# Problem Statement

# Prediction of Chronic Kidney Disease using Machine Learning

### Introduction:

Technological development, including machine learning, has a huge impact on health through an effective analysis of various chronic diseases for more accurate diagnosis and successful treatment. Kidney disease is a major chronic disease associated with aging, hypertension, and diabetes, affecting people 60 and over. Its major cause is the malfunctioning of the kidney in disposing toxins from the blood. This study analyzes chronic kidney disease using machine learning techniques based on a chronic kidney disease (CKD).

Chronic Kidney Disease (CKD) or chronic renal disease has become a major issue with a steady growth rate. A person can only survive without kidneys for an average time of 18 days, which makes a huge demand for a kidney transplant and Dialysis.

It is important to have effective methods for early prediction of CKD. Machine learning methods are effective in CKD prediction. This work proposes a workflow to predict CKD status based on clinical data, incorporating data prepossessing, a missing value handling method with collaborative filtering and attributes selection. Out of the 11 machine learning methods considered, the extra tree classifier and random forest classifier are shown to result in the highest accuracy and minimal bias to the attributes.

### Predicting Chronic Kidney Disease based on health records

Given 24 health related attributes taken in 2-month period of 400 patients, using the information of the 158 patients with complete records to predict the outcome (i.e. whether one has chronic kidney disease) of the remaining 242 patients (with missing values in their records).

Abstract: This dataset can be used to predict the chronic kidney disease and it can be collected from the hospital nearly 2 months of period.

| Data Set Characteristics: | Multivariate | Number of Instances: | 400 | Area: | N/A |
|---|---|---|---|---|---|
| Attribute Characteristics: | Real | Number of Attributes: | 25 | Date Donated | 2015-07-03 |
| Associated Tasks: | Classification | Missing Values? | Yes | Number of Web Hits: | 196331 |

**Data Set Information:**

We use the following representation to collect the dataset

- age - age
- bp - blood pressure
- sg - specific gravity
- al - albumin
- su - sugar
- rbc - red blood cells
- pc - pus cell
- pcc - pus cell clumps
- ba - bacteria
- bgr - blood glucose random
- bu - blood urea
- sc - serum creatinine
- sod - sodium
- pot - potassium
- hemo - hemoglobin
- pcv - packed cell volume
- wc - white blood cell count
- rc - red blood cell count

- htn - hypertension

- dm - diabetes mellitus

- cad - coronary artery disease

- appet - appetite

- pe - pedal edema

- ane - anemia

- class - class

**Attribute Information:**

We use 24 + class = 25 ( 11 numeric ,14 nominal)

1.Age(numerical)

age in years

2.Blood Pressure(numerical)

bp in mm/Hg

3.Specific Gravity(nominal)

sg - (1.005,1.010,1.015,1.020,1.025)

4.Albumin(nominal)

al - (0,1,2,3,4,5)

5.Sugar(nominal)

su - (0,1,2,3,4,5)

6.Red Blood Cells(nominal)

rbc - (normal,abnormal)

7.Pus Cell (nominal)

pc - (normal,abnormal)

8.Pus Cell clumps(nominal)

pcc - (present,notpresent)

9.Bacteria(nominal)

ba - (present,notpresent)

10.Blood Glucose Random(numerical)

bgr in mgs/dl

11.Blood Urea(numerical)

bu in mgs/dl

12.Serum Creatinine(numerical)

sc in mgs/dl

13.Sodium(numerical)

sod in mEq/L

14.Potassium(numerical)

pot in mEq/L

15.Hemoglobin(numerical)

hemo in gms

16.Packed Cell Volume(numerical)

17.White Blood Cell Count(numerical)

wc in cells/cumm

18.Red Blood Cell Count(numerical)

rc in millions/cmm

19.Hypertension(nominal)

htn - (yes,no)

20.Diabetes Mellitus(nominal)

dm - (yes,no)

21.Coronary Artery Disease(nominal)

cad - (yes,no)

22.Appetite(nominal)

appet - (good,poor)

23.Pedal Edema(nominal)

pe - (yes,no)

24.Anemia(nominal)

ane - (yes,no)

25.Class (nominal)

class - (ckd,notckd)

**Task to Do:**

A. Missing Value Handling: In the pre-processing step, missing values have to be handled based on their distributions to achieve reasonable accuracy
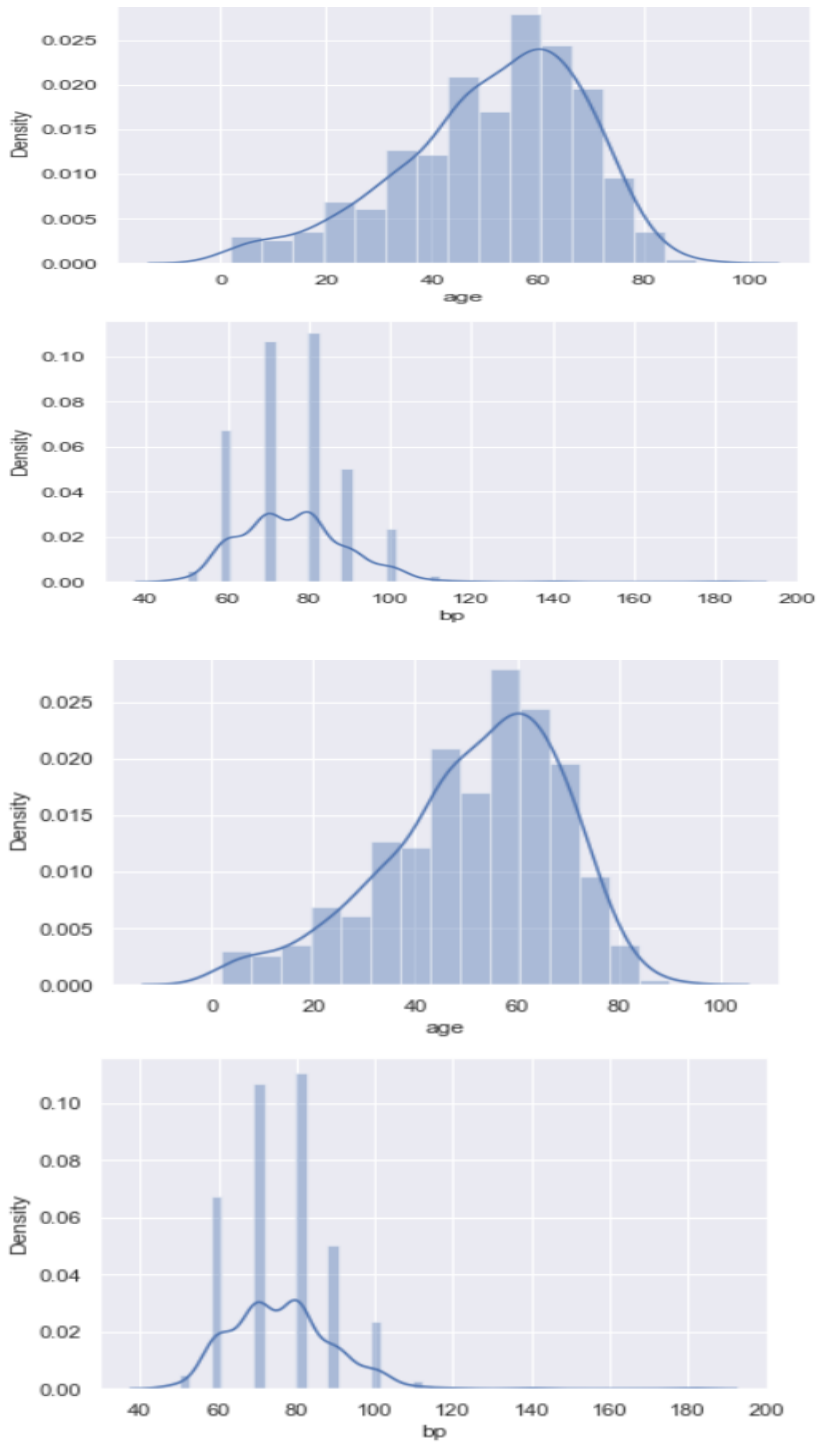
B. Data Preprocessing : Feature Selection

C. Model Training:In this work, 11 classification models can be taking in to consideration . They are logistic regression, k-Nearest Neighbors (KNN) regression, SVC with a linear kernel, SVC with RBF kernel, Gaussian NB, decision tree classifier, random forest classifier, XGB classifier, extra trees classifier, an ada  boost classifier and a classical neural network.

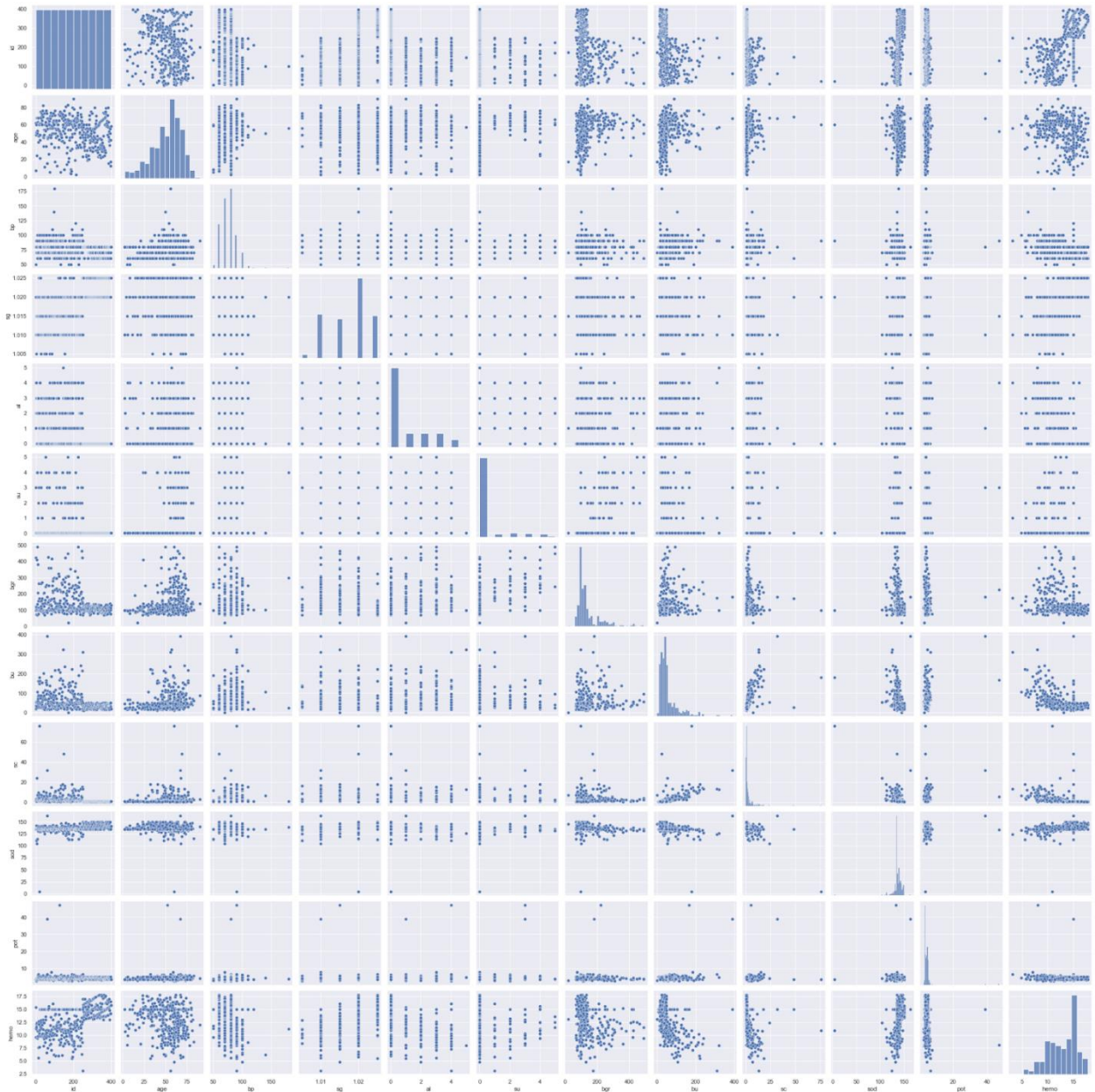D. Find out which fits the best and gives the best Accuracy to the model.

**Model Evaluation:**

**Sample Output:**
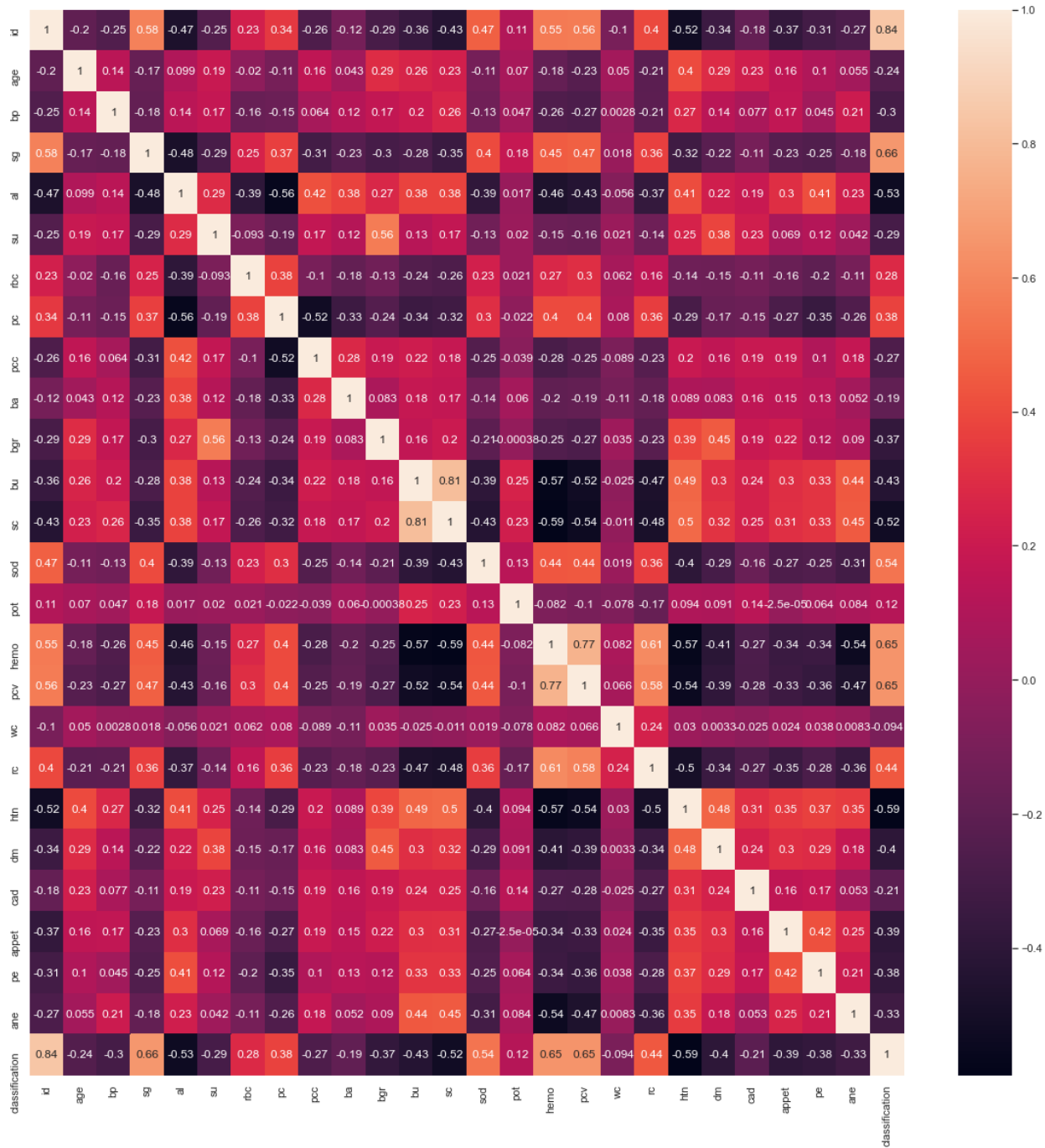
**Sample Output:**

**Sample Output:**

## Sample Output:

[ ]

```
----------------- Model Performance for  DecisionTreeClassifier(random_state=14)  with PCA ---------

                 precision    recall  f1-score   support

             0       1.00      0.98      0.99        48
             1       0.98      1.00      0.99        52

      accuracy                           0.99       100
     macro avg       0.99      0.99      0.99       100
  weighted avg       0.99      0.99      0.99       100


Confusion Matrix:
[[47  1]
 [ 0 52]]
```

```
------------------ROC FOR  Decison Tree  +PCA ----------------------------
```



Receiver Operating Characteristic