# DATA ANALYTICS

## IS 665

## "DATA MINING PROJECT"

### On

## "Classification of Iris"

- **_Rohan Uday Khambekar_**

## Description:

This is perhaps the best known database to be found in the pattern recognition literature. Fisher's paper is a classic in the field and is referenced frequently to this day. (See Duda & Hart, for example.) The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other 2; the latter are NOT linearly separable from each other.

## Question:

We are trying to predict the class of the Iris plant.

## Data Source:

http://mercury.webster.edu/aleshunas/Data%20Sets/Supplemental%20Excel%20Data%20Sets.htm

## Data Dictionary:

| Attribute | Data Type | Range of values |
|---|---|---|
| Sepal Length in cm (SL) | Numeric | 4.3 – 7.9 |
| Sepal Width in cm (SW) | Numeric | 2.0 – 4.4 |
| Petal Length in cm (PL) | Numeric | 1.0 – 6.9 |
| Petal Width in cm (PW) | Numeric | 0.1 – 2.5 |
| Classification | Varchar | Iris Setosa, Iris Virginica, Iris Versicolor |

## Sample Data:

FILE   HOME   INSERT   PAGE LAYOUT   FORMULAS   DATA   REVIEW   VIEW

A1      fx   SL

|    | A | B | C | D | E |
|----|-----|-----|-----|-----|----------------|
| 1  | SL | SW | PL | PW | Classification |
| 2  | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 3  | 4.9 | 3 | 1.4 | 0.2 | Iris-setosa |
| 4  | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 5  | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 6  | 5 | 3.6 | 1.4 | 0.2 | Iris-setosa |
| 7  | 5.4 | 3.9 | 1.7 | 0.4 | Iris-setosa |
| 8  | 4.6 | 3.4 | 1.4 | 0.3 | Iris-setosa |
| 9  | 5 | 3.4 | 1.5 | 0.2 | Iris-setosa |
| 10 | 4.4 | 2.9 | 1.4 | 0.2 | Iris-setosa |
| 11 | 4.9 | 3.1 | 1.5 | 0.1 | Iris-setosa |
| 12 | 5.4 | 3.7 | 1.5 | 0.2 | Iris-setosa |
| 13 | 4.8 | 3.4 | 1.6 | 0.2 | Iris-setosa |
| 14 | 4.8 | 3 | 1.4 | 0.1 | Iris-setosa |
| 15 | 4.3 | 3 | 1.1 | 0.1 | Iris-setosa |
| 16 | 5.8 | 4 | 1.2 | 0.2 | Iris-setosa |
| 17 | 5.7 | 4.4 | 1.5 | 0.4 | Iris-setosa |
| 18 | 5.4 | 3.9 | 1.3 | 0.4 | Iris-setosa |
| 19 | 5.1 | 3.5 | 1.4 | 0.3 | Iris-setosa |
| 20 | 5.7 | 3.8 | 1.7 | 0.3 | Iris-setosa |
| 21 | 5.1 | 3.8 | 1.5 | 0.3 | Iris-setosa |
| 22 | 5.4 | 3.4 | 1.7 | 0.2 | Iris-setosa |
| 23 | 5.1 | 3.7 | 1.5 | 0.4 | Iris-setosa |
| 24 | 4.6 | 3.6 | 1 | 0.2 | Iris-setosa |
| 25 | 5.1 | 3.3 | 1.7 | 0.5 | Iris-setosa |
| 26 | 4.8 | 3.4 | 1.9 | 0.2 | Iris-setosa |
| 27 | 5 | 3 | 1.6 | 0.2 | Iris-setosa |
| 28 | 5 | 3.4 | 1.6 | 0.4 | Iris-setosa |
| 29 | 5.2 | 3.5 | 1.5 | 0.2 | Iris-setosa |
| 30 | 5.2 | 3.4 | 1.4 | 0.2 | Iris-setosa |

# Example Set: (on adding the data)

- ## Data Screen



- ## Statistics Screen

# 1. <u>Naïve Bayes:</u>

- **Conditional Probability**

$P(A,B) = P(A|B)P(B) = P(B|A)P(A)$

**Bayes rule**

$P(B|A)P(A) = P(A|B)P(B)$

- **Naïve Bayes Classifier**

Applying Bayes rule P(Yi)

$P(X1...Xn|Yi)\ P(Yi|X1...Xn)$

$= P(X1...Xn)$

$P(Yi)\ P(X1|Yi)P(X2|Yi)...P(Xn|Yi) = P(X1...Xn)$

$Y \leftarrow \arg\max P(Yi)\ P(X1|Yi)P(X2|Yi)...P(Xn|Yi)$

$P(Y|\mathbf{X})$: posterior probability for Y P(Y): prior probability

$P(\mathbf{X}|Y)$: class-conditional probability

$P(\mathbf{X})$: evidence

Bayes theorem (Bayes rule) allows us to calculate the posterior probability $P(Y|\mathbf{X})$ using the prior probability P(Y), the class-conditional probability $P(\mathbf{X}|Y)$ and the evidence $P(\mathbf{X})$

(Which is constant and ignored).

## Process Screen:



## Description:

## Distribution Table:

| Attribute | Parameter | Iris-setosa | Iris-versicolor | Iris-virginica |
|-----------|-----------|-------------|-----------------|----------------|
| SL | mean | 5.006 | 5.936 | 6.588 |
| SL | standard deviation | 0.352 | 0.516 | 0.636 |
| SW | mean | 3.418 | 2.770 | 2.974 |
| SW | standard deviation | 0.381 | 0.314 | 0.322 |
| PL | mean | 1.464 | 4.260 | 5.552 |
| PL | standard deviation | 0.174 | 0.470 | 0.552 |
| PW | mean | 0.244 | 1.326 | 2.026 |
| PW | standard deviation | 0.107 | 0.198 | 0.275 |

Result History  ×   SimpleDistribution (Naive Bayes)  ×   ExampleSet (//Local Repository/DA/data/iris-data)  ×

Description  Charts  Distribution Table  Annotations

## Validation:

## Performance Result:



| | true Iris-setosa | true Iris-versicolor | true Iris-virginica | class precision |
|---|---|---|---|---|
| pred. Iris-setosa | 50 | 0 | 0 | 100.00% |
| pred. Iris-versicolor | 0 | 47 | 4 | 92.16% |
| pred. Iris-virginica | 0 | 3 | 46 | 93.88% |
| class recall | 100.00% | 94.00% | 92.00% | |

accuracy: 95.33% +/- 4.27% (mikro: 95.33%)

## Performance Vector:



### PerformanceVector

```
PerformanceVector:
accuracy: 95.33% +/- 4.27% (mikro: 95.33%)
ConfusionMatrix:
True:    Iris-setosa    Iris-versicolor Iris-virginica
Iris-setosa:    50        0        0
Iris-versicolor:         0        47        4
Iris-virginica: 0         3        46
kappa: 0.930 +/- 0.064 (mikro: 0.930)
ConfusionMatrix:
True:    Iris-setosa    Iris-versicolor Iris-virginica
Iris-setosa:    50        0        0
Iris-versicolor:         0        47        4
Iris-virginica: 0         3        46
```
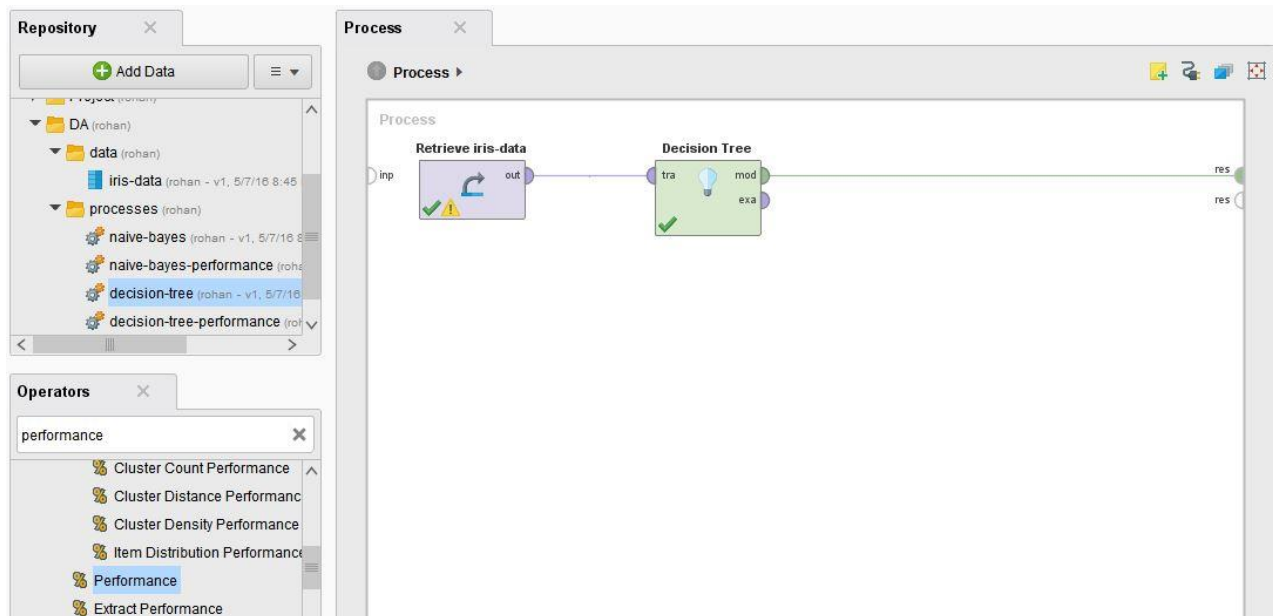
# 2. *Decision Tree*
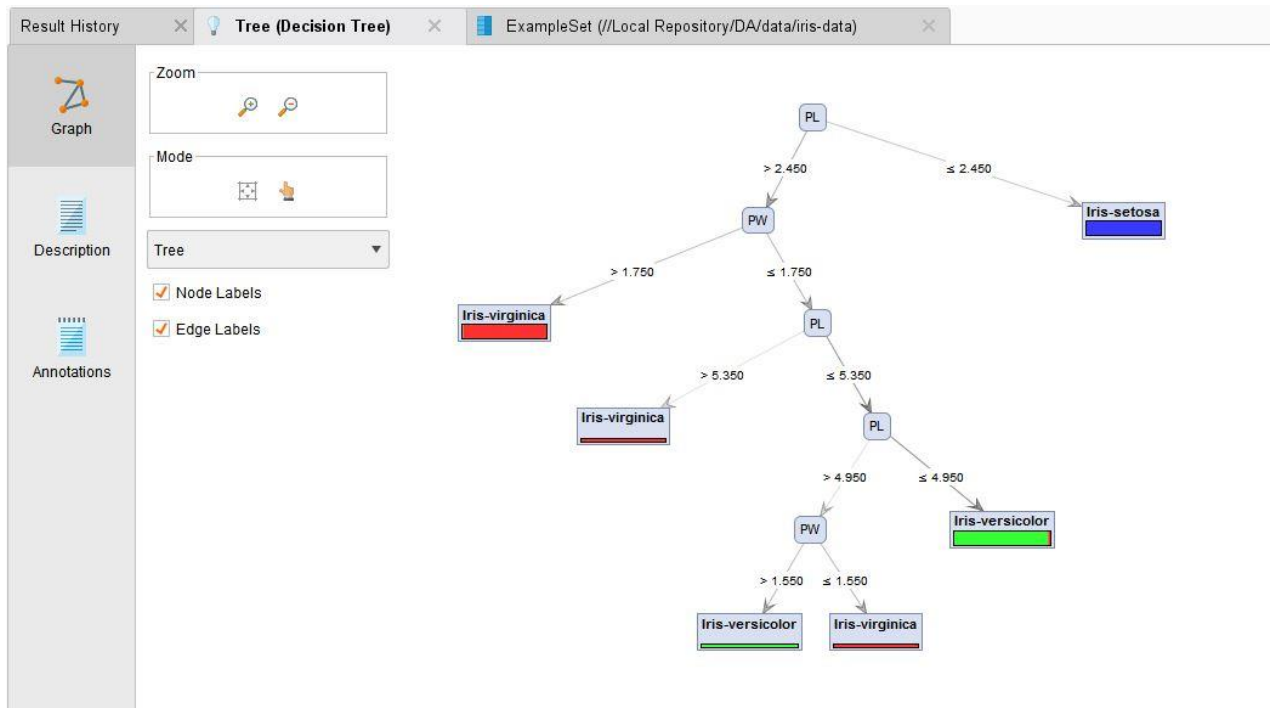
- Decision tree learning is a method commonly used in data mining.[1] The goal is to create a model that predicts the value of a target variable based on several input variables. An example is shown below. Each interior node corresponds to one of the input variables; there are edges to children for each of the possible values of that input variable. Each leaf represents a value of the target variable given the values of the input variables represented by the path from the root to the leaf.

- A decision tree is a simple representation for classifying examples. For this section, assume that all of the features have finite discrete domains, and there is a single target feature called the classification. Each element of the domain of the classification is called a class. A decision tree or a classification tree is a tree in which each internal (non-leaf) node is labeled with an input feature. The arcs coming from a node labeled with a feature are labeled with each of the possible values of the feature. Each leaf of the tree is labeled with a class or a probability distribution over the classes.

- A tree can be "learned" by splitting the source set into subsets based on an attribute value test. This process is repeated on each derived subset in a recursive manner called recursive partitioning. The recursion is completed when the subset at a node has all the same value of the target variable, or when splitting no longer adds value to the predictions. This process of *top-down induction of decision trees* (TDIDT) is an example of a greedy algorithm, and it is by far the most common strategy for learning decision trees from data.

- In data mining, decision trees can be described also as the combination of mathematical and computational techniques to aid the description, categorization and generalization of a given set of data.

- Data comes in records of the form:

$$(\mathbf{x}, Y) = (x_1, x_2, x_3, ..., x_k, Y)$$

- The dependent variable, Y, is the target variable that we are trying to understand, classify or generalize. The vector **x** is composed of the input variables, $x_1$, $x_2$, $x_3$ etc., that are used for that task.

*Source: Wikipedia*
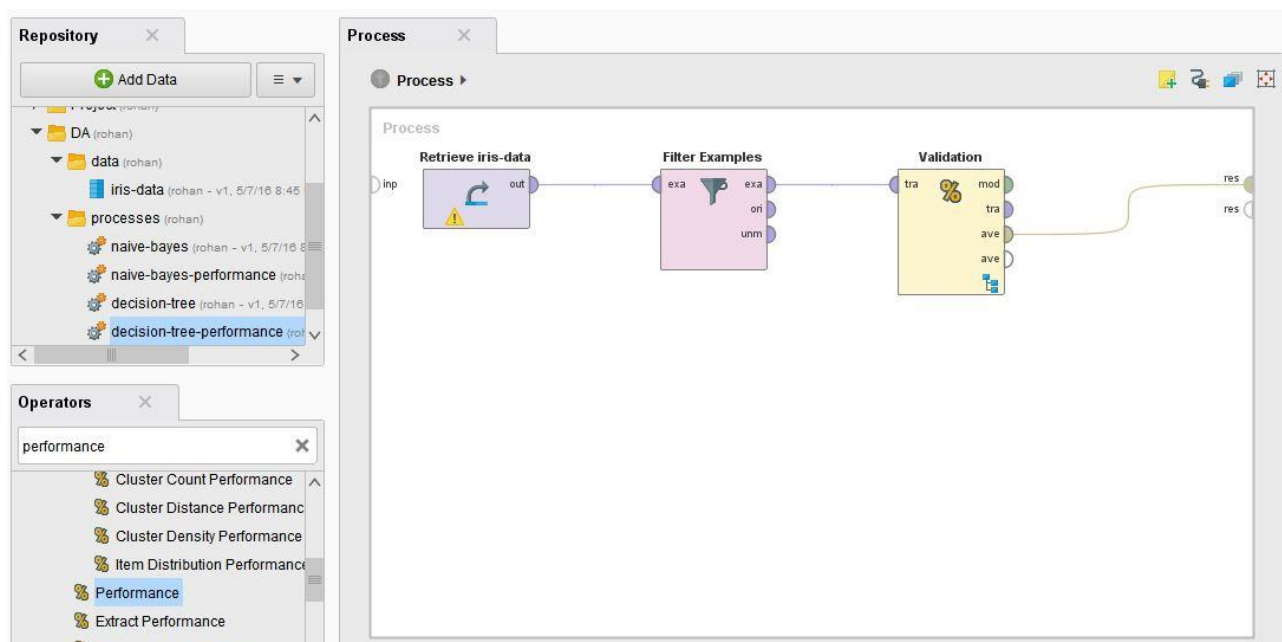
## Process Screen:



## Decision Graph:

## Tree:



```
Result History    ×   💡  Tree (Decision Tree)   ×   ▤ ExampleSet (//Local Repository/DA/data/iris-data)   ×

📈
Graph

        Tree

        PL > 2.450
        |   PW > 1.750: Iris-virginica {Iris-setosa=0, Iris-versicolor=1, Iris-virginica=45}
        |   PW ≤ 1.750
📄      |   |   PL > 5.350: Iris-virginica {Iris-setosa=0, Iris-versicolor=0, Iris-virginica=2}
Description  |   |   PL ≤ 5.350
        |   |   |   PL > 4.950
        |   |   |   |   PW > 1.550: Iris-versicolor {Iris-setosa=0, Iris-versicolor=2, Iris-virginica=0}
        |   |   |   |   PW ≤ 1.550: Iris-virginica {Iris-setosa=0, Iris-versicolor=0, Iris-virginica=2}
📄      |   |   |   PL ≤ 4.950: Iris-versicolor {Iris-setosa=0, Iris-versicolor=47, Iris-virginica=1}
Annotations  PL ≤ 2.450: Iris-setosa {Iris-setosa=50, Iris-versicolor=0, Iris-virginica=0}
```

## Validation:

## Performance:



## Performance Vector:

# Performance Comparison:

## Naïve Bayes vs Decision Trees

| Attribute | Classification | Naïve Bayes | Decision Trees |
|-----------|----------------|-------------|----------------|
| **Precision** | Iris Setosa | 100% | 100% |
| | Iris Versicolor | 92.16% | 90.20% |
| | Iris Virginica | 93.88% | 91.84% |
| | | | |
| **Recall** | Iris Setosa | 100% | 100% |
| | Iris Versicolor | 94% | 92% |
| | Iris Virginica | 92% | 90% |

We can clearly see from the above table that Naïve Bayes slightly edges out Decision Trees in terms of performance if we are to compare their Precision and Recall values.

Thus for the Iris Data Set, Naïve Bayes is better than Decision Trees.

END OF REPORT