', 'white_blood_cell_count')
```

```
catcols.add("specific_gravity")
catcols.add("albumin")
catcols.add("nugar")
print(catcols)

("hypertension", 'class", 'sliness", 'coronory_artery_timese", 'avemia', 'sugar', 'red_blood_calls', 'specific_gravity', 'bucteria', 'ped_si_edems', 'specific', 'dus_calls', 'disheteredlifus', 'pen_call_class')
```

In our data some columns some unwanted classes so we have to rectify that also for that we simply use replace()

```
data['coronavy_artery_disease'] = data.coronary_artery_disease.replace('\tno', 'no') * reviewing lims of
c(data['coronary_ortery_disease'])

cmeter(('oi': bia, 'yes': bia, 'ses': bia, non: 2))

data['diabetesmallitus'] * data.diabetesmallitus.replace(to_replace*('\tno': 'no', '\tyes': 'yes', ' yes':'
c(data['diabetesmallitus'])

downter(('pes': 527, 'no': 281, nan: 2))
```