

NAME : VARSHINI R
ROLL NO : 18BCS073
DATE : 23/06/2021

Data Warehousing and Data Mining-End Semester Practicals

1. Download a suitable dataset for classification from any Repository. List the attributes and its type in a word Doc.

This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective of the dataset is to diagnostically predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.

The datasets consist of several medical predictor variables and one target variable, Outcome . Predictor variables include the number of pregnancies the patient has had, their BMI, insulin level, age, and so on.

Attribute types	Integer, Real
Instances	768
Attributes	8

Dataset Link -

<https://www.kaggle.com/uciml/pima-indians-diabetes-database>

2. Load the dataset and set the target and feature variables. Split the dataset into training and test dataset. Build a decision tree classifier with Entropy criteria. Perform Prediction for test dataset using Entropy and print the results in the form of confusion matrix, accuracy and classification report. visualize the decision tree.

```
Decision Tree Classifier - Jupyter
localhost:8892/notebooks/Decision%20Tree%20Classifier%20.ipynb
jupyter Decision Tree Classifier Last Checkpoint a few seconds ago (autosaved)
File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3
In [7]: y_pred_class = dtree.predict(X_test)
In [7]: # calculate accuracy
from sklearn import metrics
print(metrics.accuracy_score(y_test, y_pred_class))
0.7207792207792207
In [8]: # save confusion matrix and slice into four pieces
confusion = metrics.confusion_matrix(y_test, y_pred_class)
print(confusion)
#row, column
TP = confusion[1, 1]
TN = confusion[0, 0]
FP = confusion[0, 1]
FN = confusion[1, 0]
[[82 17]
 [26 29]]
In [9]: #Classification Accuracy:
# use float to perform true division, not integer division
print((TP + TN) / float(TP + TN + FP + FN))
print(metrics.accuracy_score(y_test, y_pred_class))
0.7207792207792207
0.7207792207792207
In [10]: #Classification Error:
classification_error = (FP + FN) / float(TP + TN + FP + FN)
print(classification_error)
print(1 - metrics.accuracy_score(y_test, y_pred_class))
0.2792207792207792
```

```
Decision Tree Classifier - Jupyter
localhost:8892/notebooks/Decision%20Tree%20Classifier%20.ipynb
jupyter Decision Tree Classifier Last Checkpoint a minute ago (autosaved)
File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3
In [18]: from six import StringIO
from sklearn.tree import export_graphviz
from IPython.display import Image
import pydotplus
dot_data = StringIO()
export_graphviz(dtree, out_file=dot_data,
                filled=True, rounded=True,
                special_characters=True, feature_names = feature_cols, class_names=['0', '1'])
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph.write_png('diabetes.png')
Image(graph.create_png())
Out[18]:
In [18]: conda install -c conda-forge pydotplus
collecting package metadata (current_repodata.json): ...working... done
solving environment: ...working... done
```

3. Upload in your github account. Provide the link for access.

<https://github.com/Varshini11122000/End-Semester-Practicals-18BCS073/upload/main>