**Activity 1 –** Sentiment Analysis

In this activity, we will:
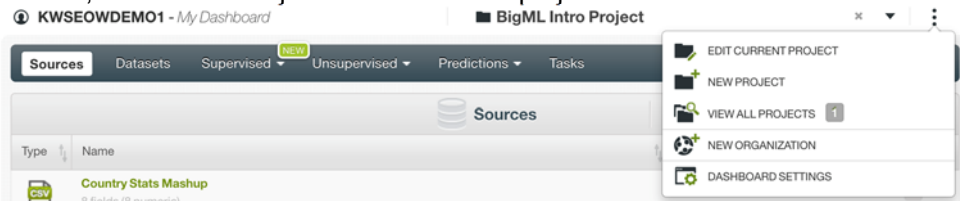❑ Create an account on BigML.com
❑ Create Project
❑ Update data source and create a dataset
❑ Create models for sentiments analysis
❑ Evaluate the model
❑ Perform predictions

1. **Setup account on BigML**
    a) Sign up for an account on https://bigml.com. The sign up process is pretty straightforward without the need for credit card information.
       * Free account allows for unlimited total tasks and storage, 16 MB dataset size, 2 parallel tasks and 1 user account with no support.
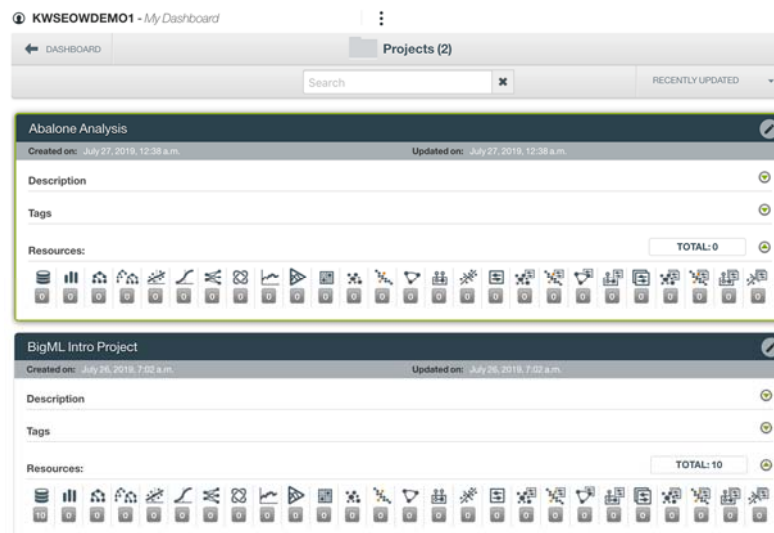
2. **Create a new Project**
    a) On the top right menu, select New Project to create a new project.



    b) Enter a project name and click Save.



    c) For new account, you should see your newly created project and the default **BigML intro Project** in your Dashboard view.
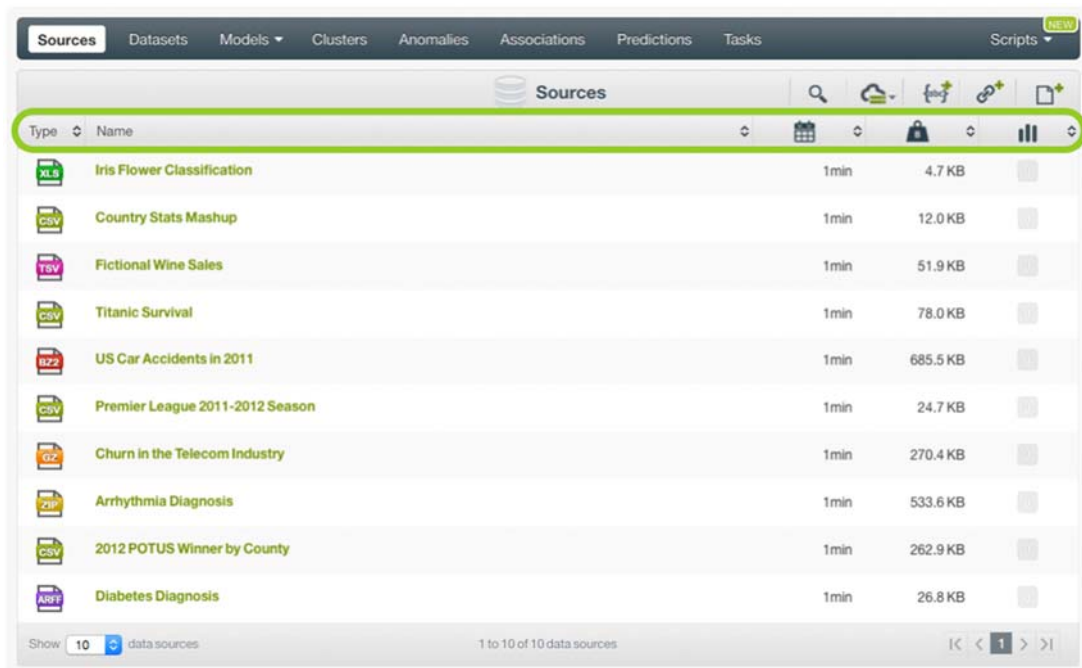


    d) Click on the project name **you used in step 2b** in the dashboard to continue with the next step.

3. **Uploading your data**
    a) The first step is to upload your data to your BigML account. BigML offers several ways to do it, you can drag and drop a local file, connect BigML to your cloud repository (e.g., S3 buckets) or copy and paste a URL. BigML automatically identifies the field types. Field types and other source parameters can alternatively be configured by clicking in the source configuration option.
    b) The first tab of the BigML Dashboard's main menu allows you to list all your available sources. You will see an empty list for your Abalone Analysis project. However, when you first create an account at BigML, you will find a list of promotional BigML sources in your **BigML intro Project**. In this source list view (see below), you will

see for each source, the **Type**, **Name**, **Age** (time since the BigML source was created), **Size**, and **Number of Datasets** that have been created using the BigML source.



c) On the top corner of the source list view, you can see the menu options show below.
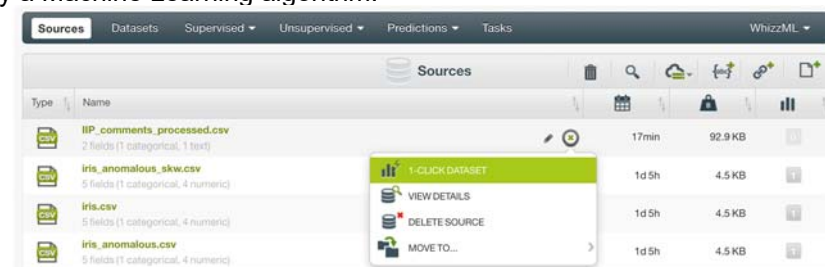


These menu options perform the following operations (from right to left):
1) **Create a source from a local source** opens a file dialog that helps you browse files in your local drives.
2) **Create a source from a URL** opens a modal window that helps you input the URL of that BigML will use to automatically download a remote source.
3) **Create a inline source** opens an editor where you can directly input or paste data into it.
4) **Cloud Storage DropDown** helps you browse through previously configured cloud storage providers.
5) **Search** searches your sources by name.

d) Download a copy of the dataset from your instructor from the github link provided by your instructor. The dataset we will be using is IIP_comments_processed.csv. The file consist of two columns, *sentiment* and *reviews*.

e) Click on **Create a source from a local source,** in the file selection dialog box, select the dataset csv file and click **Open**. The csv file will be uploaded to BigML and when completed, you should see the csv file listed on your **Sources** tab.

f) In BigML you can configure text parameters like the tokenization strategy, the case sensitivity, stemming or stop words using the Text Analysis options at the source configuration level in the Sources with the BigML Dashboard. These text configuration options define the vocabulary for all BigML models, including topic models. However, topic models are the only models that allow you to change this configuration at the topic model creation time. You can build several topic models using different text configurations without going back to the source and having to create different datasets.

Among the text configuration options for topic models you can find parameters like the language, the tokenization, whether to include or exclude stop words, the maximum size for the n-grams to be considered in your model vocabulary, the stemming, and the case sensitivity. You can also define the terms that you want to completely exclude from your model..

## 4. Create a Dataset

a) From your source view, use the **1-click dataset** option to create a dataset, a structured version of your data ready to be used by a Machine Learning algorithm.
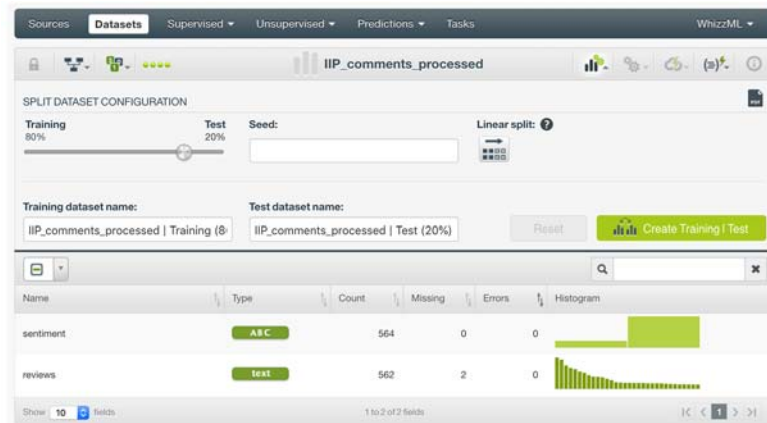


b) In the **dataset** tab, click on the newly created dataset, you will be able to see a summary of your field values, univariate statistics, and the field histograms to analyse your data distributions. This view is really useful to see any errors or irregularities in your data. You can also filter the dataset by several criteria and create new fields using different pre-defined operations as needed. However, for our simple reviews dataset, you will not be able to see too much.
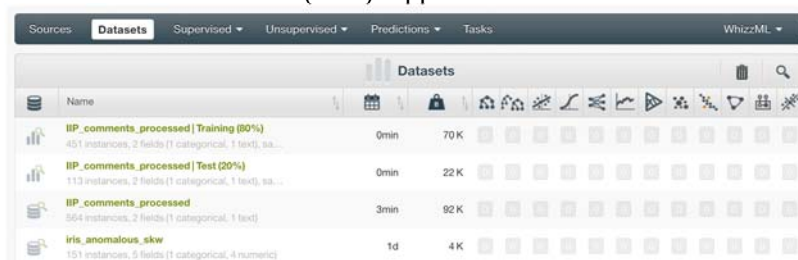


c) Once your data is clean and free of errors you can split your dataset into two different subsets: one for training your model, and the other for testing. It is crucial to train and evaluate your model with different data to ensure it generalizes well against unseen data. You can easily split your dataset using the BigML 1-click option, which randomly sets aside 80% of the instances for training and 20% for testing.

d) Can you change the split ratio. But for our activity, we will just keep it to 80-20 which is the norm. Click on **Create Training | Test**
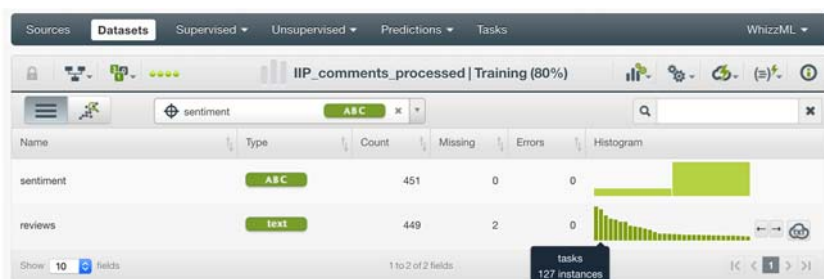


e) After the creation, return to the Datasets tab, you will see two new datasets created. One with "Training (80%)" appended to the dataset name and one "Test (20%)" appended.
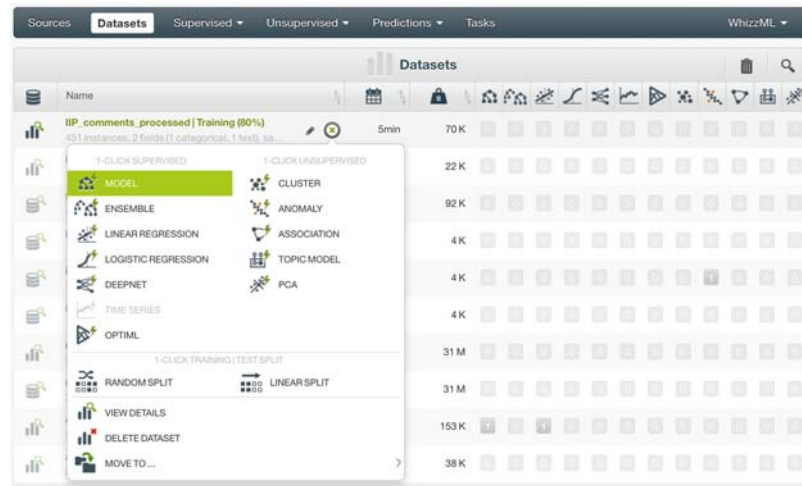


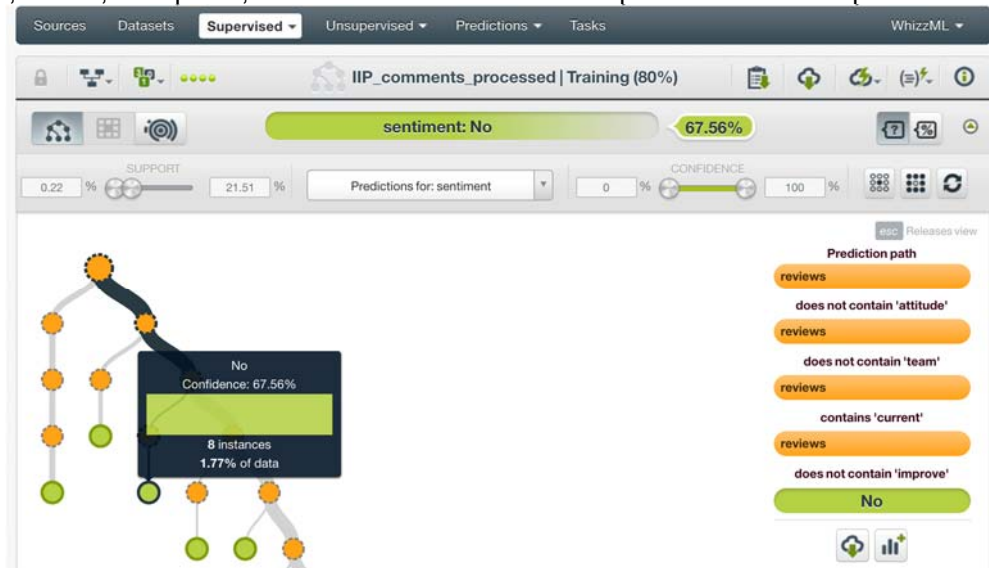## 5. Create a Model (Single Decision Tree)

a) To get a first sense of whether the text in the reviews have any power to predict the sentiment, we are going to build a single decision tree by selecting the "sentiment" as the objective field and using the reviews as input

b) When you are done creating your dataset, you can see that the reviews dataset is composed of two fields: sentiment (Yes/positive or No/negative) and reviews (the review text). Not surprisingly, the words "learn", "tasks" and "attitude" are the most frequent ones in the collection (see image below). Mouse over the histogram to see the summary.



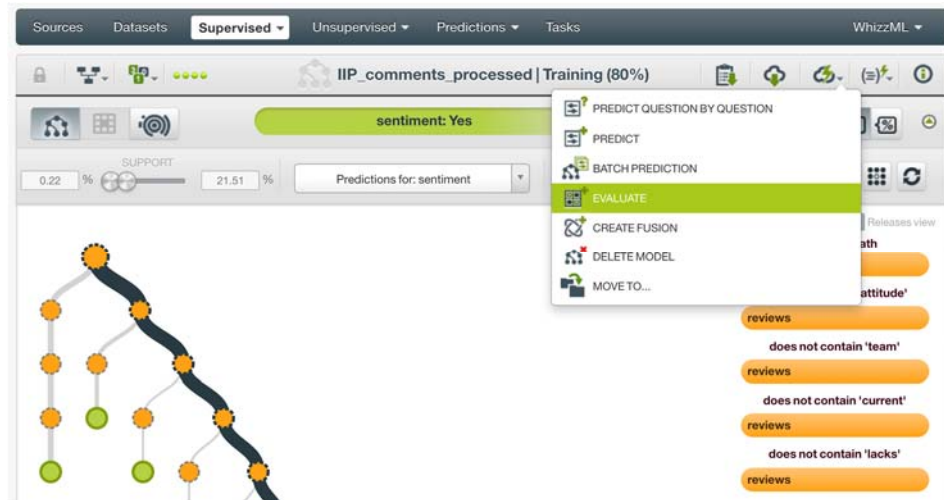c) Perform a 1-click model by clicking on the dropdown menu icon and select **MODELS**.

d) Once completed (it should take less then a minute), you can mouse over the root node and go down the tree until you reach the leaf nodes, you can see the different prediction paths. For example, in the image below you can see that if the review does not contains the terms "attitude", "team", "improve" and does contain "current" then it is a "negative" review with 67.56% of confidence. As expected, we can find words in the nodes such as "lacks", "fast", "learn", "complete", which are the terms that best split the data to best predict sentiment.
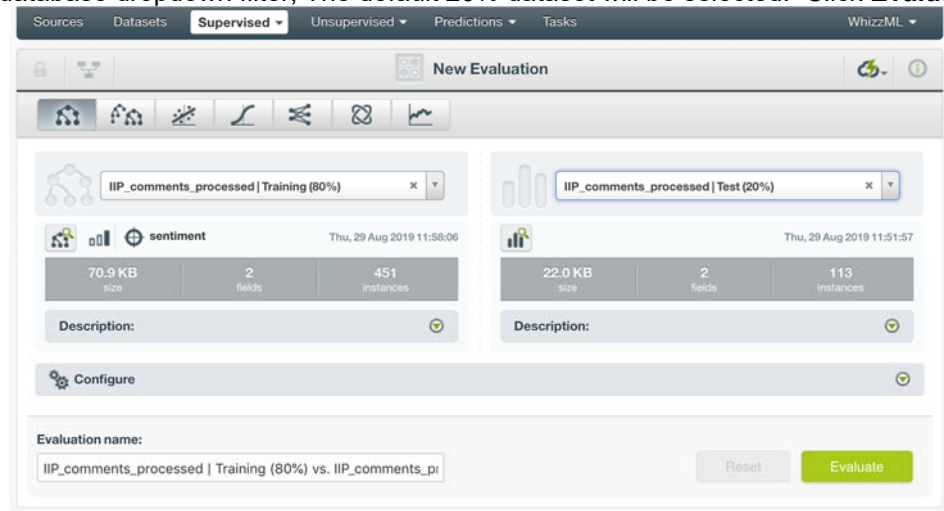


## 6. Evaluate the Model

Confidence values seem pretty high for most prediction paths in this tree, but to measure its predictive power we need to evaluate it by using data it has not seen before. We use the previously set aside 50% of the dataset, which contains the remaining 25,000 movie reviews that have not been used to train our model.
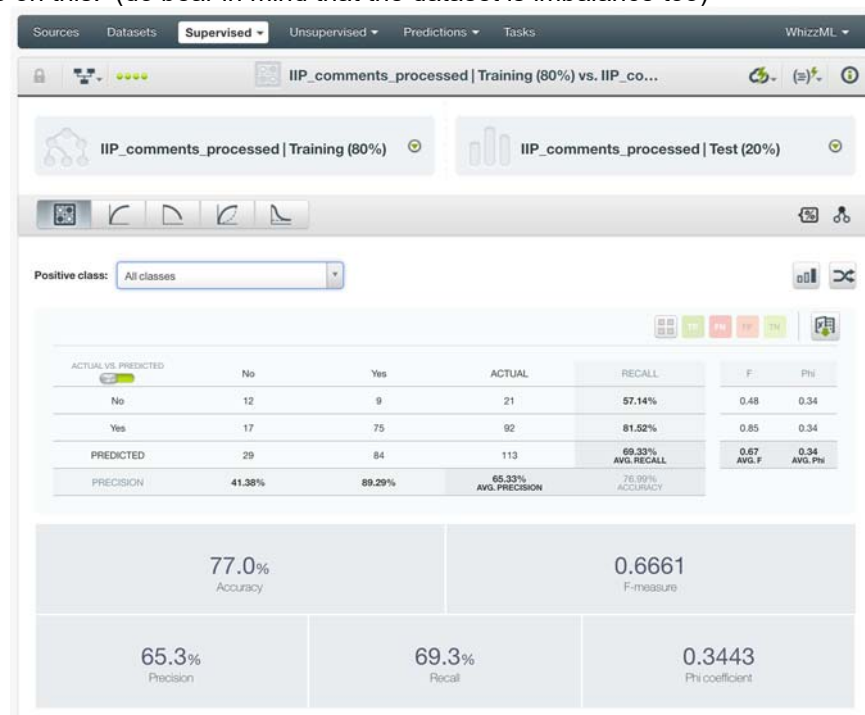
a) Navigate to the detail view of the model. Click on **Supervised**, select **MODELS** and click on the model you created. In the 1-click menu, select **EVALUATION**.

b) In the search database dropdown filter, The default 20% dataset will be selected.  Click **Evaluate**.
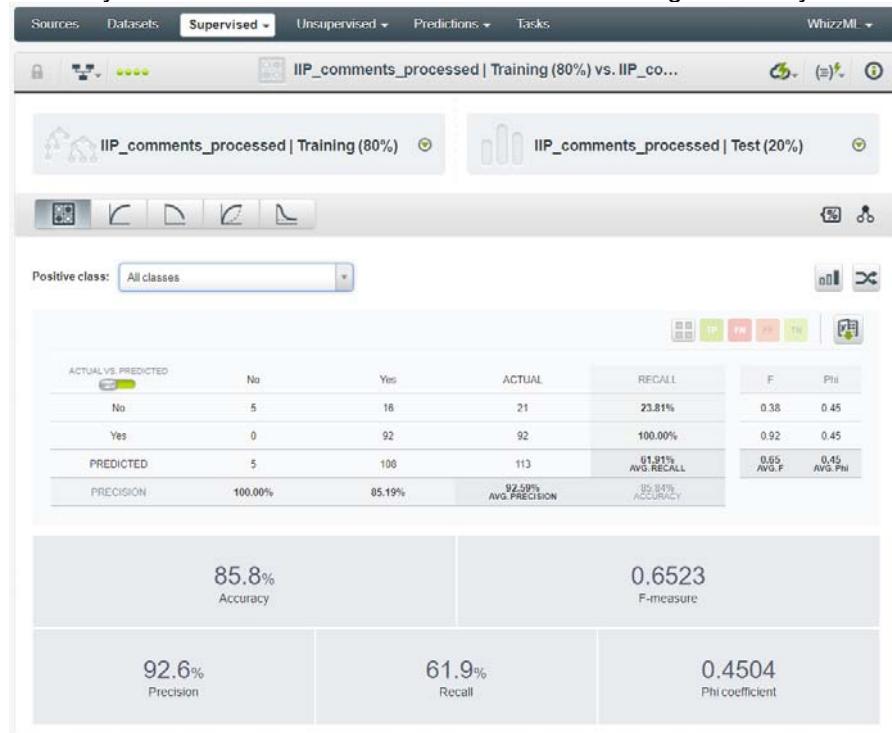


c) The evaluation yields an overall accuracy of 77.0%, which is not that bad for a 1-click model. However, we can definitely improve on this!  (do bear in mind that the dataset is imbalance too)
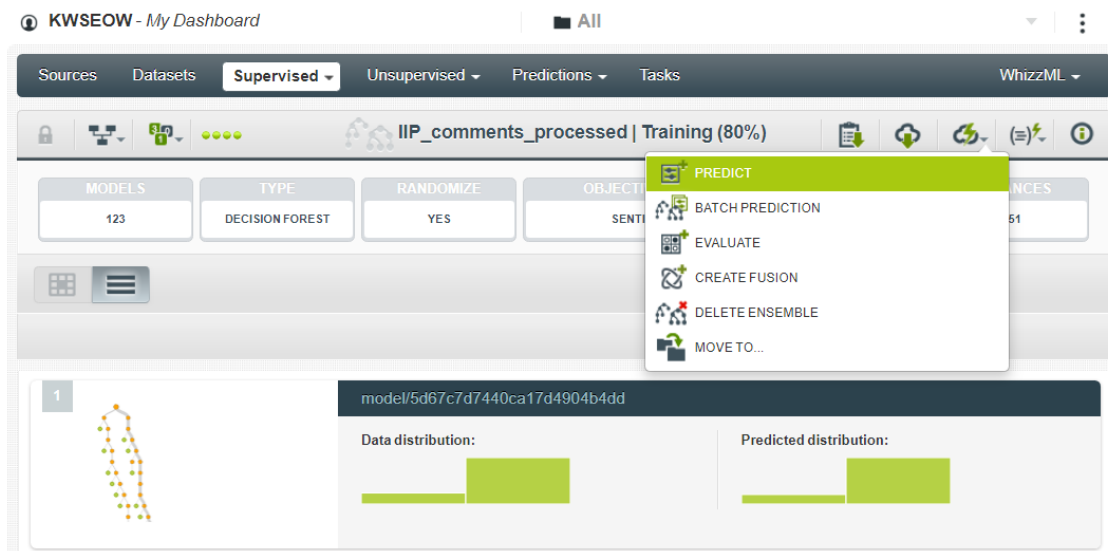
## 7. Ensemble

Go ahead to try out **Ensemble 1 click model** and check the accuracy you can get. It may take a long time 1-2 hours to build the model.

On my system, the accuracy of the default ensemble model achieves average accuracy of 85.8%.
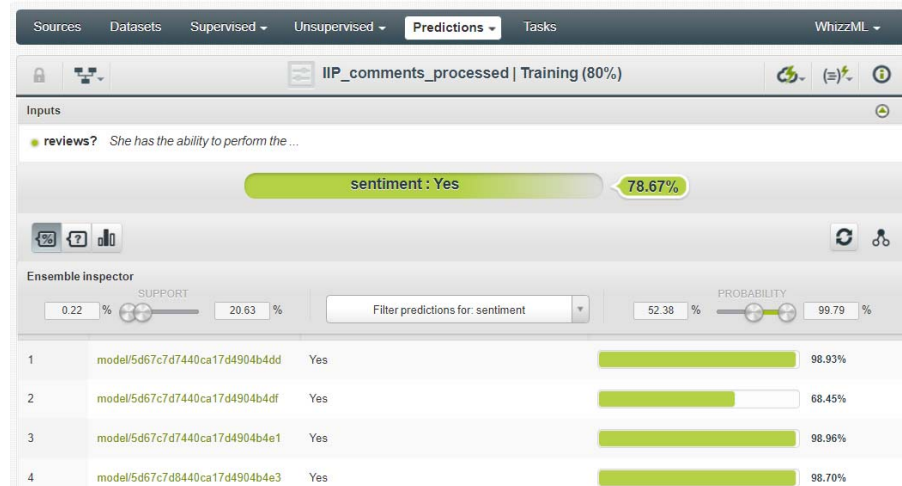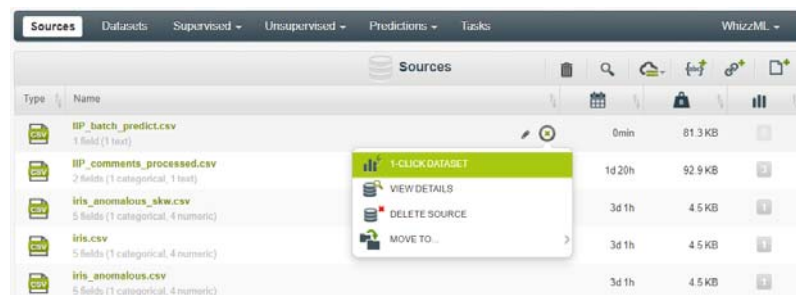


## 8. Make Predictions

a) Once you are happy with the model, you can perform a prediction. Click on **Predict**.



b) A form containing all your input fields will be displayed and you will be able to set the values for a new instance. At the top of the view, you will see the objective field prediction changing as you change your input field values, in our case, the reviews.
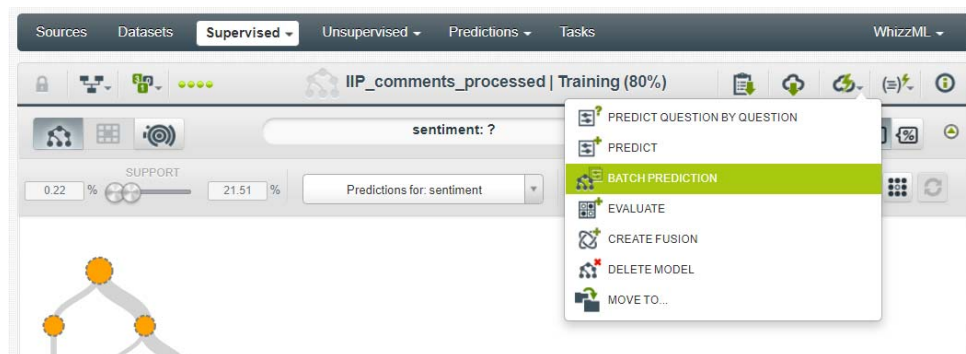
c) **Batch predictions** - Use the Batch Prediction option in the 1-click menu and select the dataset containing the instances for which you want to know the objective field value. However, before you do this, create a new source using the IIP_batch_predict.csv and then create a new dataset from this source
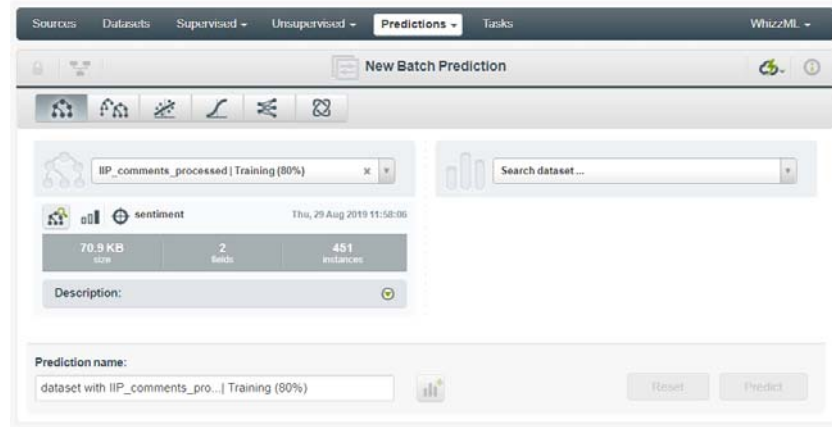


d) You can configure several parameters of your batch prediction such as the option to include both confidence interval and prediction interval in the batch prediction output dataset and file. When your batch prediction finishes you will be able to download the CSV file and see the output dataset.
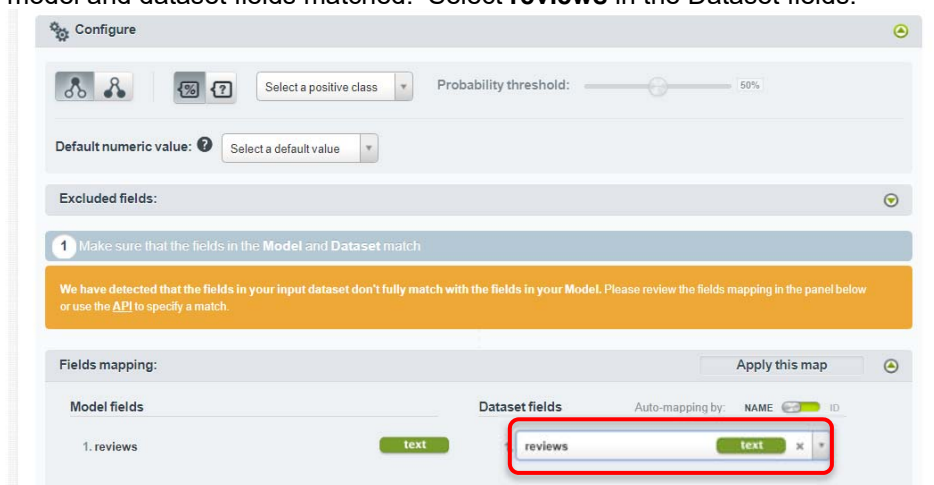
Select the model you want to use for the batch prediction. In our case, click on **Supervised**, **Models**, **Batch Predictions**.
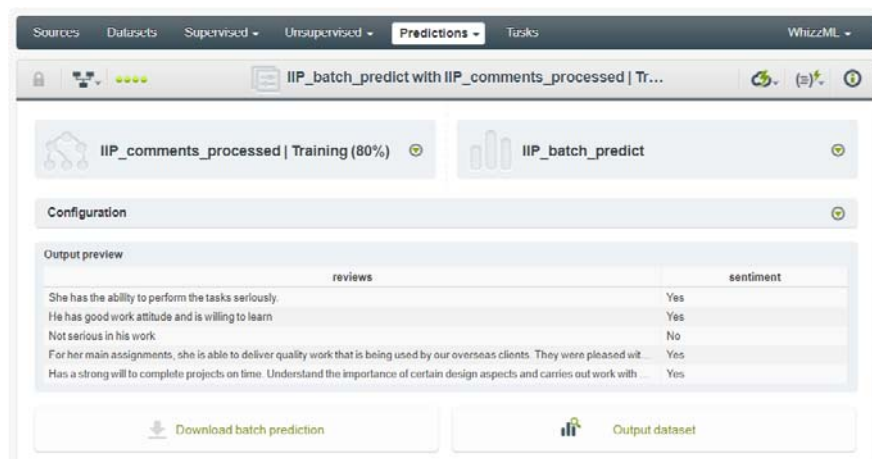


e) In the **Search Dataset…** field, search and select your dataset for batch prediction.

f) Check that the model and dataset fields matched.  Select **reviews** in the Dataset fields.



g) Click on **Predict** to start the prediction.  Usually it will just take a moment to complete.
h) Once the prediction is done, you can check out the result on the **Output preview** section or **download batch prediction** if you want to keep a copy of all the prediction.



---

Activity wrap-up:

We learn how to:
- ❑ Create an account on BigML.com
- ❑ Create Project
- ❑ Update data source and create a dataset
- ❑ Create models for sentiments analysis
- ❑ Evaluate the model
- ❑ Perform predictions