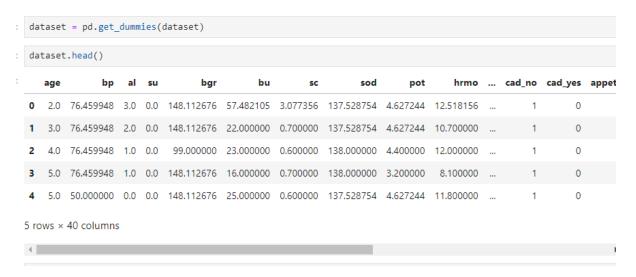
# **Machine Learning classification assignment**

- 1. From the CDK dataset we need to predict the classification column, this is a classification problem statement
- 2. Total number of rows 399

Total number of columns - 40

```
dataset.shape
(399, 40)
```

3. We are converting the string to nominal data using get\_dummies()



4. Models used

**SVM** 

```
re = grid.cv_results_
                                                                                                  □ ↑ ↓ 古 〒 🗎
 table = pd.DataFrame.from_dict(re)
 table
am_C param_gamma param_kernel
                                       params split0_test_score split1_test_score split2_test_score split3_test_score split4_t
                                        {'C': 10,
                                       'gamma':
   10
                                 rbf
                                                        0.982221
                                                                         1.000000
                                                                                           0.982051
                                                                                                            1.000000
                 auto
                                         'auto',
                                        'kernel':
                                          'rbf'}
                                        {'C': 10,
                                       'gamma':
                                                        1.000000
                                                                         1.000000
   10
                 auto
                                poly
                                         'auto',
                                                                                           0.964286
                                                                                                            1.000000
                                        'kernel':
                                         'poly'}
                                        {'C': 10,
                                       'gamma':
   10
                                                        0.982221
                                                                         1.000000
                                                                                           0.982221
                                                                                                            1.000000
                 auto
                             sigmoid
                                         'auto'.
                                        'kernel':
                                      'sigmoid'}
                                        {'C': 10,
```

The best score for SVM is 0.98

```
□ ↑ ↓ 占 〒 🗎
 from sklearn.metrics import f1_score
 f1_macro=f1_score(y_test,grid_predictions,average='weighted')
 print("The f1_macro value for best parameter {}:".format(grid.best_params_),f1_macro)
 The f1_macro value for best parameter {'C': 10, 'gamma': 'auto', 'kernel': 'sigmoid'}: 0.9834018801410106
 print("The confusion Matrix:\n",cm)
 The confusion Matrix:
  [[45 0]
  [ 2 73]]
 print("The report:\n",clf_report)
 The report:
                precision
                             recall f1-score
                                               support
                    0.96
            0
                             1.00
                                        0.98
                                                   45
            1
                    1.00
                              0.97
                                        0.99
                                                   75
                                        0.98
                                                  120
     accuracy
                    0.98
                              0.99
                                        0.98
                                                  120
    macro avg
weighted avg
                    0.98
                              0.98
                                        0.98
                                                   120
```

# **Decision tree**

F1 score is 0.942

```
回个少古早章
from sklearn.metrics import f1_score
f1_macro=f1_score(y_test,grid_predictions,average='weighted')
print("The f1_macro value for best parameter {}:".format(grid.best_params_),f1_macro)
The f1 macro value for best parameter {'criterion': 'log loss', 'max_features': 'sqrt', 'splitter': 'random'}:
0.9423437387354913
print("The confusion Matrix:\n",cm)
The confusion Matrix:
[[45 0]
[ 7 68]]
print("The report:\n",clf_report)
The report:
              precision recall f1-score support
          0
                 0.87
                         1.00
                                    0.93
                 1.00
                                              75
                          0.91
                                   0.95
          1
   accuracy
                                    0.94
                                             120
                0.93
                         0.95
                                    0.94
                                              120
  macro avg
               0.95
weighted avg
                       0.94
                                   0.94
                                             120
```

#### **KNN** classification

```
from sklearn.metrics import f1_score
f1_macro=f1_score(y_test,grid_predictions,average='weighted')
print("The f1_macro value for best parameter {}:".format(grid.best_params_),f1_macro)
The f1_macro value for best parameter {'algorithm': 'auto', 'metric': 'minkowski', 'n_neighbors': 5, 'p': 2, 'w
eights': 'distance'}: 0.9505208333333334
print("The confusion Matrix:\n",cm)
The confusion Matrix:
[[45 0]
[ 6 69]]
print("The report:\n",clf_report)
The report:
              precision recall f1-score support
          0
                  0.88
                        1.00
                                     0.94
                                                 45
                 1.00
                           0.92
                                     0.96
                                                75
          1
                                     0.95
                                               120
   accuracy
                 0.94
                           0.96
                                     0.95
  macro avg
                                                120
weighted avg
                 0.96
                           0.95
                                     0.95
                                                120
```

The F1 score is 0.950

## **Random Forest classifier**

```
from sklearn.metrics import f1_score
f1_macro=f1_score(y_test,grid_predictions,average='weighted')
print("The f1_macro value for best parameter {}:".format(grid.best_params_),f1_macro)
The f1_macro value for best parameter {'criterion': 'entropy', 'max_features': 'log2'}: 0.9916844900066377
print("The confusion Matrix:\n",cm)
The confusion Matrix:
 [[45 0]
[ 1 74]]
print("The report:\n",clf_report)
The report:
              precision
                         recall f1-score support
                  0.98
                           1.00
                                      0.99
                  1.00
                           0.99
                                     0.99
                                                 75
   accuracy
                                      0.99
                                                120
                  0.99
                            0.99
                                     0.99
  macro avg
                                                120
weighted avg
                 0.99
                            0.99
                                     0.99
                                                120
```

#### The best F1 score is 0.991

## Naïve Bayes - gaussianNB

```
from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(x_train,y_train)
y_pred = classifier.predict(x_test)
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
from sklearn.metrics import classification_report
clf_report = classification_report(y_test, y_pred)
print("The report:\n",clf_report)
print("The confusion Matrix:\n",cm)
/lib/python3.11/site-packages/sklearn/utils/validation.py:1183: DataConversionWarning: A column-vector y was
ssed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
y = column_or_1d(y, warn=True)
The report:
              precision
                          recall f1-score support
          0
                  0.94
                          1.00
                                     0.97
                                                 45
                  1.00
                           0.96
                                     0.98
                                                 75
          1
                                     0.97
                                                120
   accuracy
   macro avg
                  0.97
                           0.98
                                     0.97
                                                120
                           0.97
                                     0.98
weighted avg
                  0.98
                                                120
```

#### Bernoulli NB

```
from sklearn.model_selection import GridSearchCV
                                                                                  回↑↓古早前
from sklearn.naive_bayes import BernoulliNB
classifier = BernoulliNB()
classifier.fit(x_train,y_train)
y_pred = classifier.predict(x_test)
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
from sklearn.metrics import classification_report
clf_report = classification_report(y_test, y_pred)
print("The report:\n",clf_report)
print("The confusion Matrix:\n",cm)
/lib/python3.11/site-packages/sklearn/utils/validation.py:1183: DataConversionWarning: A column-vector y was pa
ssed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
y = column_or_1d(y, warn=True)
The report:
             precision
                         recall f1-score support
          0
                0.94 1.00 0.97
                                               45
                1.00
                          0.96
                                   0.98
                                            120
   accuracy
                                  0.97
                0.97
                          0.98
                                    0.97
  macro avg
             0.98
weighted avg
                          0.97
                                   0.98
                                              120
```

# **Logistic Regression**

```
classifier = LogisticRegression(random_state = 0)
classifier.fit(x_train,y_train)
y_pred = classifier.predict(x_test)
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
 from sklearn.metrics import classification report
clf_report = classification_report(y_test, y_pred)
print("The report:\n",clf_report)
print("The confusion Matrix:\n",cm)
/lib/py thon 3.11/s ite-packages/sklearn/utils/validation.py: 1183:\ DataConversionWarning:\ A\ column-vector\ y\ was\ packages/sklearn/utils/validation.py: 1183:\ DataConversionWarning:\ packages/sklearn/utils/validation.py: 1183:\ packag
ssed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
    y = column_or_1d(y, warn=True)
The report:
                                               precision recall f1-score support
                                                                                                                    0.99
                                   0
                                                           0.98 1.00
                                                            1.00
                                                                                           0.99
                                                                                                                           0.99
                                                                                                                                                                    75
                                                                                                                           0.99
                                                                                                                                                                120
                                                                                                                   0.99
        macro avg
                                                           0.99
                                                                                           0.99
                                                                                                                                                                 120
weighted avg
                                                            0.99
                                                                                             0.99
                                                                                                                             0.99
                                                                                                                                                                  120
The confusion Matrix:
  [[45 0]
    [ 1 74]]
```

The F1 score is 0.99

## **Best Model**

Out of these algorithms, based on the F1 score both the logistic regression and the random forest are having 0.99 score. So these 2 are considered as the best model