**ELMA EDISON**

**A00279906**

**Programming for Data Analytics – (AL_KDATA_9_1)**

**12-11-2020**

**Research on Data Analytics Applications**

## 1. ORANGE

Orange is an open source data mining tool that allows data visualization, data classification and clustering through visual programming. Extended functionalities can be used through python scripting. It has a very attractive and interactive graphical user interface.

It was easy downloading the software and I did not face any issues while installing it. I downloaded the latest version of Orange 3.27.1 from the link given below: https://orange.biolab.si/download/#windows

After downloading I watched few tutorial videos and got to know about the basic functionalities in Orange and how widgets are used in it. Given below is the screenshot of the interactive user interface in Orange (Figure 1.1).
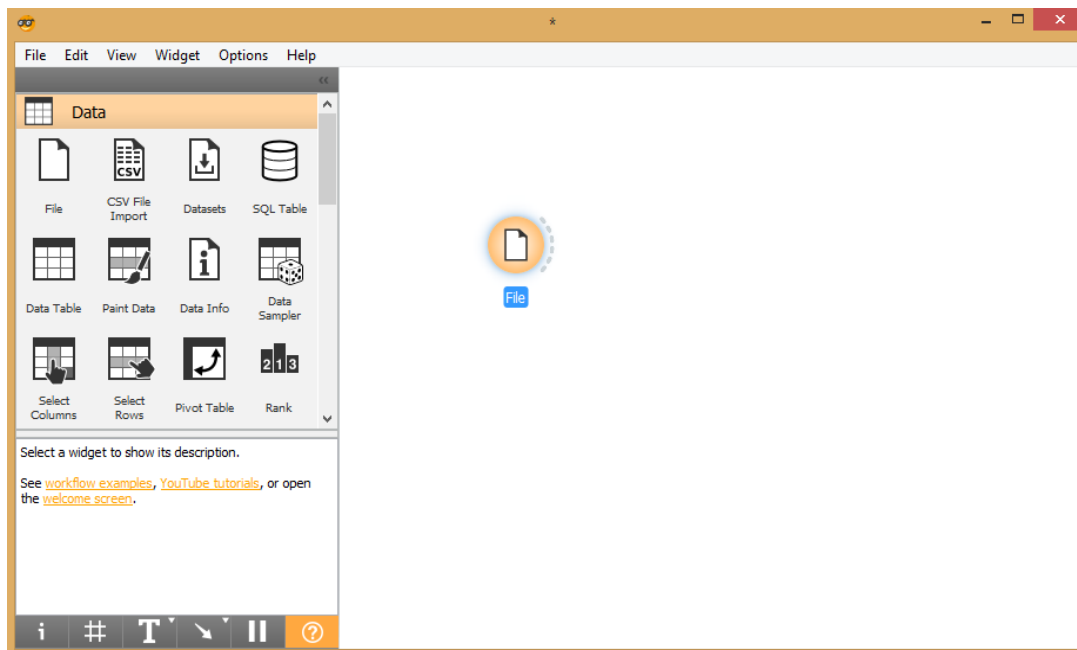


Figure 1.1

From my experience with using the tool, I think Orange is a great user-friendly data mining tool for beginners. Orange widgets are used for creating schemas for different processes such as clustering, visualization, classification etc. When you select a widget, the description about that particular widget is given in the bottom white box of the Tool Dock which is really helpful and informative. Schemas start with the File widget where in data is read and can be given to other widgets. We can also import a data table from a CSV formatted file using the CSV File Import widget and outputs the dataset. In the small example screenshot given below (Figure 1.2), the file widget is used to read the data which is then sent to both Data Table widget and Scatter Plot widget. We can select widgets from the tool dock or drag and drop or by just right clicking. It is easy to connect widgets by just dragging lines.
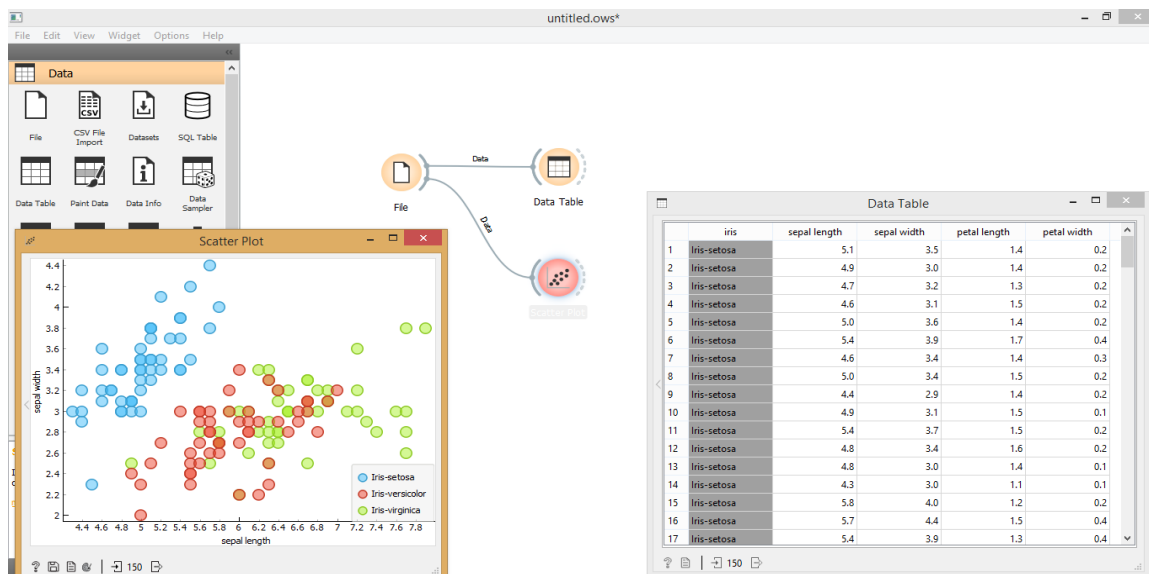
Figure 1.2

I felt that Data Sampler widget is a very useful one which can be used for various data sampling methods such as Fixed proportion of data, Fixed sample size, Cross Validation and Bootstrap. It has two outputs with sampled dataset and a dataset for remaining data (which are instances remaining in the input excluding the sampled dataset). Data Sampler can also be used for Over/Undersampling by increasing the number of instances of the sample with replacement and then by using Concatenate widget, matched and unmatched data can be appended and visualized with Distributions widget. Given below is the screenshot of an example using Data Sampler (Figure 1.3).
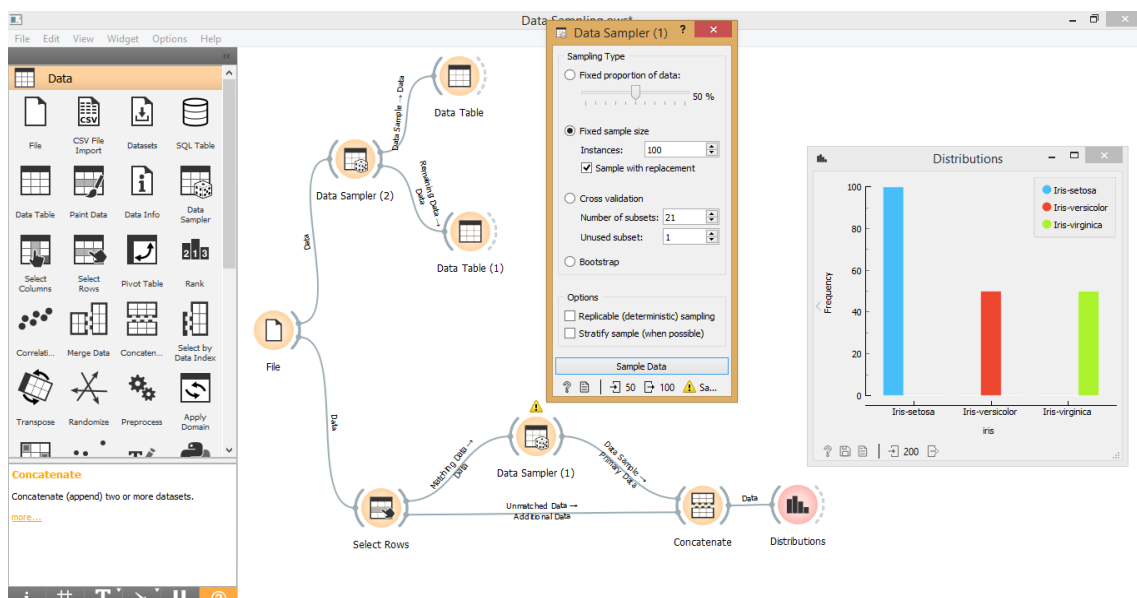


Figure 1.3

Data can be plotted with different colour, shape, size and labels for interactive data visualization. We can also change the way the plot is displayed according to the x

and y axes. Another useful point is that our data can be saved in whichever format (such as csv, xlsx, tab etc) we want using Save Data widget. Given below is the screenshot of an example (Figure 1.4).
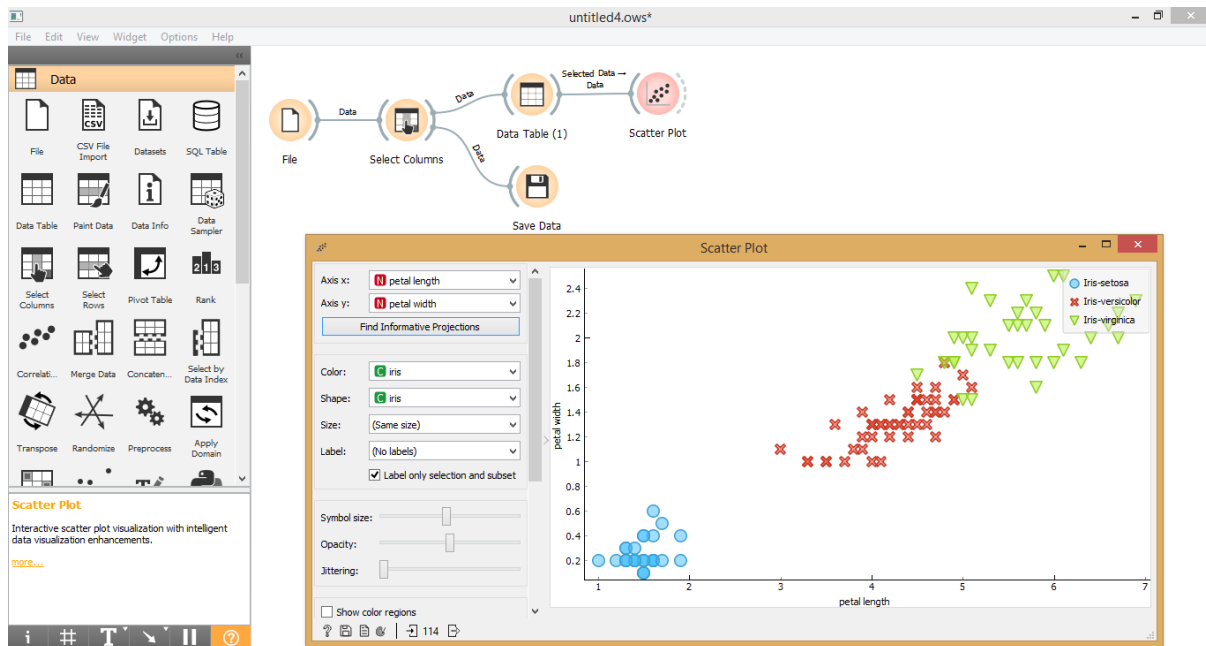


Figure 1.4

Another widget which I found interested is the Hierarchical clustering widget that can be used to group items from a matrix of distances and shows a corresponding dendogram. Given below is the screenshot of an example of the same (Figure 1.5).
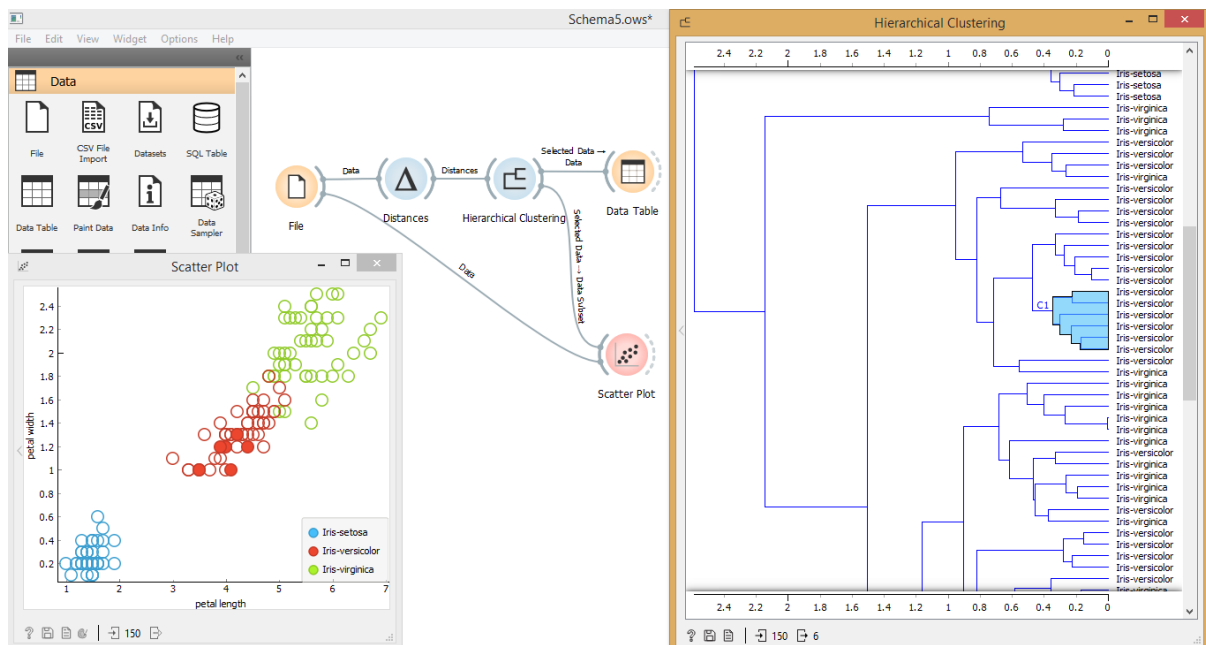


Figure 1.5

Python Script widget is where programming comes into place. It is used to extend functionalities for advanced users by running a python script that includes some input and output variables in its local namespace. Screenshot of the same is given in the Figure 1.6.
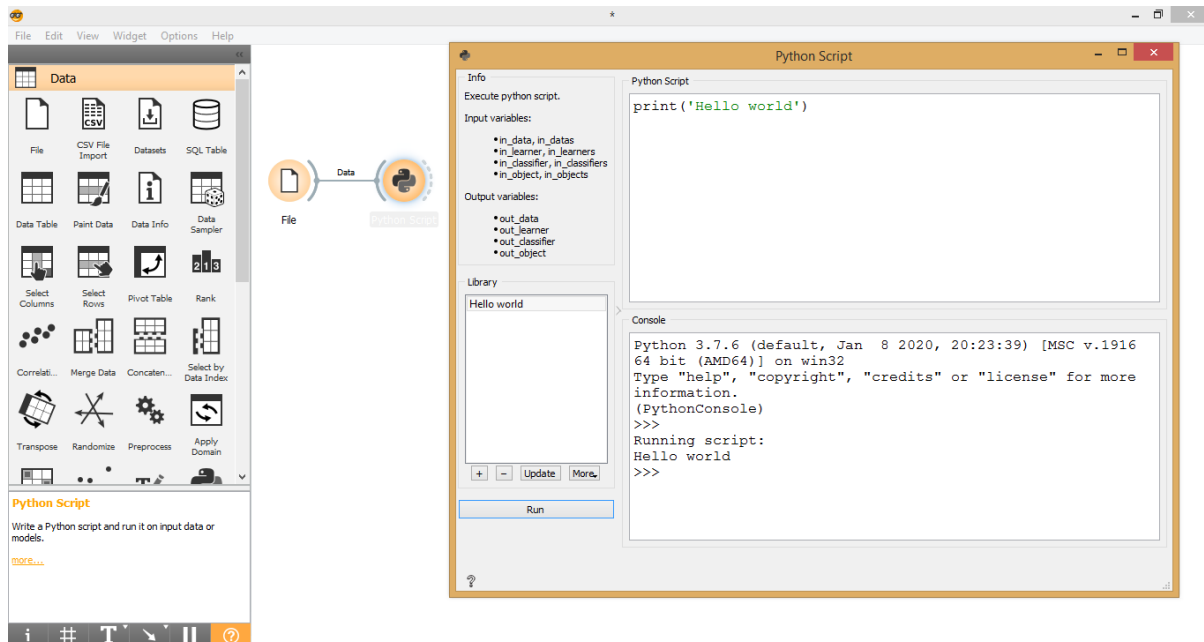


Figure 1.6

Tree widget and Predictions widget are other widgets in Orange that I learned and tried out to do predictions. Predictions widget accepts 2 inputs in which one is test data and other one is Predictors which is the output from any model widget. Here I used Tree widget as the model for classification. Tree is a very useful algorithm that can handle both discrete and continuous datasets which divides data into nodes by class purity. Screenshot of an example schema using zoo dataset for prediction is given in the screenshot (Figure 1.7). Both the test data table and model predictor data table are also shown in the Figure 1.7. Zoo dataset was the data used for model classification and I created a test data with four instances to predict the type. Here the target attribute is set as type which is the column that needs to be predicted whereas all other attributes are used for predicting.
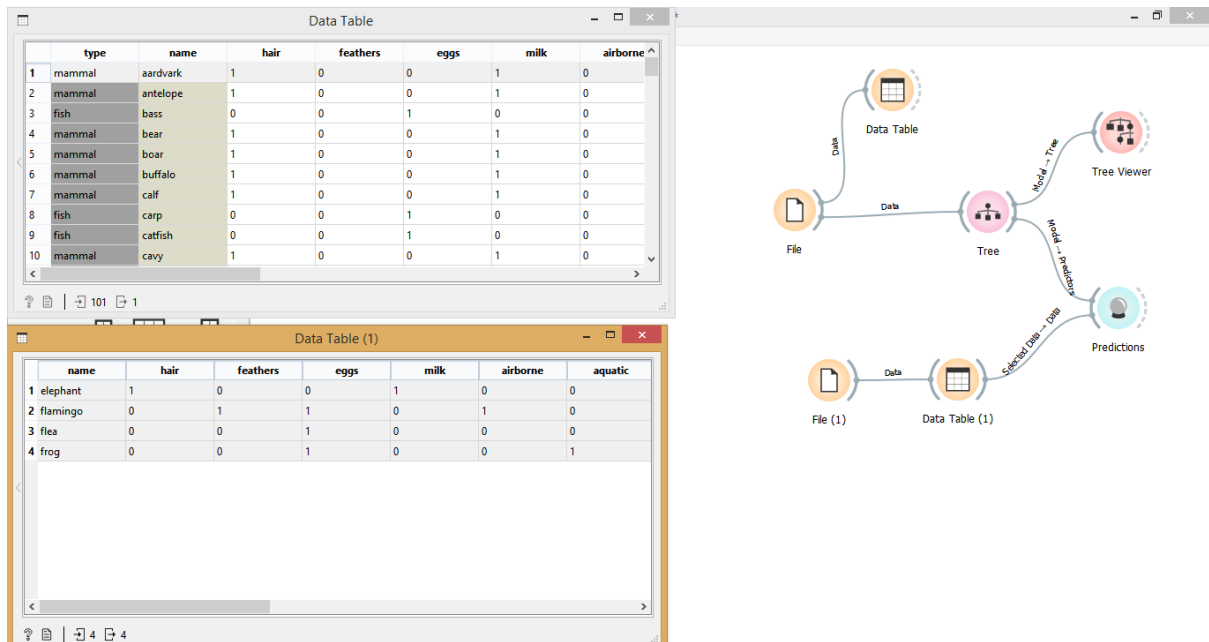
Figure 1.7

The output of Tree can be viewed using Tree Viewer (Figure1.8).
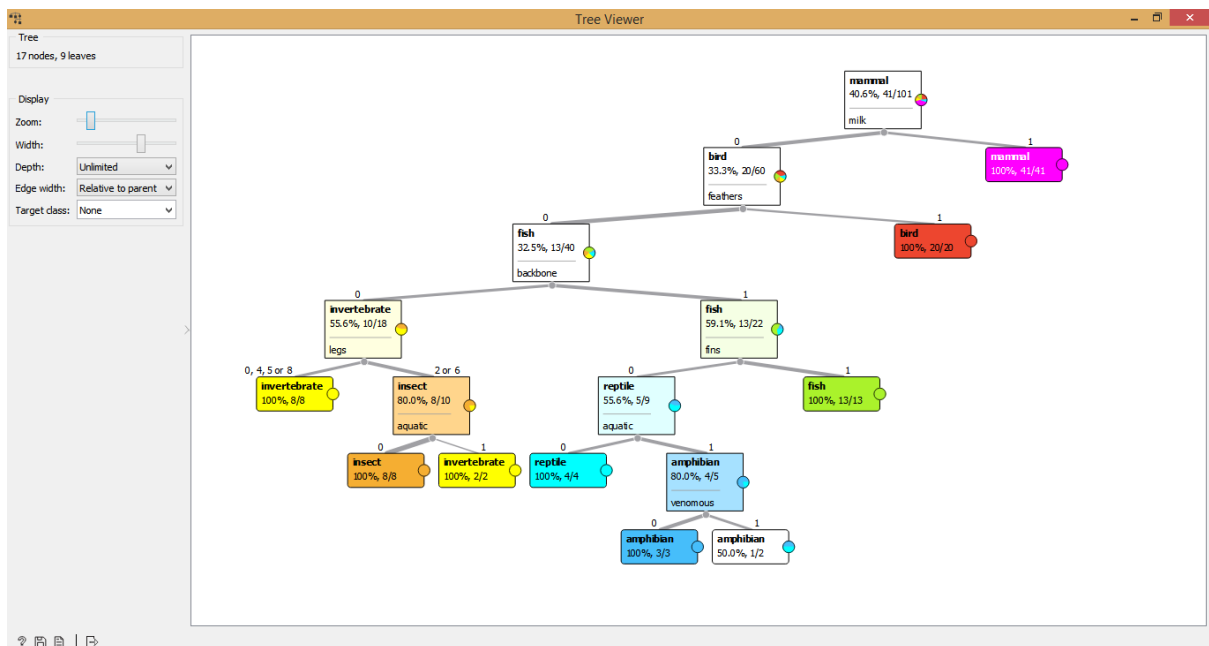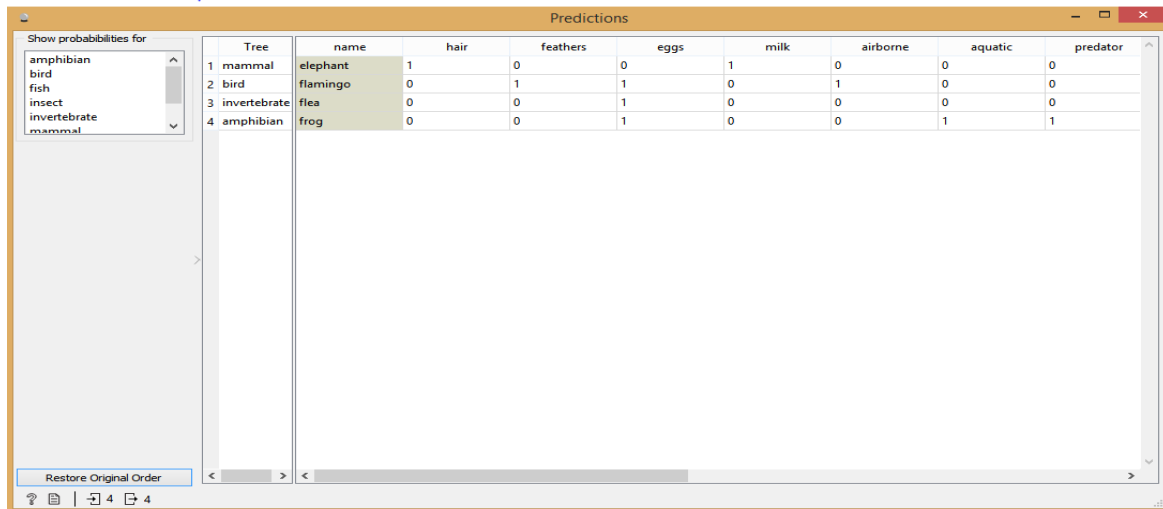


Figure 1.8

The output of Predictions is given in the screenshot below (Figure 1.9). Here the type of four instances were predicted based on other attributes such as hair, feathers, eggs, milk, airborne, aquatic, predator etc. The output of Prediction of the model is

that elephant is a mammal, flamingo is a bird, flea is an invertebrate and frog is an amphibian.



Figure 1.9

The same prediction can also be done by sampling data and using the sampled data as test data and the remaining data as the input for Tree model. Screenshot of the same schema using Data Sampler is shown in the Figure 1.10.
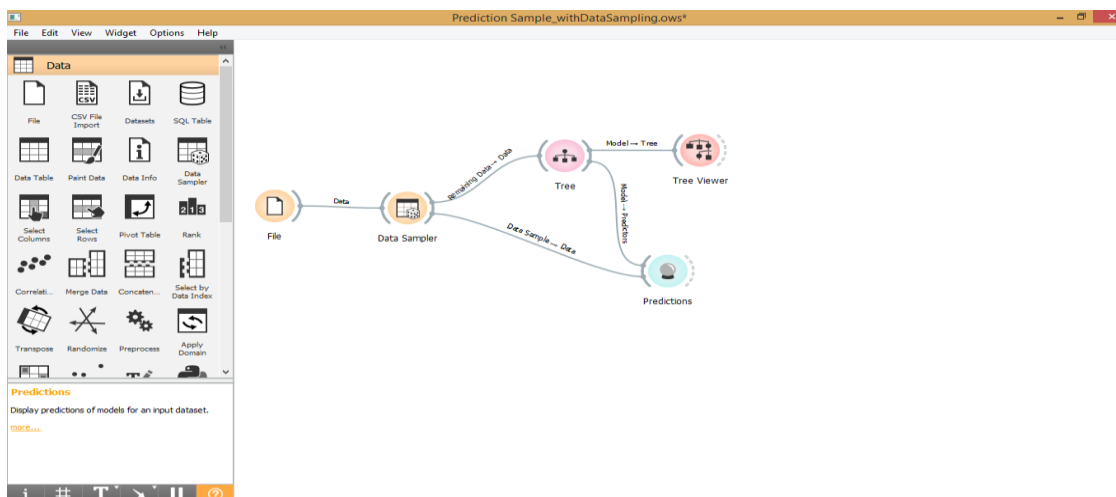


Figure 1.10

The output of Predictions of the model is that the type is being predicted for the sampled test data. The screenshot is given below in the Figure 1.11.

Figure 1.11

My experience with using Orange was really good. It has a really attractive and interactive user interface which helped me understand it easily and I was able to try out different widgets and got to know about a general idea of what the tool is and why it is used.

## 2. RAPIDMINER STUDIO

RapidMiner Studio is an extensive data science platform with visual workflow design and complete automation.

It was easy downloading the software but to install and launch, it took a huge amount of time. I downloaded the latest version of RapidMiner Studio 9.8 from the link given below:

https://rapidminer.com/get-started/

It was very slow while launching the software and asked to create an account to setup. Screenshot of the same is given below (Figure 2.1). Even after the set up and restarting the software, it takes time to open up. I also watched few youtube videos to get a basic understanding of the software.
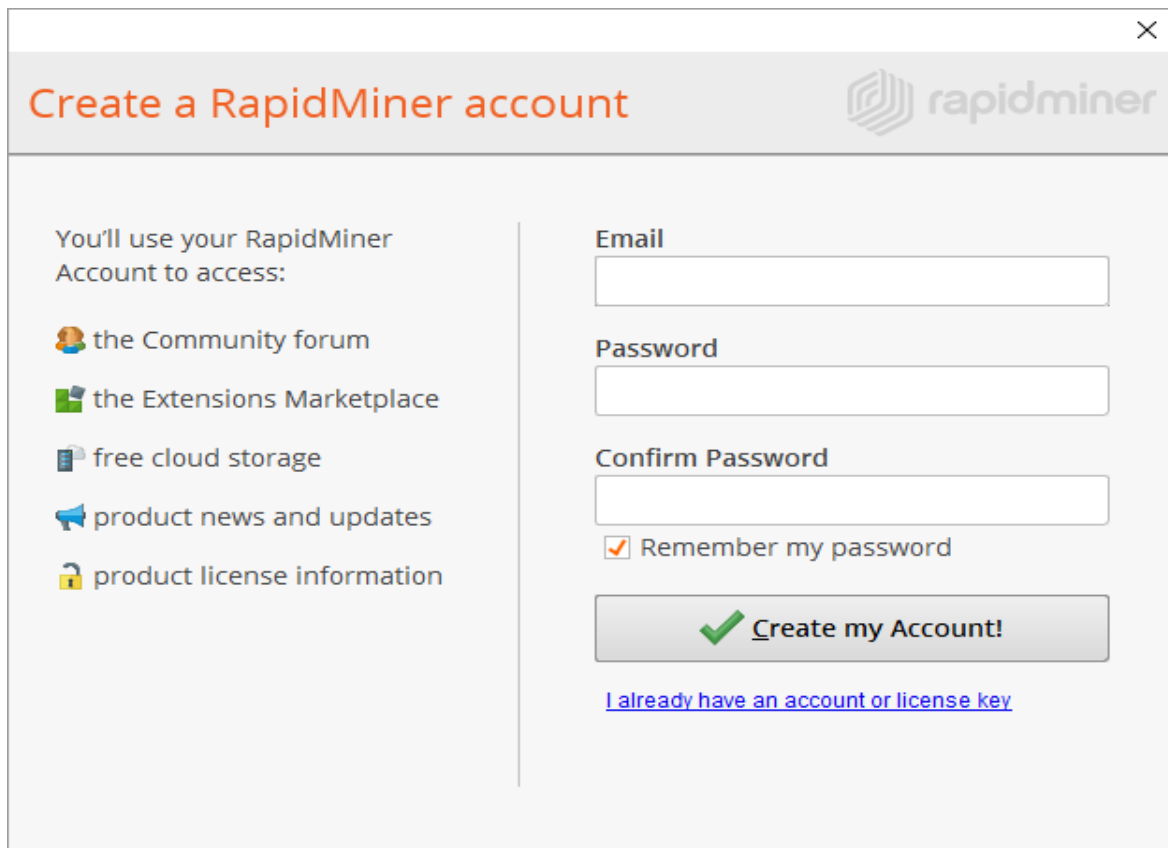


Figure 2.1

The user interface of RapidMiner Studio is somewhat user friendly but can be a bit confusing for beginners. For experts, this tool can be really helpful with a great visual workflow and automation. It can also extend functionalities using R and Python code. Its UI consists of Repository, Operators, Process panel, Views, Global search area, Ports, Parameters and Help. In Views, Design and Results can be defined as the work areas which are equipped to carry out specific tasks for data preparation, modelling etc and to access specific functionalities. The Process panel is used as a

work area to design and build any process. Repository is a storage within RapidMiner Studio for data and RapidMiner processes. Operators are building blocks that are used to create RapidMiner processes. Ports are input and output channels for operators and processes. Parameters are used to modify operator behaviour. These are few of the panels that I worked on to understand the basic functionalities. Given below is the screenshot of the RapidMiner Studio UI (Figure 2.2).
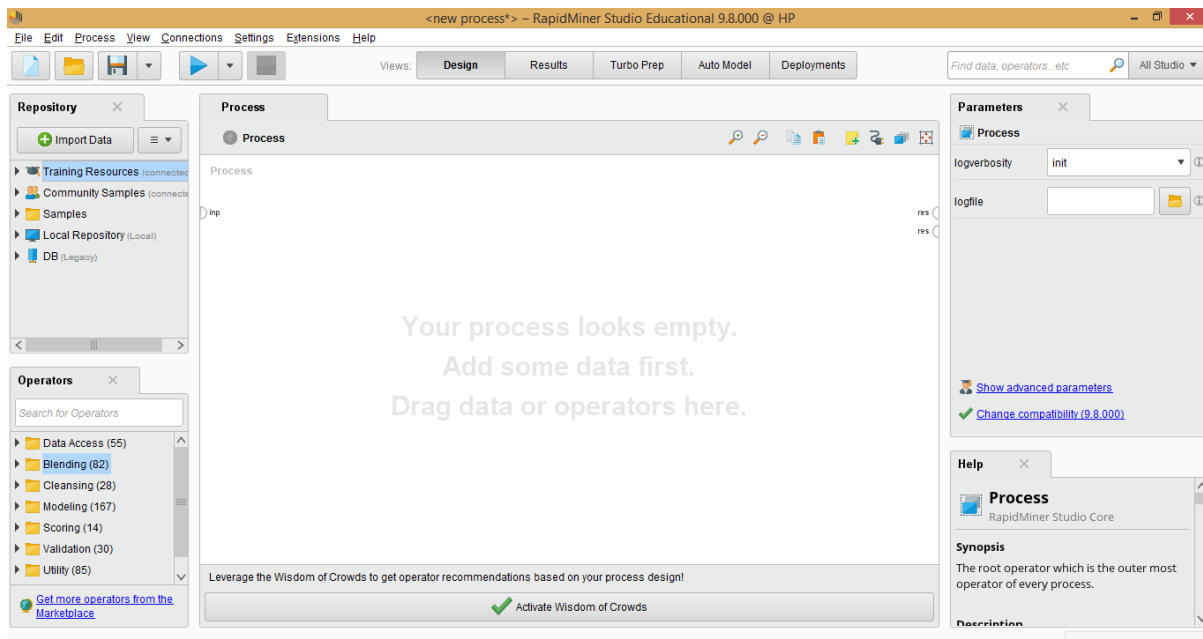


Figure 2.2

RapidMiner Studio has in-built datasets such as iris, titanic etc where the input file needs to be just dragged and dropped. The output of the file operator must be given to the output port and then we have to run the process in order to view the results. The drag + drop visual interface not only speeds up but also automates the creation of predictive models. Given below is the screenshot of the same (Figure 2.3).
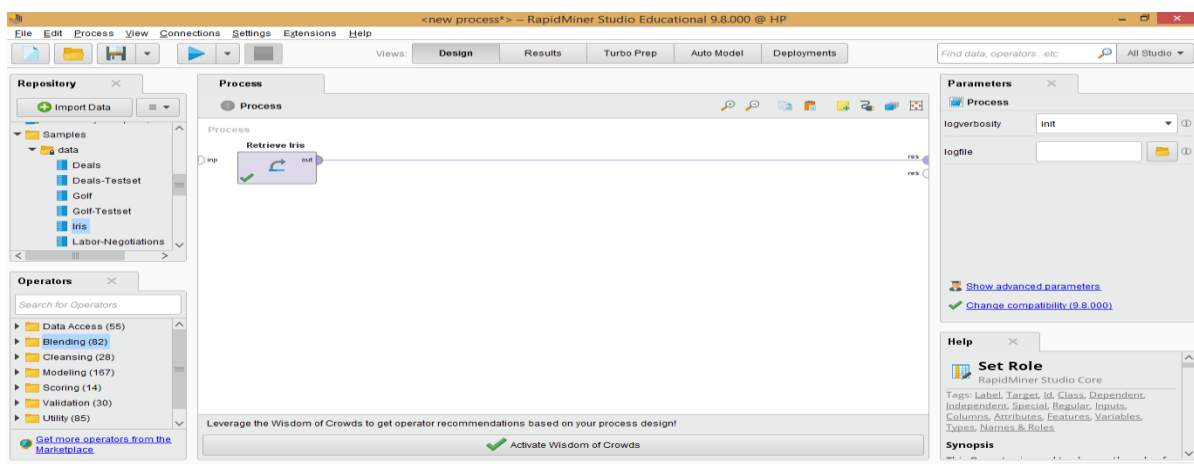


Figure 2.3

In the Results view, it is possible to view the data table (Figure 2.4), statistics (Figure 2.5) and data can be visualized using Scatterplot (Figure 2.6). Here rows are called Examples and columns are called Attributes.



Figure 2.4



Figure 2.5

There are a numerous variety of plot types under visualizations which is actually a good point about it. In the Scatterplot, x-axis columns and value columns can be changed accordingly to attain different plots. Here with a variety of plots, colour and

size can also be changed which makes it more fun to work on these visualization models. Given below is an example (Figure 2.6).



Figure 2.6

After getting a basic idea of RapidMiner Studio, I tried to predict a model by importing a dataset. Screenshot of the design view is given below (Figure 2.7). I took zoo dataset and Set Role operator was used for setting the "type" column (prediction attribute) in dataset as the label (prediction attribute role). Then the output from the Set Role that is the train data is given as input to the Decision Tree modelling operator whose output is shown in the Figure 2.8. A test data is also imported so as to predict the model using the scoring operator called Apply Model.



Figure 2.7

Figure 2.8

The result view shows the test data prediction that is the output of the Apply Model in the figure given below (Figure 2.9). It shows predictions and confidences for the inputs.



Figure 2.9

The above example was to predict a test/sampled data wherein the data was separately given as input to Apply Model operator. But in the same example given below (Figure 2.10), I used Split Data operator to split data instances with the

required ratio so that from the two outputs of Split Data operator, one is used as train data and other is used as test data. The train data is then given as input to the Split role operator where the "type" column is set as label and then given to the Decision Tree operator. The output from the Decision Tree and the test data is given as input to the Model Simulator operator.



Figure 2.10

The output of the Model Simulator is given below in the Figure 2.11 which shows the predictions, confidences and explanations for the inputs. Model Simulator is an operator that provides an easy real time modelling method so that we can make changes in output and have a much more detailed output regarding the attributes.



Figure 2.11

My experience with using RapidMiner was good. Initially I took time to grasp things but after watching some tutorials and trying out few examples, it was really interesting. RapidMiner Studio might be somewhat difficult for a beginner but for experts, I think this tool would really help them to extend functionalities and build advanced models. It has a good user interface which helped me understand the tool after few days of working on it. I was able to try out different operators and got to know about a general idea of what the tool is and why it is used.

## 3. KNIME

The Konstanz Information Miner, known as KNIME is a free open source data analytics platform which can be used for integration and reporting also.

It was easy downloading the software but to install and launch, it took a huge amount of time. I downloaded the latest version of Knime 4.2.3 from the link given below:

https://www.knime.com/download-installer/2/64bit

It was very slow while installing the software and then a dialogue box popped up saying that the startup took so long because of Antivirus tool. Screenshot of the same is given below (Figure 3.1). I also watched few youtube videos to acquire a basic knowledge of the software.



Figure 3.1

The user interface of Knime is not user friendly and it was kind of difficult for me at the beginning to understand its workflow and the working of nodes. UI of Knime consists of different panels such as Knime Explorer, Workflow Editor, Node Repository, Description, Knime console, Knime hub search and few more. Knime Explorer allows you to browse your workflows and to act upon them. Also, lot of pre-built examples by Knime are available in it. I went through some of them, but it was not that helpful for me. Local workspace is available to create either a New Knime workflow or a New workflow group. When a New Knime workflow is created, the workflow editor opens up which is a canvas for editing the currently active workflow. Node Repository is a collection of nodes which contains prewritten programs. It allows to drag and drop

multiple nodes to the workflow editor or double clicking the node will also bring nodes to create the workflow. Description panel displays the information of a node. Given below is the screenshot of the user interface of Knime analytics platform (Figure 3.2).
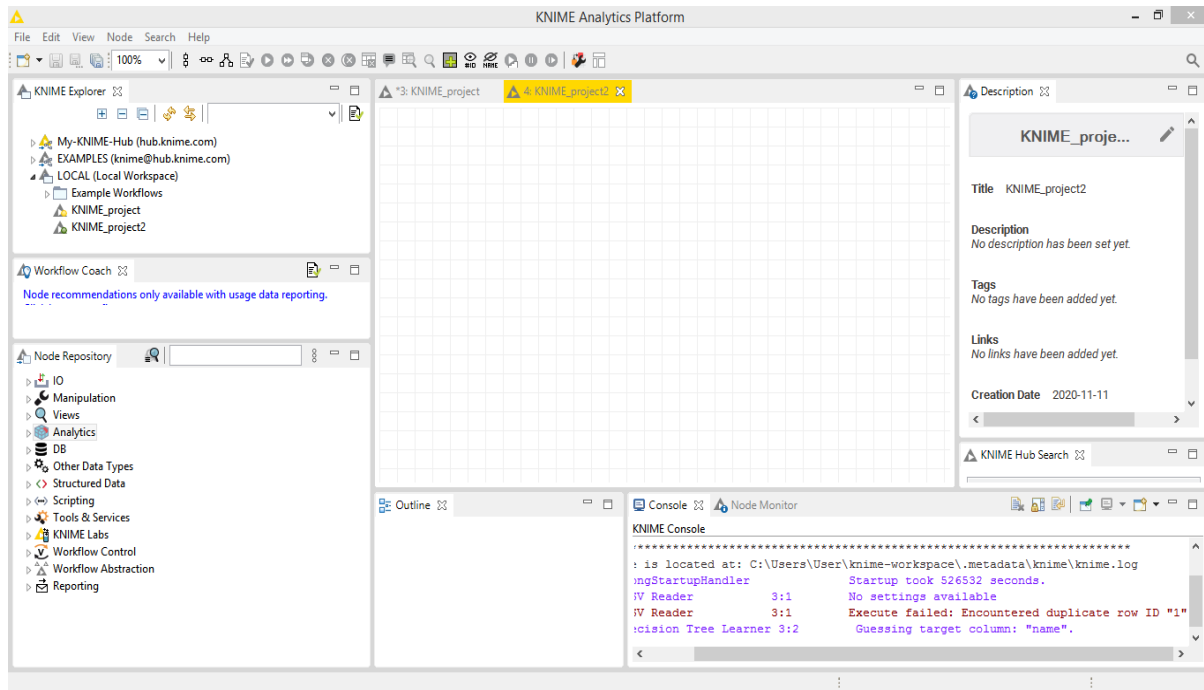


Figure 3.2

Nodes may contain one, two or more input ports on the left-hand side of the node and the output port is on the right-hand side. The node state determines the status of the node i.e., Red colour states that the node is not configured, Yellow colour states that the node is configured (ready to be executed), Green colour states that the node is executed and red cross symbol states that it is an error. Given below is an example of a CSV Reader node (Figure 3.3) which was dragged and dropped into the workflow editor. It then needs to be configured where the csv file is given as input. The node state is actually an interesting feature, but it is asked to reset, configure and re-execute everytime when changes were made. At times, I found it hard to move the nodes according to my wish. Overall, I was not satisfied with the user interface. When using Knime, my laptop got hang few times. I'm not sure whether it is my laptop's problem, but I did find it slow downing all other applications also.
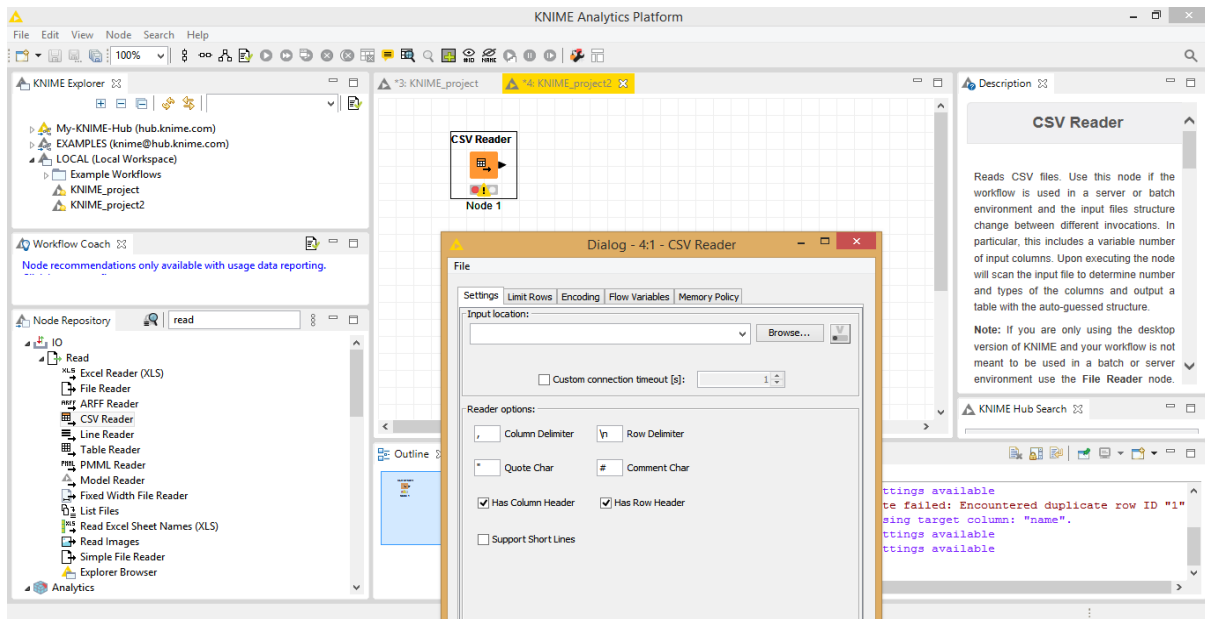
Figure 3.3

Here in an example, a CSV Reader node was used to input the iris dataset. At first, the red cross error symbol was shown because here in Knime, each row in the dataset should have a unique id. The dataset I uploaded initially did not have any unique ids and that is why it showed an error message. Screenshot of the same is given in the Figure 3.4.
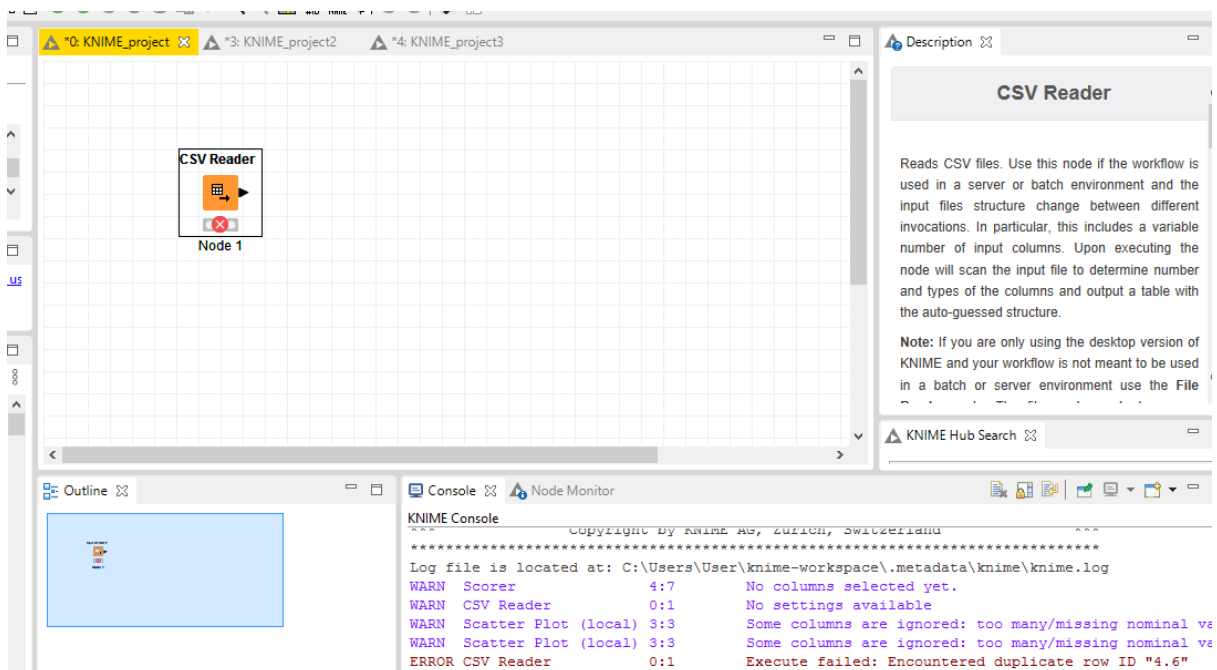


Figure 3.4

The CSV Reader is then configured with the right dataset, executed and the state is turned into Green. An example using iris dataset is shown in the Figure 3.5. After that the CSV Reader is connected to a Table View where the data table is viewed. And then

it is connected to a scatter plot to visualize the data. All the nodes have turned into Green state which means all are executed.
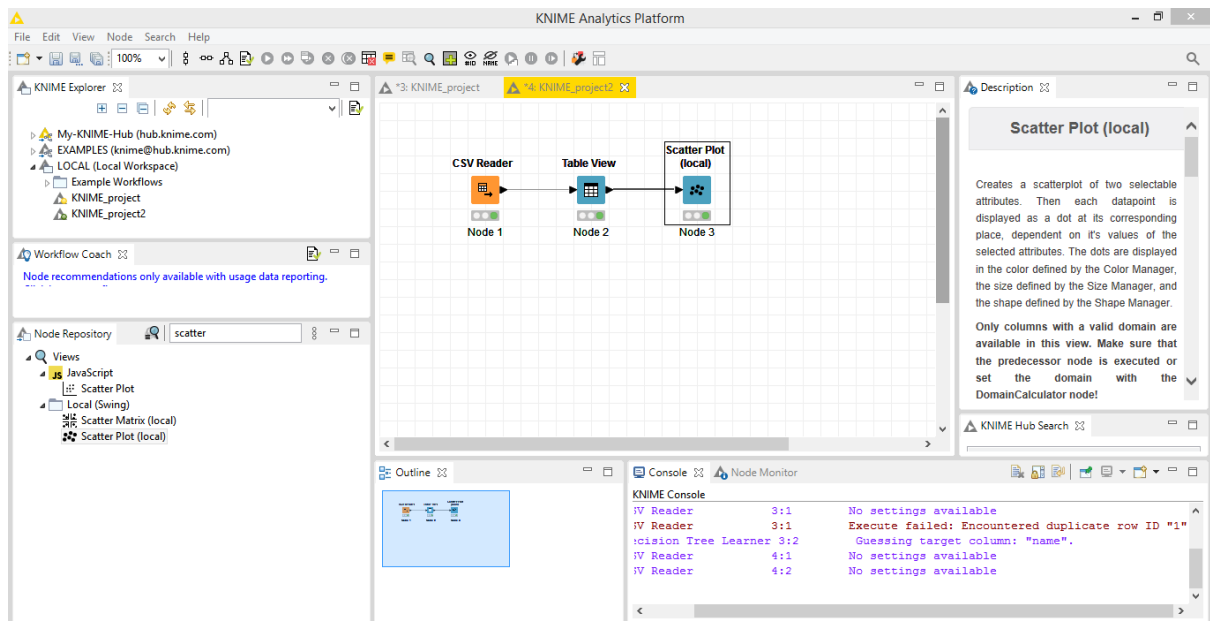


Figure 3.5

The Table view and scatterplot of the executed workflow is given in the Figure 3.6 and Figure 3.7 respectively.



| | RowID | sepal length | sepal width | petal length | petal width | iris |
|---|---|---|---|---|---|---|
| | 1 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| | 2 | 4.9 | 3 | 1.4 | 0.2 | Iris-setosa |
| | 3 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| | 4 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| | 5 | 5 | 3.6 | 1.4 | 0.2 | Iris-setosa |
| | 6 | 5.4 | 3.9 | 1.7 | 0.4 | Iris-setosa |
| | 7 | 4.6 | 3.4 | 1.4 | 0.3 | Iris-setosa |
| | 8 | 5 | 3.4 | 1.5 | 0.2 | Iris-setosa |
| | 9 | 4.4 | 2.9 | 1.4 | 0.2 | Iris-setosa |
| | 10 | 4.9 | 3.1 | 1.5 | 0.1 | Iris-setosa |

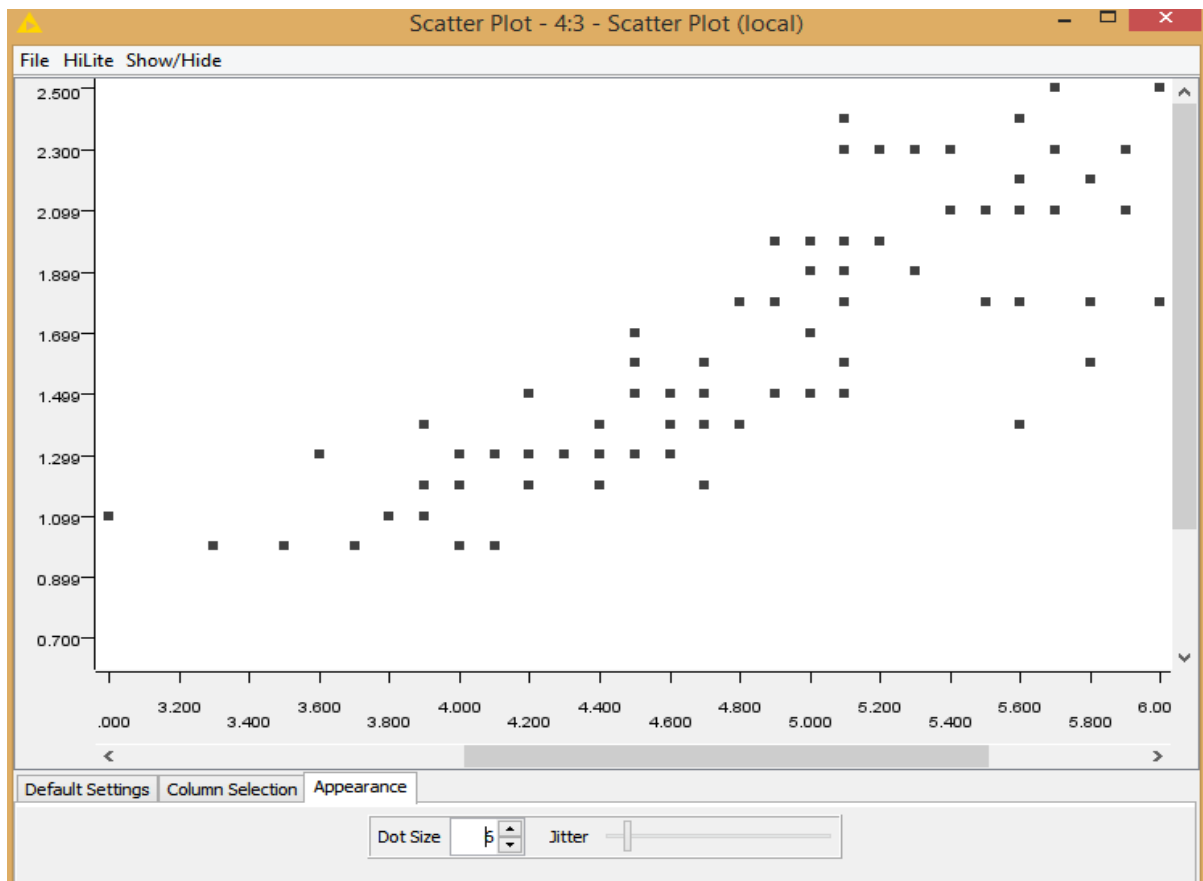Showing 1 to 10 of 150 entries

Figure 3.6

Figure 3.7

An example for prediction is shown below (Figure 3.8). Here a CSV Reader i.e., Node 1 is given the iris dataset as input (train data). It is then connected to the Decision Tree Learner (its output is given in Figure 3.9) where classification modelling is done. To get a more clear and simplified view of the output of the Decision Tree Learner, Decision Tree Viewer is used. The test data is given as input through a CSV Reader i.e., Node 5. It is connected to Table View where the data table can be viewed clearly. Both the train and test data are given to the Decision Tree Predictor in order to predict the type of iris plant species.
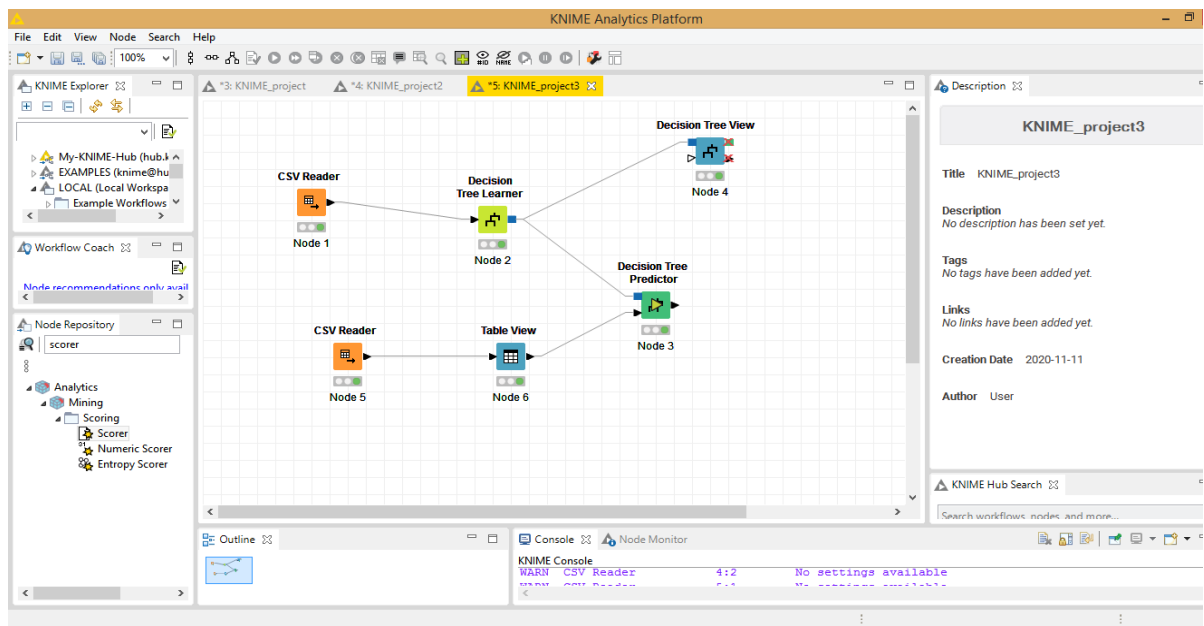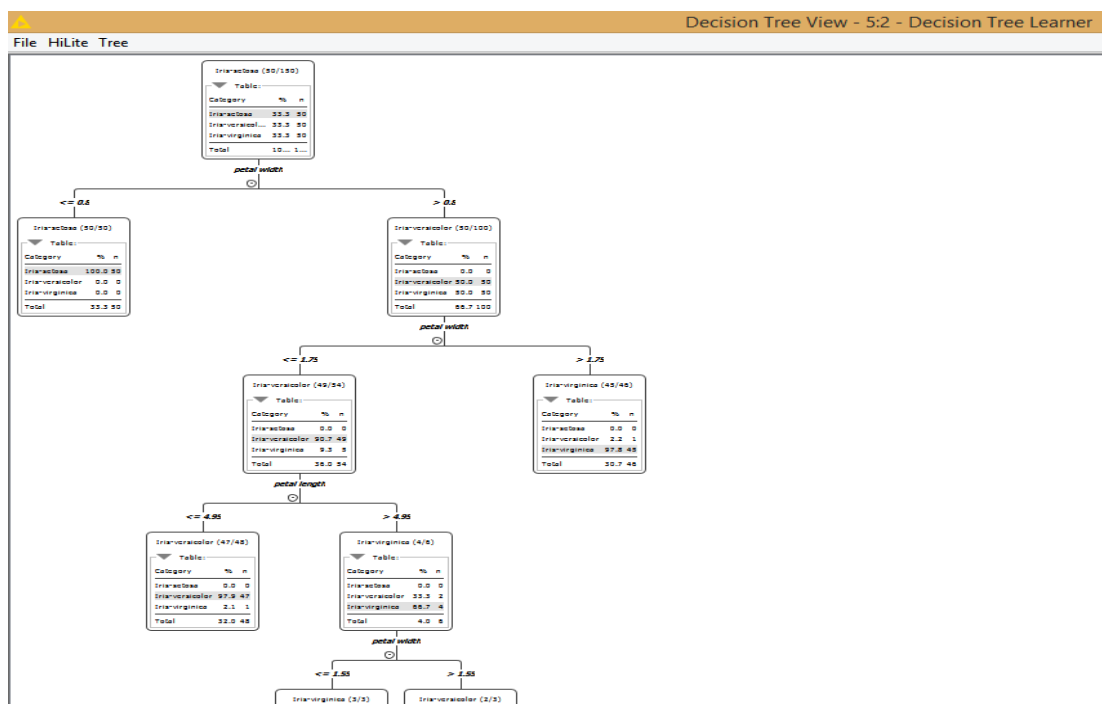
Figure 3.8



Figure 3.9

The simplified view of Decision Tree Learner is given by Decision Tree Viewer(Figure 3.10).
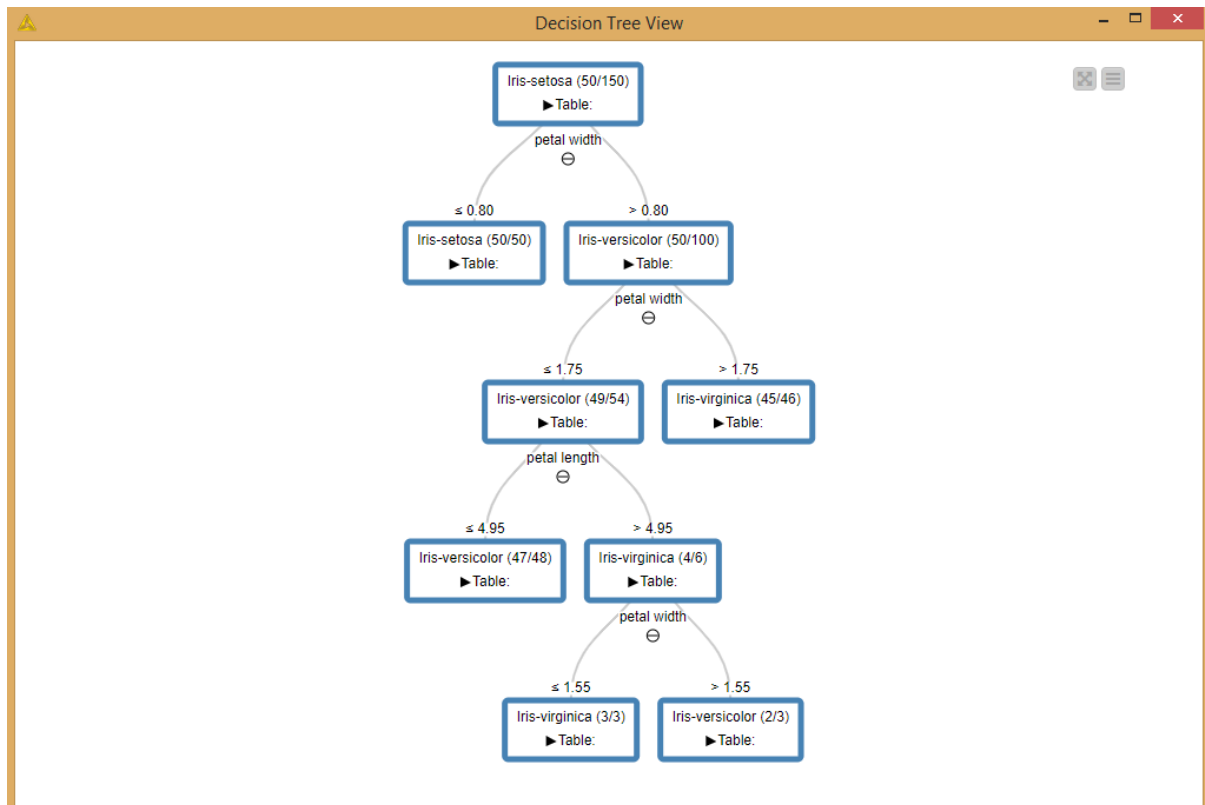
Figure 3.10

The Table View of the test data (with three instances) is given below (Figure 3.11).



Figure 3.11

The output of Decision Tree Predictor is given in the Figure 3.12 where the type of three iris species are predicted. The highlighted sections are the predicted output based on the conditions from the trained dataset. Another thing about Knime is that it has separate predictors for each models.
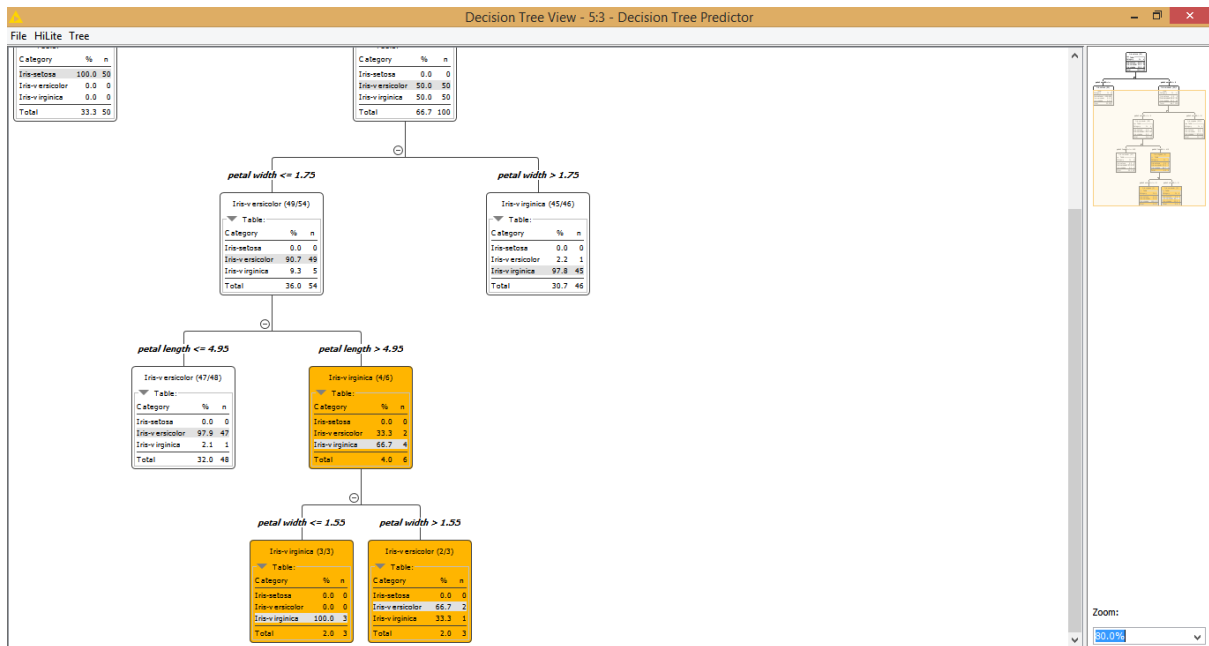
Figure 3.12

My experience with using Knime was not good. It takes a lot of time to launch the software and it slowed down all other applications. It was not user-friendly for me as I did take a lot of time to get to know about the software. After watching some tutorials and trying out few examples, I was able to do small basic examples. Knime is difficult for a beginner but for experts, I think this tool might be handy them to extend functionalities and build advanced models using Java, R and Python. It does not have a good user interface which made me spend a lot of time to get a general basic idea. After doing few examples, I was able to try out different nodes and got to know about a general idea of what the tool is and why it is used.

**REFERENCES**

1. File — Orange Visual Programming 3 documentation [WWW Document], n.d. URL https://orange3.readthedocs.io/projects/orange-visual-programming/en/latest/widgets/data/file.html (accessed 11.11.20).
2. Getting Started with Orange 01: Welcome to Orange, 2015. https://www.youtube.com/watch?v=HXjnDIgGDuI&list=PLmNPvQr9Tf-ZSDLwOzxpvY-HrE0yv-8Fy&index=1
3. Foong, N.W., 2020. Data Science Made Easy: Data Modeling and Prediction using Orange [WWW Document]. Medium. URL https://towardsdatascience.com/data-science-made-easy-data-modeling-and-prediction-using-orange-f451f17061fa (accessed 11.11.20).
4. Predictive Analytics Software | RapidMiner Studio [WWW Document], n.d. . RapidMiner. URL https://rapidminer.com/products/studio/ (accessed 11.11.20).
5. URL https://docs.rapidminer.com/latest/studio/getting-started/ui-overview.html (accessed 11.11.20).
6. Intro to the RapidMiner Studio GUI | RapidMiner, 2019. https://www.youtube.com/watch?v=Gg01mmR3-g&list=PLssWC2d9JhOZLbQNZ80uOxLypgIgWqbJA
7. KNIME, 2020. Wikipedia.https://en.wikipedia.org/wiki/KNIME
8. KNIME User Interface Walkthrough | KNIME Tutorial, 2015. https://www.youtube.com/watch?v=W0L8Feq51kE