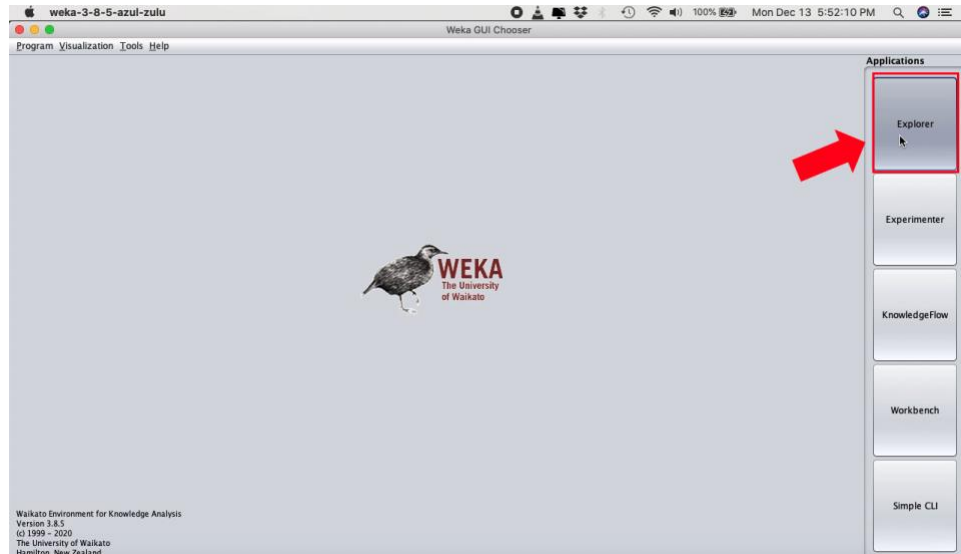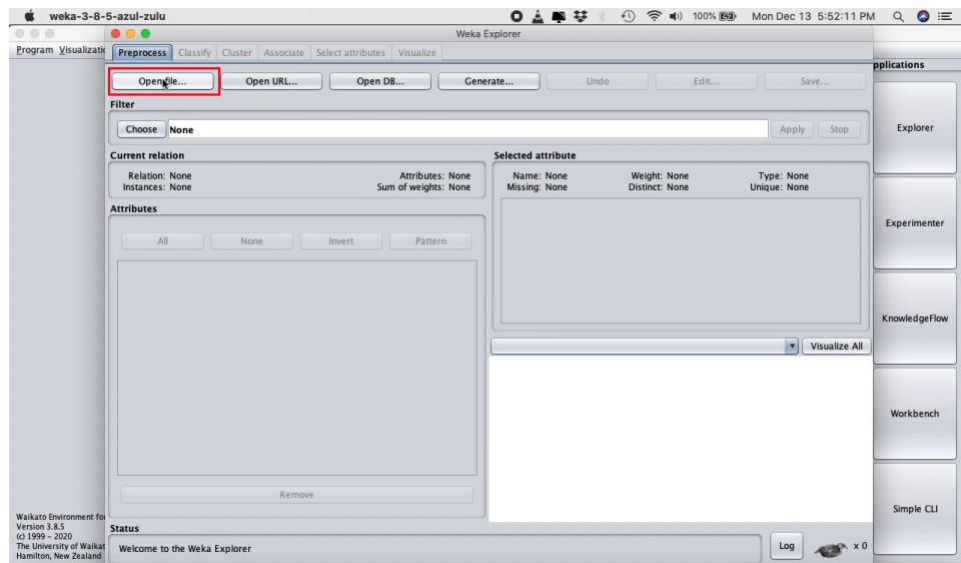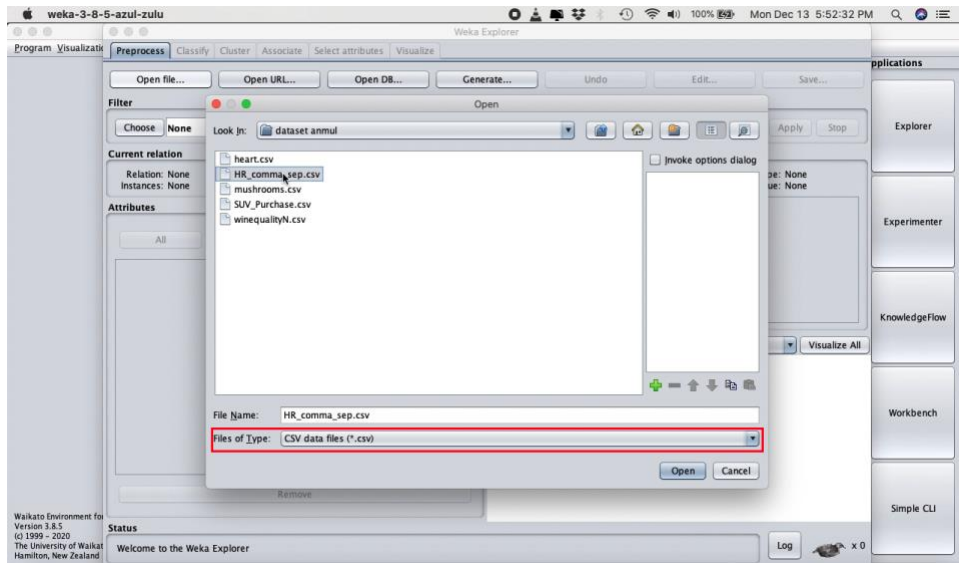**Data Processing using K-Nearest Neighbor (KNN) in WEKA**

In this part, Waikato Environment for Knowledge Analysis (WEKA) v.3.8.5 is used to create a KNN model. First, we start by choosing the Explorer menu.
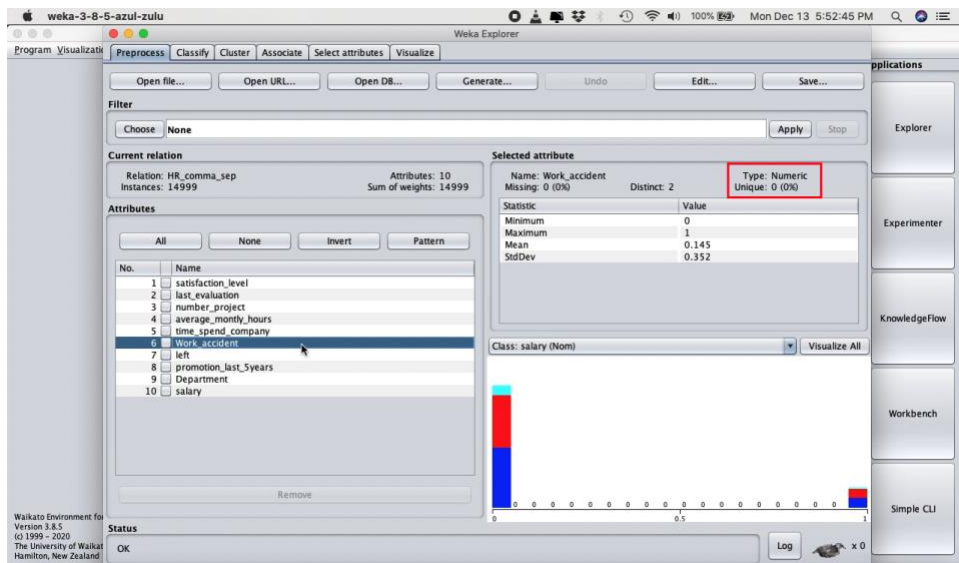


Next we will have Weka Explorer window as shown as below. Choose Open file to pick the dataset that we are going to use.
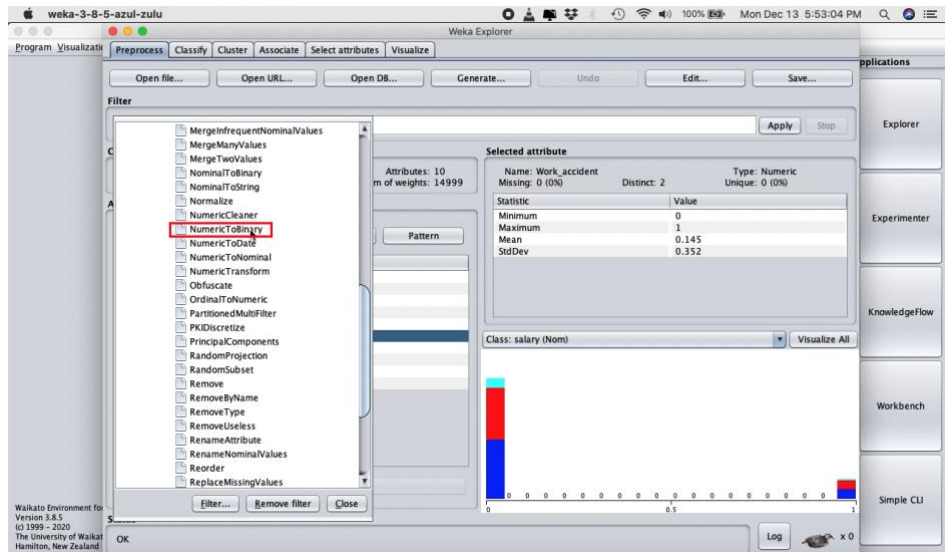


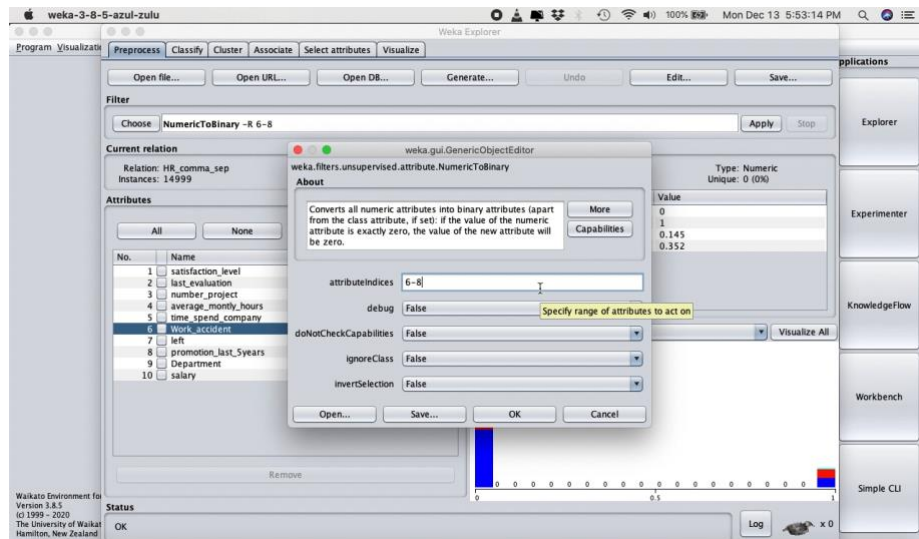Change the type of files from ARFF into CSV, which is the format of the dataset we use.

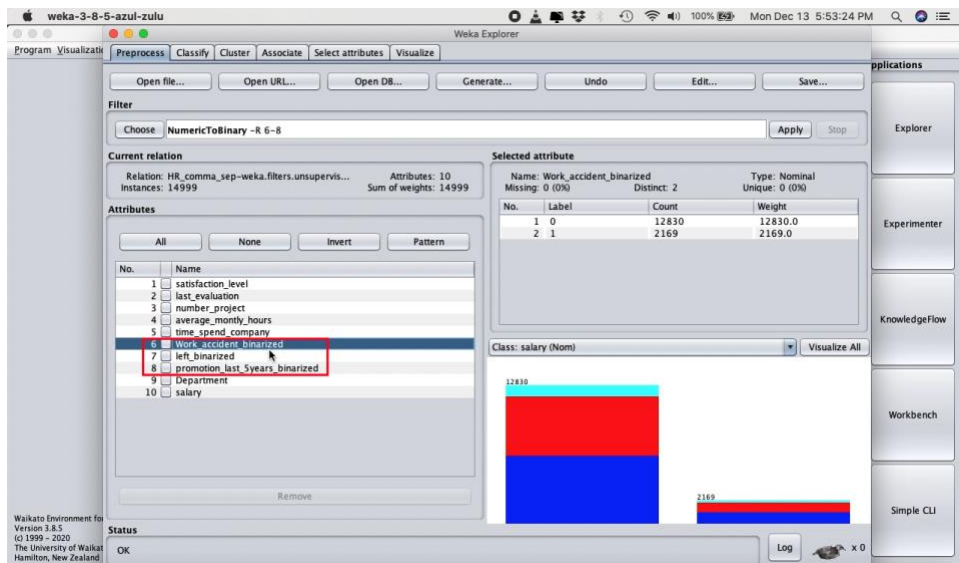Now the window will show the data with its ten attributes.



The program automatically detects categorical variable as numerical variable since it uses 0-1 scale. To change the attribute of these variables into binary, we use a NumericToBinary filter.
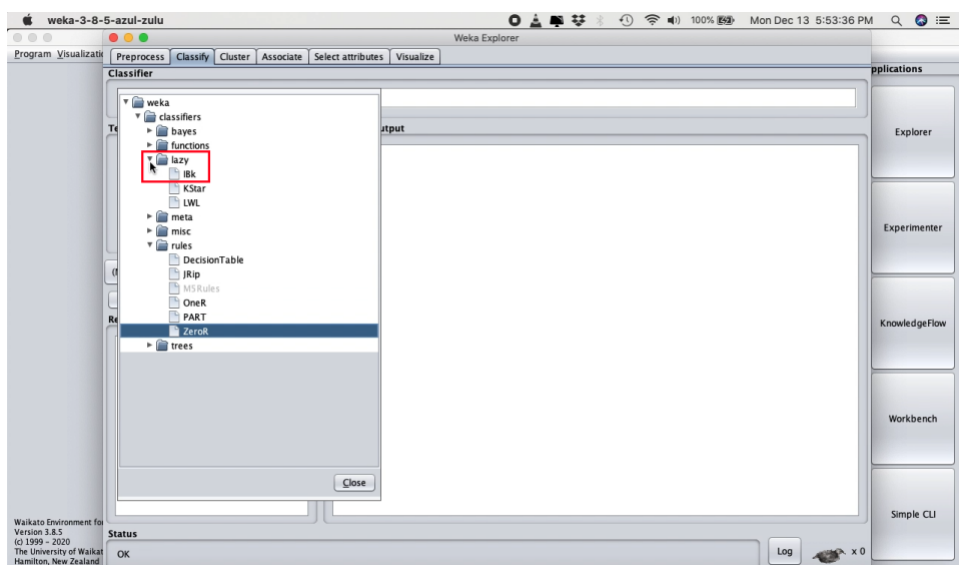
Put the number of the variables we are going to change. In this case *work_accident*, *left*, dan *promotion_last_5years* are marked as the 6[th], 7[th], and 8[th] attribute, so we put the number 6-8. Lastly, choose Apply in the Filter tab.



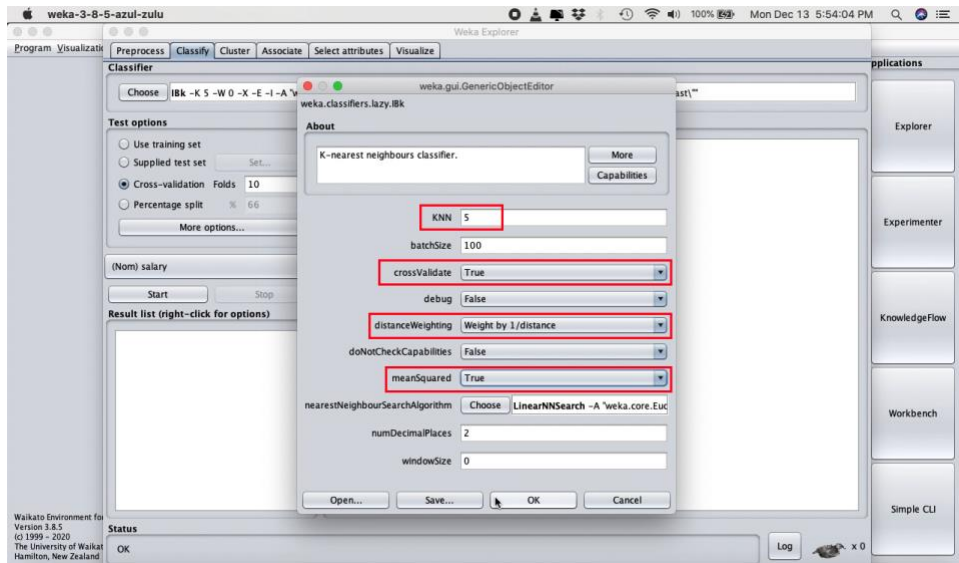The variables *work_accident*, *left*, and *promotion_last_5years* now have binarized labels behind their names.

Now we are ready to build a KNN model. In Classify tab, choose Classifier menu and pick IBk (*Instance Based Learner*), which is another name of KNN model in WEKA.



In this study we set the number of *k* to 5, set crossValidate to *True*, and picked weighted voting for distanceWeighting calculation.

Next we determine the response variable. In this case the response variable is *left* which determines whether an employee would leave the company or not. Then we can start the model building by choosing Start.



The picture below shows the output given after the model is done.

## Result

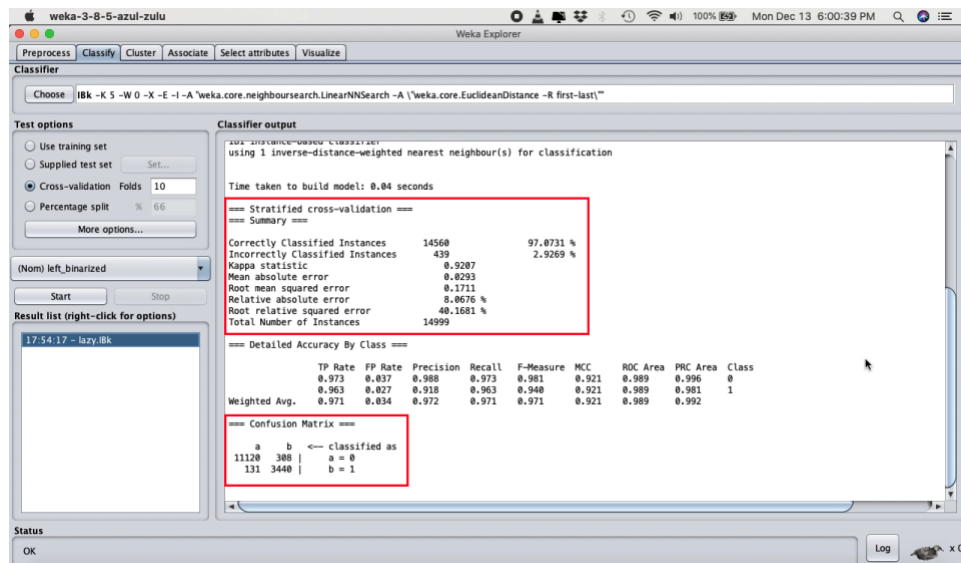From the figure it can be seen that the model successfully predicts employee turnover correctly with 97.07% accuracy. In the confusion matrix, it can be seen that the true positive, false positive, true negative, and false negative values are 3440, 308, 11120, and 131 respectively. The Kappa statistical value is a measure that compares the accuracy of observations with the expected accuracy that is formulated randomly. Landis & Koch (1977) said that the Kappa value in the range of 0.81-1 shows an almost perfect match between expectations and observations, meaning that the model is very good at predicting employee turnover.

```
Time taken to build model: 0.04 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances         14560              97.0731 %
Incorrectly Classified Instances         439               2.9269 %
Kappa statistic                          0.9207
Mean absolute error                      0.0293
Root mean squared error                  0.1711
Relative absolute error                  8.0676 %
Root relative squared error             40.1681 %
Total Number of Instances              14999

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 0.973    0.037    0.988      0.973   0.981      0.921  0.989     0.996     0
                 0.963    0.027    0.918      0.963   0.940      0.921  0.989     0.981     1
Weighted Avg.    0.971    0.034    0.972      0.971   0.971      0.921  0.989     0.992

=== Confusion Matrix ===

     a     b    <-- classified as
 11120   308 |    a = 0
   131  3440 |    b = 1
```