Studying Methods for Heat Index Forecasting in Mumbai, Colaba

Sahil Rane*

Department of Computer Science and Mathematics, Harvey Mudd College, Mumbai, Maharashtra, India

Corresponding Author*

Sahil Rane Department of Computer Science and Mathematics Harvey Mudd College Mumbai, Maharashtra, India Tel: 7506941156 Email: sahilrane18@gmail.com

Copyright: © 2022 Rane S. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Received: 09-Oct-2020, Manuscript No. JCWF-20-001-PreQc-20; Editor assigned: 13-Oct-2020, PreQC No. JCWF-20-001-PreQc-20 (PQ); Reviewed: 28-Oct-2020, QC No. JCWF-20-001-PreQc-20; Revised: 07-Jul-2022, Manuscript No. JCWF-20-001-PreQc-20 (R); Published: 15-Jul-2022

Abstract

Background: Heat Index is an important measure in determining the safety of temperature conditions for humans. Extreme heat can lead to dangerous, even deadly, health consequences, including heat stress and heatstroke. Thus, there is a need to predict the Heat Index accurately in order to warn individuals about such conditions so that they take appropriate precautions.

Methodology: In this paper, we look at weather data in Mumbai from 2008 to 2020 and attempt to come up with predictive models for the Heat Index. We carry out feature selection first in order to efficiently use a variety of algorithms to develop predictive models. We use various mathematical techniques such as: Multiple Linear Regression (MLR), Simple Exponential Smoothing (SES), Artificial Neural Networks (ANN), and Auto-Regressive Integrated Moving Average (ARIMA) models to predict the heat index. The experimental results are evaluated and compared using the Root Mean Square Error (RMSE).

Results: On experimenting with all four models, it was discovered that the ARIMA model yields the best predictive model having a RMSE of 0.354654 on testing data. This model is also concluded to be optimal as the residuals of this model are a gaussian white noise. Furthermore, the poor performance of MLR indicates that temperature cannot be accurately modelled through a linear function of the variables considered.

Keywords: Multiple linear regression • Simple exponential smoothing • Artificial neural networks • Root mean square error

Introduction

Heat Index in Celsius is often referred to as the apparent temperature or the feels like temperature. Thus, the Heat Index serves as a measure of heat risk one is exposed to. Certain precautions are necessary to be taken in order to remain safe from the adverse effects of excessive heat exposure. The Heat Index combines the temperature and the relative humidity of the atmosphere. The humidity is necessary because when the body gets hot, it begins to perspire in order to cool the body. When the relative humidity is high the rate of evaporation decreases and thus the body feels warmer and the regulatory cooling processes are unable to take place at the same rate.

This paper attempts to come up with a reliable predictive model for the Heat Index in Mumbai. Although there is several research papers based on temperature predictive models in general, there weren't any papers that focus on Mumbai specifically. Furthermore, the approach to predicting the HeatIndex instead of the temperature is unique. The models cited in several research papers, for temperature prediction in areas all around the world, were not able to improve on an RMS error of approximately 1 in most cases. Thus, this paper aims to improve the RMS error of these previous models for our dataset of Mumbai. Multiple linear regression and neural networks have been commonly used for temperature predictions with artificial neural networks being the most common. This paper in addition to these techniques explores a variety of time series analysis techniques such as exponential smoothing and auto-regressive integrated moving average models [1].

Often temperatures in Mumbai are underestimated and necessary precautions are not taken while exiting one's house. Due to the adverse effects that this heat can have on the health of individuals, it is necessary that steps be taken at an organisational level to issue warnings to citizens so that they can take the necessary precautions and steps before leaving their house. This idea was inspired by the IMD's Heat Action Plan, which has been implemented in Ahmedabad. While the temperatures in Mumbai may not be as severe, precautions are necessary in order to protect and inform individuals about the adverse effects of heat.

Methodology and Data Analysis

Feature selection

We use the Pearson Correlation Coefficient (r) as a feature selection mechanism (Table 1). In order to avoid overfitting of the data we will only be using the features that have |r|>0.5.

Table	1.	Gender	distribution	of	the stud	v subjects	(n=100)).
-------	----	--------	--------------	----	----------	------------	---------	----

Variable (lag 1)	Pearson correlation coefficient (r)
Maxtemp C	0.6266314
Mintemp C	0.645104
Avgtemp C	0.8032574
Totalprecip MM	-0.09843071
Wind speed Kmph	0.08042873
Sun hour	0.1232874
winddirdegree	0.3028105
humidity	0.1893674
Visibility Km	0.102261
Pressure MB	-0.3131072
cloud cover	-0.06016056
DewPoint C	0.6104005
WindChill C	0.7757908
WindGust Kmph	-0.1435159

The pearson correlation co-efficient is used to measure correlation between different sets of data to understand how strong the correlation is between two variables. Above we have calculated and tabulated the pearson correlation co-efficients between the HeatIndexC of the day with a variety of variables such as total precipitation, pressure, sun hours etc. with a lag of 1 day.

$$r = \frac{\sum_{i} (x_{i} - \overline{x}) (y_{i} - \overline{y})}{\sqrt{\sum_{i} (x_{i} - \overline{x})^{2}} \times \sqrt{\sum_{i} (y_{i} - \overline{y})^{2}}}$$

By the above criteria we select the variables: maxtempC (maximum temperature in Celcius), mintempC (minimum temperature in Celcius),

avgtempC (average temperature in Celcius), DewPointC (the atmospheric temperature below which water droplets begin to condense and dew can form: it depends on both the pressure and humidity), WindChillC (Combination of windspeed and temperature).

WindChillC =
$$13.12 + 0.6215 \times T - 11.37 \times (V^{0.16}) + 0.3965T (V^{0.16})$$

T = Temperature in degree Celsius

V = Wind velocity in kilometres per hour

Multivariate linear regression

First, in order to perform a multiple linear regression on our data it is necessary to carry out feature scaling as the range for each variable differs significantly. If feature scaling was not done then we variables with maximum range will dominate in the training of the regression model. We will bound all our variables in the interval (0,1). We use min-max normalisation technique on our data [2].

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

We will use Ordinary least Squares (OLS) regression

Let's call the temperature on the ith day some Ti:

$$T_{1} = \beta_{0} + \beta_{1}x_{11} + \beta_{2}x_{21} + \beta_{3}x_{31} + \beta_{4}x_{41} + \beta_{5}x_{51} + \varepsilon_{1}$$

$$T_{2} = \beta_{0} + \beta_{1}x_{12} + \beta_{2}x_{22} + \beta_{3}x_{32} + \beta_{4}x_{42} + \beta_{5}x_{52} + \varepsilon_{2}$$

$$T_{3} = \beta_{0} + \beta_{1}x_{13} + \beta_{2}x_{23} + \beta_{3}x_{33} + \beta_{4}x_{43} + \beta_{5}x_{53} + \varepsilon_{3}$$

$$\vdots$$

$$T_{n} = \beta_{0} + \beta_{1}x_{1n} + \beta_{2}x_{2n} + \beta_{3}x_{3n} + \beta_{4}x_{4n} + \beta_{5}x_{5n} + \varepsilon$$

$$x_{ij} : This is the value of the ith feature on the jth day$$

$$\varepsilon_{n} : error term on the ntn day$$

$$\beta_{i} : coefficient of the ith feature in our regression.$$

When we vectorise these equations we get:

	1	x_{11}	x_{21}	x_{31}	x_{41}	x_{51}		β_0	
$\vec{T} =$	÷	÷	÷	÷	÷	÷	×	÷	$+\vec{\varepsilon}$
	1	x_{1n}	x_{2n}	x_{3n}	x_{4n}	x_{5n}		β_5	

We can write this as:

$$\vec{T} = X\vec{\beta} + \vec{\varepsilon} \Longrightarrow \vec{\varepsilon} = \vec{T} - X\vec{\beta}\vec{T} = X\vec{\beta} + \vec{\varepsilon} \Longrightarrow \vec{\varepsilon} = \vec{T} - X\vec{\beta}$$

Since we are using OLS regression we want to minimize $\|\vec{\varepsilon}\|^2 = \sum_{i=1}^n \varepsilon_i^2$ $\sum_{i=1}^n \varepsilon_i^2 = \varepsilon' \times \varepsilon = (\vec{T} - X\vec{\beta}) \times (\vec{T} - X\vec{\beta})$

We use gradient descent with simultaneous update to minimise the error function.

$$\frac{\partial}{\partial \beta_j} \sum_{i=1}^n (x_{ji}\beta_j - T_i)^2 = \sum_{i=1}^n [(x_{ji}\beta_j - T_i) \times x_{ji}]; \forall j \in \mathbb{Z}$$

Algorithm:

Repeat using simultaneous update

$$\{\beta_j := \beta_j - \alpha \times \sum_{i=1}^n [(x_{ij}\beta_j - T_i) \times x_{ij}]\}$$

We use the R programming language to run this program and obtain the multivariate linear regression analysis.

Code for multivariate linear regression used

#First we need to read the data from the csv file with headers

data <- read.csv("weather_data_24hr_HI.csv", header = T)

#remove null row

data <- data(-4265,)

#remove columns based on feature selection

data <- data(,c(3,4,6,8,26,28))

we need to normalize the data (max-min normalization)

data\$maxtempC<-(data\$maxtempC-min(data\$maxtempC))/
(max(data\$maxtempC)-min(data\$maxtempC))</pre>

data\$avgtempC<-(data\$avgtempC-min(data\$avgtempC))/
(max(data\$avgtempC)-min(data\$avgtempC))</pre>

data\$mintempC<-(data\$mintempC-min(data\$mintempC))/
(max(data\$mintempC)-min(data\$mintempC))</pre>

data\$WindChillC<-(data\$WindChillC-min(data\$WindChillC))/
(max(data\$WindChillC)-min(data\$WindChillC))</pre>

data\$DewPointC<-(data\$DewPointC-min(data\$DewPointC))/
(max(data\$DewPointC)-min(data\$DewPointC))</pre>

Data Partition

Partitioning the Data into testing and training data

If we repeat the learning, we get the same result

set.seed(222)

ind <- sample(2,nrow(data), replace = T, prob =c(0.7,0.3))

training <- data(ind==1,)

testing <- data(ind==2,)

Training the Linear regression model

MLR <- Im(maxtempC~., data = training)

summary(MLR)

Prediction on training

output <- predict(MLR, training(,-1))

df <- data.frame(output, training(,1))

df

sum((df(,1)-df(,2))^2)/nrow(training)

Prediction on testing

output <- predict(MLR, testing(,-1))</pre>

df <- data.frame(output, testing(,1))

df

sum((df(,1)-df(,2))^2)/nrow(testing)

Explanation of code: We first import the comma-splitted values containing the raw data for our analysis. We then keep only the columns that were deemed statistically significant by the pearson correlation test. We then carry out the max-min normalisation of the data as specified above. We then partition 70% of the data for training our regression model and 30% of the data to test the regression model. We create the multivariate ordinary least squares linear regression model using gradient descent as outlined above. We then calculate the sum of squared error for both the training and the testing data using the regression model [3].

Results: The regression equation obtained by our program is:

h_ β (x)=15.910+ 2.323×x_1-1.906×x_2+4.246×x_3+10.469×x_4+12.8 61×x_5

- x_1:normalised maxtempC lag 1
- x_2: normalised mintempC lag 1
- x_3: normalised avgtempC lag 1
- x_4:normalised DewPointC lag 1
- x_5:normalised WindChillC lag 1

Metric of analysis for results: We will now be analysing the results produced by the multivariate linear regression model. For analysis we will be finding the square root of mean sum of squared errors (RMS errors) returned by the program. We define the function of our prediction model (hypothesis function) as $h_{-\beta}(x)$. The formula for the RMS error is.

RMS error= $\sqrt{(\sum_{i=1}^{n} (h_{\beta}(x)-T_{i})^{2})/n}$

For the training data the sum of squared errors was: 2.261942

For the testing data the sum of squared errors was: 2.540919

Since the behaviour of our model is similar for the training and testing data we can conclude that our model is not over or underfitting the dataset. However, the RMS error is still quite high.

Simple exponential smoothing

We will also attempt to forecast the HeatIndex using time series analysis techniques. We will begin with simple exponential smoothing as a technique because there is no clear trend or seasonality in the data when the time series is plotted (Figure 1). We don't use Holt's Exponential smoothing or Holt-Winter Exponential Smoothing as our data does not have clear increasing/decreasing trends or any visible seasonality. In simple exponential smoothing, we say the predicted observation is a weighted average of previous observations [4]. While calculating the weighted averages, the weights decrease exponentially as lag increases, that is, the smallest weights are associated with the oldest observations:



Figure 1. Time series plot.

Notation:

l_o:start value

y,:actual temperature at time t

 $\overline{y_t}$: forecasted temperature at time t

 $0 \le \alpha \le 1$

Weighted average form of exponential smoothing:

 $\overline{y_{t+1}} = \alpha y_t + (1 - \alpha) \overline{y_t}$

Let's try to arrive at the above generalisation by considering smaller examples.

 $v_{1} = l_{1}$ $\overline{y_{2}} = \alpha y_{1} + (1 - \alpha) l_{0}$ $\overline{y_{3}} = \alpha y_{2} + (1 - \alpha) \overline{y_{2}}$ $\overline{y_{4}} = \alpha y_{3} + (1 - \alpha) \overline{y_{3}}$ \vdots $\overline{y_{t}} = \alpha y_{t-1} + (1 - \alpha) \overline{y_{t-1}}$ $\overline{y_{t+1}} = \alpha y_{t} + (1 - \alpha) \overline{y_{t}}$

When we substitute each equation in the next we get:

$$\overline{y_3} = \alpha y_2 + (1 - \alpha) (\alpha y_1 + (1 - \alpha) l_0)$$

= $\alpha y_2 + (1 - \alpha) (\alpha) y_1 + (1 - \alpha)^2 l_0$
 $\overline{y_4} = \alpha y_3 + (1 - \alpha) (\alpha y_2 + (1 - \alpha) (\alpha) y_1 + (1 - \alpha)^2 l_0)$

$$= \alpha y_{3} + (1-\alpha)(\alpha) y_{2} + (1-\alpha)^{2} (\alpha) y_{1} + (1-\alpha)^{3} l_{0}$$

$$\vdots$$

$$\overline{y_{t+1}} = \sum_{i=0}^{t-1} (\alpha (1-\alpha)^{i} y_{t-i}) + (1-\alpha)^{t} l_{0}$$

This is the weighted average form of the simple exponential smoothing model. We want to find the value of α that minimises the Sum Of Squared Errors (SSE), therefore, we use the gradient descent algorithm to minimise the squared error.

$$SSE = \sum_{i=1}^{t} (y_i - \overline{y_i})^2 = \sum_{i=1}^{t} e_i^2$$

We use the R programming language to run gradient descent in order to minimise error.

Code for simple exponential smoothing used

#first read the comma splitted values containing the dataset

data <- read.csv("weather_data_24hr_master1.csv", header = T)

#filter only maxtempC out which is relevant to our timeseries

data <- data(,18)

#remove NA values

data <- data(-4265)

data

#Let's partition the data into testing and training data.

training <- data(1:2990)

testing <- data(2991:4264)

#Create a timeseries using the maxtempC data

timeseries <- ts(training, frequency = 365, start = c(2008,183))

timeseries

#plot the timeseries

plot.ts(timeseries)

install.packages("TTR")

library("TTR")

#Simple exponential Smoothing

timeseriesforecasts <- HoltWinters(timeseries, beta=FALSE, gamma=FALSE, I.start = 31)

timeseriesforecasts

timeseriesforecasts\$fitted

plot(timeseriesforecasts)

#calculate the sum of squared errors for our forecasts

a<-(timeseriesforecasts\$SSE)

#calculate the mean of the SSE

MSE<- a/length(training)

#calculate the RMS error

sqrt(MSE)

testmodel<-HoltWinters(timeseriestest, alpha = 0.7159879, beta = FALSE, gamma = FALSE)

#calculate the sum of squared errors for our forecasts

a<-testmodel\$SSE

#calculate the mean of the SSE

MSE<-a/length(testing)

#calculate the RMS error

sqrt(MSE)

Explanation of code: We first import the comma-splitted values containing the raw data for our analysis. We then remove all columns except the HeatIndexC for our timeseries. We partition our dataset into training data and testing data. We train our simple exponential smoothing model on the first 70% of our data. We test this model using the remaining 30% of the data. We then convert the vector into a timeseries. We then plot the timeseries and run simple exponential smoothing on our timeseries. This gives us the minimum value of a. We also find the RMS error of our model on the training data. We then use the same value of α to run our simple exponential smoothing model on the testing data and calculate the RMS error for the same.

Results: By running gradient descent we get the result that the error is minimised when the learning parameter α = 0.7159879

This makes out weighted average exponential smoothing equation to be as follows:

$$\overline{y_{t+1}} = \sum_{i=0}^{t-1} ((0.7159879)(0.2840121)^i y_{t-i}) + (0.2840121)^i .31$$

Metric of analysis for results: We will now be analysing the results produced by the simple exponential smoothing model. One of the main issues that arises in the use of this model is that it does not perform well on sudden fluctuations in the data whereas it performs very well for gradual increases or decreases in our data. The high alpha value indicates high dependence on previous day values of the max temperature. The high alpha also indicates that as lag increases data beyond lag=3 $(0.7159879)(0.2840121)^3 = 0.01640273287$ become too small to have relevance to our forecasted value. To determine how well our model fits the data we calculate the RMSE for the testing data and the training data [5]. We define our prediction model function for simple exponential smoothing to be: $\overline{y_{\alpha}}(t)$ (t). The temperature on the tth day is defined as M(t).

Therefore we get that the RMS error is:

 $RMS \, error = \sqrt{\frac{\sum_{t=1}^{T} \left(\overline{y_a}(t) - M(t)\right)^2}{T}}$ For the training data the RMS error was: 1.326023

For the testing data the RMS error was: 1.423472

The model gives very similar values for the RMS training and testing data, thus overfitting of data is not an issue for our model.

Artificial neural networks

A neural network is inspired from the model of a human brain consisting of neurons. When neural networks learn they independently find a variety of connections in the data which helps with complicated predictions when we have large datasets with several variables. Each neuron receives the values of the variables (features) from the training set and calculates a weighted average of these values. The result of this calculation is passed through a non-linear activation function.

For the mth neuron we supply the vector x of training examples and it calculates the weighted average and returns a zm value. In this scenario b is a bias constant that is added to the outcome of each neuron.

$$z = \sum_{i=1}^{n} (w_i x_i) + b \; (\forall j = 1, 2, 3, \dots, m)$$

The activation function g(z) is applied on each z value to give our forecast value (\hat{y}) . Thus the output of each neuron is passed to the activation function.

Without an activation function, the neural network would simply return a linear function thus not being able to model complex data while having small error.

The way a neural network is trained is by the value of the loss function. We use the sum of squared error as our error function thus during the training our model minimises the error of the neural network. In our learning the values of the weights and the bias parameter are changed in order to minimise the error function. We calculate the partial derivative of the loss function in order to arrive at a minimum value for the error. We use the backpropagation algorithm in order to train our neural network [6].

For our regression problem the loss function would be Mean Squared Error, which squares the difference between actual (y_i) and predicted value (ŷ_i).

Sum of squared error
$$(C) = \sum_{i=1}^{m} (y_i - \widehat{y_i})$$

Using chain rule we get the following:

$$\frac{\partial C}{\partial w_i} = \frac{\partial C}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial W_i}$$

First lets calculate $\overline{\partial \hat{v}}$:

$$\frac{\partial C}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} \times \left(\sum_{i=1}^{n} \left(y_i - \widehat{y}_i \right)^2 \right) = 2 \times \sum_{i=1}^{m} \left(y_i - \widehat{y}_i \right)$$

Now let's calculate the $\frac{\partial \hat{y}}{\partial r}$:

$$\frac{\partial \hat{y}}{\partial z} = \frac{\partial}{\partial z} \times g(z)$$

$$= \frac{\partial}{\partial z} \left(\frac{1}{1+e^{-z}}\right)$$

$$= \frac{e^{-z}}{\left(1+e^{-z}\right)^2}$$

$$= \frac{1}{\left(1+e^{-z}\right)} \times \left(1-\frac{1}{1+e^{-z}}\right)$$

 $=g(z)\times(1-g(z))$

Finally let's calculate the $\frac{\partial z}{\partial w}$:

$$\frac{\partial z}{\partial w_i} = \frac{\partial}{\partial w_i} \times (z)$$
$$= \frac{\partial}{\partial w_i} \sum_{i=1}^n (w_i x_i) + b$$
$$= x_i$$

Therefore, we get:

$$\frac{\partial C}{\partial w_i} = x_i \times g(z) \times (1 - g(z)) \times 2 \times \sum_{i=1}^m (y_i - \widehat{y_i})$$

We also calculate $\frac{\partial C}{\partial b}$ which we can get from the above formula as the input for the bias operator is 1.

$$\frac{\partial C}{\partial b} = g(z) \times (1 - g(z)) \times 2 \times \sum_{i=1}^{m} (y_i - \widehat{y_i})$$

After backpropagation is carried out we focus on optimisation which is done through gradient descent. For this we define our learning rate as α .

Repeat until convergence with simultaneous update:

$$\begin{cases} w_i := w_i - \left(\alpha \times \frac{\partial C}{\partial w_i}\right) \\ b := b - \left(\alpha \times \frac{\partial C}{\partial b}\right) \end{cases}$$

We now implement this algorithm using R.

Code for artificial neural networks used

#First we need to read the data from the csv file with headers

data <- read.csv("weather_data_24hr_master1.csv", header = T)</pre>

#remove null row

data <- data(-4265,)

#remove columns based on feature selection

data <- data(,c(4,7,9,17,18,20))

we need to normalize the data (max-min normalization)

data\$maxlag1<-(data\$maxlag1-min(data\$maxlag1))/ (max(data\$maxlag1)-min(data\$maxlag1))

d a t a \$ a v g l a g 1 < - (d a t a \$ a v g l a g 1 - min(d a t a \$ a v g l a g 1))/ (max(data\$avglag1)-min(data\$avglag1))

data\$cloudcover<-(data\$cloudcover-min(data\$cloudcover))/
(max(data\$cloudcover)-min(data\$cloudcover))</pre>

data\$HeatIndexC<-(data\$HeatIndexC-min(data\$HeatIndexC))/
(max(data\$HeatIndexC)-min(data\$HeatIndexC))</pre>

data\$WindChillC<-(data\$WindChillC-min(data\$WindChillC))/ (max(data\$WindChillC)-min(data\$WindChillC))

#Data partition to divide our data into training anad testing data (70% training 30% testing)

#we set seed in order to be able to repeat the learning

set.seed(222)

ind <- sample(2,nrow(data), replace = T, prob =c(0.7,0.3))

training <- data(ind==1,)</pre>

testing <- data(ind==2,)

#install the neural network packages in R

install.packages("neuralnet")

library(neuralnet)

#We create a neural network n trained on the training data

#This neural network has the error function as the sum of squared error

#It has the activation function as the sigmoid function

#It has 2 neurons

n <- neuralnet(maxtempC~.,

data = training,

hidden = 2,

stepmax = 9999999,

err.fct ='sse',

act.fct = 'logistic',

linear.output = T)

n

#We plot our trained neural network

plot(n)

We calculate the RMS error for our neural network on the training and testing dataset

output <-compute(n, training)</pre>

p1 <- output\$net.result

sqrt(sum((training\$maxtempC-p1)^2)/nrow(training))

max((training\$maxtempC-p1))

output <- compute(n, testing)

p2 <- output\$net.result

sqrt(sum((testing\$maxtempC-p2)^2)/nrow(testing))

Explanation of code: We first import the csv file containing the raw data for our analysis. We then retain all the feature columns from our feature selection and remove all NA values. We first carry out feature scaling of our data through min-max normalisation. We then partition our dataset into training data (70%) and testing data (30%). We train an artificial neural network with 2 neurons to fit our training data. We find the RMS error of our model on the training data. We then run our model on the testing data and calculate the RMS error [7].

Results: The neural network plot has been included below with the weights on each "wire" connecting our input to the neurons and our neurons to the output layer (Figure 2).



Figure 2. Neural network plot.

The value of the bias operator for the hidden layer is -1.95219 and -5.32014 while the bias operator for the output layer is 23.67867.

Metric of analysis for results: We will now analyse the results of the neural network model. The use of neural networks allows us to predict the data using a non-linear hypothesis function. The high weights on the WindChillC, maxtempC, and DewPointC indicate that these are the dominant features in our learning. A variety of number of neurons was tried out. 2 neurons were chosen as larger values seemed to overfit our data and would have an overly complicated hypothesis function. To determine how well our model fits the data we calculate the RMSE for the testing data and the training data. We define our prediction model function for our neural network to be: $\hat{y}_i(x)$. The actual temperature is represented by y_i.

Therefore we get the formula for the RMS error to be:

$$RMS\,error = \sqrt{\frac{\sum_{i=1}^{n} \left(\hat{y}_{i}\left(x\right) - y_{i}\right)^{2}}{n}}$$

For the training data the RMS error was: 1.304715

For the testing data the RMS error was: 1.379394

Since the RMS error for the training and testing data using our model is similar we can conclude that our model does not overfit our data. It also performs slightly better than the exponential smoothing model and significantly better than the multivariate regression model.

Auto-Regressive Integrated Moving Average (ARIMA)

In order to apply the ARIMA technique we require our timeseries data to be stationary. This means that our data has a constant mean, constant variance (deviation from the mean), and no seasonality. It is necessary to confirm whether our time series is stationary using the Augmented Dickey-Fuller test (ADF test) and the Kwiatkowski-Phillips-Schmidt-Shin test (KPSS test) [8].

The ADF test is a unit root test. The mathematics and proofs related to the ADF test will not be described in this paper but can be found in [1]. For our purposes we only need to know the null and alternative hypothesis as defined by the test.

H_a: (Null hypothesis) The given data is non-stationary

 H_{1} : (Alternative hypothesis) The given data is stationary

On running the ADF test on our time series we get a p-value ≤ 0.01

Since the p-value is less than the significance level of 0.05, we can safely reject the null hypothesis and conclude that our data is stationary. However, for large datasets the adf test can be erroneous rejecting the null hypothesis a vast majority of the times. Hence, we use the KPSS test to confirm that our data is stationary. We use the KPSS test for level stationarity.

For the KPSS test the hypotheses are as follows:

 H_{0} : (Null hypothesis) The data is level stationary

 H_{1} : (Alternative hypothesis) The given data is level non-stationary

On running the KPSS test we get the p-value ≥ 0.1

Since the p-value is greater than the significance level of 0.05 we cannot reject the null hypothesis which supports our previous conclusion that our data is stationary. Since our data is stationary we can proceed by using auto-regressive integrated moving average processes on our data to come up with an acc urate forecasting model for the data.

Since our data is stationary we do not require an ARIMA model and we can directly use an ARMA model (Auto-Regressive Moving Average).

Before we can define an ARMA model we will define white noise.

Definition 4.1: White noise: It is an identically and independently distributed stochastic process {X_t , t \in Z} with mean zero such that there is no serial correlation between values of stochastic process in the present and past.

A gaussian white noise is when our stochastic process X, is:

 $X_t \sim \mathcal{N}(0,\sigma^2)$

White noise timeseries cannot be predicted as they are a sequence of random numbers. If the series of forecast errors are white noise it suggests that our model cannot be improved.

We used the ARIMA model on our timeseries and the residuals of our model were a gaussian white noise with mean 0 and variance 1 (Figure 3). However, we want to improve our model despite this white noise. To do this we will attempt to smooth our timeseries in order to reduce the white noise in our timeseries. We will use a triangular moving average smoothing technique in order to minimise the error caused by this white noise. For this we will defined our new smooth timeseries to be ⁻y and our original time series to be y [9].



Figure 3. Plotting the residuals of our ARIMA (2,0,1) model.

Mathematically our smoothed moving average time series of order 2 will be defined as:

$$\overline{y_t} = \frac{y_t + y_{t-1}}{2}$$

To get the triangular moving average (\mathcal{Y}_t) to smooth our time series we apply the moving average of order 2 again.

 $\overline{\overline{y_t}} = \frac{\overline{y_t} + \overline{y_{t-1}}}{2}$

An auto-regressive model is one in which we assume that the values of the time series in the future depends on past values of the time series (Figure 4). Essentially, it is a linear model relating values of the time series to past values of the time series. Let this time series be y_t . Then, for this timeseries the kth order auto-regression model (AR(k)) is:



Figure 3. Plotting the residuals of our ARIMA (2,0,1) model.

$$y_t = \beta_0 + \beta_1 y_{t-1} + \beta_2 y_{t-2} + \ldots + \beta_k y_{t-k} + \dot{\mathbf{Q}}_{t-k}$$

$$y_t = \dot{\mathbf{Q}}_t + \beta_0 + \sum_{j=1}^{k} \beta_j y_{t-i}$$

A moving average model is one in which we assume that the values of the time series in the future depend on the previous residual terms of the time series. It is a linear model relating the value of the time series to past values of the error. Therefore, the kth order moving average model (MA(k)) is:

$$y_t = \dot{\mathbf{o}}_t + \theta_0 + \theta_1 \dot{\mathbf{o}}_{t-1} + \theta_2 \dot{\mathbf{o}}_{t-2} + \theta_3 \dot{\mathbf{o}}_{t-3} + \dots + \theta_k \dot{\mathbf{o}}_{t-k}$$

$$y_t = \dot{\mathbf{o}}_t + \theta_0 + \sum_{i=1}^k \theta_i \dot{\mathbf{o}}_{t-i}$$

An ARMA model simply combines an auto-regressive model and a moving average model. Thus, an ARMA model is such that a value of the time series can be predicted based on previous values of the time series. Thus an ARMA (p,q) model is:

$$y_{t} = \dot{\mathbf{o}}_{t} + \beta_{0} + \beta_{1}y_{t-1} + \beta_{2}y_{t-2} + \dots + \beta_{k}y_{t-p} + \theta_{0} + \theta_{1}\dot{\mathbf{o}}_{t-1} + \theta_{2}\dot{\mathbf{o}}_{t-2} + \dots + \theta_{k}\dot{\mathbf{o}}_{t-q}$$
$$y_{t} = \dot{\mathbf{o}}_{t} + \beta_{0} + \sum_{i=1}^{p} \beta_{i}y_{t-i} + \theta_{0} + \sum_{i=1}^{q} \theta_{i}\dot{\mathbf{o}}_{t-i}$$

The algorithm for choosing the best ARMA (p,q) model is not outlined in this paper. However, we use maximum likelihood estimation to arrive at the best ARMA model. We use the R programming language along with the Hyndman-Khandakar algorithm [2] to minimise Akaike's Information Criterion (AIC) [10].

$$AIC = -2\log(L) + 2(p+q+k+1)$$

L : Likelihood of the data

$$k = 1 i f \beta_0 + \theta_0 \neq 0$$

 $k = 0 if \beta_0 + \theta_0 = 0$

Code for ARIMA used

#first read the comma splitted values containing the dataset

data <- read.csv("weather_data_24hr_HI.csv", header = T)

#filter only maxtempC out which is relevant to our timeseries

data <- data(,3)

#remove NA values

data <- data[-4265]

#looking at the first few values of our vector

head(data)

#sequential data partition into training (70%) and testing data (30%)

training <- data(1:2990)

testing <- data(2991:4264)

#convert the data into a time series

ts <- ts(data, frequency = 365, start = c(2008,183))

#smooth the time series using the triangular moving average

tst2 <- ts(rollmean(rollmean(ts,2),2),frequency = 365, start = c(2008,183))

#plotting the original time series and the smoothed time series

plot(ts, col = 'green')

lines(tst2, col = 'blue')

#carrying out the augmented dickey-fuller test on our time series for stationarity

adf.test(ts)

 $\# carrying \ out \ the \ kpss \ test \ to \ further \ confirm \ level \ stationarity \ of \ the \ data$

kpss.test(ts, null = "Level")

tstrain <- ts(training, frequency = 365, start = c(2008,183))

#smoothing the training time series using triangular moving average

ts1 <- ts(rollmean(rollmean(tstrain,2),2))

#converting the testing data into a time series

tstest <- ts(testing, frequency = 365, start = c(2016,252))

#smoothing the testing time series using triangular moving average

ts2 <- ts(rollmean(rollmean(tstest,2),2))

#we fit the training data to the arma(2,2) model

fit<-arima(ts1, c(2,0,2))

#finding out the RMS error of the arma(2,2) model

accuracy(fit)

fitted(fit)

#applying our previously derived model on the testing data

refit<-Arima(ts2, model = fit)

#calculating the accuracy of our model on the testing data

accuracy(refit)

#plotting the residuals of our model

plot(residuals(fit))

Explanation of code: We first import the comma splitted values and filter only the relevant HeatIndexC data from the dataset. We then create multiple time series: one containing the entire dataset, one containing the training data, and one containing the testing data. We then use triangular moving averages to smooth these time series. We also run the ADF test and KPSS test on the data to test for stationarity. We then fit the ARMA (2,2) model to our data and calculate the RMS error on the training and testing data. We plot our residuals to observe whether it is a white noise time series to ensure optimality of our time series.

Results: According to the ADF and KPSS test we get that our time series is stationary. We get that an ARMA (2,2) model best fits our data. Our model is as follows:

 $y_t = \dot{\mathbf{o}}_t + 32.6958 + 0.8550 y_{t-1} + 0.1253 y_{t-2} + 1.9714 \dot{\mathbf{o}}_{t-1} + 0.9715 \dot{\mathbf{o}}_{t-2}$

Metric of analysis for results: We will now analyse the results of the ARMA (2,2) model. To determine how well our model fits the data we calculate the RMS error for the testing data and the training data. We define our prediction model function for our ARMA model to be: $(y_i) \hat{(}x)$. The actual temperature is represented by y_i .

$$RMS \, error = \sqrt{\frac{\sum_{i=1}^{n} (\hat{y}_i(x) - y_i)^2}{\sum_{i=1}^{n} (\hat{y}_i(x) - y_i)^2}}$$

The plot of our residuals is also a white noise time series. Thus, we can conclude that our ARMA model is optimal as white noise cannot be predicted [11].

For the training data the RMS error was: 0.335002

For the testing data the RMS error was: 0.354654

Conclusion

Since the RMS error for the training and testing data using our model is similar we can conclude that our model does not overfit our data and is a reliable predictor even on unseen testing data. It also performs significantly better than the exponential smoothing model, multivariate regression model as well as the neural network model.

Declaration

Availability of data and materials

The data used during this paper includes both data from the Indian Meterological Department as well as data from https://www. worldweatheronline.com regarding a variety of weather phenomena between 2008 and 2020. The R programming language was also used to run the programs for this paper.

Competing interests

No competing interests

Funding

Not applicable

Acknowledgements

I would like to thank Dr. Shubhangi Bhute (Indian Meterological Department) for guiding me with the project, helping me ideate, and find data for the research paper. Furthermore, I am grateful to Dr. Mahendra Mehta for helping me explore various areas in statistical prediction providing the knowledge base I needed to work on this paper.

References

- 1. Hyndman, R.J., & Khandakar, Y. "Automatic time series forecasting: the forecast package for R". J Stat Softw 27 (2008):1-22.
- 2. Reilly, D., "The AUTOBOX system". Int J Forecast 16.4 (2000):531-533.
- Shenstone, L., & Hyndman, R.J. "Stochastic models underlying Croston's method for intermittent demand forecasting". In J Forecast 6 (2005):389-402.
- Smith, J., & Yadav, S. "Forecasting costs incurred from unit differencing fractionally integrated processes". Int J Forecast 10.4 (1994):507-514.
- Taylor, J.W. "Exponential smoothing with a damped multiplicative trend". Int J Forecast 19.4 (2003):715-725.
- 6. Wallis, K.F. "Asymmetric density forecasts of inflation and the Bank of England's fan chart". Natl Inst Econ Rev 167 (1999):106-112.
- So, M.K., & Chung, R.S. "Dynamic seasonality in time series". Comput Stat Data Anal 70 (2014):212-226..
- John, S. "The three-parameter two-piece normal family of distributions and its fitting." Commun Stat-Theory Methods 11.8 (1982):879-885.
- Williams, W.H., & Goodman, M. "A simple method for the construction of empirical confidence limits for economic forecasts." J Am Stat Assoc 66.336 (1971):752-754.
- Aitchison, J., & Dunsmore, I.R. "Statistical Prediction Analysis." Cambridge University Press, Cambridge, 1975.
- Thombs, L.A., & Schucany, W.R. "Bootstrap prediction intervals for autoregression." J Am Stat Assoc 85.410 (1990):486-492.