## Assignment 3.2 - ML development using Python

#### **SIT 788**

### 1. Introduction

The dataset "Acute Inflammations" [1], has been used for this case study. This data set has been prepared by a medical expert for presumptive diagnosis of two diseases of urinary system, namely: "Inflammation of urinary bladder" and "Nephritis of renal pelvis origin". These two diseases constitute the two output class labels for each feature row in this dataset. The features of this dataset include 6 variables:

1. Temperature of patient [continuous floating-point values, ranging from 35C – 42C]

Apart from temperature, 5 categorical variables which take binary values [yes/no] are the remaining features of this dataset:

- 2. Occurrence of nausea
- 3. Lumbar pain
- 4. Urine pushing
- 5. Micturition pains
- 6. Burning of urethra, itch, swelling of urethra outlet

The dataset has a total of 119 datapoints, if we split the dataset on "Inflammation of urinary bladder" then the data has 60 negative and 59 positive class samples. Whereas, if we split the dataset on "Nephritis of renal pelvis origin", the data has 69 negative and 50 positive samples. The dataset is slightly imbalanced, also, due to two labels for each feature vector, this dataset is an example of multilabel dataset. This dataset requires supervised learning and makes the problem to be that of a multilabel classification. In the following sections, Decision Tree and Random Forest Classifier have been used to train a classifier and perform comparison across different metrics.

# 2. Creating ML Pipeline

The ML pipeline built for this consists of following steps:

1 – importing the required python packages

```
In [1]: #!pip install scikit-multilearn
import pandas as pd
import zipfile as zf
import pickle
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
import codecs
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.metrics import multilabel_confusion_matrix
from sklearn.metrics import accuracy_score
from skmultilearn.problem_transform import LabelPowerset
from time import time
```

Fig 1: imported libs

2 – Loading the data from source, the data is tab separated hence tab is used as the delimiter, also because the data is not in csv format by default therefore *codecs.open()* is used to read the data in

required pandas dataframe. The dataframe is then assigned meaningful column names denoting feature variable names.

#### Load the dataset

Fig 2: Loading dataset

3 – Exploratory data analysis and data preparation is done, as the dataset contains categorical strig values, hence they are transformed to numerical values as they are better suited for classification tasks (although in theory string categorical values can be used in case of decision trees, but sklearn does not offer support for such values].

#### Print 1st 5 rows and data tyes of column In [4]: display(df.head(5)) display(df.dtypes) **0** 35.9 no no 1 35.9 no no **2** 36.0 no no yes yes yes yes no **3** 36.0 no no no yes no no no **4** 36.0 no yes no float64 nausea object lumbar pain urine pushing micturian pains urethra pain object object object object urinary\_bladder\_inflammation renal\_pelvis\_nephritis object object dtype: object

Fig 3: Data pre-processing

# Dataframe after data pre-processing:

Split data into train and test

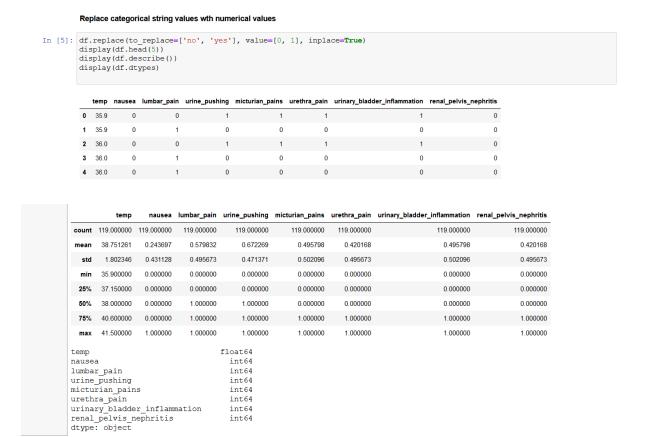


Fig 4: Data visualization

4 – After pre-processing is done, the data is then split into training and testing set. The training set contains 60% of the data whereas 40% data belongs to the test set (Lower than this fraction in testing could not differentiate the two classifiers). The 6 feature variables belong to the feature vector X and the two decision variables constitute the output label vector Y. Random state is used for reproducibility of the results

Fig 5: Data splitting into train and testing set.

5 – The first model trained is **Decision Tree**, no hyper parameter tuning is done for any model and default parameters are used as they yield high accuracy for this dataset. The model is trained on training set and prediction is done on test set. Decision tree takes around 0.23 milliseconds for a single prediction and 5 milliseconds to train.

Fig 6: Decision Tree classifier

6 – Similarly, a second classifier: **Random Forest** is trained on the same dataset. As Random Forest is an example of ensemble classifier, and has more complexity to it, it takes slightly longer (around 0.69 milliseconds) to predict a single sample and 283 milliseconds to train.

```
Training Random Forest

In [8]: clf_rf = LabelPowerset(RandomForestClassifier())

millis_a = int(time() * 1000)
    clf_rf = clf_rf.fit(x_train, y_train)
    millis_b = int(time() * 1000)
    dif = millis_b - millis_a
    print ("RF train time: ", dif)

millis_a = int(time() * 1000)
    y_pred_rf = clf_rf.predict(x_test)
    millis_b = int(time() * 1000)
    dif = millis_b - millis_a
    avg_t = dif_/len(x_test)
    print ("RF time per prediction: ", avg_t)

RF train time: 283
    RF time per prediction: 0.6875
```

Fig 7: Random Forest classifier. Note the increase in prediction and training time.

The **confusion matrix** of both the classifiers [Fig 8], denote high values across all metrics. Both the classifiers are used with Label Powerset, this technique is used to transform a multi-label problem into a multi-class problem. This involves mapping combinations of all labels to a unique id which is treated as a label, this is what transforms a multi-label from to a multi-class scenario.

7 – Comparison of both models, both the models are compared across various metrics: accuracy, precision, recall and f1-score. Confusion matrix for multilabel classification is also used.

```
print (classification report(y test, y pred_dt, zero_division=0))
print("DT Accuracy: ", accuracy_score(y_test, y_pred_dt))
print(multilabel_confusion_matrix(y_test, y_pred_dt))
      print (classification_report(y_test, y_pred_rf, zero_division=0))
print("RF Accuracy:", accuracy_score(y_test, y_pred_rf))
print(multilabel_confusion_matrix(y_test, y_pred_rf))
      precision
                             recall f1-score
                                               support
                      0.95 1.00
1.00 1.00
                                        0.97
                0
                                        1.00
                     0.98 1.00
                                        0.99
                                                   41
        micro avg
                    0.97 1.00
0.98 1.00
0.73 0.73
        macro avg
                                        0.99
      weighted avg
                                        0.99
                                                    41
       samples avg
                                        0.73
                                                   41
      DT Accuracy: 0.979166666666666
        [ 0 19]]
       [[26 0]
      precision recall f1-score support
                      1.00 1.00
1.00 1.00
                                        1.00 19
                0
                                        1.00
      micro avg 1.00 1.00
macro avg 1.00 1.00
weighted avg 1.00 1.00
samples avg 0.73 0.73
                     1.00 1.00
        micro avg
                                        1.00
                                                   41
                                       1.00
      weighted avg
                                                    41
                                        0.73
                                                   41
      RF Accuracy: 1.0
        [ 0 19]]
       [[26 0]
        [ 0 22]]]
```

Fig 8: Confusion Matrix, Accuracy, Precision, Recall and f1-score.

The comparison shows that the Random Forest is the best performing model (discussed in detail in later section).

8 – Finally, the best model (Random Forest) is exported as a zip file.

Exporting and Saving zipped model

# In [10]: model\_f = open('best\_model.pkl', 'wb') pickle.dump(clf\_rf, model\_f) model\_f.close() zf.ZipFile('best\_model.zip', mode = 'w').write('best\_model.pkl')

Fig 9: Exporting the model

### 3. Results

The two models trained were:

**Decision Tree**: Decision Trees are a supervised learning algorithms which can be used for both classification and regression problems. Here, a decision tree classifier is used. While only being trained on 60% of the data, decision trees reach an **accuracy of 97.9%**, **98% precision and 100% recall** while taking only **0.23 milliseconds to predict and 5 milliseconds to train**. This shows the strength of decision trees as they can be used to create a base pipeline very quickly for benchmarking purposes. Another advantage of decision trees is that they provide insight into what features are important during model training and provides under the hood knowledge of the final model trained.

Random Forest: Random Forest come from the family of ensemble classifiers. Being a collection of multiple decision trees, they tend to be more robust to changes and require lesser data to generalize. Here, Random Forest reach a 100% test accuracy, precision and recall while being trained on 60% data. This shows the strength of ensemble models over standalone models. Another advantage of Random Forest is that they are less prone to overfitting and hence can be used in real world scenarios. Because of the complexity and large number of trees involved, the training time is huge (283 milliseconds) as compared to Decision Tree.

# 4. Model Comparison and conclusion

On the given Dataset, the best model is Random Forest because of **higher accuracy and precision** than Decision Trees. Another key observation is, even when the size of training data is reduced, Random Forest maintain their higher accuracy and precision. The only downside of using Random Forest over Decision Tree is higher inference time, but as the difference is 0.4 milliseconds between both the prediction time, therefore the trade-off between time and higher accuracy-precision is acceptable. Also, due to the nature of the dataset being a medical dataset, accuracy-precision takes precedence over time as a misclassification is costlier than a delayed classification.

Moreover, the generalisability offers protection against overfitting scenarios which Decision Tree does not, this makes Random Forest the better performing model on this dataset.

**Note**: Both the classifiers reach **100%** accuracy, precision, recall once the training size is increased to **70%**, even in this scenario the Random Forest would be the classifier to go with because of its robustness and ability to handle unseen data better, as the training data is not large in volume and the high metrics could be a result of overfitting.

## References

[1] Dua, D. and Graff, C. (2019). UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science.