# American International University Bangladesh
## Department of Computer Science
## Project Report

Semester: Spring'21-22
Section: D
Course Name: Data Warehousing And Data Mining

Submitted To,
*Dr. Akinul Islam Jony*
*Assistant Professor*
*Department of Computer Science*
*American International University-Bangladesh*

TITLE: WEATHER FORECASTING USING DATA MINING APPLICATION

## GROUP MEMBERS INFORMATION

| NAME | ID |
|------|-----|
| NAJMUL UDDIN | 18-38293-2 |
| SUMONA ISLAM | 18-37288-1 |
| MD. MILKAN ISLAM | 18-36958-1 |
| MAHZABIN MOSTARY | 18-38606-2 |

# Introduction

## Objective

Data mining is the process of uncovering patterns and finding anomalies and relationships in large datasets that can be used to make predictions about future trends. The main objective of data mining is to extract valuable information from available data. This generally involves the use of database techniques such as spatial indexes. Thus, these patterns can be seen as a kind of summary of the input data.In addition to being able to be used in additional analysis or, for example, in machine learning and predictive analysis. One of the examples we can give is data mining. This could identify several groups in the data, which can then be used to obtain more accurate results being able to predict problems through a decision support system. Neither data collection, data preparation,nor interpretation of results and information is part of the data mining stage.

## Description

Data mining involves the use of complicated data analysis tools to discover previously unknown, interesting patterns and relationships in large data set. These tools can include statistical models, mathematical algorithm and machine learning methods. The most significant analysis phase in the knowledge discovery in database (KDD) process is data mining. The basic purpose of data mining is to extract meaningful information from large amounts of raw data and convert it into a format that can be used effectively and efficiently. Data mining jobs are typically classified into two types: descriptive and predictive classification techniques. Data classification is the process of organizing data into categories/groups in such a way that data objects of same group are more similar and data objects from different groups are very dissimilar. The classification algorithm assigns a classification each instance to a certain class so that classification is possible. There will be the least amount of mistake. It's utilized to create models that are correct. Classification techniques can handle processing of large volume of data. It can predict categorical class labels and classifies data based on model built by using training set and associated class labels and then can be used for classifying newly available test data. Thus, it is outlined as an integral part of data analysis and is gaining more popularity.

## Outcome

1. To fully understand standard data mining methods and techniques such as association rules, data clustering and classification.

2. Learn new, advanced techniques for emerging applications (social network analysis, stream data mining).

3. Gain practical intuition about how to apply these techniques on datasets of realistic sizes using modern data analysis frameworks.

# Dataset URL

## Data processing Steps

The data processing and filtering consists of the following stages.

### Data Source

A data source is the location where data that is being used originates from. A data source may be the initial location where data is born or where physical information is first digitized, however even the most refined data may serve as a source, as long as another process accesses and utilizes it.

### Data Store

A data store is a repository for persistently storing collections of data, such as a database, a file system or a directory. In an information technology context, data stored can be of any type that can be rendered in digital format and placed in electronic media.

### Prepared Data

Once data store is completed then data prepared for selecting and preprocessing.

### Patterns

When data prepared for selecting and preprocessing after that data make patterns for interpretation and assimilation.

**Dataset Description**

| | date | precipitation | temp_max | temp_min | wind | weather |
|---|---|---|---|---|---|---|
| 1 | date | precipitation | temp_max | temp_min | wind | weather |
| 2 | 4/2/2001 | 0 | 12.8 | 5 | 4.7 | drizzle |
| 3 | 4/3/2001 | 10.9 | 10.6 | 2.8 | 4.5 | rain |
| 4 | 4/4/2001 | 0.8 | 11.7 | 7.2 | 2.3 | rain |
| 5 | 4/5/2001 | 20.3 | 12.2 | 5.6 | 4.7 | rain |
| 6 | 4/6/2001 | 1.3 | 8.9 | 2.8 | 6.1 | rain |
| 7 | 4/7/2001 | 2.5 | 4.4 | 2.2 | 2.2 | rain |
| 8 | 4/8/2001 | 0 | 7.2 | 2.8 | 2.3 | rain |
| 9 | 4/9/2001 | 0 | 10 | 2.8 | 2 | sun |
| 10 | 4/10/2001 | 4.3 | 9.4 | 5 | 3.4 | rain |
| 11 | 4/11/2001 | 1 | 6.1 | 0.6 | 3.4 | rain |
| 12 | 4/12/2001 | 0 | 6.1 | -1.1 | 5.1 | sun |
| 13 | 4/13/2001 | 0 | 6.1 | -1.7 | 1.9 | sun |
| 14 | 4/14/2001 | 0 | 5 | -2.8 | 1.3 | sun |
| 15 | 4/15/2001 | 4.1 | 4.4 | 0.6 | 5.3 | snow |
| 16 | 4/16/2001 | 5.3 | 1.1 | -3.3 | 3.2 | snow |
| 17 | 4/17/2001 | 2.5 | 1.7 | -2.8 | 5 | snow |
| 18 | 4/18/2001 | 8.1 | 3.3 | 0 | 5.6 | snow |
| 19 | 4/19/2001 | 19.8 | 0 | -2.8 | 5 | snow |
| 20 | 4/20/2001 | 15.2 | -1.1 | -2.8 | 1.6 | snow |

*FIGURE : SAMPLE OF DATASET*

To predict the weather, we need to measure the weather. If you want to know the weather is like tomorrow, it's pretty important to know what the weather was like today and yesterday. Knowing the average weather on a particular day of the year is also useful. Collecting data every day can show you patterns and trends, and help you figure out how our atmosphere works.

In here, it's about 5236 instances in 5236 days and for each of these days we have 6 values of 6 attributes. We are going to be doing is predicting the weather attribute from the other attributes.

## Dataset Exploration

Data exploration refers to the initial step in data analysis in which data analysts use data visualization and statistical techniques to describe dataset characterizations, such as size, quantity, and accuracy, in order to better understand the nature of the data. Classification problem sometimes called supervised learning problem. We take as input a data set of classified examples. These examples are independent example with a class value attached and the idea is to produce automatically a model that is some kind of model that can classify your examples. That's the classification problem. An instance with the different attribute values a fixed set of feature and then we add to that a class to get a classified example. That's what we have in our training data set. These attributes can be discrete of continuous. What we looked at in the weather data was discrete or we call them normal attribute values where they belong to a certain fixed set well. And also the class can be discrete or continuous.

## How We Deal With Missing Value in Our Datasets

Incomplete data is an unavoidable problem in dealing with most of the real world data sources. There are two most used methods for dealing with missing values. In our project we used repeated by most frequent/average value. To estimate each of the missing values is a less cautious method, using the values is present in the datasets. For a categorical attribute, the most commonly occurring (non-missing) value is a simple yet effective technique to do this. This is easy to justify if the attribute values are very unbalanced. For example, if attribute X has possible values a, b and c which occur in proportions 80%, 15% and 5% respectively, it seems reasonable to estimate any missing values of attribute X by the value a. If the values are more evenly distributed, say in proportions 40%, 30% and 30%, the validity of this approach is much less clear

## 1. Implementation of Decision Tree

**About Decision Tree**

 Decision tree is the most powerful and popular tool for classification and prediction. A Decision tree is a flowchart like tree structure, where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (terminal node) holds a class label. Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features. A tree can be seen as a piecewise constant approximation. A decision tree is a very specific type of probability tree that enables you to make a decision about some kind of process.

**Representation of Decision Tree**

Decision trees classify instances by sorting them down the tree from the root to some leaf node, which provides the classification of the instance. An instance is classified by starting at the root node of the tree, testing the attribute specified by this node,  then moving down the tree branch corresponding to the value of the attribute as shown in the above figure. This process is then repeated for the subtree rooted at the new node. The decision tree in above figure classifies a particular morning according to whether it is suitable for playing tennis and returning the classification associated with the particular leaf.(in this case Yes or No). For example, the instance (Outlook = Rain, Temperature = Hot, Humidity = High, Wind = Strong ) would be sorted down the leftmost branch of this decision tree and would therefore be classified as a negative instance. In other words we can say that decision tree represent a disjunction of conjunctions of constraints on the attribute values of instances.

## The Strengths of Decision Tree

- Decision trees are able to generate understandable rules.
- Decision trees perform classification without requiring much computation.
-  Decision trees are able to handle both continuous and categorical variables.
- Decision trees provide a clear indication of which fields are most important for prediction or classification.

## The weaknesses of Decision Tree Methods

- Decision trees are less appropriate for estimation tasks where the goal is to predict the value of a continuous attribute.
- Decision trees are prone to errors in classification problems with many class and relatively small number of training examples.
- Decision tree can be computationally expensive to train. The process of growing a decision tree is computationally expensive must be formed and compared.

## MODEL DEVELOPMENT BY WEKA TOOLS

Apply the weather data in Weka.  Explain the steps on how to implement   decision Tree in Weka Tool.



*FIGURE 01*

AFTER SELECTING DATA ,FOR DECISSION TREE WE SELECTED J48.AFTER MINING THE DATA OR DECISSION TREE WE GET 96.65% ACCURACY INSTEAD WE GET 3.35% INCORRECTLY CLASSFIED INSTANCE

```
=== Confusion Matrix ===

      a      b      c      d      e    <-- classified as
    200      7     22      0      4 |     a = drizzle
      2   2408     12      2      0 |     b = rain
     16     32   2087      0     23 |     c = sun
      0      3      0    137      0 |     d = snow
      3     11     64      0    200 |     e = fog
```

*FIGURE 02*

HERE IS OUR CONFUSION MATRIX. WHICH IS COMPARE ACTUAL VALUE WITH PREDICTED VALUE



*FIGURE 03*

HERE HAS BEEN DESCRIBED, SIZE OF THE TREE IS 591 .WE GOT 0.9365 MEAN ABSOLUTE ERROR.HERE WE HAVE DISPLAY,ALL THE CLASSIFIED OF DECISSION TREE

**FIGURE 04**

FLOW CHART TO ATTRIBUTE OF WEATHER. TEMPERATURE VS WIND, WIND VS WEATHER.



**FIRGURE 05**

HERE IS THE DECISSION TREE OF MY DATASET.NUMBER OF THE LEAVE IS 296 AND SIZE OF THE TREE IS 591

## 2. Implementation Of Naïve Bayes Algorithm

### About Naïve Bayes Algorithm
The Naive Bayes algorithm is called "naive" because it makes the assumption that the occurrence of a certain feature is independent of the occurrence of other features. Bayes It refers to the statistician and philosopher Thomas Bayes and the theorem named after him, Bayes' theorem, which is the base for the Naive Bayes Algorithm

### Advantage of Naïve Bayes Algorithm
- It is simple and easy to implement
- It doesn't require as much training data
- It handles both continuous and discrete data
- It is highly scalable with the number of predictors and data points
- It is fast and can be used to make real-time predictions
- It is not sensitive to irrelevant features

### Disadvantage of Naïve Bayes Algorithm
- This limits the applicability of this algorithm in real-world use cases.
- This algorithm faces the 'zero-frequency problem' where it assigns zero probability to a categorical variable whose category in the test data set wasn't available in the training dataset
- Its estimations can be wrong in some cases, so you shouldn't take its probability outputs very seriously.
- Its estimations can be wrong in some cases, so you shouldn't take its probability outputs very seriously.

# MODEL DEVELOPMENT BY WEKA TOOLS

*FIGURE : 06*

HERE IS MY DATASET DESCRIPTION



*FIGURE : 07*

AFTER APPLIFED 10 FOLD CROSS VALIDATION I GOT 71.77% ACCURACY AND 29.21% INCORRECT CLASSFIED INSTANCE
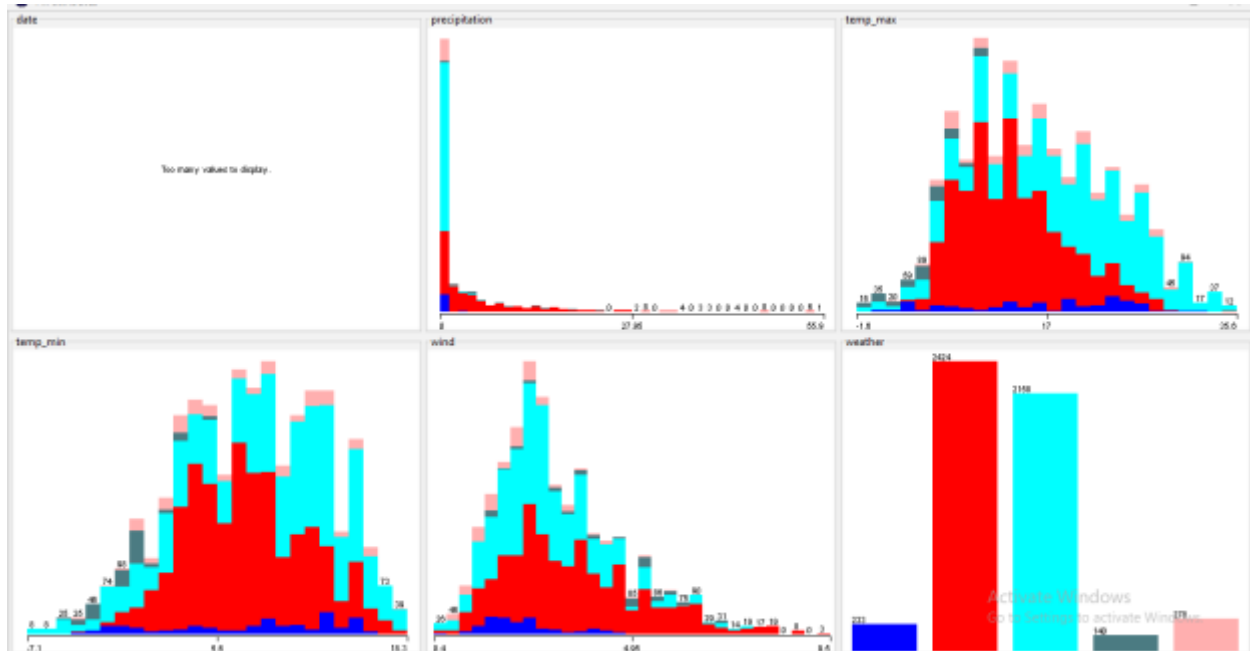
*FIGURE : 08*
FLOW CHART  OF NAÏVE BAYES ALGORITHM APPLIED ON DATASET

## 3.  Implementation of   Nearest Neighbor Classification.

### About Nearest Neighbor Classification

K-NN Classification is used to develop this model. The k-nearest neighbor (k-NN) method is one of the data mining techniques considered to be among the top 10 techniques for data mining. The k-NN method tries to classify an unknown sample based on the known classification of its neighbors. In our dataset, the set of samples with known classification is available, which is the training set. Intuitively, each sample should be classified similarly to its surrounding samples. Therefore, if the classification of a sample is unknown, then it could be predicted by considering the classification of its nearest neighbor samples.

## ADVANTAGES

- **K-NN is pretty intuitive and simple**: K-NN algorithm is very simple to understand and equally easy to implement.
- **K-NN has no assumptions**: K-NN is a non-parametric algorithm which means there are assumptions to be met to implement K-NN.

- **No Training Step:** K-NN does not explicitly build any model, it simply tags the new data entry based learning from historical data. New data entry would be tagged with majority class in the nearest neighbor.

# DISADVANTAGES

- **K-NN slow algorithm**: K-NN might be very easy to implement but as dataset grows efficiency or speed of algorithm declines very fast.
- **Curse of Dimensionality:** KNN works well with small number of input variables but as the numbers of variables grow K-NN algorithm struggles to predict the output of new data point.
- **K-NN needs homogeneous features**: Deciding to build k-NN using a common distance, like Euclidean or Manhattan distances, it is completely necessary that features have the same scale, since absolute differences in features weight the same.

- 

- ## MODEL DEVELOPMENT BY WEKA TOOLS



*FIGURE : 09*

In this picture we can see that the data set is loaded on the Weka tool to perform the model development process. In the current relation section of this software, it is visible that, the name

of the relation is weather. Also there are 6 attributes, instances of 5235 and the sum of weights are the same as instances



H
or
at



*FIGURE : 11*

Here on the classifier section, we have chosen lazy -> IBK on the classifier and it is using the Euclidian distance method to determine the accuracy of the k nearest neighbor, by using a 10 folds cross validation technique. The value of the K is 1, which finally it shows that the model has 97.78% of correctly classified instances and 2.216% of incorrectly classified instances.



changed, now it's Manhattan distance here.

Here the value of K has been significantly increased to 4000. After that it's visible, the accuracy level has been massively reduced to a 46.53% and the percentage of error has been increased. The incorrectly classified instances now 53.47%

# MODEL COMPARISON:

Comparison of used classification methods in this project; Decision Tree (DT), k-Nearest Neighbor (k-NN), Naïve Bayes (NB), has been compared. The effects of parameters including size of the dataset, kind of the independent attributes, and the number of the discrete and continuous attributes are the key parameters of this comparison. Based on the results, it can be concluded that in the datasets with few numbers of records, the comparison between classifiers may not do correctly. When the number of the records and the number of the attributes in each record are increased, the results become more stable. Decision Tree, k-NN obtain higher AUC than Naïve Bayes.

Nearest neighbor classifiers are based on learning by analogy, that is by comparing a given test tuple with training tuples which are similar to it. The training tuples are described by n attributes. Each tuple represents a point in an n-dimensional space. In this way, all of the training tuples are stored in an n-dimensional pattern space. When given an unknown tuple, a k-nearest neighbor (k-NN) classifier searches the pattern space for the k training tuples which are closest to the unknown tuple. These k training tuples are the k-nearest neighbors of the unknown tuple.

 We have developed KNN model and got pretty good accuracy. But the K value plays a significant role to differ the error and incorrect instances percentage. Yet KNN classification is more optimizable and easier to use in weka tools.  Considering as many as instances to develop the model leads to some errors and less correct results. On that account Decision tree is more accurate on the results.

A decision tree (DT) is a flowchart-like tree structure, where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (or terminal node) holds a class label . The topmost node in a tree is the root node.   During tree construction, attribute selection measures are used to select the attribute which best partitions the tuples into distinct classes. Three popular attribute selection measures are Information Gain, Gain Ratio, and Gini Index. When DTs are built, many of the branches may reflect noise or outliers in the training data. Tree pruning attempts to identify and remove such branches, with the goal of improving classification accuracy on unseen data. By using decision tree we have achieved the best possible outcome comparing with the other two models.

Naïve Bayes has given the lowest accuracy in comparison with other models. As the basic idea in NB approaches is to use the joint probabilities of words and categories to estimate the probabilities of categories given a document. The naïve part of NB methods is the assumption of word independence, i.e., the conditional probability of a word given a category is assumed to be independent from the conditional probabilities of other words given that category. This assumption makes the computation of the NB classifiers far more efficient but in case of our dataset this model was not efficient enough.

## DISCUSSION AND ANALYISIS

From this project, we have gained knowledge about the efficacy of DT, k-NN, NB has been investigated and these three methods is compared to each other with respect to different circumstances of attributes deployment. In this model development and comparison process, the size of datasets, types of the attributes and the number of discrete and continuous attributes have been considered. From our analysis we have shown that for small datasets, in all cases deviation is such high that the comparison between classifiers may not do correctly. For larger datasets results become more stable and the comparisons can be done. That is the reason for choosing our large dataset to perform the model development and comparison with maximum accuracy. Decision Tree as an implementation of that, show an efficient performance in all datasets. Two classifiers including k-NN and Naïve Bayes are also high efficacy in the current work. Following the used process of evaluation of our project, we can choose the best classifier according to data type and continuous or discrete attributes existing in datasets. Therefore, using the above conclusion our most efficient classifier with respect to dataset characteristics is Decision Tree. Lastly, it should be mentioned that our current project is based on simulation data and the generated datasets that are not dependent to a special problem. Consequently, the above results can be extended to a wide range of problems and real datasets made from data collected from actual real-world problems are suitable for model developing and comparing the mentioned methods. In future studies, we will be using the experience gained from this project, taking larger real datasets and using other evaluation criteria and applying new classifiers on datasets with more variables could be as open problems in this field.