# CSCU9T6 – Data Mining Assignment 2020

Student ID: 2626628

# Contents

## 1.0 Introduction

The aim of this report is to produce a data analysis for world of bargains, a supply chain company with 100 shops in the UK. The objective is to identify why some shops are doing better than other shops and identify what factors are driving the profits. This report should also help Ivor Buquetlowd, the owner of world of bargains to build a computer program that takes input and predicts how much money each shop would make.

Using the Weka software, I have predicted the revenue by using Correlation Based Feature Selection with the 'CorrelationAttributeEval' technique. This technique requires Ranker search method. Correlation is more formally referred to as Pearson's correlation coefficient in statistics. To predict the revenue of the shops I have used 'Multilayer perceptron' technique. For classification task I have used 'decision trees – J48' technique. I also used 'cfsSubset' technique for profit prediction and this technique use 'GreedyStepwise' search method.

Data mining is used to discover the patterns in large data sets and database system. Data mining technique can be used to analyse the world bargains data and predict shop revenues. For this task, data mining helps to develop a computer program which helps to manage the choice of new shop location and provide an estimation through prediction or classification. So, data mining process is the best and well suitable for this task.
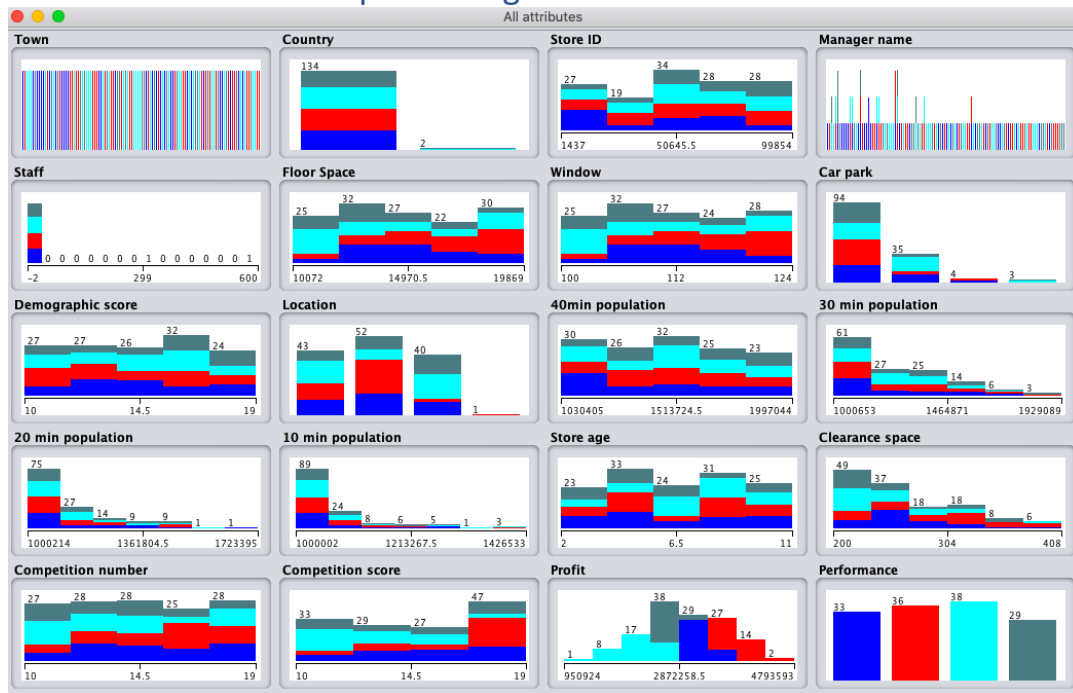
## Data mining process

## 1.1  Data Summary

 The given store data is in "comma-separated values"(csv) containing 136 observations and 20 attributes. So, in the data many attributes are self-explanatory, but some attributes are hard to understand like:

| Attributes | Data Type | Description |
|---|---|---|
| Town | Nominal | It is a place, where stores are available. |
| Country | Nominal | It's a geographic location |
| Store ID | Numeric | It is a unique number of your store account |
| Manger name | Nominal | Name of the manger for each store |
| Staff | Numeric | Number of the staff |
| Floor Space | Numeric | The space within the shop e.g. square meters or square foots of total shop area. |
| Window | Numeric | The showcasing space of the products. |
| Car park | Nominal | Weather a store has a car park or not. |
| Demographic | Numeric | It is the scale of the score e.g. 0 to 100 or 0 to 5 related to the structure of population. |
| Location | Nominal | It is a place of the shop |
| 40min population | Numeric | From the shop this is the number of populations within the radius of particular shop location. |
| 30min population | Numeric | From the shop this is the number of populations within the radius of particular shop location. |
| 20min population | Numeric | From the shop this is the number of populations within the radius of particular shop location. |
| 10min population | Numeric | From the shop this is the number of populations within the radius of particular shop location. |
| Store age | Numeric | Time since the store established |
| Clearance space | Numeric | It is a sale, which a larger number of items are discounted or reduced price e.g. square meters or square foot to place all sale items in that area. |
| Competition number | Numeric | Number of competitors |
| Competition score | Numeric | It's the scale of the competitor's e.g. o to 100 or 0 to 10 |
| Profit | Numeric | Revenue of the stores |
| Performance | Nominal | Performance of the store |

## 2.0 Data Preparation

The steps below describe the pre-processing of the data and Step 1 is done manually using Microsoft excel.

### Pre-processing Variable Values



## 2.1 Step 1 – Problems Found and Fixed

I have found 4 problems in the "store.csv" file.

### 1. Country

As per the requirements it was mentioned that there are shops in the UK but in the given data there are 2 stores in the France, this might be a data entry issue. So, I assumed these 2 stores are also belongs to UK. So, to fix this I amended France with UK.



### 2. Location

There are 4 locations in the data.
1. High street
2. Retail park
3. Shopping centre
4. Village

- There are 40 shops in High street, in the Retail park there are 43 shops, 52 shops in Shopping centre and 1 shop in village. The shop entry was confirmed geographically

to be located on a High Street in a small fishing village. So, I amended the village to High Street.

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Town | Country | Store ID | Manager | Staff | Floor Spa | Window | Car park | Demogra | Location | 40min po | 30 min po | 20 min po | 10 min po | Store age | Clearance | Competiti | Competiti | Profit | Performa |
| 37 | Southwick | UK | 26307 | Hannah | 9 | 16915 | 117 | Yes | 13 | Village | 1697206 | 1222492 | 1162941 | 1018063 | 4 | 368 | 10 | 18 | 4185306 | Excellent |

## 3. Staff

In the staff I found 3 errors, one with the value of -2 and I have corrected this with the positive value 2, one with the value of 300 and one with a value of 600 these are clear errors and I have deleted these rows from the data.

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Town | Country | Store ID | Manager | Staff | Floor Spa | Window | Car park | Demogra | Location | 40min po | 30 min po | 20 min po | 10 min po | Store age | Clearance | Competiti | Competiti | Profit | Performa |
| 4 | Skipton | UK | 2039 | Valentina | -2 | 12288 | 105 | No | 12 | Retail Park | 1595638 | 1281661 | 1104490 | 1011395 | 11 | 219 | 13 | 18 | 2297810 | Poor |
| 54 | Sherborne | UK | 44722 | Ethan | 300 | 15053 | 112 | Yes | 11 | Shopping Cer | 1067570 | 1025791 | 1001489 | 1000793 | 3 | 310 | 14 | 18 | 3758014 | Excellent |
| 109 | South Pethe | UK | 82709 | Mariana | 600 | 17744 | 119 | Yes | 14 | Shopping Cer | 1343476 | 1296509 | 1093566 | 1090240 | 2 | 312 | 16 | 18 | 3895318 | Excellent |

## 4. Car park

The given data for car park is inconsistent, sometimes it says Yes/No and sometimes it is Y/N. So, to make it consistent across all the data fields I replaced all Yes to Y and all No to N.
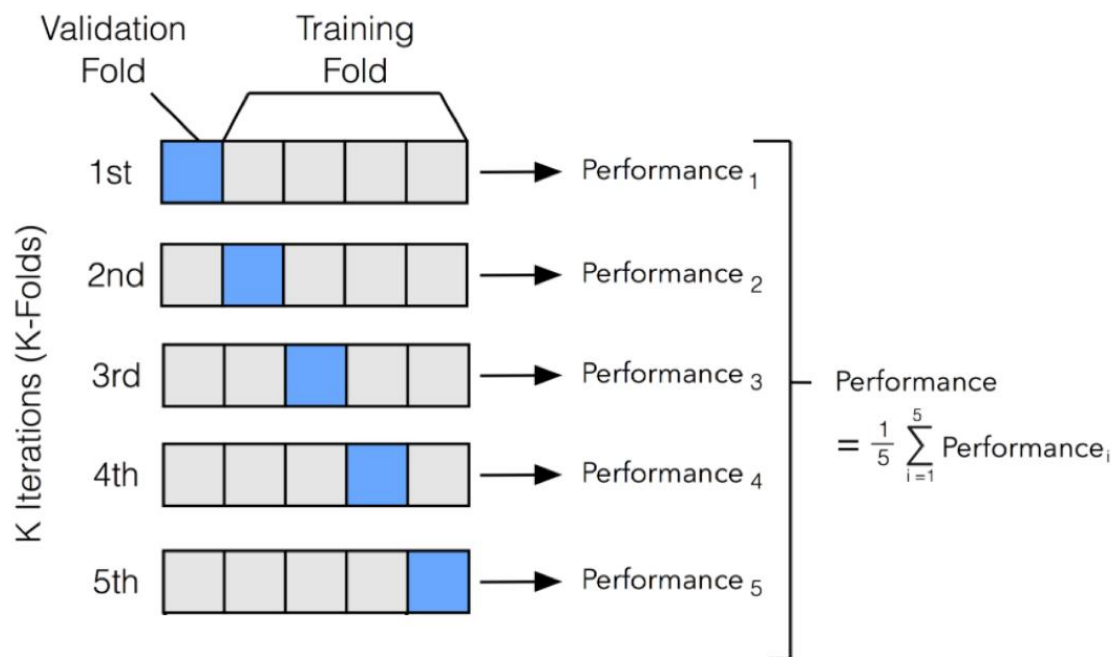
## 2.2 Step 2 – Splitting the data into 2 parts Training and Test

Using Weka Cross-validation I have divided the data into 2 sets: training set and test set. This is to make sure the data built not only to execute well on training data but also on unseen data. Basically, training data is an initial set of data that is used to understand how a program apply technologies like neutral networks(NN) and produce some advanced results.

"K-Fold Cross Validation" is one of the most common techniques we use for model evaluation and model selection. For larger training data, k is commonly chosen as 10 and for smaller training data, k is chosen as 5. As you can see the figure below which illustrates the process of 5-fold cross validation.

### K-fold cross-validation

# 3.0  Model Training

## Variable I chose to Investigate

I have chosen these variables using Weka software/Select attributes. The Select attributes are divided in two parts:

1. Attribute Evaluator
2. Search Method

So, I have chosen a technique called Correlation Based Feature Selection. This technique gives the ranked attributes according to profit.

| | |
|---|---|
| 1. Performance<br><br>2. Competition score<br><br>3. Window<br><br>4. Clearance space<br><br>5. Floor space<br><br>6. Car park<br><br>7. Location<br><br>8. Competition number<br><br>9. Staff |  |

# 4.0 Technique Description

## Multi-layer perceptron(MLP)

Multi-layer perceptron (MLP) can be used for prediction and classification tasks, also MLP is called as neutral networks(NN). An MLP consist of 3 layers of nodes:

1. input layer
2. hidden layer
3. output layer

Each node is neutron which uses a non-linear function except the input node. Each node connects with a certain weight to every node. As you can see figure below, it's an example of MLP. There are 4 input layer, 5 hidden layers and 1 output layer. So, Input node receives a value for a given instance in the data and passes the value through the node as output.

## Multilayer Perceptron Network

Decision trees

A decision tree is a structure that includes a root node, branches and leaf nodes, where a node represents a single variable. Every Internal node denote a test on an attribute, branch node denotes the outcome of a test and leaf nodes represent the classified/predicted variable and the top most node is root node. As you can see figure below, is the structure of decision trees.

### Structure of decision tree



Many decision trees use ID3 algorithm which use entropy to calculate the homogeneity of a sample.

## 5.0 Predicted results(Profit)

In the figure below, the variable values distribution post-processing for the profit prediction.

### Chosen Variable values for Prediction



Matrix X represents all the attributes and our targeted variable Y which is Profit.

Example:

X = [Staff, Floor Space, Window, Car Park, Demographic Score, Location, 40min population, 30min population, 20min population, 10min population, Store age, Clearance Space, Competition Score ]

y = [ Profit ]

For prediction task the training data was loaded via "Select Attributes" tab In the Attribute Evaluator of Weka software. I have chosen 'cfcSubsetEval' technique identifies the most effective attributes by considering the individual predictive and with the degree of redundancy between them. If we select 'cfcSubsetEval', Weka will give a pop-up alert window saying in the search method should use 'Greedy Stepwise' method. So, Weka detects automatically which is relevant to the attribute. I have used 'Greedy Stepwise' as full training set containing all attributes in it, results can be seen below:

## Prediction results

```
=== Run information ===

Evaluator:    weka.attributeSelection.CfsSubsetEval -P 1 -E 1
Search:       weka.attributeSelection.GreedyStepwise -T -1.7976931348623157E308 -N -1
-num-slots 1
Relation:     fina
Instances:    134
Attributes:   10
              Staff
              Floor Space
              Window
              Car park
              Location
              Clearance space
              Competition number
              Competition score
              Profit
              Performance
Evaluation mode:    evaluate on all training data




=== Attribute Selection on all input data ===

Search Method:
        Greedy Stepwise (forwards).
        Start set: no attributes
        Merit of best subset found:    0.627

Attribute Subset Evaluator (supervised, Class (numeric): 9 Profit):
        CFS Subset Evaluator
        Including locally predictive attributes

Selected attributes: 1,3,4,5,7,8,10 : 7
                      Staff
                      Window
                      Car park
                      Location
                      Competition number
                      Competition score
                      Performance
```

## 5.1 Multilayer Perceptron(MLP)

I have loaded the data in training set in the Weka classify. I have chosen Multilayer Perceptron technique to predict the profit and these are the results.

You can see in the 2 new columns in the figure saying 'Predicted' and 'Error'

In the predicted column it shows the predicted profit.

In the Error column it shows the difference between the given profit and what are actually predicted.



The hyperparameters evaluated the MLP are 'Number of hidden layers ', 'Learning rate' and 'Momentum'

**Number of hidden layers:** This is a process; how many hidden layers are used to find relationships in data. A single hidden layer is more a simple relationship in the data, it's better to have multiple hidden layers to have more complex relationships in the data.

**Learning rate:** This value impact how quickly the network learns, so a low learning rate will result slow learning of weights. A high learning rate makes the weight and results in the progressive learning.

**Momentum:** This influence, the model which can escape local minima during backpropagation.

Below are the results MLP model.

**Multilayer Perceptron Model Results**

| Model Number | Hidden Layers | Learining Rate | Momemtum | Correlation Coefficient | RMS Error(£) |
|---|---|---|---|---|---|
| 1 | 1 | 0.1 | 0.1 | 0.7377 | 485552.837 |
| 2 | 1 | 0.2 | 0.1 | 0.7348 | 492137.988 |
| 3 | 1 | 0.1 | 0.2 | 0.7365 | 487259.022 |
| 4 | 1 | 0.2 | 0.2 | 0.7315 | 497500.431 |
| 5 | 2 | 0.1 | 0.1 | 0.7343 | 492494.27 |
| 6 | 2 | 0.2 | 0.1 | 0.6781 | 552056.947 |
| 7 | 2 | 0.1 | 0.2 | 0.7326 | 494866.834 |
| 8 | 2 | 0.2 | 0.2 | 0.6976 | 542345.231 |

# Prediction model

```
=== Run information ===

Scheme:        weka.classifiers.functions.MultilayerPerceptron -L 0.1 -M 0.2 -N 500 -V 0 -S 0 -E 20 -H 1
Relation:      fina-weka.filters.unsupervised.attribute.Remove-R10
Instances:     134
Attributes:    9
               Staff
               Floor Space
               Window
               Car park
               Location
               Clearance space
               Competition number
               Competition score
               Profit
Test mode:     5-fold cross-validation

=== Classifier model (full training set) ===

Linear Node 0
    Inputs    Weights
    Threshold    -0.4994089530104155
    Node 1    1.1654338314242292
Sigmoid Node 1
    Inputs    Weights
    Threshold    -0.4267816585248884
    Attrib Staff    1.3316692079955184
    Attrib Floor Space    0.2986526735471699
    Attrib Window    0.37531010288023464
    Attrib Car park=N    -0.47348576474389586
    Attrib Location=Retail Park    0.015167051248528641
    Attrib Location=Shopping Centre    0.4330957826014408
    Attrib Location=High Street    -0.13775304445303158
    Attrib Clearance space    0.044641616717041935
    Attrib Competition number    0.4329481497889247
    Attrib Competition score    1.036987297484464
Class
    Input
    Node 0


Time taken to build model: 0.2 seconds

=== Cross-validation ===
=== Summary ===

Correlation coefficient             0.7365
Mean absolute error             376174.6116
Root mean squared error         487259.0219
Relative absolute error            63.7883 %
Root relative squared error        68.1245 %
Total Number of Instances          134
```

# 6.0 Performance Classification

In the figure below, the variable values distribution for the performance classification data.

Example:

X = [ Staff, Floor Space, Window, Car Park, Demographic Score, Location, 40min population, 30min population, 20min population, 10min population, Store age, Clearance Space, Competition Score ]

y = [ Performance ] ∈ {'Poor', 'Reasonable', 'Good', 'Excellent'}

## Variable values for classification

In the prediction task above, I used 'cfcSubsetEval' and 'Greedy Stepwise' method. To perform a selection for classification task I used 'CorrealtionAttributeEval' and 'Ranker' method, below are the results of the model performance using all the attributes available to the model.

## Selection algorithm results of classification

```
=== Run information ===

Evaluator:      weka.attributeSelection.CorrelationAttributeEval
Search:         weka.attributeSelection.Ranker -T -1.7976931348623157E308 -N -1
Relation:       fina-weka.filters.unsupervised.attribute.Remove-R9
Instances:      134
Attributes:     9
                Staff
                Floor Space
                Window
                Car park
                Location
                Clearance space
                Competition number
                Competition score
                Performance
Evaluation mode:    evaluate on all training data




=== Attribute Selection on all input data ===

Search Method:
        Attribute ranking.

Attribute Evaluator (supervised, Class (nominal): 9 Performance):
        Correlation Ranking Filter
Ranked attributes:
 0.221    1 Staff
 0.208    8 Competition score
 0.174    4 Car park
 0.156    3 Window
 0.151    2 Floor Space
 0.144    6 Clearance space
 0.133    5 Location
 0.12     7 Competition number

Selected attributes: 1,8,4,3,2,6,5,7 : 8
```

## 6.1 Decision Trees

Below table are the results from the decision tree algorithm. I used Weka's J48 decision trees, it is used to build a classification model.

There are 2 main hyperparameters which helped to get the results:

**1. confidenceFactor** – a small confidence factors will force a high degree of pruning. It helps to reduce the size of the decision tree y removing parts of the tree.

**2. minNumObj** – The minimum number of instances per leaf.

Decision Tree Model Results

| Model Number | confidenceFactor | minNumObj | Classificatio n Accuracy (%) | Total Misclassified Errors (%) |
|---|---|---|---|---|
| 1 | 0.1 | 1 | 40.2985 | 59.7015 |
| 2 | 0.2 | 1 | 40.2985 | 59.7015 |
| 3 | 0.1 | 2 | 41.0448 | 58.9552 |
| 4 | 0.2 | 2 | 43.2836 | 56.7164 |
| 5 | 0.1 | 3 | 37.3134 | 62.6866 |
| 6 | 0.2 | 3 | 38.0597 | 61.9403 |
| 7 | 0.1 | 4 | 32.8358 | 67.1642 |
| 8 | 0.2 | 4 | 33.5821 | 66.4179 |

# Decision tree visualise

# Classification model

```
=== Run information ===

Scheme:       weka.classifiers.trees.J48 -C 0.1 -M 2
Relation:     fina-weka.filters.unsupervised.attribute.Remove-R9
Instances:    134
Attributes:   9
              Staff
              Floor Space
              Window
              Car park
              Location
              Clearance space
              Competition number
              Competition score
              Performance
Test mode:    5-fold cross-validation

=== Classifier model (full training set) ===

J48 pruned tree
------------------

Competition score <= 16
|   Window <= 121
|   |   Car park = Y
|   |   |   Clearance space <= 237
|   |   |   |   Staff <= 7
|   |   |   |   |   Staff <= 6
|   |   |   |   |   |   Floor Space <= 10944: Poor (2.0)
|   |   |   |   |   |   Floor Space > 10944: Reasonable (10.0/1.0)
|   |   |   |   |   Staff > 6: Poor (4.0/1.0)
|   |   |   |   Staff > 7
|   |   |   |   |   Competition number <= 14: Reasonable (3.0/1.0)
|   |   |   |   |   Competition number > 14: Excellent (4.0/1.0)
|   |   |   Clearance space > 237
|   |   |   |   Location = Retail Park: Good (10.0/4.0)
|   |   |   |   Location = Shopping Centre
|   |   |   |   |   Staff <= 8: Good (10.0/4.0)
|   |   |   |   |   Staff > 8: Excellent (2.0)
|   |   |   |   Location = High Street
|   |   |   |   |   Staff <= 7: Poor (5.0)
|   |   |   |   |   Staff > 7: Good (5.0/1.0)
|   |   Car park = N: Poor (23.0/8.0)
|   Window > 121
|   |   Car park = Y: Excellent (8.0/2.0)
|   |   Car park = N: Poor (3.0)
Competition score > 16
|   Staff <= 6
|   |   Competition number <= 15: Reasonable (13.0/6.0)
|   |   Competition number > 15: Excellent (8.0/2.0)
|   Staff > 6
|   |   Competition score <= 18
|   |   |   Car park = Y: Excellent (10.0/3.0)
|   |   |   Car park = N: Good (5.0/1.0)
|   |   Competition score > 18: Excellent (9.0/1.0)

Number of Leaves  :     18

Size of the tree :      34


Time taken to build model: 0.04 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances          55               41.0448 %
Incorrectly Classified Instances        79               58.9552 %
Kappa statistic                          0.2081
Mean absolute error                      0.3196
Root mean squared error                  0.4834
Relative absolute error                 85.4414 %
Root relative squared error            111.772  %
Total Number of Instances              134

=== Detailed Accuracy By Class ===

               TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC
Area  Class
               0.303    0.208    0.323      0.303   0.313      0.097  0.492     0.268
Good
               0.618    0.190    0.525      0.618   0.568      0.407  0.684     0.393
Excellent
               0.500    0.240    0.452      0.500   0.475      0.253  0.663     0.446
Poor
               0.172    0.152    0.238      0.172   0.200      0.023  0.548     0.247
Reasonable
Weighted Avg.  0.410    0.200    0.392      0.410   0.399      0.204  0.602     0.346

=== Confusion Matrix ===

  a  b  c  d   <-- classified as
 10  7 11  5 |  a = Good
  8 21  2  3 |  b = Excellent
  6  5 19  8 |  c = Poor
  7  7 10  5 |  d = Reasonable
```

# Recommendation

Based on the results produced by the Multilayer perceptron(MLP) and Decision tree(J48). It is strongly recommended for Ivor Buquetlowd that there is no necessity to collect the population data.

From the results produced by the MLP predictive model it is highly recommended that these variables Staff, Window, Car park, Location, Competition number, Competition score and Performance are sufficient to achieve Ivor Buquetlowd objectives.

Predictive modelling is more of an approach than a process. It is typically a machine learning algorithm. These algorithms perform the data mining statistical analysis, determining trends and patterns in the data. Weka software have built in algorithms that can be used to make predictive models. One of the predictive model that had been used in this report is **Multilayer Perceptron(MLP).**

As per given scenario, I strongly recommend MLP for prediction task. It fulfils the Ivor Buquetlowd's requirements such as:

1. Why some shops are doing better than others?
2. How much money that shop should make?

In the predictive model the data is classified in a simple but powerful form of multiple variable analysis called **Decision trees**. These are produced by algorithms that are identified various ways of splitting data into branch like segments. Decision trees partition data into subset based on the categories of input variables like:

1. Performance
2. Competition score
3. Window
4. Clearance space
5. Floor space
6. Car park
7. Location
8. Competition number
9. Staff

It helps Ivor Buquetlowd to make decisions for improving better profits and forecasting the revenue of each shop in the supply chain.

I conclude that predictive analytics solutions help organisations to turn their data into timely insides for better and faster decision making.