



2020 - STUDY

MONGOLIA'S RAINFALL FORECAST USING RSTUDIO

A study to help small farmers in emerging countries
fight climate change

Mongolia's Rainfall Forecast using Rstudio

Morais, Manoela N. Morais/jnn@hotmail.com

Abstract—One of the essential sectors of Mongolia's economy is agriculture, which is sensitive to climate variation. The most important climatic element which impacts agriculture productivity is the rainfall. Therefore, rainfall prediction becomes an important issue in Mongolia, especially for small farmers that lack technological resources. In this paper, I propose to model the daily rainfall prediction over Dornod Province to help farmers to forecast rain over 1 year. I applied several modeling methodologies, which include time series regression and Machine learning. The method that presented the lower MAE and RMSE was the Arima Model, followed by the default Simple Times series.

Keywords—Time Series Regression, Rainfall Forecasting, Statistical forecasting, Mongolia weather, Farming Technology

I. INTRODUCTION

The agricultural industry is one of the most vulnerable to climate change as it directly depends on rainfall and temperature. Rainfall information is normally used by farmers for crop selection and water resource management in agriculture. The occurrence of prolonged dry periods or heavy rain at the critical stages of crop growth and development may lead to a significant reduction of the farmer's crop yield. Although climate change affects differently different crops and regions, it is expected that a decrease in farming productivity (Lobell, 2007). In fact, some decline can already be seen. In the last years, developed countries with a good farm infrastructure had more than 17% yield decrease because of climate changes (Nelson et al, 2014).

Climate variability is a major source of risk for agriculture and food systems and its changes are making it harder for the farmers to guess/predict the weather patterns by just looking for last year's results. The climate change impact is even more accentuated to small-low-income farmers in emerging countries, given their small scale, the burden of economic exclusion, and lack of access to critical resources. Around 500 million smallholder farms in the developing

world are supporting almost 2 billion people, and 70 % depending partly or completely on agriculture for their livelihoods (IFAD, 2011). Those farmers do not have adequate networking or tools from which they can learn and easily exchange know-how. Therefore, temperature and rainfall prediction become a significant factor in agriculture to reach the best crop performance.

To this study, I selected the Dornod province, in Mongolia, and an emerging country, which its economy is largely based upon agriculture. The focus of this study is to provide to the small-house-hold-farmers the rainfall forecast with a good prediction and low resources.

1.2 Research Goal

The main goal of this research is to propose a solution to help farmers selecting their crop and its start farming period by giving them useful rainfall prediction. Therefore, I am going to test several simples' models to identify one that can predict with lower errors the rainfall forecast for Dornod Province.

| Dornod Province | | References: |
|------------------------------|-----------------|-------------|
| Soum (sum of rice) | Climate | North |
| 1 Chuluut Khoroos (Khoroskh) | wet | Central |
| 2 Durbidulime | wet | South |
| 3 Bayanbulag | wet | |
| 4 Bayanbulag (Bayanbulag) | wet | |
| 5 Gersan (Gersan) | wet | |
| 6 Tsagaan Ovoo | dry | |
| 7 Sengelen | dry | |
| 8 Chodolun | dry | |
| 9 Bayanbulag (Bayanbulag) | dry | |
| 10 Kherlen | dry | |
| 11 Bulgan | dry | |
| 12 Halanbulag | moderately mild | |
| 13 Mutud | moderately mild | |
| 14 Khulbulag | wet | |

Figure 2. Dornod Provinces used for weather model study

DATA

I am using 5 years of meteorological data between January 2015 and May 2020 from the Nasa website from Dornod Province (see all Soums in Figure 01). From the Nasa available indices, a set of 14 indices of climate variables were selected for this study (Figure 2). For long-range forecast (LFR), which is more than 2 weeks and up to 2 years, generally, the rain forecast is consistent with fundamental variables such as previously temperature and precipitation (Doblas-Reyes & all, pp. 12). However, we are using other variables to identify the interference and importance in the rainfall prediction results.

Of the selected variables, seven indices refer to temperature, one to precipitation, and 4 to other variables. The temperature indices describe cold extremes as well as warm extremes. The precipitation indices describe wet extremes.

| Parameter | Units | Description |
|--------------------|----------|--------------------------------|
| T2M | °Celsius | Temperature at 2 Meters |
| T5 | °Celsius | Earth's Temperature |
| rs | mm | Surface Rainfall |
| T2M_NH | °Celsius | Temperature at 2 Meters |
| SLP | hPa | Sea Level Pressure at 2 Meters |
| WSW | m/s | Wind Speed at 2 Meters |
| ALBID_TGA_2M_2M_NH | W/m² | Top-of-atmosphere Radiation |
| ALBID_SF_2M_2M_NH | W/m² | Surface Radiation at 2 Meters |
| PRECIP | mm/day | Precipitation |
| T2M_NHGB | °Celsius | Temperature Range at 2 Meters |
| T2M_NH | °Celsius | Temperature at 2 Meters |
| T2M_NH | °Celsius | Temperature at 2 Meters |
| ALBID_SF_2M_2M_NH | W/m² | Surface Radiation at 2 Meters |
| T2M_NH | °Celsius | Temperature at 2 Meters |

Figure 1. 14 Indices selected from the Nasa climatology website available indices.

II. Modeling

The methodology used will be first, get historical data and random from the Nasa website. If the model is not times series related, I will divide it in 2 categories in 80% for training (build the models) and 20% for test the models (Figure 04). If the model contains time-series I will divided for the training 2015-2019 data and 2020 results will be used for testing. Then, I am going to run several regression models for rain forecast for autoregressive and multivariable regressions. Finally compare RMSE and MAE results and select the one, that present lower results. For illustration of the models I going to represent initially Dashbalhar results, for the others soum please vide table at attachment.

A quantitative forecast of rainfall is extremely difficult and realizable. Generally, it is done only a couple hours of their occurrence with a Doppler. However, for agriculture operations, a quantitative forecast of rain is not as important as a forecast of the (i) non-occurrence of rains and (ii) type of rain spell that can be expected. Therefore, after modeling and get the results in "millimeters of rainfall", I am going to classify the rainfall based on the USG definition of raining (Figure 03). USG defines that Absent of raining ≤ 0 mm/hh, Slight rain ≤ 0.5 mm/hh, Moderate

$\leq 4\text{mm/hh}$ and heavy rain (High) $> 8\text{mm/hh}$ rainfall. Although the information is given by USG in mm/hour, I am going to consider the same value for the whole day (mm/day) to classify. The main goal of doing this classification is to give farmers the information they need and decrease forecast errors. The main goal of doing this classification is to give to farmers the information they need and decrease the forecast errors.

| Classification from website | | |
|-----------------------------|-------------|---------|
| No rain | 0 - 0.001 | mm/hour |
| Slight rain: | 0.001 - 0.5 | mm/hour |
| Moderate rain: | 0.5 - 4 | mm/hour |
| Heavy rain | 4 - 8 | mm/hour |

Figure 3. Rain data classification extracted from the USG, Gov Resource: <https://www.usgs.gov/mission-areas/water-resources>

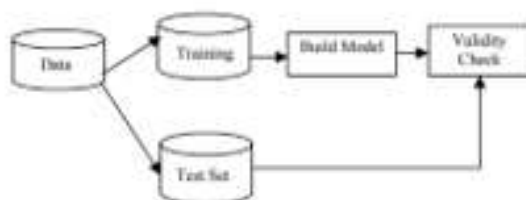


Figure 4. Overview of the forecasting Planning methodology

III. Single Regression Models

The autoregressive process is a regression on itself. Therefore, Y_t is a linear combination of the p most recent past values of itself plus the term “ e_t ” that incorporates the error (Equation 01). We are going to start this stage by establishing a Baseline model, which is basically the mean of the historical data, to be used as a reference and compared with the linear Times series model, Times series *Seasonal naïve Method*, and *Arima model*. The $Y(t)$ will be rainfall (PRECTOT).

$$Y_t = \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + e_t$$

Baseline Model

Assuming...that all variables are independent and there is no rainfall prediction (no time series), I constructed a dependent guess (not random) based on the mean of the value itself, using the training data. The Results obtained were RMSE= 2.187456 and MAE= 1.190749, these results will be used as a reference to compare with other models.

```
> Base Model <- mean(trainSPRECTOT)
```

Times Series linear model (TS)

The rainfall observations collected through daily data is sequential over time. Therefore, we are going to use time-series to model the stochastic mechanism that gives rise to an observed series. I am going to predict and compare it with the test set. If one of those models has the best fit, I will forecast it. We will start by analyzing the data extracted from NASA, if the variables' time series are (non) stationarity and possess a unit root.

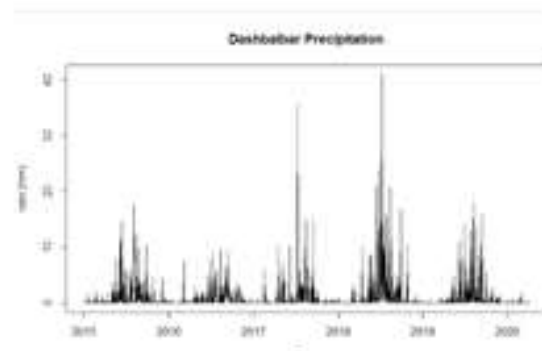


Figure 5. Precipitation for ~~Qashbelbar~~ through the last 5 years. Extracted from RStudio. We can see times series and a seasonal trend.

Checking UnitRoot

To make statistical inferences about the structure of a stochastic process on the basis of an observed record of that process,

we need to verify the assumption that the data is stationarity and therefore, the probability laws that govern the behavior of the process do not change over time. In a sense we formulate the null and alternative hypotheses for the unit root. Considering the equation for AR(p):

$$\Delta Y_t = \alpha + \delta t + \rho Y_{t-1} + \gamma_1 \Delta Y_{t-1} + \dots + \gamma_{p-1} \Delta Y_{t-p+1} + \varepsilon_t$$

H0: $\rho=0$ (non-stationary, has UNIT ROOT)

Ha: $\rho < 0$ (stationarity)

When Y has no unit root and stationarity condition $2 < \rho < 0$ holds after incorporating a time trend, we call these series trend stationary. We tested the unit roots with Augmented Dickey-Fuller Test (ADF).

```
## Unit Root Test
library(FarTimeGTS)

#new data
adfTest(dailySDM,type="c") # type="nc" no const
adfTest(dailyS,type="c") # type="nc" no const
adfTest(dailyPRECOT,type="c") # type="nc" no c
adfTest(dailyRH2M,type="c") # type="nc" no const
adfTest(dailyT2M,type="c") # type="nc" no c
adfTest(dailyT2M_MAX,type="c") # type="nc" no c
adfTest(dailyT2M_MIN,type="c") # type="nc" no c
adfTest(dailyT2M,type="c") # type="nc" no const
adfTest(dailyT2M_RANGE,type="c") # type="nc" no const
adfTest(dailySPD,type="c") # type="nc" no const
adfTest(dailyT2MNET,type="c") # type="nc" no c
adfTest(dailyALLSKY_TOA_SW_DWN,type="c") # type="nc" no c
adfTest(dailyALLSKY_SFC_SW_DWN,type="c") # type="nc" no c
adfTest(dailyALLSKY_SFC_LW_DWN,type="c") # type="nc" no c
```

The Augmented Dickey-Fuller (ADF) test uses the t-statistic from the regression and compares it to a DF-t critical value that can be found in statistical packages. If the unit root hypothesis $\rho=0$ has not been rejected, the conclusion is that at least one-unit root exists in the process and we need to test for possible second unit root for the differenced series ΔY until we get a stationary process. From results of all 14 indices we see that for H0: unit root hypothesis is rejected since p-value < 0.01 $<$ significance level of 10% (or 5%). And therefore, we can assume that the indices are stationary and can be used as it is, without the need of differentiation.

After creating a simple linear time series with a simple seasonality and an increasing trend and forecast for a year value we obtained a the RMSE=1.784812 and a MAE =0.9783169, both lower than the base model.

```
##### Time Series -Linear model #####

ts.train= ts(train$PRECOT,frequency=365.25,start=c(2015))
ts.test=ts(test$PRECOT,frequency=365.25,start=c(2016))

#TEST PERFORMANCE
test.pred.ts= predict(ts.train,test$PRECOT, h=365)
RMSE.ts =- sqrt(mean((test.pred.ts$mean)^2))
MAE.ts =- mean(abs(test.pred.ts$mean))
MAE.ts
```

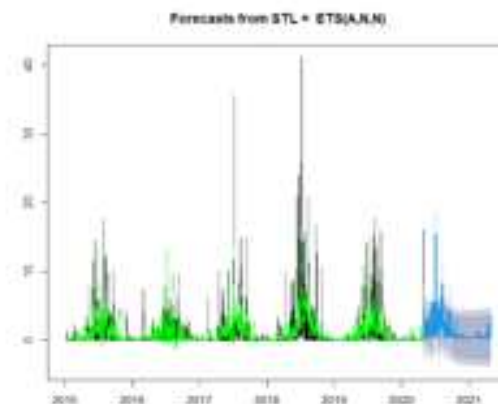


Figure 6. Forecast from Datalark. PRECOT for 365 days using default Times series method. Where we can see in green the forecast model in top of the past historical values.

After some days I tested again this model adding more recently data. The plot gave me Figure 07. Although, RMSE and MAE slight changed, we can see that the mean in non-Zero (there is no period without raining!). Which shows that for long periods (365 days) we might have a problem forecasting values.

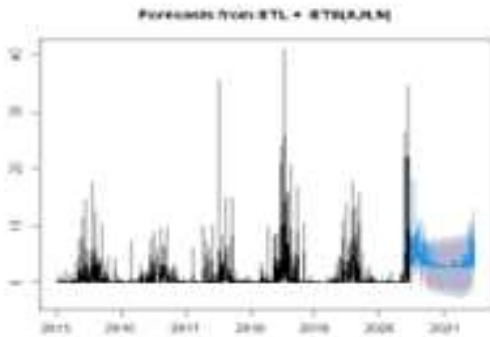


Figure 7. Forecast from Dataskar. PRECTOT for 365 days using default Times series method. After recently shower in the place.

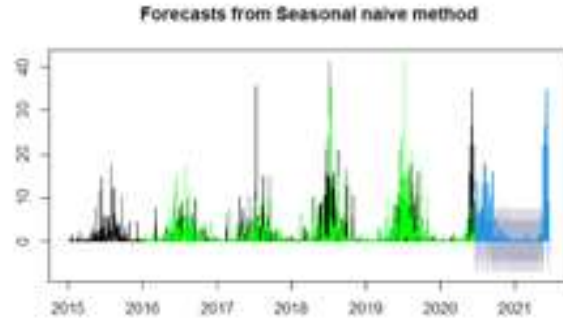


Figure 8. Forecast PRECTOT for 365 days using Times series with exponential smoothing method.

TS + Seasonal naïve Method

The seasonal Naïve (snaive) model verify if rainfall in the place are highly seasonal. The seasonal naive model makes the forecast using the most recently observation from the same season (Equation 2). The $Z[t]$ is normal error. The Exponential smoothing generally is good for short term forecast and refer to error, trend and seasonality. The model uses the exponentially weighted moving average (EWMA) to “smooth” a time series and trying to eliminate the random effect (RPubs).

$$Y[t] = Y[t-m] + Z[t] \quad \text{Equation 2}$$

```
##### method 02: TS- naive (Random walk forecast)
fit_ts_train <- snaive(ts_train,frequency=12*36.5373)
fit_ts <- snaive(ts_reg,frequency=12*36.5373)

#RATAPRICE
test_pred_sml<- predict(fit_ts_train,test$PRECTOT, n=365,35)
mse_sml<- sqrt(mean((test_pred_smltest)^2,na.rm=TRUE))
RMSE_sml
MAE_sml <- mean(abs(test_pred_smlmean),na.rm=TRUE)
MAE_sml
```

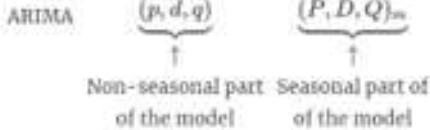
I got for this model a RMSE = 3.373169, which increased compared with than the base model, however the MAE= 1.110575 was lower. Which means the ET+ Snaive model fit better, however the error when happens is bigger the base model.

Arima Method

Arima also known as AutoRegressive Integrated Moving Average models is widely used approaches for time series weather forecasting. Arima Model is a combination of three mathematical models, using autoregressive(p), integrated(d), moving-average (q) (ARIMA) models for time series data. An ARIMA (p, d, q) model can account for temporal dependence in several ways. Firstly, the time series is d-differenced to render it stationary. If $d = 0$, the observations are modelled directly, and if $d = 1$, the differences between consecutive observations are modelled.

For this paper I started using the auto.arima function in Rstudio, the result was Arima (3,1,2) with non-Zero mean. This shows that Rstudio could not identify the seasonal part of the zero. This happened because Daily data is challenging as it often involves multiple seasonal patterns, and so we need to use a method that handles such complex seasonality. Therefore, I also runned

arma function with
n 1 to 25 and and



```

start = time.time()
y = y[train][:NPROCTEST, frequency=100, start=0.0201, end=0.0201]
y_train = train[NPROCTEST, frequency=100, start=0.0201, end=0.0201]
test_y, test_t = get_y, test_t

# fit the model
# using the logistic regression
# Fit the logistic model
bestFit = fitLogistic(y_train)
# Fit the model
fit = LogisticRegression()
fit = fit.fit(y_train, x_train)
# Predict the test set
y_test_fit = fit.predict(x_test)
# Evaluate the model
bestFit = fit
# Print the results
print('BestFit = %s' % bestFit)

```

Both the RMSE =1.497 and MAE=1.098 decreased compared with the base model. The autocorrelation plot - ACF graph in figure 08 shows that for the first 500 lags, almost all sample autocorrelations fall inside the 95 % confidence bounds indicating the residuals appear to be random.

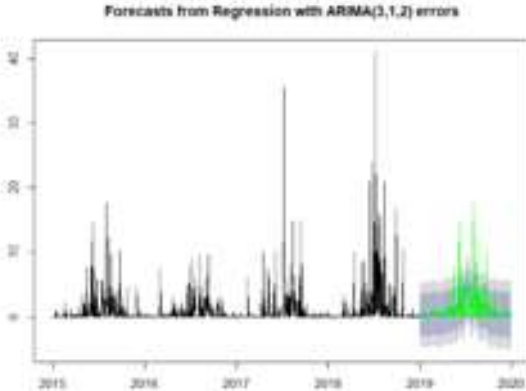


Figure 9. Forecast PRECOTOT for 365 days using Times series with auto-regressive method. Green the real value and in blue the forecast.

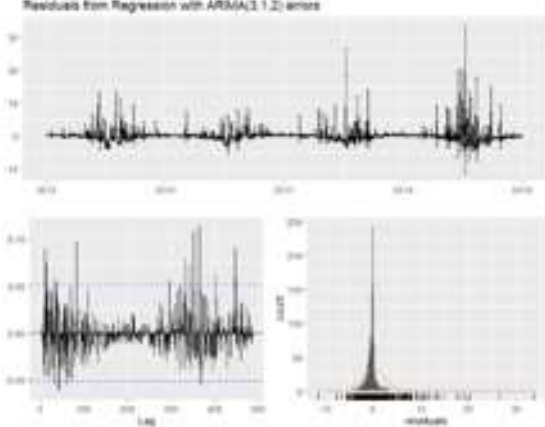


Figure 10. Residual from the forecast of PRECTOT for 365 days using Times series with ~~auto arima~~ method.

IV. Multivariable linear Regression

For a Multiple linear regression, the dependent variable is assumed to be a linear function of several independent variables (predictors), where each of them has a weight (regression coefficient) that is expected to be statistically significant in the final model.

Multiple Log- linear regression

Using a log-linear regression model to the dependent variable ($y_x = \text{PRECTOT}$) in Equation x. We started with the all potential variables (Figure 06) and then eliminated from the model, those that were not statistically significant $p < 0.05$ (Figure 10). The adjusted R-squared for the log-linear model with all 12 variables is 0.59 which means that 59% of the variance in our dependent variable mm in rainfall (PRECTOT) can be explained by the set of predictors in the model; Although not all variables are significant in this first model, after we try to used just the significant values (figure 12), we can observe that the results of the Adjusted R-square for all variables $>$ Adjusted R-square for significant ones. This happens, probably because the variables

might have some degree of dependent
between each other.

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im} + \epsilon_i \quad \text{Equation X}$$

$i = 1, 2, \dots, n$, where, $\varepsilon_i \sim \text{i.i.d. } N(0, \sigma^2)$

4) <http://regression.wustl.edu> using the training data.

```
ltn.reg = log(oddsratio[i]) - d0[m0] - d0[m1, x0] - d0[m2] - d0[m3, x0] + y0[i]
summary(ltn.reg) # report the model
accuracy(ltn.reg)
```

The result obtained from the log linear model with all the variables was a RMSE= 2.791962 that have slight increases when compared with the base model, however the MAE= 1.186184 was lower. This shows that the model fits better than the base model. This model makes some assumptions which includes linearity, constant variance, normality and independence between the parameters. We can see on figure 13 the normal QQ plot of the residuals, the values are not normal after x-axis greater than 1.5 where points are away from the line. In figure 14, we can see the plot of the forecast of the variable PRECTOT against the past value, using linear regression.

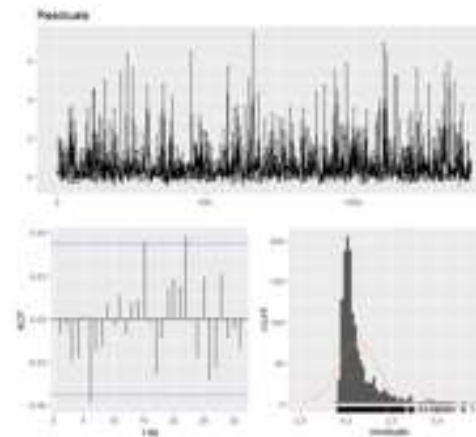
[illegible]

Figure 11. Results(summary) and residuals from the log linear regression for the variable PRECOTOT against all others 12 variables.

```
Call:
lm(formula = log(disp/sector ~ 13 + dissem2 + dissem3 + disps +
  dist2m_range + zsa11307_TDA_3m_0m + zsa11307_SFC_3m_0m +
  zsa11307_SAC_3m_0m, data = train)

Residuals:
    min       1Q   median       3Q      max
-2.12188  -0.21921  -0.01038  0.16247  2.39918

Coefficients:
(Intercept)      3.4662050      1.5787414      2.193  -0.03643  +
dissem2         0.0125607      0.0007868      17.232  + 2e-16 ***
dissem3         0.0130484      0.0070887      4.521  6.32e-06 ***
disps           -0.0532381      0.0168223      -1.201  -0.00189 +
dist2m_range    -0.0310493      0.0679151    -20.400  + 7e-16 ***
zsa11307_TDA_3m_0m  0.0230638      0.0018510      12.967  + 1e-16 ***
zsa11307_SFC_3m_0m  -0.0042215      0.0018771      32.603  + 1e-18 ***
zsa11307_SAC_3m_0m  0.0442682      0.0027572      16.054  + 1e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9878 on 1028 degrees of freedom
Multiple R-squared:  0.3815. Adjusted R-squared:  0.1803
F-statistic: 181.4 on 7 and 1028 df, p-value: 3.2e-16
```

Figure 12. Results(summary) and residuals from the log linear regression for the variable PRECTOT against significant variables.

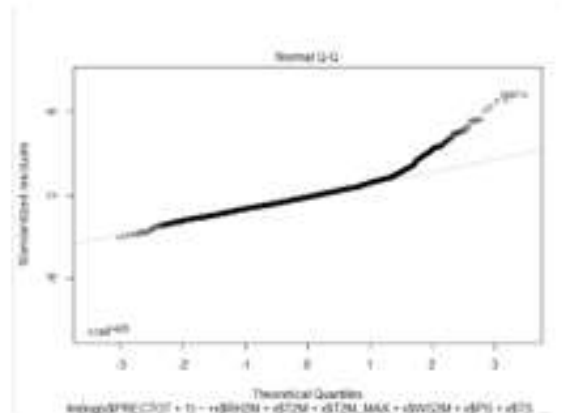


Figure 13. QQ plot of the residuals of the log-linear models for $y(\text{PRECTOT})$ against the 14 variables.

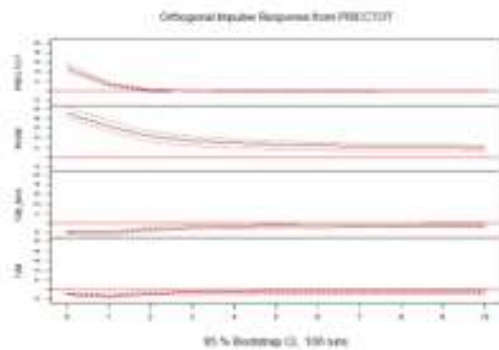


Figure 16. Impulse response for VAR Model

V. Machine Learning Models

Neural Network Model – library ([nnetar](#))

The Artificial Neural Networks (ANN) are a set of algorithms, designed to recognize patterns. They interpret sensory data through a kind of machine perception, labeling or clustering raw input. Artificial ANNs have become very popular, and prediction using ANN is one of the most widely used techniques for rainfall forecasting.

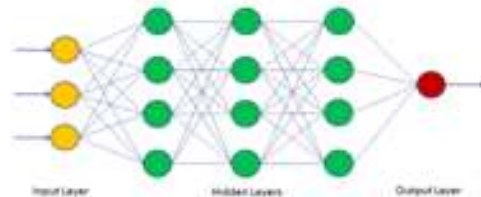


Figure 17. Artificial Neural Network: [structure](#) Source: <https://www.datacamp.com/community/tutorials/neural-network-models-r>

To perform the Neural Network forecast I used the [nnetar](#) function in the [forecast](#) package for R, that fits a neural network model to a non-linear time series. Used with big data sets. The NN model is organized in multiples layers, the simplest networks contain no hidden layers and are

equivalent to linear regressions. The coefficients attached to these predictors are called “weights”. The forecasts are obtained by a linear combination of the inputs.

```
forecast method: nnar(17,3,100[100])
Model Information:
Average of 10 networks, each of which is
a 18-10-5 network with 201 weights
options were: n.linear.output.units
Error measures:
Training set: MSE 0.0213389 RMSE 0.462780 MAE 0.314708 MAPE 0.1001360 MAE2 0.0800773
```

Regression trees

Regression tree for continuous outcome variables, is a simple and popular machine learning algorithm. In contrast with previous linear models it makes no assumptions about the relation between the outcome and predictors. It is the basis of a very powerful method that we will also use in this tutorial, called random forest



We can interpret that as when the downward Thermal Infrared (longwave) is lower than 26 we have lower than 0.23mm of precipitation (71% of the days). In the same way when the downward Thermal Infrared is between 26-31 and the humidity is lower than 65% (16% of the days), predict a wet day of 0.69mm.

.... For more please e-mail moraismn@hotmail.com....

ATTACH 02

RMSE and MAE Results from the models for Darbarr *(need to add for other places)*

| Type of Model | Model | RMSE | | MAE | |
|----------------------------------|--|----------|------|-----------|------|
| | | Value | Rank | Value | Rank |
| Autoregressive Regression Models | Base Model | 2.729586 | 4 | 1.302 | 6 |
| Autoregressive Regression Models | Time Series | 1.784812 | 2 | 0.9783169 | 2 |
| Autoregressive Regression Models | Times Series with exponential smoothing method | 3.373169 | 7 | 1.110575 | 4 |
| Autoregressive Regression Models | Auto ARIMA | 1.496647 | 1 | 1.098219 | 3 |
| Multivariate linear Regression | Log-Linear Regression | 2.791962 | 5 | 1.186184 | 5 |
| Multivariate linear Regression | Vector autoregression (VAR) | 5.10147 | 9 | 4.632723 | 9 |
| Machine Learning | Neural network models | 2.22929 | 3 | 0.9534483 | 1 |
| Machine Learning | Tree Model | 3.511225 | 8 | 1.410056 | 7 |
| Machine Learning | Random Forest | 3.305167 | 6 | 1.414745 | 8 |

Entry

Nasa Data .Retrieved from: <https://power.larc.nasa.gov/data-access-viewer/>

Alexy, V. (2018). USA advanced retail sales 24 month forecasting analysis. *RPubs*. Retrieved from: <https://www.rpubs.com/alev2301/421553>

Das, H. P., Doblas-Reyes F. J., Garcia, A., Hansen, J., Mariani, L., Nain, A., Ramesh, K., Rathore L. S. & Venkataraman, R. Weather and Climate Forecasts for Agriculture. *Agrometeorology*. Retrieved from: http://www.agrometeorology.org/files-folder/repository/gamp_chapt4.pdf

Pedro M. (2015). Modelling – predicting the amount of rain. *R-bloggers*. Retrieved from: <https://www.r-bloggers.com/part-4a-modelling-predicting-the-amount-of-rain/>

Chai, T. and Draxler1, R. R. (2014). Root mean square error (RMSE) or mean absolute error (MAE)? – Arguments against avoiding RMSE in the literature. *Geoscientific Model Development*. Retrieved from: <https://www.geosci-model-dev.net/7/1247/2014/gmd-7-1247-2014.pdf>

FAO. retrieved from: <http://www.fao.org/3/y2722e/y2722e0y.htm>

Lobell D.B., Field C.B. Global scale climate–crop yield relationships and the impacts of recent warming. *Environ. Res. Lett.* 2007;2:014002. doi: 10.1088/1748-9326/2/1/014002. Retrieved from: <https://iopscience.iop.org/article/10.1088/1748-9326/2/1/014002>

IFAD. 2011. Rural poverty report 2011. New realities, new challenges: new opportunities for tomorrow's *génération* (available at <http://www.ifad.org/rpr2011/report/e/rpr2011.pdf>).

Nelson, G., van der Mensbrugghe, D., Abammad, H., Blanc, E., Calvin, K., Hasegawa, T., Havlik, P., Heyhoe, E., Kyle, P., Lotze-Campen, H., von Lampe, M., Mason d'Croz, D., van Meijl, H., Müller, C., Reilly, J., Robertson, R., Sands, R., Schmitz, C., Tabeau, A., Takahashi, K., Valin, H. & Willenbockel, D. 2014b. Agriculture and climate change in global scenarios: why don't the models agree? *Agricultural Economics*. 45(1): 85–101.