

Forecasting Weather Using Machine Learning

Md. Azwad Hasan Chowdhury
Electrical & Computer Engineering
North South University
Dhaka, Bangladesh
azwad.kenpachi@gmail.com

Sadman Kabir Soumik
Electrical & Computer Engineering
North South University
Dhaka, Bangladesh
sadmanks@gmail.com

Rifat Arefin Badhon
Electrical & Computer Engineering
North South University
Dhaka, Bangladesh
rifat.arefin@northsouth.edu

Sarwat Islam Dipanzan
Electrical & Computer Engineering
North South University
Dhaka, Bangladesh
dipanzan@live.com

Abstract— *The question of predicting weather has baffled mankind for centuries but with machine learning techniques it is now possible to predict weather with good accuracy. This project proposes to predict weather more accurately within the context of Bangladesh. As predicting weather is a challenging task, the application of machine learning in this sector has promising results. The people of Bangladesh suffer a lot due to bad weather pattern and it is an on-going problem. This study hopes to achieve an insight on how machine learning techniques can prove to be helpful using classification algorithms to predict cloud patterns based on past data of Dhaka. Using state of the art classifiers from statistical models such as Naïve Bayes, Decision Tree, Logistic Regression and non-probabilistic models like Support Vector Machine, gives accurate results up to 70% and more when predicting weather patterns.*

Keywords—*weather, forecasting, machine learning, logistic regression, support vector machine, naïve bayes, gaussian, multinomial, decision tree.*

I. INTRODUCTION

The recent climate is quite unpredictable in the context of Bangladesh and as a result a lot of people are faced with destruction and calamity. This affects a large population of farmers that are faced with tough decisions as crops are destroyed at the onset of bad weather. Due to unpredictable patterns, most outdoor activities have to consider the effects of weather before an event.

People are faced with the difficult task of planning ahead and plans have to be changed at the last moment if an unforeseen change of events occurs such as a heavy downpour. This inconvenience forces alternate measures and costs people valuable time and resources. Bangladesh is also very prone to floods and rain is a constant threat to the infrastructure of cities like Dhaka which are unplanned and densely populated.

Weather data from the past is an invaluable source for finding patterns and co-relations among certain weather variables. The weather data is constantly monitored by government stations installed at certain parts of the country which collect various different kinds of data such as amount of rainfall, temperature, cloud, humidity. These data can then be used to find very useful metrics

The current model for predicting weather works on the basis of Numerical Weather Prediction (NWP), which is tasked to predict future weather characteristics using present conditions. [6]

The machine learning techniques applied focus on the previous data collected from authorized government stations situated in Dhaka. Using various implementations of already available popular Python packages like scikit-learn library to work on the dataset and prepare a model.

The authenticity of the dataset is crucial in building upon the foundation of this project as erroneous data greatly alters outcomes of the algorithms. The efficiency of the algorithm greatly depends on a good dataset that is varied and not heavily skewed. The assumption that the data available also includes the possibility of anomalies due to faulty measurements and bad records. The advent of such conditions has to be dealt with in such a manner that does not disrupt the structure of the dataset so that the algorithms can work better even when generalizing new samples of collected data.

II. RELATED WORK

Over the last few decades, there has been a substantial amount of research done using numerical weather data. The data was primarily used for applying machine learning techniques. Other techniques included fuzzy logic and data mining. The majority of weather forecasting relies on generative approaches and the underlying principle are based on numerical methods.

Some people applied comprehensive tree-based learning algorithm. For example, N. Hasan and M.T Uddin used tree-based algorithm, namely C4.5 and their output result had 96% accuracy. They also used Naïve Bayes, but the accuracy of C4.5 was much better in terms of f-score. [1] Lin and Chen [2] worked on typhoon rainfall forecasting model using ANN and their result shows that excessive spatial rainfall information may not increase the generalization of forecasting model. Awan and Awais [3] also tried to predict weather events based on fuzzy RBS method for Lahore, Pakistan. They used two different datasets of 365 examples of with only 4 features, and 2500

examples with 17 features. They mentioned their finding that fuzzy RBS method was sensitive to random sampling with replacement technique that was applied to produce. Another article reviews we found useful is S. B. Kotsiantis [4] mentioned few statistical classifiers to build classification tress. Using information entropy from a set of training examples of pre-classified samples, where each sample comprises of N-dimensional vector. H. Zhang, X. Zhao and S.Zou [5] proposed a classification algorithm naming Neuron Classification Algorithm (NCA). The algorithm has higher approximation function. They introduced the law of attraction here which increases the accuracy of weather forecast. It has been known to classify the test samples more accurately than Euclidian distance. As dataset, they chose forecast of abnormal megathermal weather in North of Zhejiang province. Combining Neural Networks and ARIMA Models for Hourly Temperature Forecast H. S. Hippert , C. E. Pedreira and R. C. Souza in their works [7] on hourly temperature forecasting, proposed to use the combination of Neural Networks and ARIMA models. The forecast is done on the basis of previous temperature records, maximum and minimum temperature data supplied by weather service. The previous temperatures are used as input to Artificial Neural Network (ANN) where there is only one output node, the predicted temperature on a particular time. Their results show that hybrid system based on ARMA model produces more accuracy than auto progressive models.

Simple classifiers such as Support Vector Machines (SVM) or Artificial Neural Networks (ANN) are also widely used to classify certain parameters such as rainfall and cloud states.

III. DATASET

The dataset was obtained from Bangladesh Meteoritical Department (www.bmd.gov.bd). The data set is comprised of the following variables:

Temperature(°C): The temperature variable has great impact on precipitation and is greatly related to humidity (%). The dataset has both maximum and minimum temperatures along with the mean values.

Cloud (okta): The cloud variable measured in Okta impacts the Earth's surface by reflecting incoming sunlight. It is also responsible for absorbing the heat emitted from surface and radiating to space. The dataset comprises of daily cloud data measured in a range from 0-8 Okta.

Wind Speed (knots): The wind variable measured in knots shows how quick the air is moving. The wind speed also has a direction and has various impacts on surface water and evaporation. The dataset comprises of daily prevailing wind speed.

Rainfall (mm): The rainfall variable is a very important metric in weather forecasting. It helps the environment to continue to stay in its position the way it should be. Agriculture of Bangladesh mostly depends on rainfall. The dataset has daily basis of total rainfall data in millimeters.

Sunshine (hour): The sunshine variable measures the amount of sunshine at a particular place. The Sun is the basic

cause of our changing weather. The day-night cycles in the weather have obvious causes and effects on weather.

Humidity (%): The humidity variable measured in percentage helps to calculate the amount of moisture in the air at a give time on a given day, which is simply the ratio of water vapor and dry air.

Sea Level Pressure (millibars): The sea level pressure variable plays a significant role in the formation of weather condition in a certain area. It is a component that has mass and weight. This means vast ocean of air inserts huge amount of pressure. So, it is natural that the air will affect the Earth's weather.

TABLE 1: SAMPLE TRAINING DATA

Date	Humidity	Max Temp	Min Temp	Sunshine	Cloud
01/01/1988	77	26.7	12.9	8.4	0
02/01/1988	76	26	12.9	8.4	0
03/01/1988	73	27.5	14.5	7.8	1
04/01/1988	71	27.2	15.8	6.4	2
05/01/1988	75	27.8	15.4	8.0	1

TABLE 2: TOTAL AVAILABLE PARAMETERS

Parameter Number	Parameter Name
1	Day
2	Month
3	Year
4	Humidity (%)
5	Max Temp (in °C)
6	Min Temp (in °C)
7	Rainfall (in mm)
8	Sea Level Pressure (in mb)
9	Sunshine (hours)
10	Wind Speed(knot)
11	Cloud (in okta)

The dataset from BMD that was used for the project has past weather data from for the last 30 years (1988-2017) from Dhaka. The data collected is from a daily basis and this roughly equates to 109 samples for the entire set. The origin of the data is from Bangladesh Meteorological Department's Dhaka station with a resolution of 25 kilometers. Absolute coordinates: (Lat 23 Deg 46 Mts.N & Long 90 Deg 23 Mts.E).

The training and test set is divided into two segments having 70% and 30% data split across the two categories.

IV. METHODOLOGY

The purpose of our project is to predict the cloud states of various days using other features from the dataset. The initial cloud data from the data set had a range of 0-8. This number signifies the amount of cloud on a given day, with 0 being the least cloudy and 8 being completely cloudy. The numbers in between are cloud states in the intermediate range. The cloud ranges were then compressed to reduce the original range given in the data set. This is done so that classifiers have to deal with lower number of classes when tasked with predicting. The ranges were aggregated to 0 representing the cloud states 0, 1 and 2 labelling it “clear skies”. The class 1 then represented the cloud states 3, 4 and 5 giving it the label “half cloudy”. The class 2 then defined the cloud states 6, 7 and 8 labelling it “fully cloudy”. The cloud states are defined by the numbers ranging from 0 – 2, so this is a classification problem, namely multi-class classification. The dataset is labelled with cloud values in that range, so the learning algorithms are tasked with a supervised approach. We used several machine learning techniques which is constructed using the Python scikit-learn library. Parameter tuning of the models has been done by 5-fold cross validation using sklearn GridSearchCV.

IV. MODEL SELECTION

1. **Multinomial Logistic Regression:** Logistic regression is known as a binary classifier but can also act as a multi-class classifier or otherwise known as multinomial logistic regression which can identify more than 2 classes using methods such as One vs All. The logistic regression calculates the probability of a class based on the hypothesis that works using the sigmoid function.

$$h_{\theta}(X) = \frac{1}{1 + e^{-\theta^T X}}$$

The multinomial logistic regression function is using “l1” penalization and “liblinear” with a C value of 0.01.

2. **Decision Tree:** Decision Tree Classifier is an algorithm suited for classification tasks and falls under the supervised criteria. Decision trees works under the basic principle that it has to predict target values and it does so by forming trees from input nodes. The decision tree used in this implementation works using the “gini” criterion.

$$Gini(E) = 1 - \sum c_j = 1 - p_j^2$$

Here, c is the number of classes and p is the fraction of records.

3. **Naïve Bayes Multinomial:** Naïve Bayes Multinomial falls under the category of probabilistic classifiers. The algorithm is devised using Bayes theorem and works on the principal that there is a strong(naïve) independence between the features. Naïve Bayes classifiers are beneficial

in supervised learning setup because they can be trained very efficiently.

From the family of Naïve Bayes classifiers, we have implemented Gaussian Naïve Bayes, which is characterized by the following equation:

$$p(x = v | C_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(v-\mu_k)^2}{2\sigma_k^2}}$$

4. **Support Vector Machine:** Support Vector Machines are a group of algorithms that is used for both classification and regression and are associated with supervised learning models. The most prominent feature of SVMs is that these class of algorithms fall into the “non-probabilistic” criteria. SVM works on the basis that it maximizes the Euclidean distance between the points that fall on the furthest lines and the median which is the decision boundary. A hinge loss function calculates the false classifications. When tuning for hyperparameters using the built-in GridSearchCV function using scikit-learn it is found that the algorithm works best with kernel “rbf” along with a C value of 0.01 and cv value of 5.

V. RESULTS

The results after applying the different kinds of models are discusses along with their performance measures. Models tested are: Multinomial Logistic Regression, Decision Tree, Naïve Bayes Multinomial and Support Vector Machine. The training data fed to these algorithms were 70% of the original dataset and the rest 30% for testing. The confusion matrix generated for these multi-class classifications delves a bit more on the evaluation criteria and hence are given below.

TABLE 3: CONFUSION MATRIX FOR LOGISTIC REGRESSION

ACTUAL CLASS	PREDICTED CLASS			
		0	1	2
	0	1079	98	69
	1	148	389	343
	2	23	91	1048

TABLE 4: CONFUSION MATRIX FOR DECISION TREE

ACTUAL CLASS	PREDICTED CLASS			
		0	1	2
	0	927	309	10
	1	107	585	188
	2	14	217	931

TABLE 5: CONFUSION MATRIX FOR MULTINOMIAL NAÏVE BAYES

ACTUAL CLASS	PREDICTED CLASS			
		0	1	2
	0	1222	15	9
	1	626	134	120
	2	370	361	431

TABLE 6: CONFUSION MATRIX FOR SUPPORT VECTOR MACHINE

ACTUAL CLASS	PREDICTED CLASS			
		0	1	2
	0	1145	164	37
	1	126	493	261
	2	24	127	1011

The confusion matrix of all the models are then used to find the average precision, recall and fit scores. The precision score is giving an overall general sense of the times the model is able to correctly predict and how often. The recall score is the actual relevant results correctly predicted by the model.

TABLE 7: AVERAGE VALUES OF PRECISION, RECALL AND FIT-SCORE

Model	Precision	Recall	Fit Score	Support
Logistic Regression	0.76	0.77	0.76	3288
Decision Tree	0.76	0.74	0.75	3288
Multinomial NB	0.55	0.54	0.50	3288
SVM	0.77	0.78	0.77	3288

After running the models using our training set and testing it against the data that was split at the beginning for testing, it is evident that Support Vector Machines is more accurate against the test set compared to the other models. The nature of SVM allows it to run much better with a kernel trick but it becomes a lot slower when testing against the whole training set. The values of c have to be modified in order to gain better performance.

TABLE 8: CHECKING ACCURACY USING CROSS VALIDATION 10

Model	Mean Accuracy	Standard Deviation
Logistic Regression	0.7392	0.010
Decision Tree	0.723	0.0219
Multinomial NB	0.5418	0.012
SVM	0.759	0.0199

Using cross validation value of 10, the mean accuracy and standard deviations are also collected. This shows that the results are well within range with expected errors.

TABLE 9: TRAIN AND TEST ACCURACY OF THE MODELS TESTED

Model	Training Accuracy (%)	Testing Accuracy (%)
Logistic Regression	74.2	76.9
Decision Tree	76.8	74.05
Multinomial NB	54.27	54.34
SVM	76.42	77.52

VI. CONCLUSION

After running all tests with the proposed models with the dataset provided after initial setup, the optimized parameters after tuning could predict reasonably well. One exception to this was the model Multinomial Naïve Bayes receiving just 54% in both training and test accuracy. The other models received average scores but the difference between testing and training accuracy proved to be less than 5%. This ensures that the models did not have an over fit and were performing well within ranges with errors. For future work we hope to bring in more models to test against our already optimized models and see if they generate better results. The project revolved around classification but there are other critical weather variables which are continuous and regression models can be used to see if they perform better or worse than their classification counterparts. A better dataset with even more variables that govern weather patterns would prove invaluable for testing out models that are better and more advanced in their prediction capabilities.

REFERENCES

- [1] N. Hasan, M.T Uddin and N.K. Chowdhury "Automated Weather Event Ananysis with Machine Learning"
- [2] G.-F. Lin and L.-H. Chen, "Application of an artificial neural network to typhoon rainfall forecasting," *Hydrological Processes*, vol. 19, no. 9, pp. 1825–1837, 2005.
- [3] M. S. K. Awan and M. M. Awais, "Predicting weather events using fuzzy rule based system," *Applied Soft Computing*, vol. 11, no. 1, pp.56–63, 2011.
- [4] S.B. Kotsiantis, "Supervised machine learning: A review of classification techniques," in *Proceeding of the 2007 Conference on Emerging Artificial Intelligence Applications in Computer Engineering: Real World AI Systems with Applications in Computer Engineering: Real World AI systems with Applications in eHealth, HCI Amsterdam, The Netherlands, The Netherlands: IOS Press, 2007*, pp. 3–24.
- [5] H. Zhang, X. Zhao and S.Zou, "Neuron Classification Algorithm and Magthermal Weather Forecast" Beijing: Meteorology Press, 2002.
- [6] Dirmeyer, Paul A.; Schlosser, C. Adam; Brubaker, Kaye L. "Precipitation, Recycling, and Land Memory: An Integrated Analysis." <https://journals.ametsoc.org/doi/pdf/10.1175/2008JHM1016.1> . Dec. 2016.
- [7] H. S. Hippert , C. E. Pedreira and R. C. Souza, "Combining Neural Network and ARIMA models for Hourly Temperature Forecast", *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks*, vol 12, no. 1, pp. 57-28, 2014.