# Titanic Dataset (CRISP-DM)

**Business Understanding**: Business Understanding is the first step in the process. In this step, we figure out what problem or question we want to solve or answer.

**Data Understanding:** The second step concerns the data used to resolve the business problem. This data must be collected, and its significance must be comprehended. For instance, if the data is stored in a table, it is crucial to understand the precise meaning of each column. The data comprehension may necessitate a return to the first step, i.e. business comprehension.

**Data Preparation:** In this step, the data get prepared for the next step, i.e. the modelling. Preparing the data includes, for example, a strategy to deal with missing values. In general, the data could be cleaner than it is. But since data of low quality can lead to wrong answers to business questions, the data preparation step also includes data cleaning.

**Modelling:** This step is one of the core steps where one tries to find a model which describes the data. In terms of our business problem of predicting the price of a car based on its features, modelling means finding a function that maps the features of a car to its price. It could be possible that the modelling step requires us to turn back to the data preparation step (for example, if we realise in the modelling step that the data still needs to be sufficiently cleaned).

**Evaluation:** Once the model is constructed, its performance must be evaluated. This occurs during the evaluation phase. For instance, we could predict the prices of automobiles based on a model for which we know the actual prices and then compare how closely the predicted values match the actual values.

**Deployment:** In the last step, if we are satisfied with the model's performance, we need to productionize it. For the car price prediction example, that could mean that we build an application that allows entering certain car features in a form and returns the price prediction to the user.

# CRISP-DM Example Walkthrough for The Titanic Dataset

## Business Problem

As mentioned above, the CRISP-DM process starts with understanding the business problem. There are some business questions regarding this dataset:

- What factors are crucial to surviving?
- Who has a higher chance to survive: Jack, as a young male, travelling alone as the passenger of $3^{rd}$ class, or Rose, as a young female, travelling with family and passenger of $1^{st}$ class?

## Data Understanding and Data Preparation

The data is taken from the module. Below is data looks like:

| Out[2]: | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked | ship |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3.0 | Braund, Mr. Owen Harris | male | 22 | 1 | 0 | A/5 21171 | 7,25 | NaN | S | Titanic |
| 1 | 2 | 1 | 1.0 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38 | 1 | 0 | PC 17599 | 71,2833 | C85 | C | Titanic |
| 2 | 3 | 1 | 3.0 | Heikkinen, Miss. Laina | female | 26 | 0 | 0 | STON/O2. 3101282 | 7,925 | NaN | S | Titanic |

The data above contains several problems, such as missing values, incorrect labelling data, outlier detection, and data duplication. This process is described in detail in the Jupyter logbook file. The resulting output is a new clean dataset in csv form. [Link to Data Cleaning]
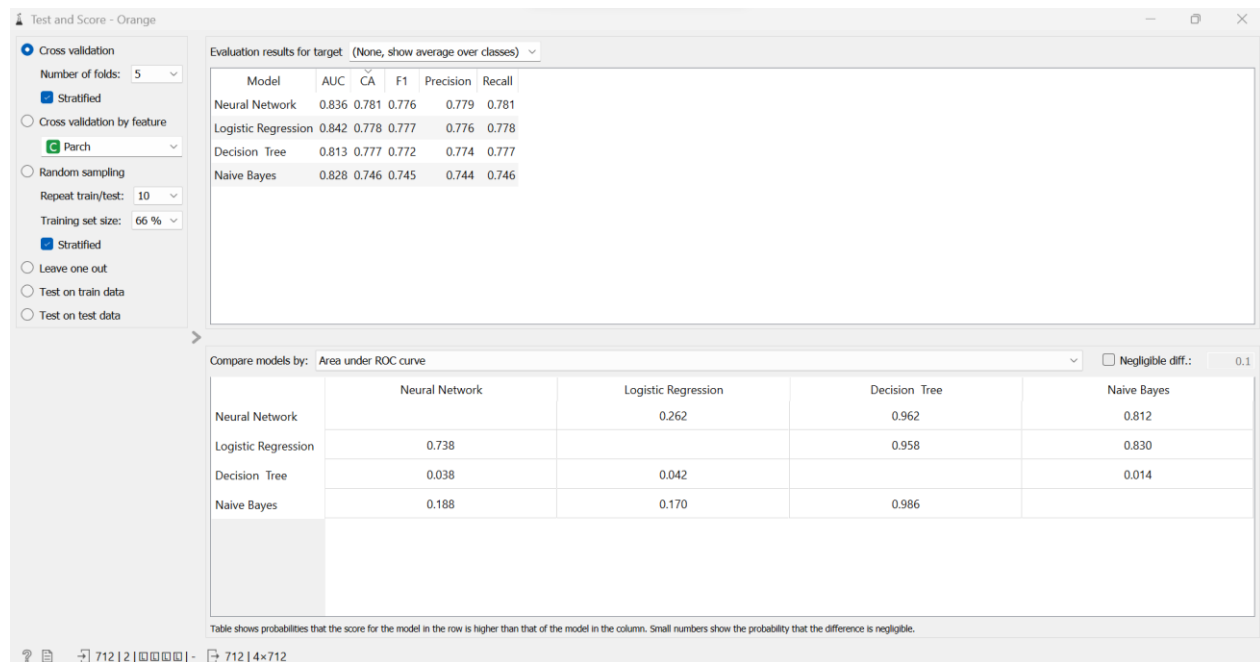
## Modelling

To generate a model to predict who has a high chance of survival and the determining factors in this case, we used several algorithms such as Neural Network, Logistic Regression, Naive Bayes, and Decision Tree.

- **Neural networks** classify and cluster. These add clustering and classification to our data. This classifies data when trained on a labelled dataset and groups unlabeled data by example input similarities. Deep neural networks can extract features fed to other algorithms for clustering and classification to be used as components of larger machine-learning applications involving 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.
- **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.
- **Decision Trees** (DTs) is a non-parametric supervised classification and regression learning method. The objective is to develop a model capable of predicting the value of a target variable using simple decision rules inferred from the data features. A tree can be viewed as a piecewise constant approximation.

**Evaluation**

To evaluate the performance of these models, we use classification accuracy. Classification accuracy is a metric that summarises a classification model's performance by dividing the number of correct predictions by the total number of predictions. It is simple to calculate and intuitive to comprehend, making it the most popular metric for classifier model evaluation. This intuition only works when the distribution of examples among classes is grossly asymmetrical.
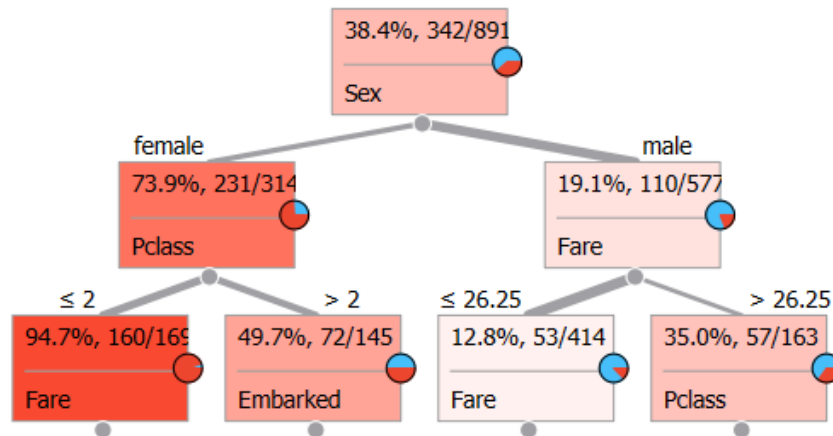
Below are the results:



Regarding the results, Neural Network has a higher classification accuracy of 78.1%.

**The answer to the Business Model:**

1. 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 the class of passengers; 94.7% of women in the first and second classes stayed than those in the third. For men, the fare has the most significant impact after gender, there 35% of the men passengers survived if they paid > 26.25.



2. According to the model created by the author, Rose is predicted to survive in four algorithms.

The following picture shows the survival percentage for two passengers (Rose and Jack).



The above picture illustrates that Rose has a bigger chance of surviving than Jack, with 84% for Logistic Regression, 92% for Neural Network, 91% for Naïve Bayes, and 75% for Decision tree.

## Conclusion

To summarise, the factors that can improve survival in this problem are gender, passenger class, and fare. The female gender has a greater chance of survival than the male, and passengers in 1$^{st}$ and 2$^{nd}$ classes have a greater chance of survival than passengers in 3$^{rd}$ class. Thus, Jane is predicted to survive in all four methods used. Of all the methods, the method with the highest accuracy is the Neural network, with a percentage of 78.1%, and Naive bayes is in the bottom position with an accuracy rate of 74.6%.

# TUTORIAL

### 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 labelling 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]
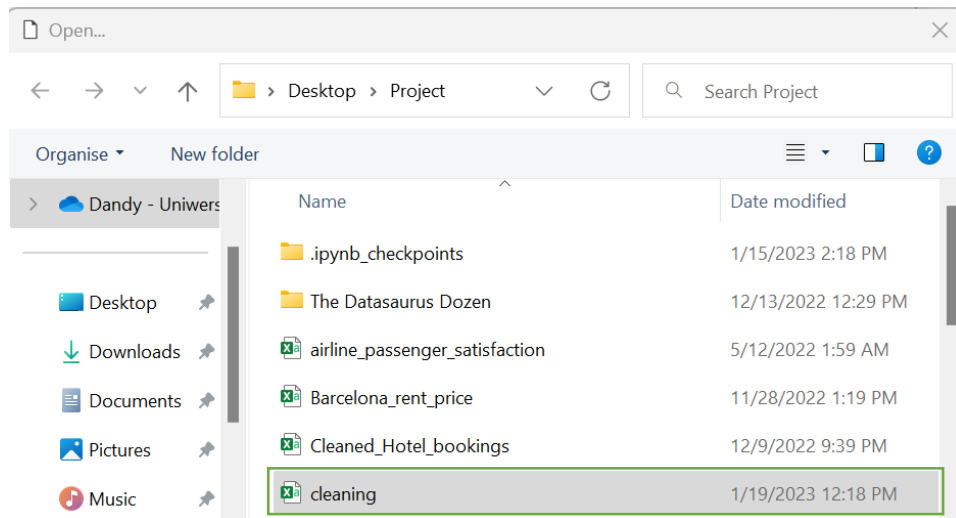
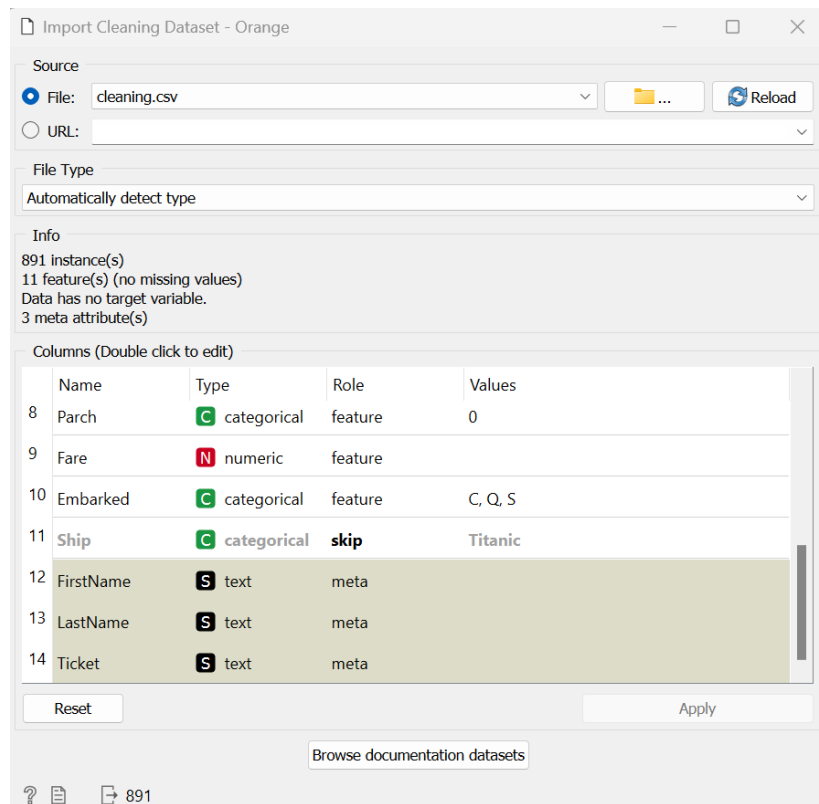### b. Preprocess data using Orange Data Mining Software

#### i. Start the orange software.



The picture above is the orange project for the titanic dataset. [Link to Download Orange Modelling]

#### 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.



### 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 categorised as young, and more than 26 years old in the adult category. We have a new dataset as below.

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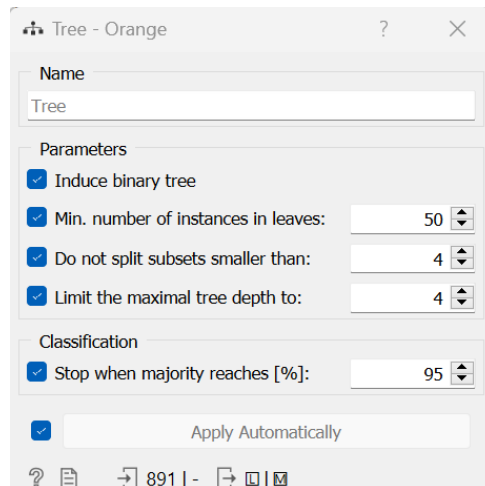


We filtered some data that will be used in the analysis. Because we already have an "AgeGroup" column, we will use that column instead of age. Hence, the columns used include Pclass (Passenger Class), Sex (Passenger Gender), Sibsp (Number of siblings/Spouses abroad), Fare (Ticket price), and Embarked (Passenger Embarkation) in the features space. Then, the "survived" column is put as a target.
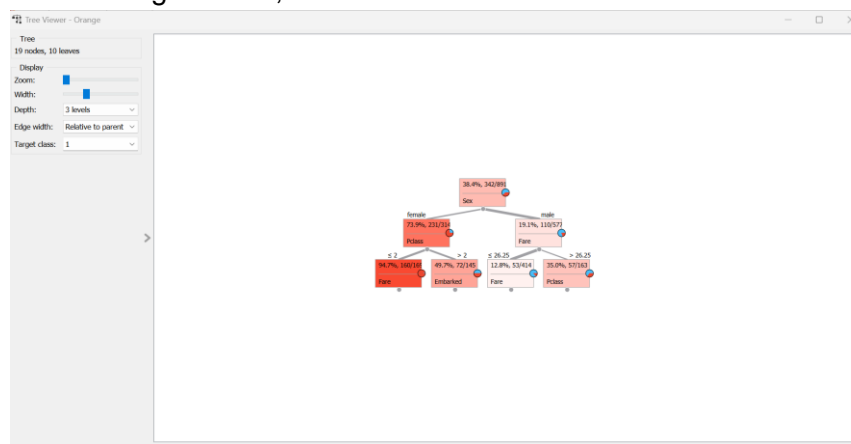
**iv.    Tree and Tree Viewer**
In this step, we can know what factors affect the possibility of passengers surviving. We set the limit maximal for the tree as four, as below.
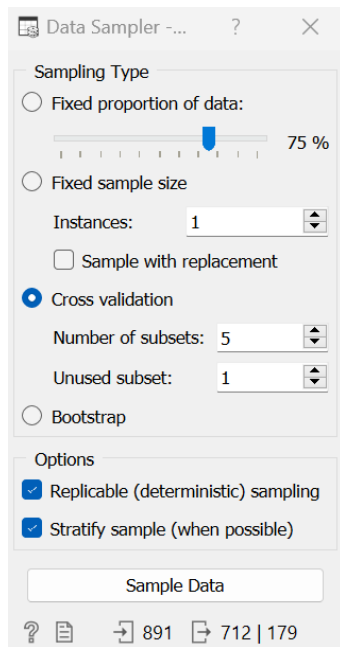
After setting the tree, we can know the tree in Tree Viewer as below:



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.
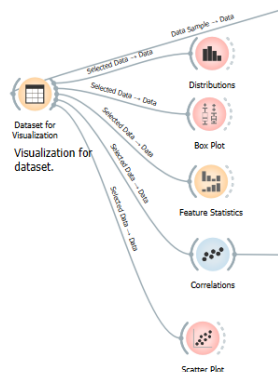
v.    **Data Sampler**

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. Cross-validation, also known as rotation estimation or out-of-sample testing, is several similar model validation techniques used to evaluate how the results of a statistical analysis will generalise to an independent data set. Cross-validation is a resampling technique that employs distinct data subsets to test and train a model across multiple iterations. It is primarily employed in situations where prediction is objective, and one wishes to estimate how well a predictive model will perform in practice. In a typical prediction problem, a model is provided with a dataset of known data for training (training dataset) and a dataset of unknown data (or first-seen data) for testing (called the validation dataset or testing set).

## vi.   Data Visualization

We can choose these visualisations for visualisation data, and the authors provide the visualisation in Jupyter Notebook. [Link to Data Visualization]
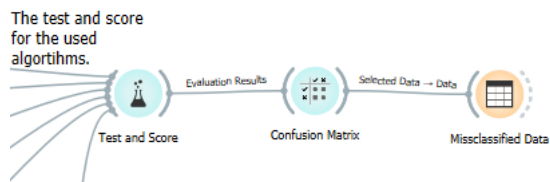
### vii. Algorithms

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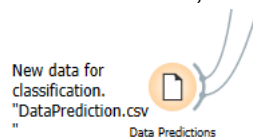


### viii. The Test, Result, and Misclassification Data.

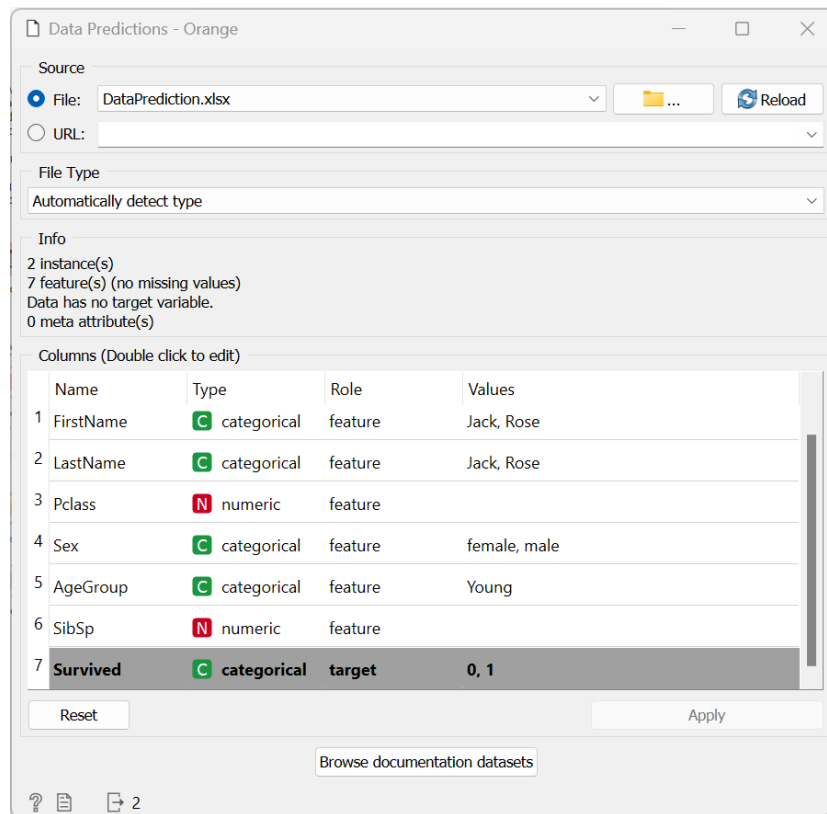This step lets us know the test and score, confusion matrix, and Misclassification data.



### 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.



In the next step, we set the dataprediction.csv like the following picture:

Then we check the results in the predictions.



Check the result
for predictions

Predictions