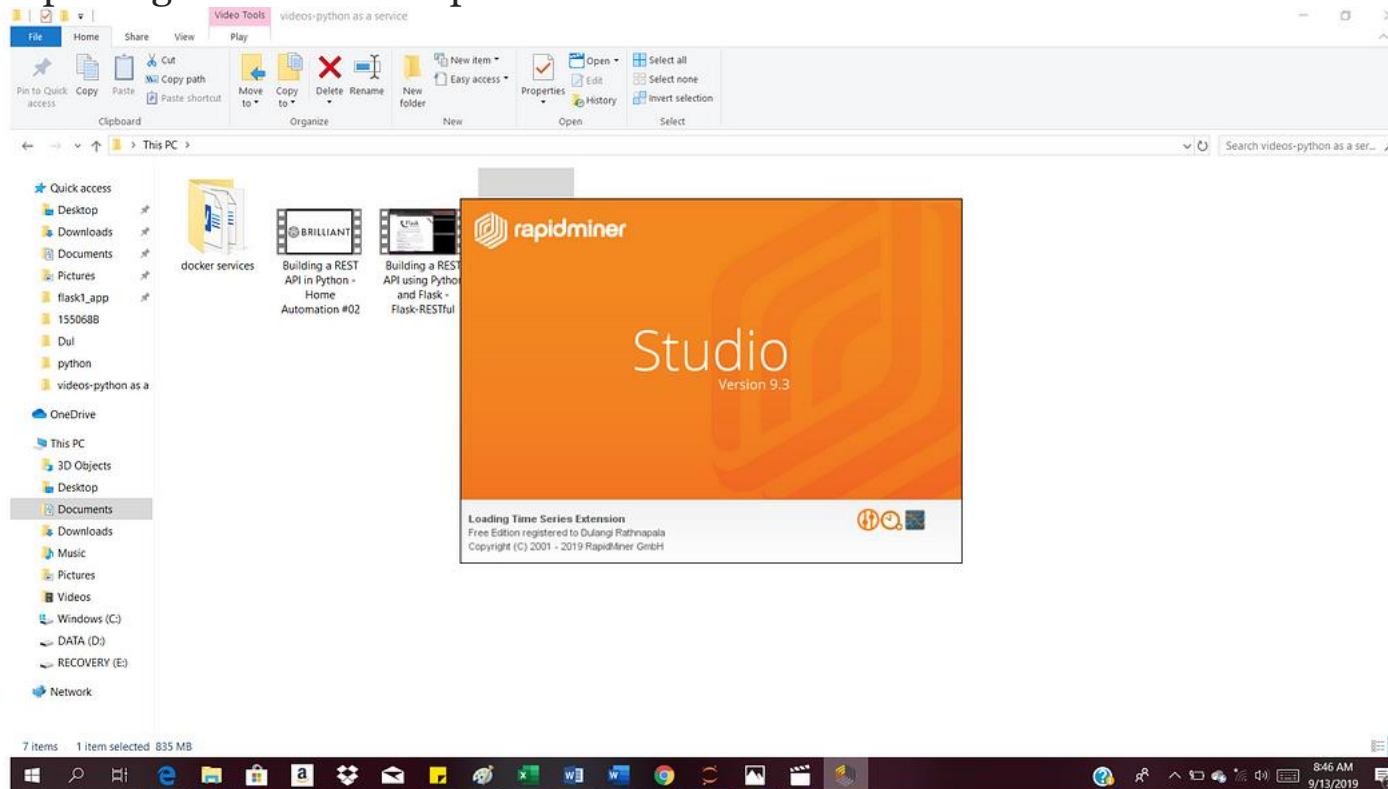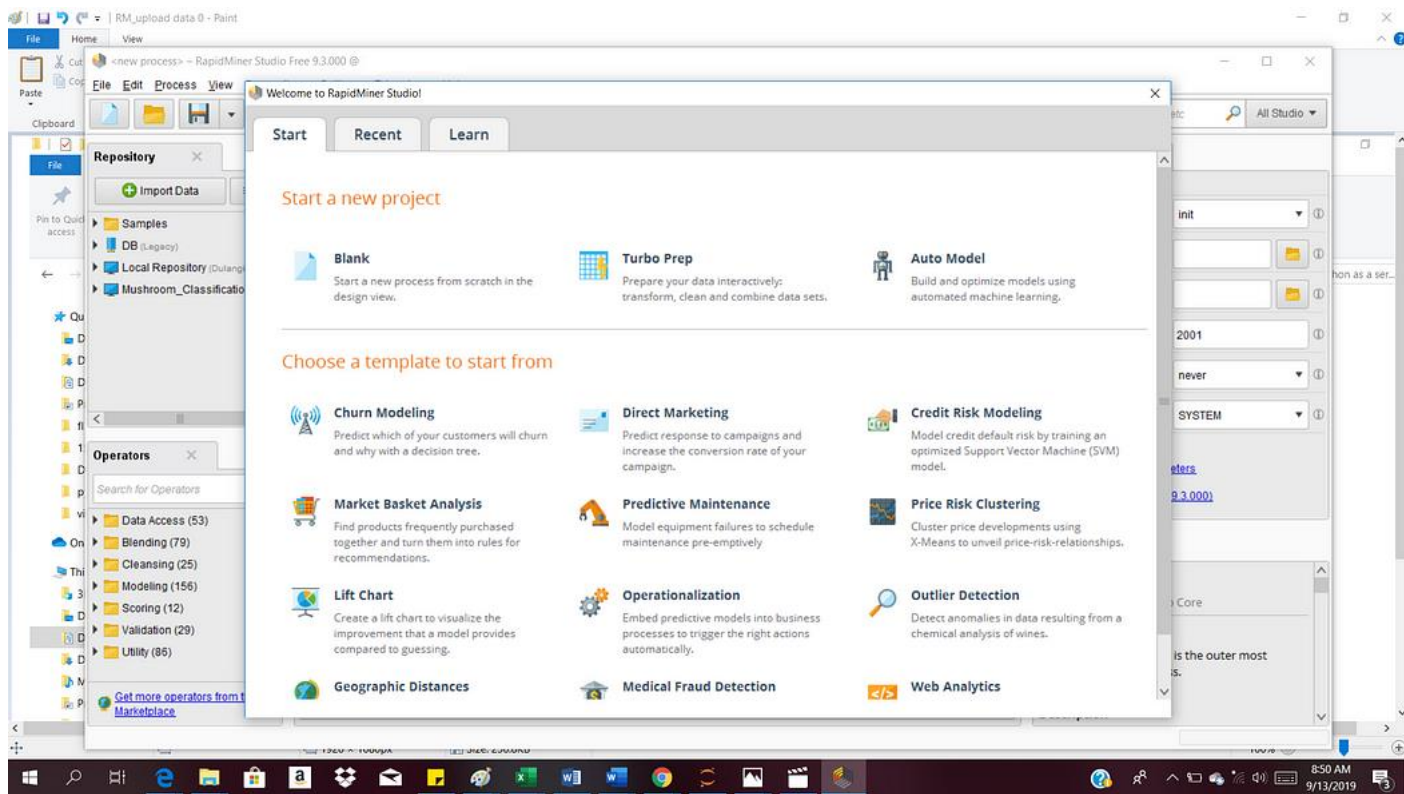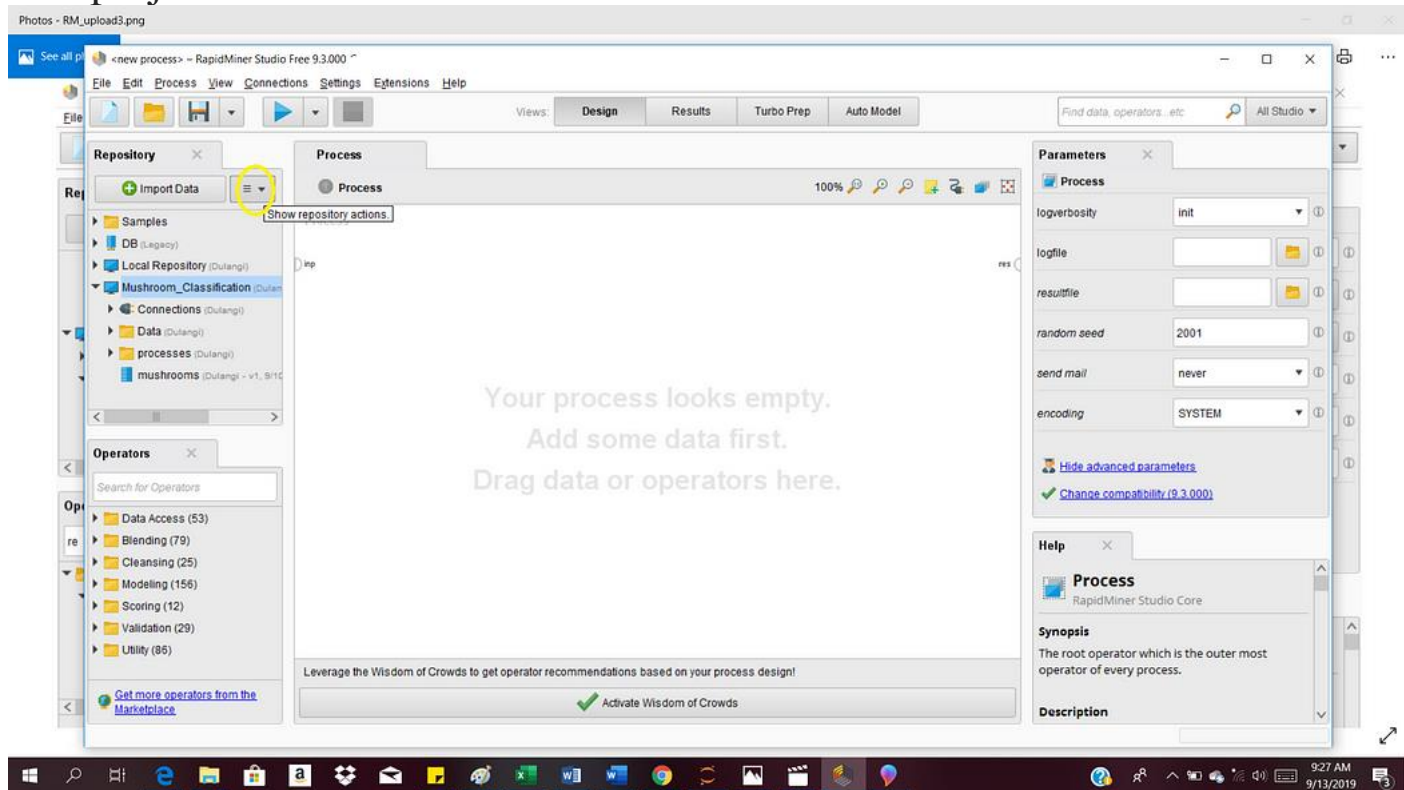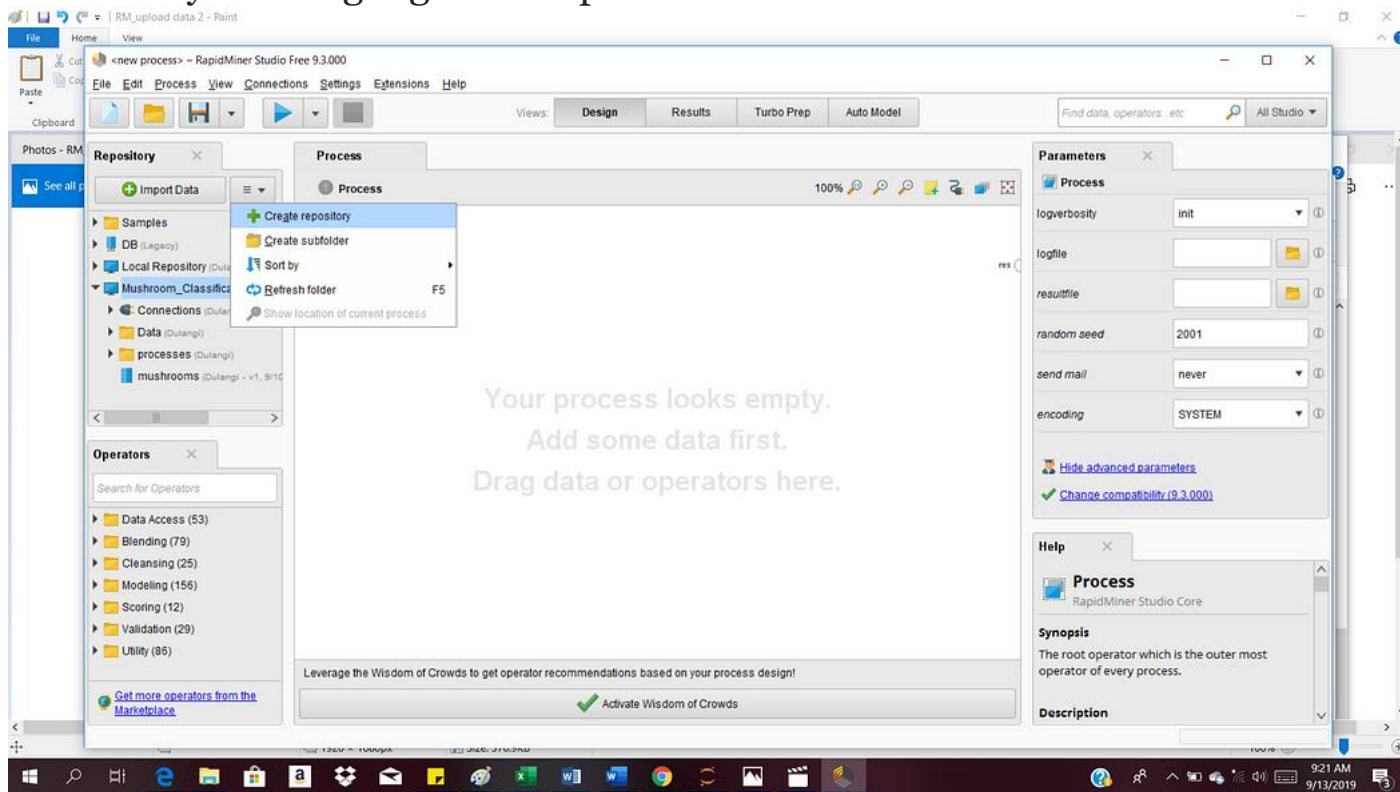# Importing a data set to Rapid Miner



When you first download RapidMiner using this url
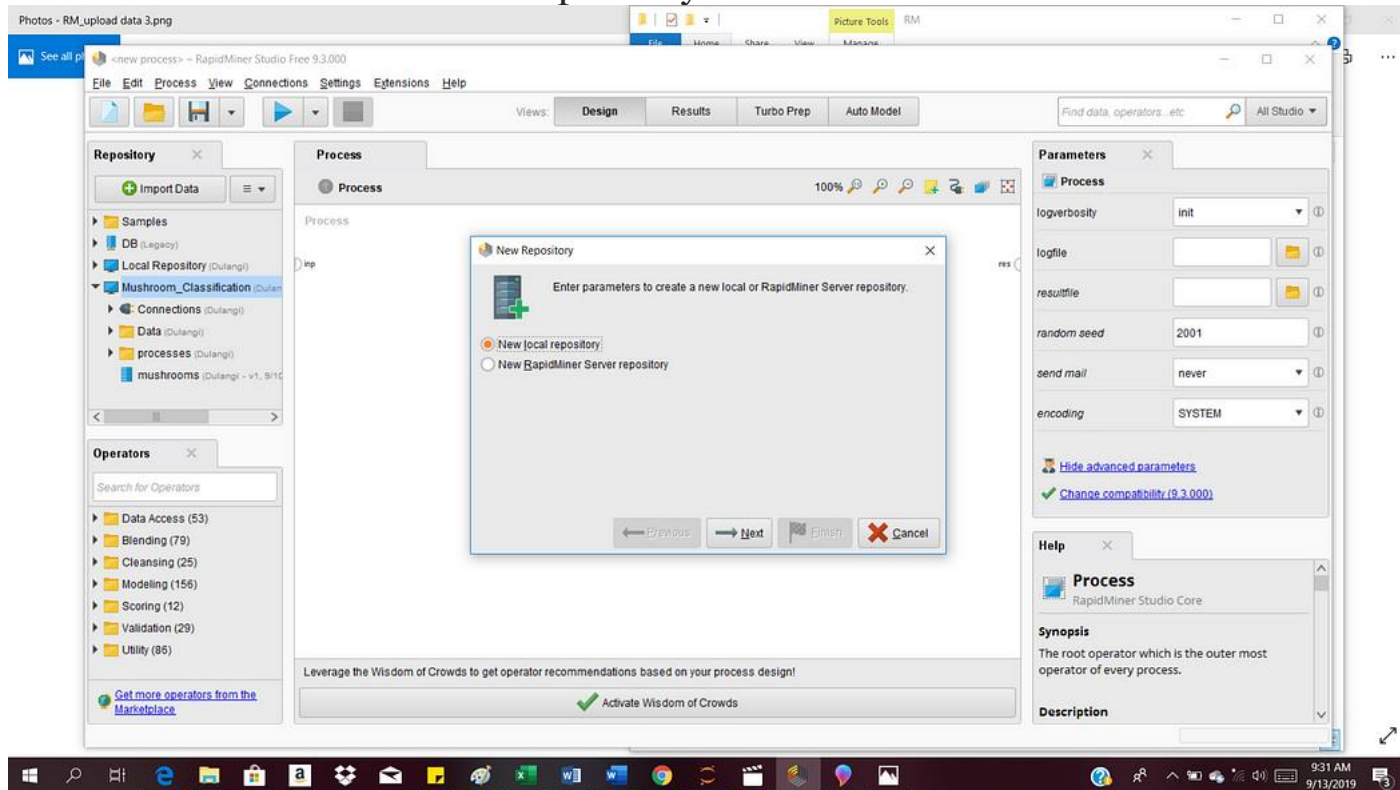:https://rapidminer.com/get-started/, and opening of RapidMiner .

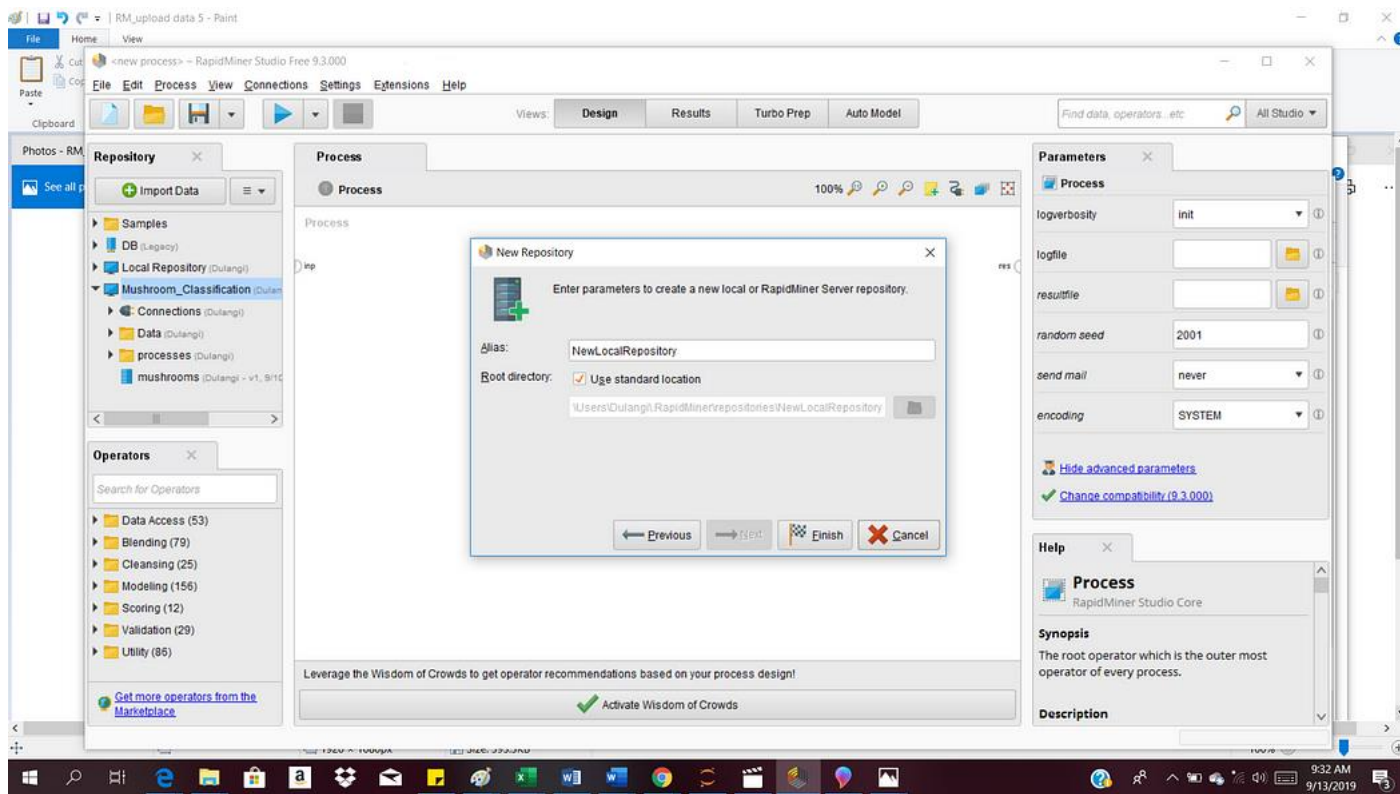Now when this page appears, select blank project, to start work on a new project.

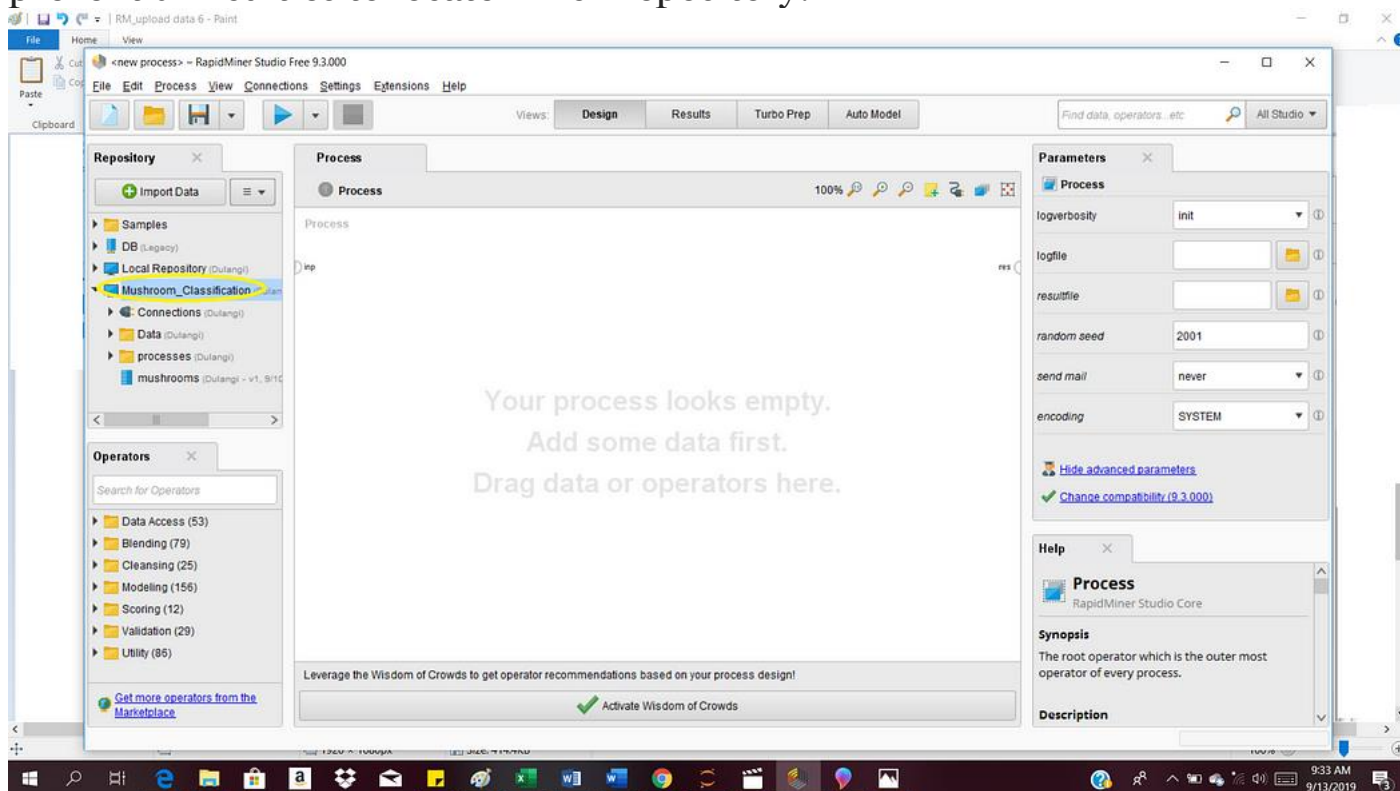Select the yellow highlighted drop down arrow.



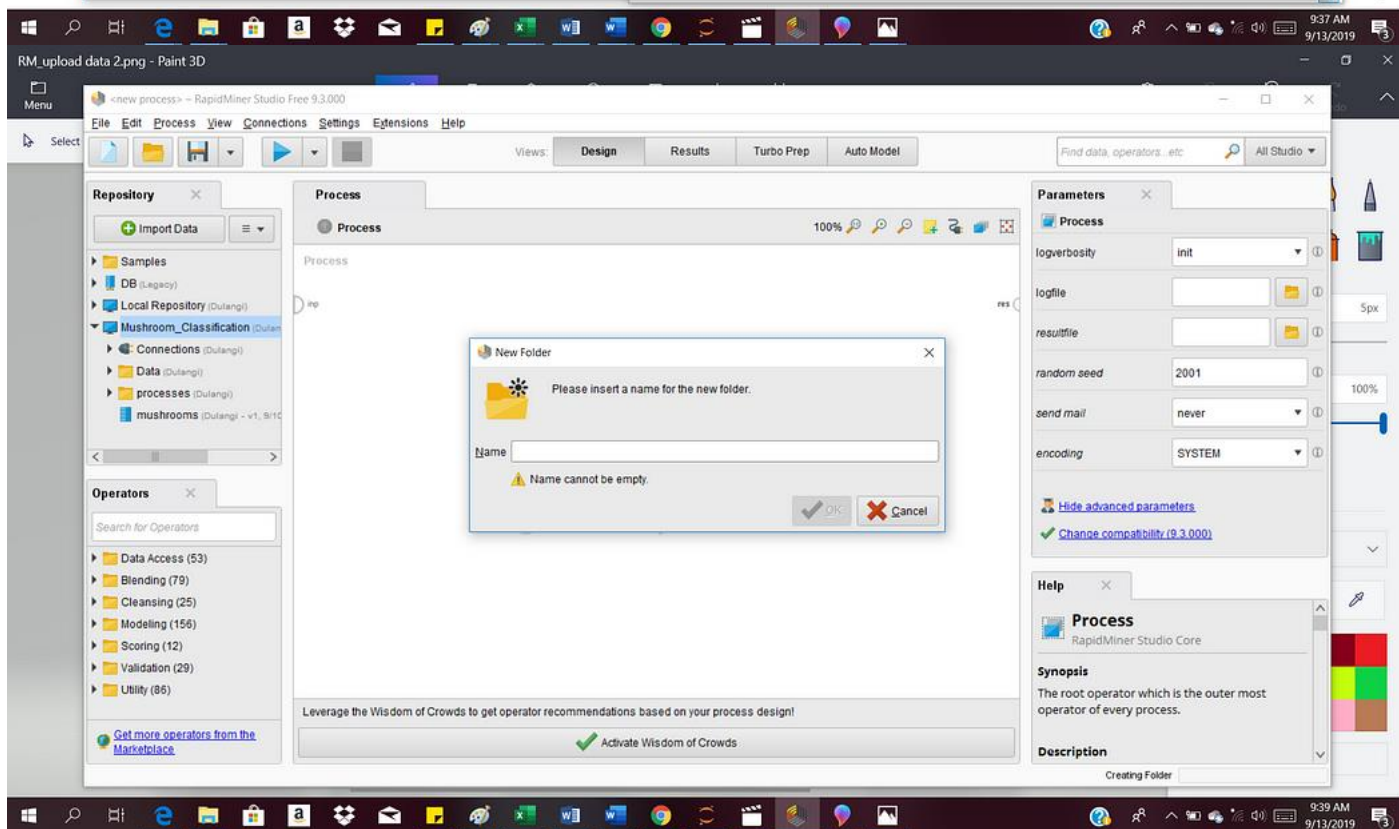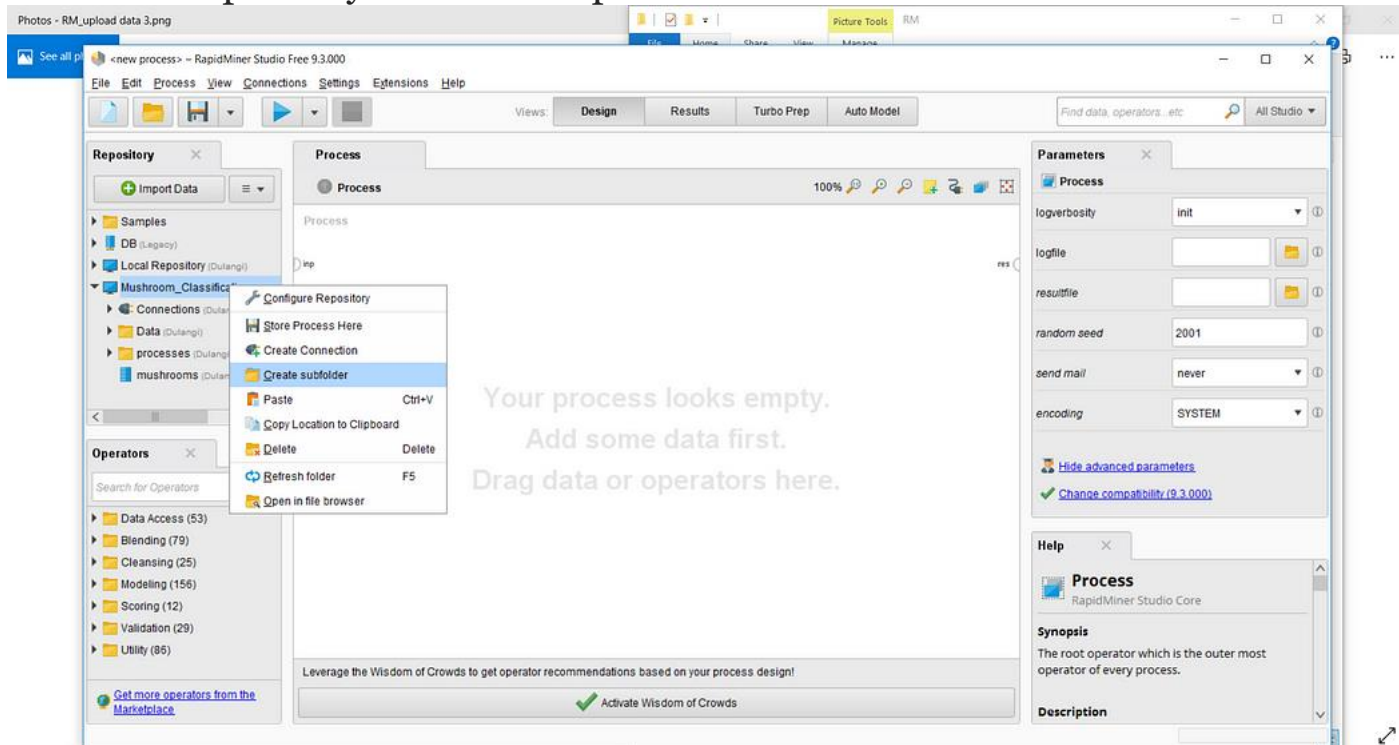First we have to create a new repository.

Give a name to the repository, its advisable give the same name as name of your data set, as you work with many projects later to prevent difficulties to locate which repository.
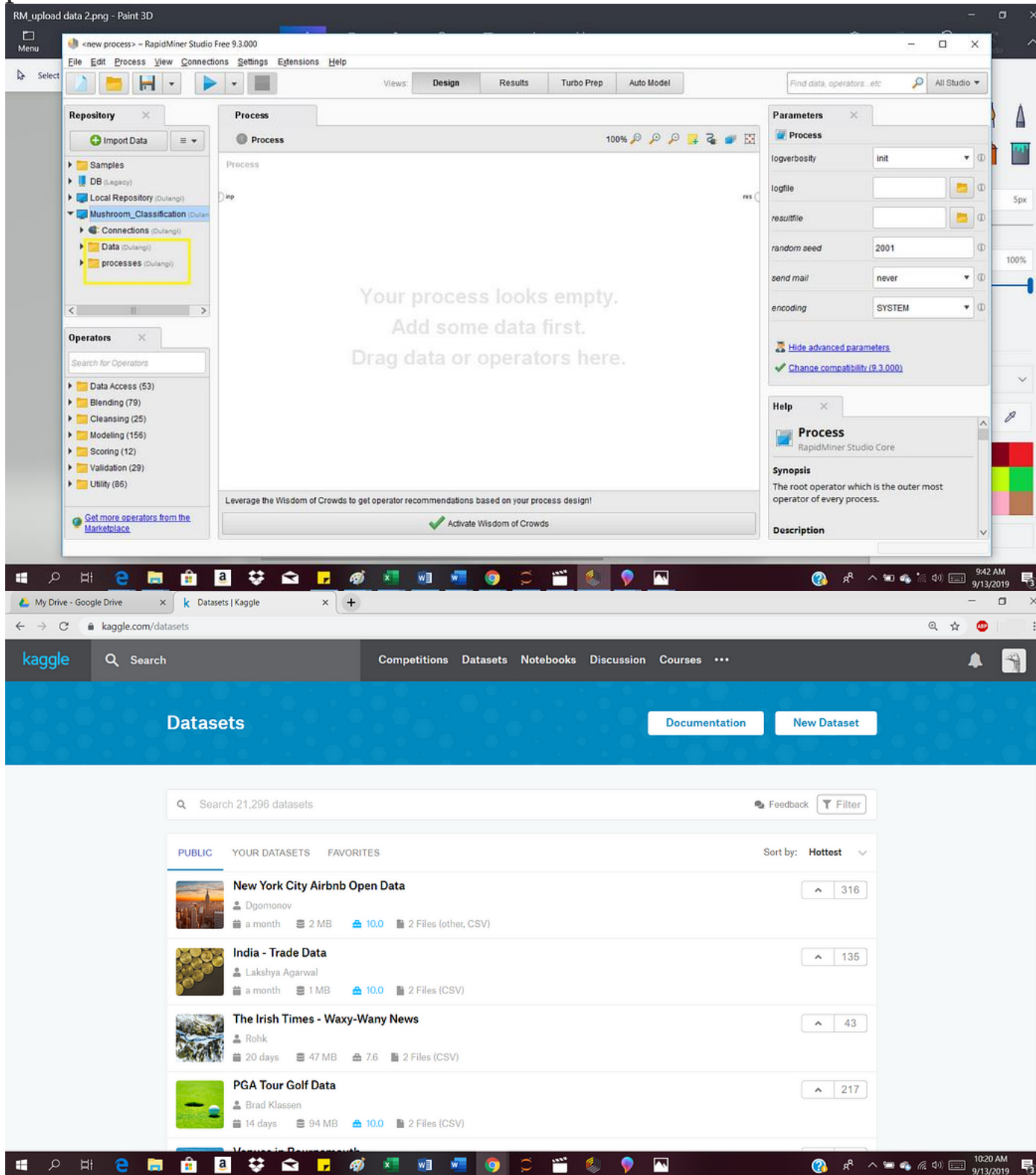
Now select the repository you created, and create two sub folders under the repository as Data and processes.
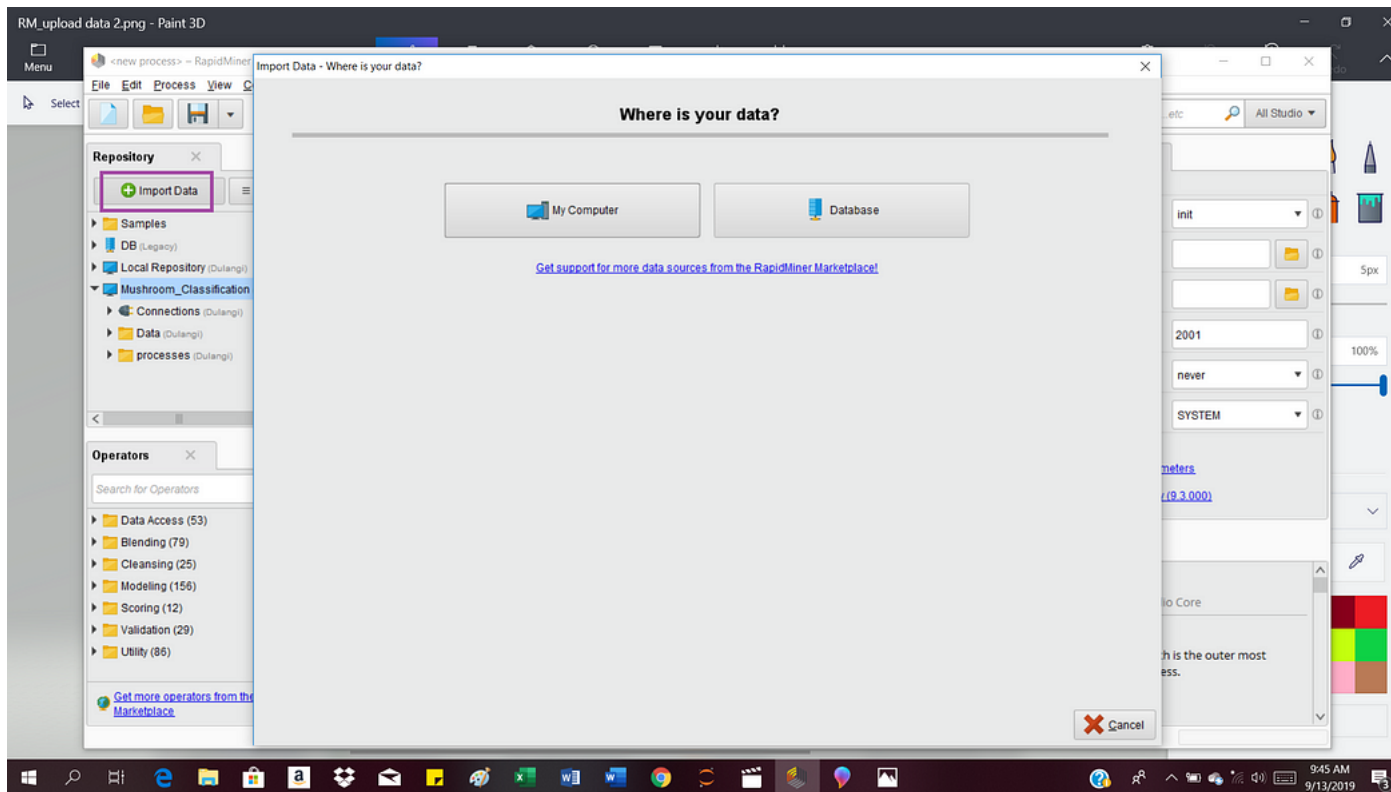
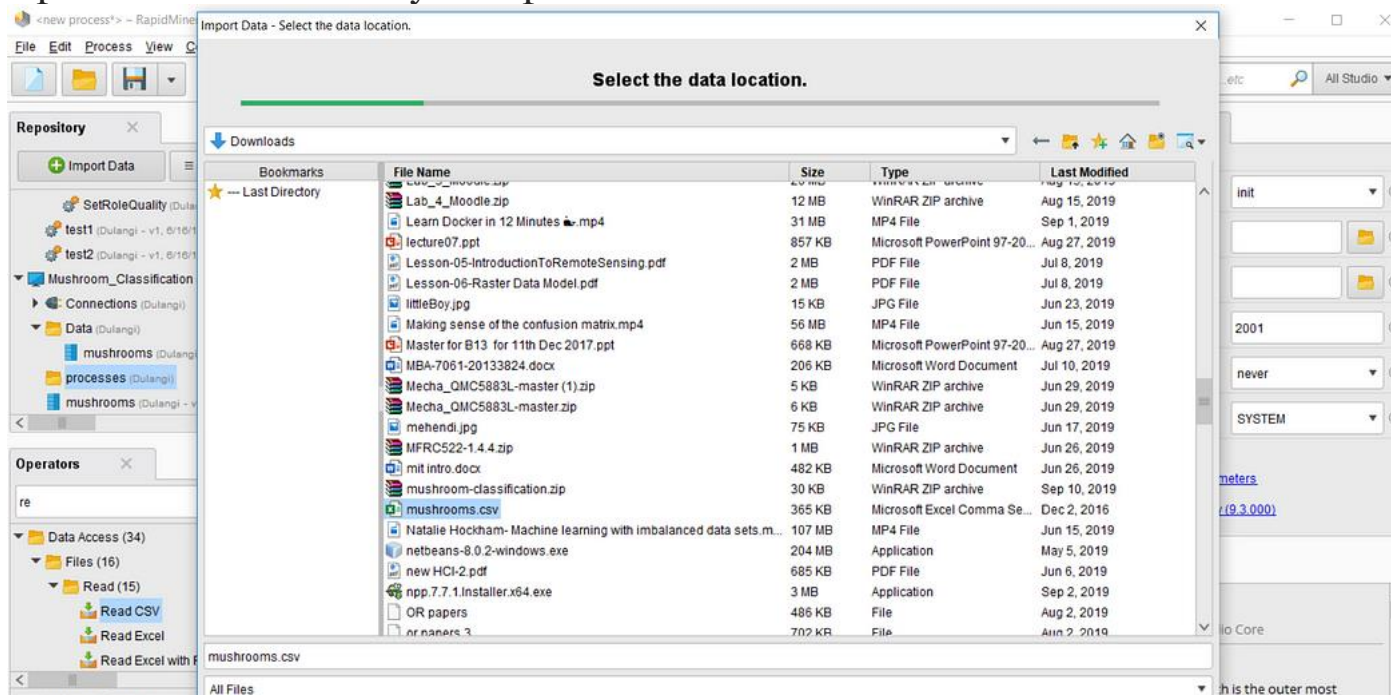Type in Data and repeat same method for creating sub folder for process.

Get a data set for free from Kaggle. A good data set has at least four columns, and at least 1000 rows. Select an appropriate data set, there are two generally used types of classification as:

1. *For supervised learning, we need to chose a data set with a label class, and algorithms such as : Linear Regression, Decision Tree, Random Forest, Naive Bayes. For example lets say we are studying the factors responsible for heart attack via a data set, and in the data set studied are gender, age, occupation, income levels, cholesterol levels in blood and so on, and also there is a class label, column called has got heart attack, hence we can train using this data, and predict for a number of other people having similar factors the probability of them getting a heart attack.*

2. *In unsupervised learning, we don't have a class label for the training data set, methods such as clustering is used. For example lets say the amount of students enrolling for a particular course module, there is no classifier in that as class label.*

Now store first select the purple highlighted box, Import Data, and normally the data downloaded from kaggle will be in Downloads as a zip file. Hence select My Computer. Extract files to download.

Select the mushroom.csv file.

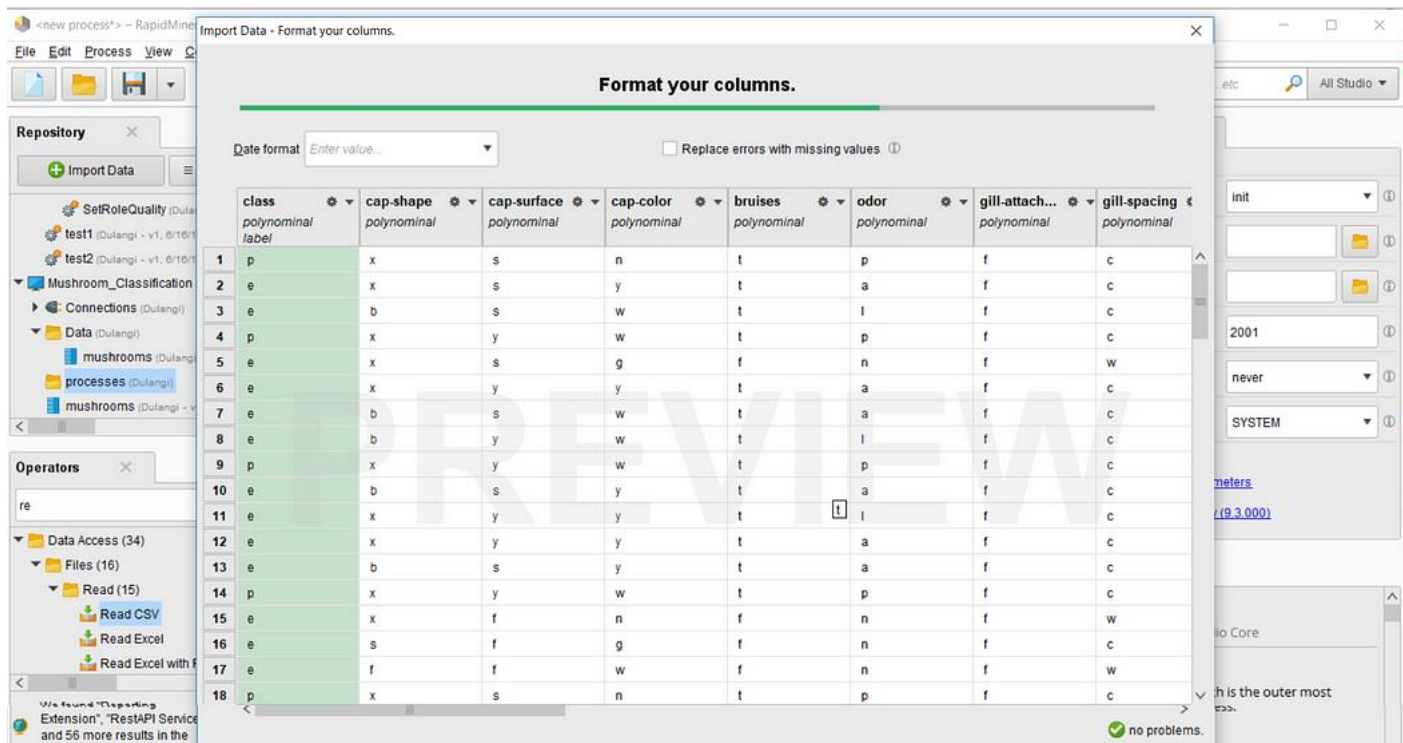If the project you are working on is a supervised classification, then define class label by selecting the green highlighted drop down arrow.



Chose label from drop down menu.

The class label column will be highlighted in light green as shown.



Now save the mushroom.csv file in Data sub folder. Do each of changes to the data such as remove duplicates, fill missing values

with mean, and store the output data files under Data sub folder. Each of the processes done on data store in process sub folder.