MACHINE LEARNING BASED RAINFALL PREDICTION CH.MOUNIKA¹, S.K.ALISHA²

¹MCA Student, B V Raju College, Kovvada, Andhra Pradesh, India.

²Associate Professor, B V Raju College, Kovvada, Andhra Pradesh, India.

ABSTRACT:

Rainfall prediction is one of the challenging and uncertain tasks which has a significant impact on human society. Timely and accurate predictions can help to proactively reduce human and financial loss. This study presents a set of experiments which involve the use of prevalent machine learning techniques to build models to predict whether it is going to rain tomorrow or not based on weather data for that particular day in major cities of Australia. This comparative study is conducted concentrating on three aspects: modeling inputs, modeling methods, and pre-processing techniques. The results provide a comparison of various evaluation metrics of these machine learning techniques and their reliability to predict the rainfall by analyzing the weather data.

Keywords: Rain fall, weather, ML.

1. INTRODUCTION:

India's welfare is agriculture. The achievement of agriculture is dependent on rainfall. It also helps with water resources. Rainfall information in the past helps farmers better manage their crops, leading to economic growth in the country. Prediction of precipitation is beneficial to prevent flooding that saves people's lives and property. Fluctuation in the timing of precipitation and its amount makes forecasting of rainfall a problem for meteorological scientists. Forecasting is one of the utmost challenges for researchers from a variety of fields, such as weather data mining,

environmental machine learning, functional hydrology, and numerical forecasting, to create a predictive model for accurate rainfall. In these problems, a common question is how to infer the past predictions and make use of future predictions. A variety of sub-processes are typically composed of the substantial process in rainfall. It is at times not promis ing to predict the precipitation correctly by on its global system. Climate forecasting stands out for all countries around the globe in all the benefits and services provided by the meteorological department. The job is very complicated because it needs specific

numbers and all signals are intimated without any assurance. Accurate precipitation forecasting has been an important issue in hydrological science as early notice of stern weather can help avoid natural disaster injuries and damage if prompt and accurate forecasts are made. The theory of the modular model and the integrati2on of different models has recently gained more interest in rainfall forecasting to address this challenge. A huge range of rainfall prediction methodologies is available in India. In India, there are two primary methods of forecasting rainfall. Regression, Artificial Neural Network (ANN), Decision Tree algorithm, Fuzzy logic and team process of data handling are the majority frequently used computational methods used for weather forecasting The basic goal is to follow information rules and relationships while gaining intangible and potentially expensive knowledge. Artificial NN is a promising part of this wide field.

Rainfall prediction remains a serious concern and has attracted the attention of governments, industries, risk management entities, as well as the scientific community. Rainfall is a climatic factor that affects many human activities like agricultural production, construction, power generation, forestry and tourism, among others [1]. To this extent, rainfall prediction is essential since this variable is the one with the highest correlation with adverse natural events such as landslides, flooding, mass movements and avalanches. These incidents have affected society for years [2]. Therefore, having an appropriate approach for rainfall prediction makes it possible to take preventive and mitigation measures for these natural phenomena

To solve this uncertainty, we used various machine learning techniques and models to make accurate and timely predictions. These paper aims to provide end to end machine learning life cycle right from Data preprocessing to implementing models to evaluating them. Data Preprocessing steps include imputing missing values, feature transformation, encoding categorical features, feature scaling and feature selection. We implemented models such as Logistic Regression, Decision Tree, K Nearest Neighbour, Rule-based and Ensembles.

2. LITERATURE SURVEY

1.Climate Change and Human Health: Risks and Responses

The long-term good health of populations depends on the continued stability and functioning of the biosphere's ecological and physical systems, often referred to as lifesupport systems. We ignore this longestablished historical truth at our peril: yet it

is all too easy to overlook this dependency, particularly at a time when the human species is becoming increasingly urbanized and distanced from these natural systems. The world's climate system is an integral part of this complex of life-supporting processes, one of many large natural systems that are now coming under pressure from the increasing weight of human numbers and economic activities.

By inadvertently increasing the concentration of energy-trapping gases in the lower atmosphere, human actions have begun to amplify Earth's natural greenhouse effect. The primary challenge facing the world community is to achieve sufficient reduction in greenhouse gas emissions so as to avoid dangerous interference in the climate system. National governments, via the UN Framework Convention on Climate Change (UNFCC), are committed in principle to seeking this outcome. In practice, it is proving difficult to find a politically acceptable course of action-often because of apprehensions about possible short-term economic consequences.

This volume seeks to describe the context and process of global climate change, its actual or

likely impacts on health, and how human societies should respond, via both adaptation strategies to lessen impacts and collective action to reduce greenhouse gas emissions. As shown later, much of the resultant risk to human populations and the ecosystems upon which they depend comes from the projected extremely rapid rate of change in climatic conditions. Indeed, the prospect of such change has stimulated a great deal of new scientific research over the past decade, much of which is elucidating the complex ecological disturbances that can impact on human well-being and health—as in the following example.

The US Global Change Research Program (Alaska Regional Group) Assessment recently documented how the various effects of climate change on aquatic ecosystems can interact and ripple through trophic levels in unpredictable ways. For example, warming in the Arctic region has reduced the amount of sea ice, impairing survival rates for walrus and seal pups that spend part of their life cycle on the ice. With fewer seal pups, sea otters have become the alternative food source for whales. Sea otters feed on sea urchins, and with fewer sea otters sea urchin

populations are expanding and consuming more of the kelp that provides breeding grounds for fish. Fewer fish exacerbate the declines in walrus and seal populations. Overall, there is less food available for the Yupik Eskimos of the Arctic who rely on all of these species.

Global climate change is thus a significant addition to the spectrum of environmental health hazards faced by humankind. The global scale makes for unfamiliarityalthough most of its health impacts comprise increases (or decreases) in familiar effects of climatic variation on human biology and Traditional environmental health health. have been focused concerns long on toxicological or microbiological risks to health from local environmental exposures. However, in the early years of the twenty-first century, as the burgeoning human impact on the environment continues to alter the planet's geological, biological and ecological systems, a range of larger-scale environmental hazards to human health has emerged. In addition to global climate change, these include: the health risks posed by stratospheric ozone depletion; loss of biodiversity; stresses on terrestrial and ocean food-producing systems; changes in hydrological systems and the supplies of freshwater; and the global dissemination of persistent organic pollutants.

Climate change and stratospheric ozone depletion are the best known of these various global environmental changes. Human societies, however, have had long experience of the vicissitudes of climate: climatic cycles have left great imprints and scars on the history of humankind. Civilisations such as those of ancient Egypt, Mesopotamia, the Mayans, the Vikings in Greenland and populations European during the four centuries of Little Ice Age, all have both benefited and suffered from nature's great climatic cycles. Historical analyses also reveal widespread disasters, social disruption and disease outbreaks in response to the more acute, inter-annual, quasi-periodic ENSO (El Niño Southern Oscillation) cycle (1). The depletion of soil fertility and freshwater supplies, and the mismanagement of water catchment basins via excessive deforestation, also have contributed to the decline of various regional populations over the millennia.

2. Geomorphology, natural hazards, vulnerability and prevention of natural disasters in developing countries

The significance of the prevention of natural disasters is made evident by the commemoration of the International Decade for Natural Disaster Reduction (IDNDR). This paper focuses the role of on geomorphology in the prevention of natural disasters in developing countries, where their devastating impact has consequences. Concepts such as natural hazards, natural disasters and vulnerability have a broad range of definitions; however, the most significant elements are associated with the vulnerability concept. The latter is further explored and considered as a key factor in understanding the occurrence of natural disasters, and consequently, in developing and applying adequate strategies for prevention. Terms such as natural and human vulnerabilities are introduce and explained as target aspects to be taken into account in the reduction of vulnerability and for prevention and The of mitigation natural disasters. importance of the incorporation not only of geomorphological research, but also of geomorphologists in risk assessment and management programs in the poorest countries is emphasized.

Atmospheric and climatic hazards:
 Improved monitoring and prediction for disaster mitigation

The last few years have seen enormous damage and lossof life from climate and weather phenomena. The mostdamaging events have included the severe 1997/98 ElNiño (with its near-global impacts), HurricaneMitch, and floods in China in mid-1998. What have welearnt regarding the causes, variability, and predictability, of these phenomena? Can we predict theoccurrence of these extreme events, and therebymitigate their damage? This paper reviews what we havelearnt in the last decade or so regarding thepredictability of these climate and weather extremes. The view starts with the largest (El Niño) scales, and works towards the scale of individualthunderstorms. It focuses on the ofour improved practical outcomes knowledge with regard to decreasing theimpact of natural disasters, rather than describing indetail the scientific knowledge underlying theseoutcomes. The paper

concludes with a discussion of some of the factors that still restrict our ability tomitigate the deleterious effects of atmospheric andclimatic hazards.

4. Exploratory Data Analysis: the Best way to Start a Data Science Project

Exploratory Data Analysis is a set of techniques that were developed by Tukey, John Wilder in 1970. The philosophy behind this approach was to examine the data before building a model. John Tukey encouraged statisticians to explore the data, and possibly formulate hypotheses that could lead to new data collection and experiments. Today Data scientists and analysts spend most of their time in Data Wrangling and Exploratory Data Analysis also known as EDA. But what is this EDA and why it is so important? This article explains what is EDA and how to apply EDA techniques to a dataset.

EXICITING SYSTEM:

Rainfall prediction is important as heavy rainfall can lead to many disasters. The prediction helps people to take preventive measures and moreover the prediction should be accurate. There are two types of prediction short term rainfall prediction and long term Prediction mostly rainfall. short term prediction can gives us the accurate result. The main challenge is to build a model for long term rainfall prediction. Heavy precipitation prediction could be a major drawback for earth science department because it is closely associated with the economy and lifetime of human.

Dis Advantages

We can just do it by having the historical data analysis of rainfall and can predict the rainfall for future seasons. We can apply many techniques like classification, regression according to the requirements and also we can calculate the error between the actual and prediction and also the accuracy. Different techniques produce different accuracies so it is important to choose the right algorithm and model it according to the requirements

Proposed System

It's a cause for natural disasters like flood and drought that square measure encountered by individuals across the world each year. Accuracy of rainfall statement has nice importance for countries like India whose economy is basically dependent on

The agriculture. dynamic nature of atmosphere, applied mathematics techniques to provide sensible fail accuracy for precipitation statement. The prediction of precipitation using machine learning techniques may use regression. Intention of this project is to offer non-experts easy access to the techniques, approaches utilized in the sector of precipitation prediction and provide a comparative study among the various machine learning techniques.

Advantages

1. It is a powerful technique for testing relationship between one dependent variable and many independent variables.

2. It allows researchers to control extraneous factors. 3. Regression asses the cumulative effect of multiple factors.

4. It also helps to attain the measure of error using the regression line as a base for estimations.

3. METHODOLOGY

The predictive model is used to prediction of the precipitation. The first step is converting data in to the correct format to conduct experiments then make a good analysis of data and observe variation in the patterns of rainfall. We predict the rainfall by separating the dataset into training set and testing set then we apply different machine learning approaches (MLR, SVR, etc.) and statistical techniques and compare and draw analysis over various approaches used. With the help of numerous approaches we attempt to minimize the error.

MODULES:

1.Add Product Details

To build project I used some sample products image to train product identification models

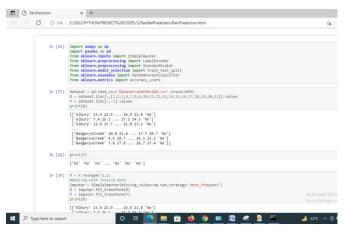
2.Train Model

In this Module screen train model generated with 100% accuracy and now show product to web cam.

3. Add/Remove Product from basket

To allow application to identify product image and then show in text area and if we again show same product then application will remove from text area.

OPERATION: Packeages

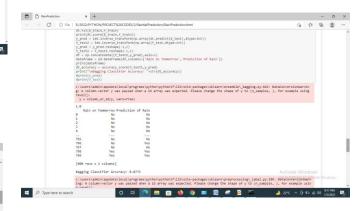


C O File	E/2022/P1THON/PROJECT%20CODES/2/RainfalPrediction/RainPrediction.html	16	th.	6	
1+ [136]:	<pre>fronting.mail. Listerier = two devections for (entering the state of the state Listerier = frite (state), vision (preset = int, invest presentation area (Listerier and Lister), and presentation preset = int, invest presentation area (Listerier and Lister), and presentation preset = interpresentation area (Listerier and Lister), and presentation preset presentation area (Listerier and Lister), and Lister and Lister preset presentation area (Listerier and Lister), and Lister and Lister preset preset area (Listerier and Lister), and Lister and Lister and Lister preset preset preset preset preset preset presentation area (Lister), and Lister and Lister preset preset preset preset preset pres</pre>				
	c:\users\admin\appdata\local\programs\python\python\T\lD\site-packages\ipython lawscher.py:3: DetaConversions n/vector y was passed when a id array was expected. Plaus change the shape of y to (m_samples_), for example us This is separate from the lpython(python(appdate so we can world doing theoris will	arning: A c ing ravel()	malam 		
P Type here to search	о н 💽 🚍 前 🍁 🏟 💷 🖻 🚅 🖉 📨 🌙	2210 ~ 0	\$5.0	100	ę

Random Forest



Bagging Classifer

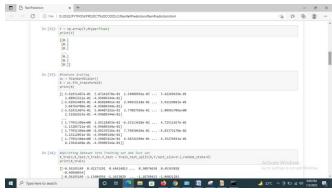


Gradient Boosting

Analysis

- → 0	iction	+ E/2022/P/THON/PR0/ECT%20C0DE5/2/RainfallPrediction/RainPrediction.html	16 P	⊕ (8
		$\begin{aligned} & \textit{Hiscoling bitset} \\ & \textit{ist} - \textit{ishifticode}() \\ & \textit{ist} - \textit{ishifticode}() \\ & i$			
		$ \begin{array}{c} [23,12,23,2,150,21,4,9] \\ [67,24,25,,12,24,4,9] \\ [67,24,25,,12,24,4,9] \\ [67,24,25,,12,24,24,3] \\ [16,42,24,,12,27,24,7] \\ [16,42,24,,12,27,24,7] \\ [16,42,24,,12,27,24,7] \\ [16,42,24,,12,27,44] \\ [16,42,,12,27,44] \\ [16,42,,12,,12,12,,12,12,,12,,12,,12,,12,,12,,12,,12,,12,,12,,12,,12,,$			
	In [31]:	print(Y)			
		([e] [e] [e] [e] [e] [e]			
	In [32]:	Y = np.array(Y,dtype-float) print(Y)	Activate Window Go to Settings to active		10/5.
		[[0,]	🔤 🌛 22°C ^ ĝ 40 @		

Training



Algorithms

C	③ File	E:/2022/PY1	THON/PROJECT%20COD	ES/2/RainfallPredic	tion/RainPrediction.8	ntml			10	£≞	۲
	In [143]:	dt = Grad dt.fit(X_ print(dt. y_pred = Y_test3 = y_pred = Y_test3 = df = np.c dataframe print(dat dt_accura	cy = accuracy_score(Gradient Boosting Ac pred)	<pre>in(n_estimators) (n)(n_array(dt)) rm(np_array(Y_1) (1) (1) y_pred),axis=1 columns=["Rain ((Y_test1,y_pred))</pre>	s=170,max_depth=1 predict(X_test),d test,dtype=int))) on Tommorrow','Pr)	type-int))	u,])				
		c:\users\ column-ve ().		ien a 1d array w				:1454: DataConversion ples,), for example			
		c:\users\ column-ve (). y = col 0.878125 Rain 0 1 2 3 4 4 795 796 797 798 799	ctor y was passed wh	ien a 1d array v rue)), for example		wel	

Xgboost

C V C I I EXCEPTION/PROCENS/CONSTRUCT/Absolution/Index Constants V I I I EXCEPTION/PROCENS/CONSTRUCT/Absolution/Index Constants I I I I I I I I I I I I I I I I I I I		a	() Dia	E/2022/EVTHO	N/RRO IECTS 20CODI	15/2/Rainfal Prediction /S	PainPrediction html			~	Gh	0
<pre>Ing i A colume-writer y was passed when is 14 arrays was reported. Plane that maps of y"to (c_imbles,), for example with</pre>			() Hit	Y_test4 = let #print(y_pres #print(Y_test y_pred = y_pt Y_test4 = Y_t dataframe = t print(datafrated d_accuracy	s.inverse_transfo) t) red.reshape(-1,1) test4.reshape(-1,2) atenate((Y_test4, od.DataFrame(df,c me) = accuracy_score(rm(np.array(Y_test,) y_pred),axis=1) olumns=['Rain on To y_test1,y_pred)	dtype=int))	Rain'])	10	μ		*
Kain on Tommerrae Prediction of Kain 1 No 2 No 3 No 4 No 755 No 756 No 708 No 799 No 799 Yes 700 No				<pre>ing: A column g ravel(). y = column c:\users\adm ing: A column g ravel().</pre>	or_id(y, warn=Tr in\appdata\local\ orvector y was pa	ssed when a 1d arra we) programs\python\pyt ssed when a 1d arra	y was expected. Please o	change the shape of y to (n	_samples,), for example	usin		
[860 rows x 2 columns] Activate Windows				Rain on 1 9 1 2 3 4 795 796 797 798	No No No No No No No Yes	No No No No No No No Ves						
X68oost Accuracy: 0.875 Go to Settings to activate Wi												

CONCLUSION

In this paper, we explored and applied several preprocessing steps and learned there impact on the overall performance of our classifiers. We also carried a comparative study of all the classifiers with different input data and observed how the input data can affect the model predictions.

We can conclude that Australian weather is uncertain and there is no such correlation among rainfall and the respective region and time. We figured certain patterns and relationships among data which helped in determining important features. Refer to the appendix section. As we have a huge amount of data, we can apply Deep Learning models such as Multilayer Perceptron, Convolutional Neural Network, and others.

It would be great to perform a comparative study between the Machine learning classifiers and Deep learning models.

REFERANCES

1. World Health Organization: Climate Change and Human Health: Risks and Responses. World Health Organization, January 2003

2. Alcntara-Ayala, I.: Geomorphology, natural hazards, vulnerability and prevention of natural disasters in developing countries. Geomorphology 47(24), 107124 (2002)

3. Nicholls, N.: Atmospheric and climatic hazards: Improved monitoring and prediction for disaster mitigation. Natural Hazards 23(23), 137155 (2001)

4. [Online] InDataLabs, Exploratory Data Analysis: the Best way to Start a Data Science Project. Available: https://medium.com/@InDataLabs/ why-start-adata-science-project-with-exploratory-dataanalysis-f90c0efcbe49

5. [Online] Pandas Documentation. Available: https://pandas.pydata.org/ pandasdocs/stable/reference/api/pandas.get_dummies.ht ml

6. [Online] Sckit-Learn Documentation Available: https://scikit-learn.org/ stable/modules/generated/sklearn.feature_extract ion.FeatureHasher. html

7. [Online] Sckit-Learn Documentation Available: https://scikit-learn.org/ stable/modules/generated/sklearn.preprocessing. MinMaxScaler.html

8. [Online] Sckit Learn Documentation Available: https://scikit-learn.org/ stable/modules/generated/sklearn.feature_selectio n.SelectKBest.html

9. [Online] Raheel Shaikh, Feature Selection Techniques in Machine Learning with Python Available: https://towardsdatascience.com/ feature-selection-techniques-in-machine-learningwith-python-f24e7da3f36e

10. [Online] Imbalanced-learn DocumentationAvailable:https://imbalanced-learn.readthedocs.io/en/stable/introduction.html

11. V. Veeralakshmi and D. Ramyachitra, Ripple Down Rule learner (RIDOR) Classifier for IRIS Dataset. Issues, vol 1, p. 79-85.

12. [Online] Aditya Mishra, Metrics to Evaluate your Machine Learning Algorithm Available: https://towardsdatascience.com/ metrics-toevaluate-your-machine-learning-algorithmf10ba6e38234