# DATA MINING- PROJECT

Submitted to: DR. MD. MAHBUB CHOWDHURY MISHU

Submitted By: MD. SAMIUL HAUQUE  CHOWDHURY

MARCH 10, 2021

Section: A

**Project Title:** Application of KNN to a data-set on weka and analyzing the accuracy by applying different algorithms.
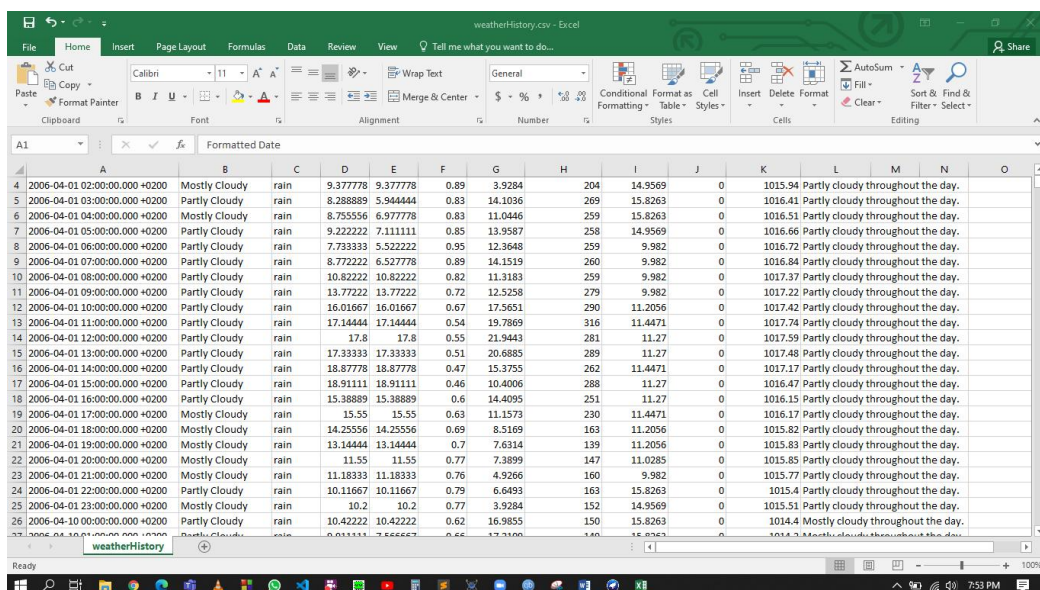
## Introduction:

One of the most common data mining concept is classification techniques. The objective of classification is to specifically predict the target class for each case in the data. In this project, I used a widely utilized classification procedure called **"K-nearest neighbors (KNN)"** which is also known as lazy learner algorithm. By utilizing K-NN, it will recognize the category or class of a specific dataset. In this project, I used a weather data set and name of the data set is therHistory.csv**"**. It has some data records of an area that recorded temperature, apparent temperature, humidity, wind speed, wind bearing, visibility, and pressure type data and by analyzing we need to predict next day weather.

## Methods of KNN:

1. Assigning the value of K in KNN   K=1, 3, 5... (I used K=1)
2. Now Applying Euclidean distance d= sqrt(x2-x1)$^2$-(y2-y1)$^2$
3. Now we need to sort the values of distances.
4. Now we need to assign the new data points to that category for which   number of the data the number of neighbor is maximum.
5. Among these k neighbors, we need to count the points in each category.

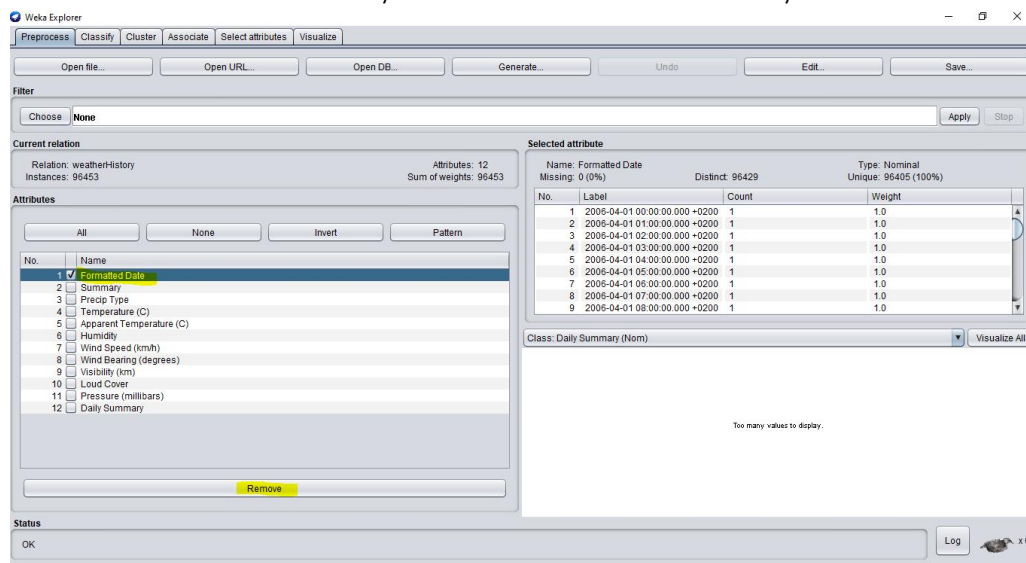## DataSet:

## Dataset Description:

Name of this dataset is "weatherHistory.csv". This dataset contain 12 attributes and 96453 instances. It has 12 columns and 96454 rows.

## Procedure in weka:

First I need to upload my dataset to weka from the open file section. Now next procedure is to delete less important attribute to increase significant of accuracy.
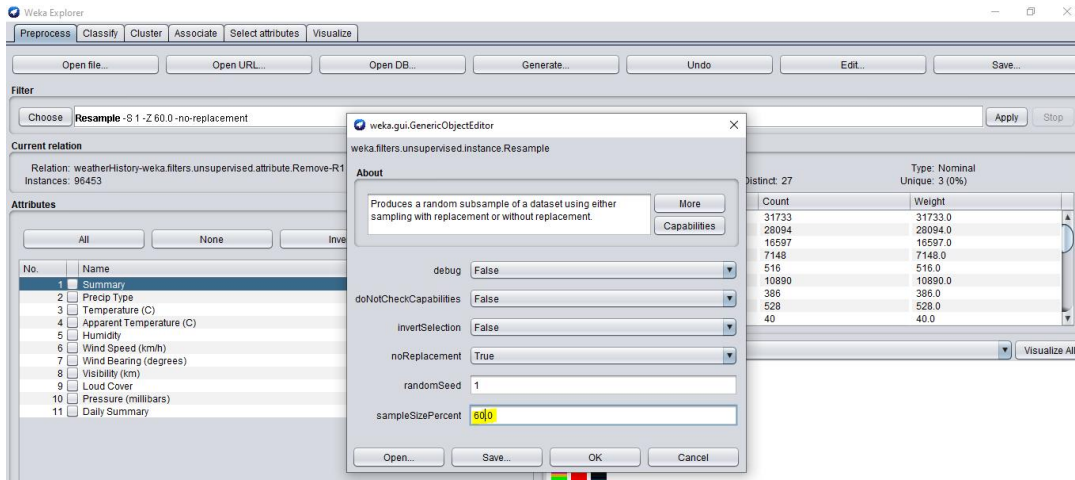
1. ### Deleting Unnecessary attributes:
   Here "formatted date" is an unnecessary attribute that has no significant relation with the entire dataset. That's why I have deleted this attribute by remove function.
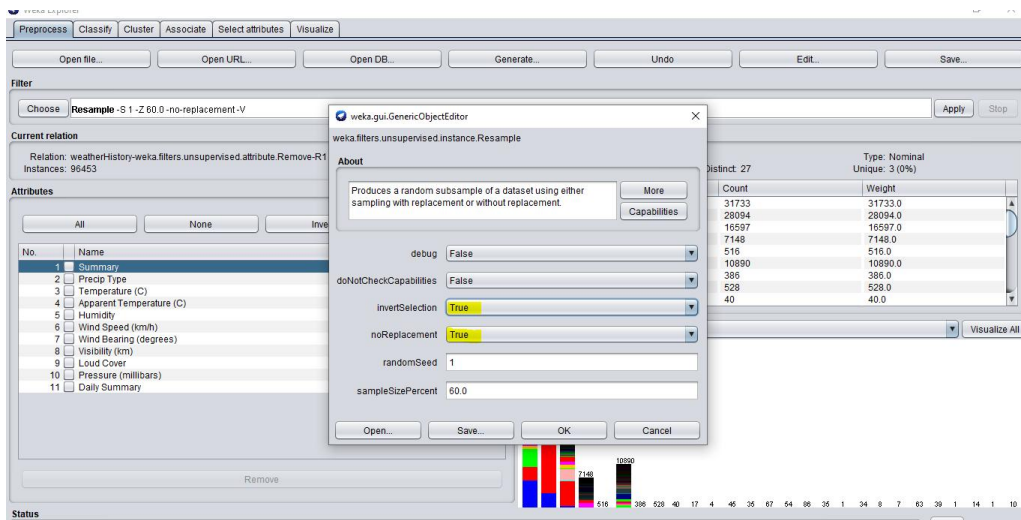


2. ### Splitting data set:

Now we need to divide data into two set one is training set another is test set.

Here first one is we have divided into training set by taking 60% of data set.

And the next picture we have saved for test set by taking the 40% of the rest.



3.  Applying Algorithms:
    Now we need to apply KNN, Naïve Bayes and decision tree algorithms to our train set and test set.

For KNN: classify>choose>lazy>IBk

For Naïve Bayes: classify>choose> bayes >NaiveBayes

For KNN: classify>choose>trees>j48

**Here are some screenshots of test and train set accuracy result in different algorithms.**

## KNN (Test set and Train set):

## Naïve Bayes (train set and test set):

## Decision Tree (Training Set and Test Set):

<u>Results:</u>

<span style="color:red">C = Correctly Classified Instances</span>

<span style="color:red">In = Incorrectly Classified Instances</span>

By analyzing the data table of accuracy of different algorithms we can say that KNN is best algorithm for this specific problem. Accuracy of KNN is 99.96% and 99.94% for training set and test set where accuracy of Decision Tree is 79.08% and 88.44%. At the mean time the accuracy of Naïve Bayes for this data set is only 22.37% that means it can correctly classify only 22.37% instances. So the accuracy of Naïve Bayes is very poor for this data set.

|  | KNN | Naïve Bayes | Decision Tree |
|---|---|---|---|
| Training set | [C: 99.96%, In: 0.03%] | [C:22.37%, In:77.56%] | [C:79.89%, In:20.10%] |
| Test Set | [C:99.94%, In: 0.05%] | [C:22.46%, In:77.54%] | [C:88.44%, In:19.55%] |

<u>Why KNN is best for this dataset:</u>

According to my results that I got from weka I can say that KNN is definitely suitable for this type of dataset.

According my data set there are 12 columns and our targeted column name is Daily Summary which gives us what is the weather of today. By analyzing this data set we need to find what will be the weather of tomorrow. This problem is a classification problem because it has numeric and discrete values and it has multiple classes. So it is a classification and supervised algorithm. Here we can't predict what our output will be so our dataset need unsupervised algorithms. Our data set is also unlabeled dataset.

I have executed this data set with k=1. So we only use nearest neighbor to define the category. As I have separated this data set into two part that are training set and test and after executing

this data set I got 99.96% and 99.94% accuracy respectively. The second best accuracy value I got from decision tree. Decision tree uses several decision tree for several output and it need to re-execute when two output results are same. That's why Decision tree is not a suitable for this type of dataset. On the other hand Naïve Bayes is for basically sentiment analysis and text classification so it is not a superior algorithm for this data set and we got very poor accuracy for Naïve Bayes that is correct instances are only 22.37% where incorrect instances are 77.54%.

When we implemented KNN it works by searching the distances (Euclidean distance) between a query and all the possible examples of dataset and after that we assign a value for K (I used k=1) then it select nearest label. Our data set need to measure the distance among of **temperature**, **wind speed**, **pressure**, **visibility**, **Precip type**, **wind bearing** to give a proper decision that whether the weather of next day will be Partly cloudy throughout the day or Mostly cloudy throughout the day or other state. That's why KNN give us the best accuracy.

To recapitulate, based on above decision and analyzation undoubtedly we can say that KNN is best for this dataset.

**Data-Set Reference:** https://www.kaggle.com/abhishek20182/performing-analysis-of-meteorological-data