

Tutorial (Titanic Dataset)

Step 1: Define and understand the purpose – predict or classify? Step 2: Obtain data (may involve random sampling)

We know that we use the titanic dataset. In this case, we want to know the passengers who survived the titanic accident. Thus, the classification method is the best method for solving this problem.

Step 3: Explore, clean, and pre-process data.

a. The cleaning Dataset – Titanic Dataset

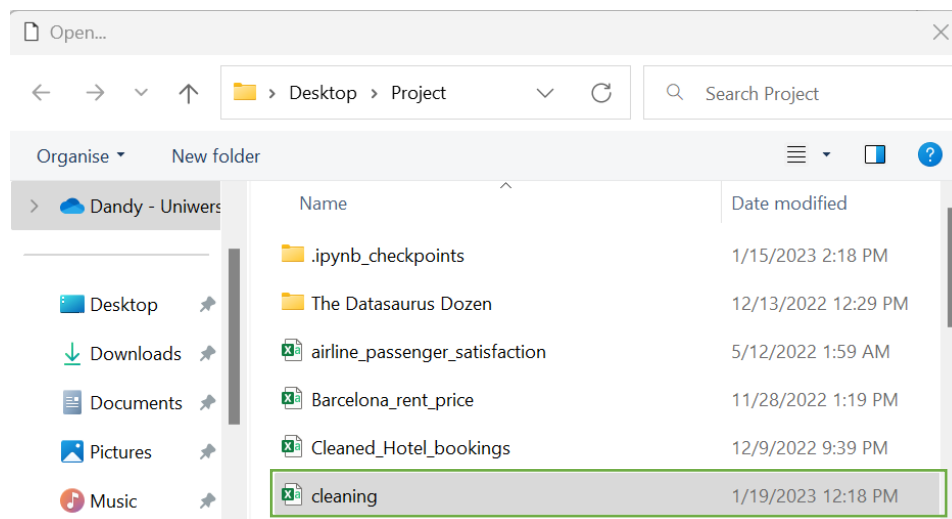
In this step, we used Jupyter Notebook to deal with some issues. There are several problems found in this dataset, such as missing values, incorrect labeling data, outlier detection, and data duplication. This process is explained in detail in the Jupyter notebook file. The resulting output is a new dataset that is clean. [\[Link to Data Cleaning\]](#)

b. Preprocess data using Orange Data Mining Software

i. Start the orange software.

ii. Import Cleaning Dataset in the file

For the first step, we choose “import cleaning dataset.” Then we import the “cleaning.csv” file.



For the next step, we close or skip it.

Import Cleaning Dataset - Orange

Source

☒ File: cleaning.csv ☐ URL:

File Type

Automatically detect type

Info

891 instance(s)
11 feature(s) (no missing values)
Data has no target variable.
3 meta attribute(s)

Columns (Double click to edit)

	Name	Type	Role	Values
8	Parch	categorical	feature	0
9	Fare	numeric	feature	
10	Embarked	categorical	feature	C, Q, S
11	Ship	categorical	skip	Titanic
12	FirstName	text	meta	
13	LastName	text	meta	
14	Ticket	text	meta	

Reset Apply

Browse documentation datasets

891

iii. Python Script and New Dataset

A Python script was used to create a new dataset with the addition of the AgeGroup column. Age less than 26 years old is categorized as young, and more than 26 years old in the adult category. We have a new dataset as below.

New Dataset - Orange

Info

891 instances (no missing data)
No features
No target variable.
15 meta attributes

Variables

☒ Show variable labels (if present)
☒ Visualize numeric values
☒ Color by instance classes

Selection

☒ Select full rows

Restore Original Order

Send Automatically

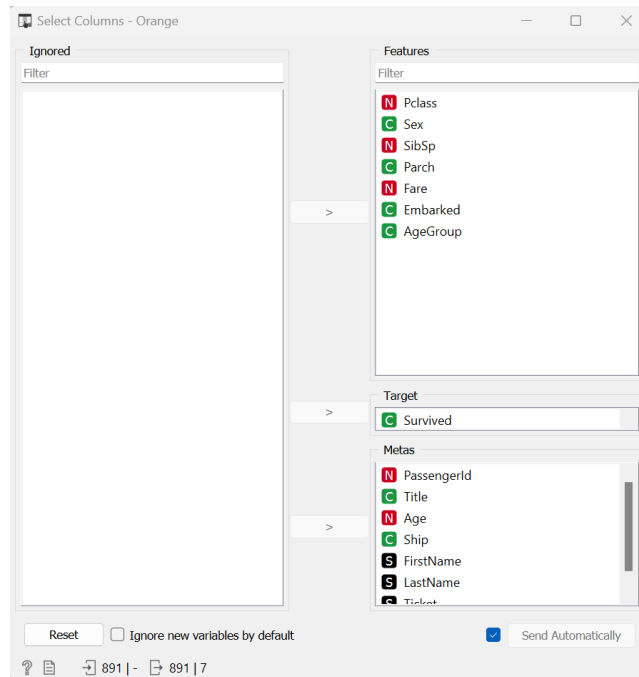
	PassengerId	Survived	Pclass	Title	Sex	Age	SibSp	Parch	Fare	Embarked
1	1	0	3	Mr.	male	22	1	0	7.90	S
2	2	1	1	Mrs.	female	35	1	0	31.00	C
3	3	1	3	Miss.	female	26	0	0	7.92	S
4	4	1	1	Mrs.	female	35	1	0	31.00	S
5	5	0	3	Mr.	male	35	0	0	8.05	S
6	6	0	3	Mr.	male	28	0	0	8.46	Q
7	7	0	1	Mr.	male	35	0	0	31.00	S
8	8	0	3	Master.	male	22	2	0	21.08	S
9	9	1	3	Mrs.	female	27	0	0	11.13	S
10	10	1	2	Mrs.	female	22	1	0	30.07	C
11	11	1	3	Miss.	female	22	1	0	16.70	S
12	12	1	1	Miss.	female	35	0	0	26.55	S
13	13	0	3	Mr.	male	22	0	0	8.05	S
14	15	0	3	Miss.	female	22	0	0	7.90	S
15	16	1	2	Mrs.	female	35	0	0	16.00	S
16	17	0	3	Master.	male	22	2	0	29.12	Q
17	18	1	2	Mr.	male	28	0	0	13.00	S
18	19	0	3	Planko.	female	31	1	0	18.00	S
19	20	1	3	Mrs.	female	28	0	0	7.90	C
20	21	0	2	Mr.	male	35	0	0	26.00	S
21	22	1	2	Mr.	male	34	0	0	13.00	S
22	23	1	3	Miss.	female	22	0	0	8.03	Q
23	25	0	3	Miss.	female	22	2	0	21.08	S
24	26	1	3	Mrs.	female	35	1	0	31.00	S
25	27	0	3	Mr.	male	28	0	0	7.90	C
26	28	0	1	Mr.	male	22	2	0	31.00	S
27	29	1	3	Miss.	female	28	0	0	7.90	Q
28	30	0	3	Mr.	male	28	0	0	7.90	S

891 | 891

0°C Rain/snow

10:27 AM 1/19/2023

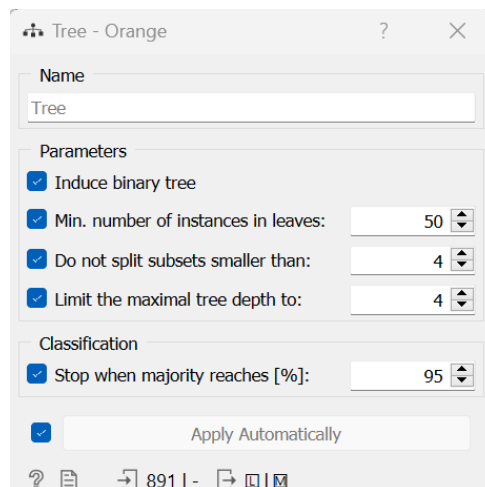
There are 891 instances without missing values and a new “Age group” column. Then we filter a new data set as below.



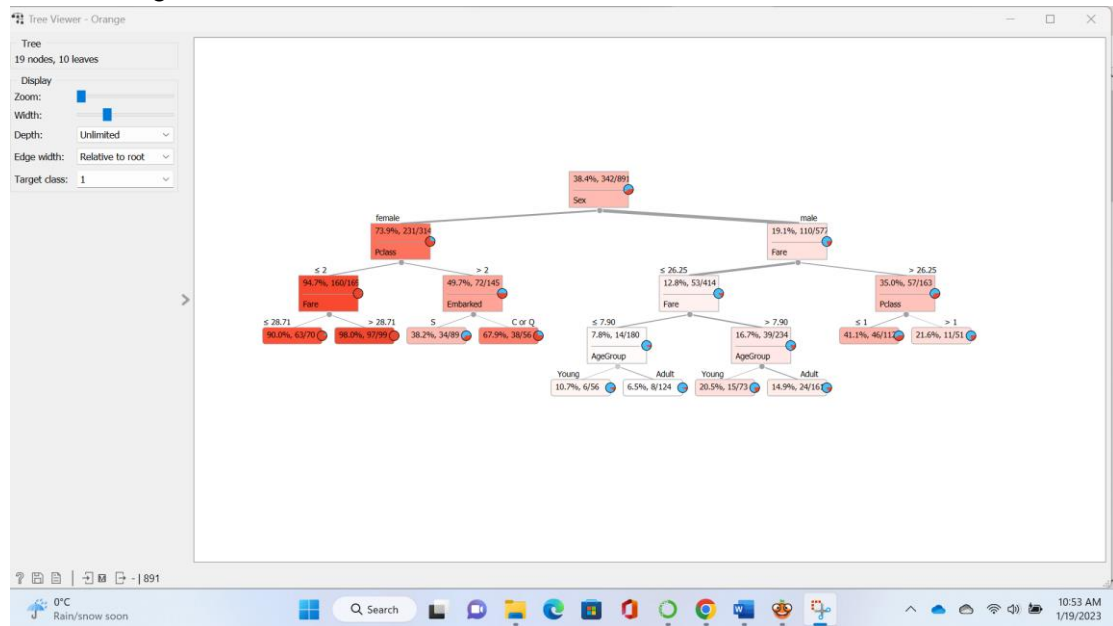
We filtered some data that will be used in the analysis. Because we already have an “AgeGroup” column, we will use that column instead of age. Hence, the columns used include Pclass (Passenger Class), Sex (Passenger Gender), Sibsp (Number of siblings/Spouses abroad), Fare (Ticket price), and Embarked (Passenger Embarkation) in the features space. Then, the “survived” column is put as a target.

iv. Tree and Tree Viewer

In this step, we can know what factors affect the possibility of passengers surviving. We set the limit maximal for the tree as four, as below.



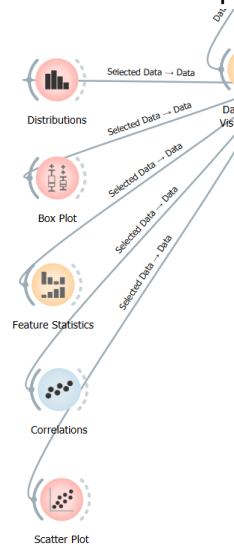
After setting the tree, we can know the tree in Tree Viewer as below:



The highest factors that affect survival rates are gender. Women have a higher chance of survival than men. 73.9% of women survive, and men only 19.1%. The second factor is class; 94.7% of women in the first and second classes stayed than those in the third.

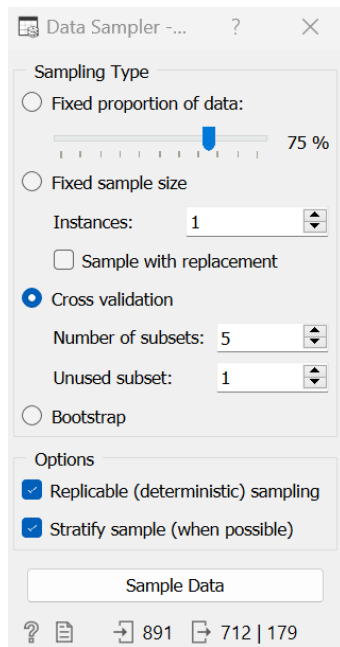
v. Data Visualization

We can choose these visualizations for visualization data, and the authors provide the visualization in Jupyter Notebook. [\[Link to Data Visualization\]](#)



vi. Data Sampler

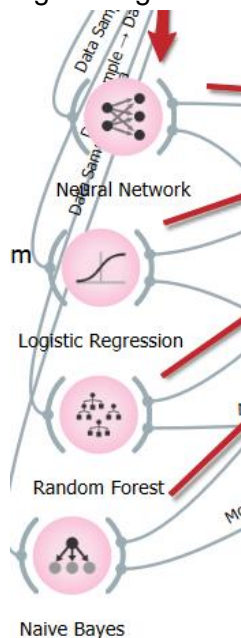
In this step, we choose a cross-validation sampling type and set your data sampler as below.



We usually use cross-validation to tune the hyperparameters of a given machine-learning algorithm to get good performance according to some suitable metric.

vii. Algorithms

In this step, the author provides four classification algorithms: the neural network, logistic regression, random forest, and naïve bayes.



Neural networks help us cluster and classify. You can think of them as a clustering and classification layer on top of the data you store and manage. They help to group unlabeled data according to similarities among the example inputs, and they classify

data when they have a labeled dataset to train on. (Neural networks can also extract features fed to other algorithms for clustering and classification so that you can think of deep neural networks as components of larger machine-learning applications involving algorithms for reinforcement learning, type, and regression.)

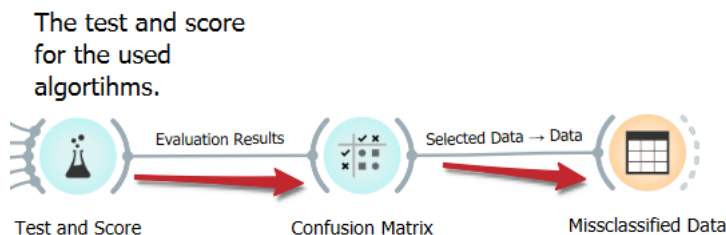
Logistic Regression is a classification technique used in machine learning. It uses a logistic function to model the dependent variable. The dependent variable is dichotomous, i.e., there could only be two possible classes (e.g., either the cancer is malignant or not). As a result, this technique is used while dealing with binary data.

Among all the available classification methods, random forests provide the highest accuracy. **The random forest** technique can also handle big data with numerous variables running into thousands. It can automatically balance data sets when a class is more infrequent than other classes in the data. The method also handles variables fast, making it suitable for complicated tasks.

The Naive Bayes is fast and easy to implement, but their most significant disadvantage is the independent predictor requirement. In most real-life cases, the predictors are dependent, hindering the classifier's performance.

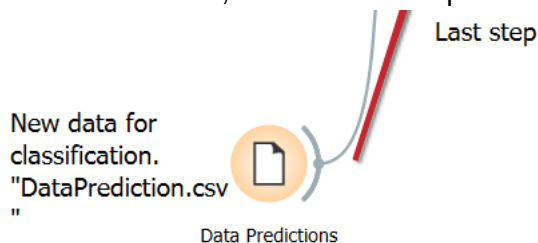
viii. **The Test, Result, and Misclassification Data.**

This step lets us know the test and score, confusion matrix, and Misclassification data.



ix. **The Predictions**

We provide new data for classifying them. We import a new data set, "DataPrediction.csv," into the first step's data predictions.



Then we check the results in the predictions.



Check the result
for predictions

Predictions

