

# Store Performance Classification

ITNPBD6 Assignment 1.Student Number 3142459

## 1.Introduction

The aim is to build classification model to classify a store' performance good, or bad using data about stores in the UK.

## 2. Project Methodology

In particular, we will go through CRISP-DM as:

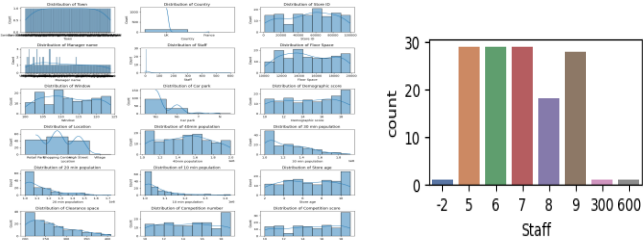
- Environment setup import required libraries
- Preparing Data
- Splitting Dataset into Training Set and Test set
- Training Model
- Evaluating Model using Confusion Matrix

## 3. Preparing Data

There are 136 examples and 19 variables.

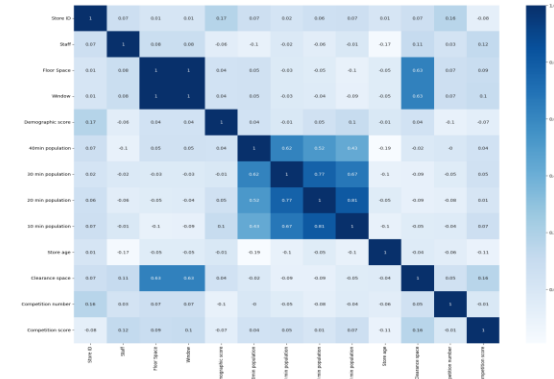
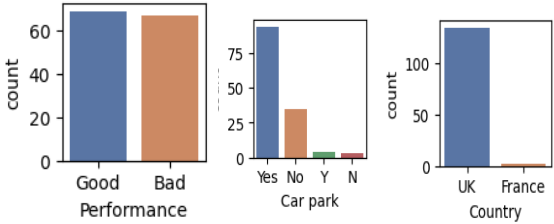
Variable	Type
Categorical	Nominal: Town, Country, Car park, Location, Manager name Ordinal: Performance
Numerical	Discrete: StoreID, Staff, Competition score, Competition number, Store age, Window, Demographic score, 40min population, 30 min population, 20 min population, 10 min population Continuous: FloorSpace, Clearance space

- There is no missing values.
- Performance as either a good or bad. It is a binary classification problem (balanced).
- Car park has inconsistency.
- Country has wrong data entry (France)
- Staff has outliers, replace them with median.



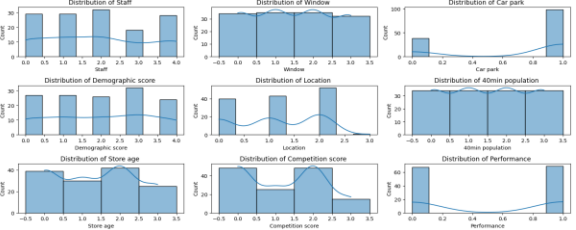
## 7. References

- Logistic Regression Classifier Tutorial | Kaggle
- Scikit-learn DecisionTreeClassifier
- tf.keras.Model | TensorFlow v2.12.0
- Hyperparameter tuning - GeeksforGeeks



In the heatmap, there is high correlation like between Floor Space and Window as in the dark color. Competition number has strong negative correlation with 40min population.

- Convert categorical variables ["Car park", "Location", "Performance"] to numerical using one-hot encoding in preprocessing module in [Sklearn](#)
- classifying Staff, Store age, Competition score, 40 min population, Demographic score and Window like in the histogram below.
- Drop the variables: Town, Country, Store ID, each have a unique value. Floor Space, Competition number, 30 min, 20 min, and 10 min population, as these are correlated. Manager name has no impact.
- Split the data into an input X of 8 features and output y 1 binary target variable (Performance)



## 4. Splitting Dataset into Training Set and Test set

The split ratio 70:30

## 5. Training Model

Training Logistic regression with respect to type of solver, etc. to find the best combination of parameters using GridSearchCV. This approach searches for the best set of hyperparameters from a grid of parameters

```
param_grid = [
    {'penalty': ['l1', 'l2', 'elasticnet', 'none'],
     'C': np.logspace(-4, 4, 20),
     'solver': ['lbfgs', 'newton-cg', 'liblinear', 'sag', 'saga'],
     'max_iter': [100, 200, 500]}
]
```

Tuned Logistic Regression Parameters: {'C': 0.23357214698981212, 'max\_iter': 30, 'penalty': 'l1', 'solver': 'saga'}  
Best score is 0.7894736842105263

**Decision Tree** consider the function to measure the quality of a split. The maximum depth of the tree. Number of features to consider, etc. RandomizedSearchCV approach is applied.

```
# Creating the hyperparameter grid
param_dist_dt = {'criterion': ['gini', 'entropy'],
                 'max_depth': [3, 5, None],
                 'max_features': randint(1, 6),
                 'min_samples_leaf': randint(1, 6)}
}
```

Tuned Decision Tree Parameters: {'criterion': 'gini', 'max\_depth': None, 'max\_features': 5, 'min\_samples\_leaf': 5}  
Best score is 0.6631205673758864

## Neural Network

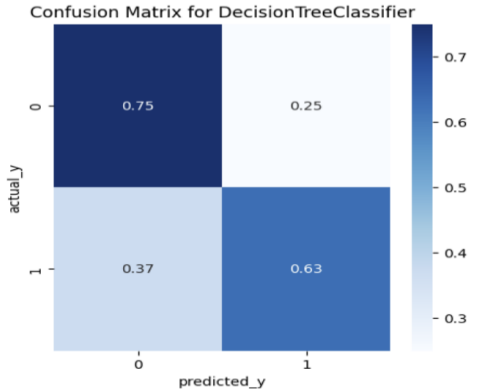
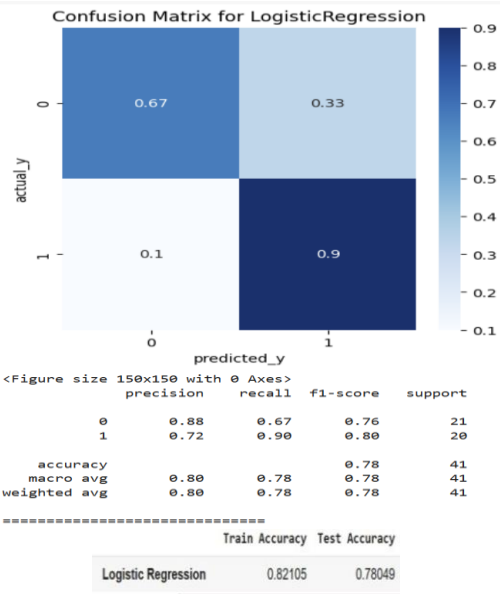
```
n_features = 8
W_model = models.Sequential(name='DeepNN', layers=[
    # Hidden layer 1
    layers.Dense(name='l1', input_dim=n_features,
                 activation='relu'),
    # Hidden layer 2
    layers.Dense(name='l2', units=units,
                 activation='relu'),
    # Layer output
    layers.Dense(name='output', units=1, activation='sigmoid')])
W_model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

params = {'batch_size': [100, 20, 50, 25, 32],
          'nb_epoch': [1, 20, 10, 30, 40],
          'unit': [5, 4, 3],
          }
```

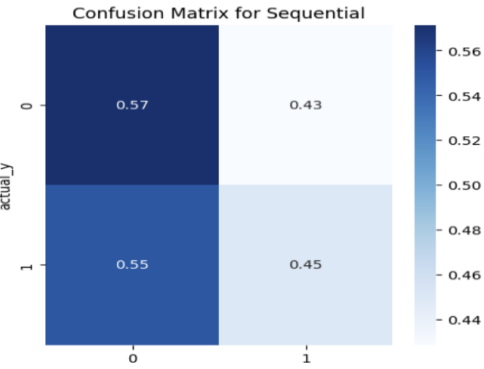
Tuned Neural Network Parameters: {'batch\_size': 25, 'nb\_epoch': 40, 'unit': 3}  
Best score is 0.5894736826419831

## 6. Evaluating Model

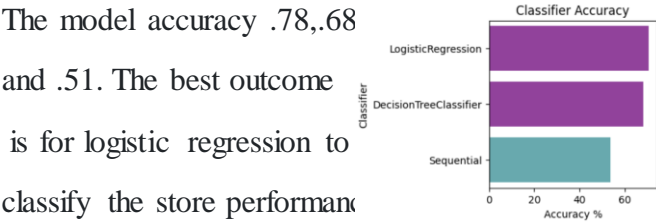
Confusion matrix to evaluate accuracy classification model.



<Figure size 150x150 with 0 Axes>



<Figure size 150x150 with 0 Axes>



The model accuracy .78, .68 and .51. The best outcome is for logistic regression to classify the store performance

with .9 correctly predicted as good (1 class) out of actual, .67 to predict the bad performance class. The F-score is .80 which is good.

