

A PROJECT ON

**3 YEARS FORCAST OF TIME SERIES HOSPITAL INPATIENT
WAITING LIST DATASET BASED ON **A**UTO **R**EGRESSIVE
INTEGRATED **M**OVING **A**VERAGE (**ARIMA**) MODEL**

BY

PAVAN M SUNDER

ABSTRACT

The aim of this project is to obtain 5 years future sales forecast for a time series data based on **Auto Regressive Integrated Moving Average (ARIMA)** model.

The data set used for this project is monthly data of hospital inpatient waiting list

A stationary data is generally desirable to work on as it facilitates easier and reliable implementation of predictive models. Stationarity means that the statistical properties of a process generating a time series data do not change over time.

However, most often the real-life data will be non-stationary and we might have to apply a certain transformation to make the data stationary. The data set used for this project was initially non – stationary with decreasing trend and small seasonal variations. A transformation was applied to make the data stationary.

ARIMA model is then fitted on transformed data set. Then future 3 years forecast was obtained based on fitted model.

CONTENTS

| | |
|-------------------------------|---|
| 1. INTRODUCTION | 1 |
| 2. DATA | 2 |
| 3. DATA PREPROCESSING..... | 3 |
| 4. ARIMA MODEL..... | 4 |
| 5. FORECAST FOR 3 YEARS | 7 |
| 6. CONCLUSION | 8 |

1. INTRODUCTION

There are so many prediction problems that involve a time component. In this project I have considered a time series dataset for analysis and prediction. A time series data is a series of data points indexed (or listed or graphed) in time order. Time series analysis comprises methods for analysing time series data in order to extract meaningful statistics and other characteristics of the data. **Time series forecasting** is the use of a model to predict future values based on previously observed values.

Autoregressive Integrated Moving Average (ARIMA) model is one of the models we can deploy to better understand the time series data or to predict future points in the series (forecasting). ARIMA is a generalization of an Auto Regressive (AR) Moving Average (MA) model. ARIMA models can be applied in some cases where data show evidence of non-stationarity. An initial differencing step (corresponding to the "Integrated" part of the model) can be applied one or more times to eliminate the non-stationarity.

The AR part of ARIMA indicates that the evolving variable of interest is regressed on its own lagged (prior) values. The MA part indicates that the regression error is actually a linear combination of error terms whose values occurred contemporaneously and at various times in the past. The I (for "integrated") indicates that the data values have been replaced with the difference between their values and the previous values (and this differencing process may have been performed more than once). The purpose of each of these features is to make the model fit the data as well as possible.

Non-seasonal ARIMA models are generally denoted $ARIMA(p,d,q)$ where parameters p , d , and q are non-negative integers, p is the order (number of time lags) of the autoregressive model, d is the degree of differencing (the number of times the data have had past values subtracted), and q is the order of the moving-average model.

In this project I have used ARIMA model to forecast 3 years inpatient waiting list numbers for a non-stationary time series hospital inpatient waiting list dataset.

2. DATA

The data set used for this project is monthly hospital inpatient data. There are 145 observations over a 12-year period. The original dataset is credited to Department of Health, UK.

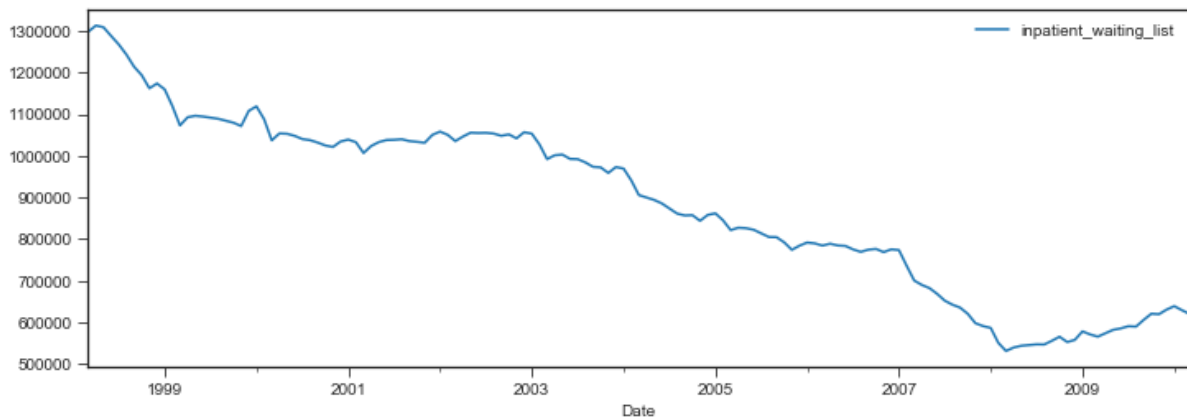


Figure 1. Original hospital inpatient data

The data apparently shows non – stationarity with varying trend and seasonal variations.

The non – stationarity can further be confirmed by using rolling statistics plot and augmented dickey fuller test as shown in Figure 2.

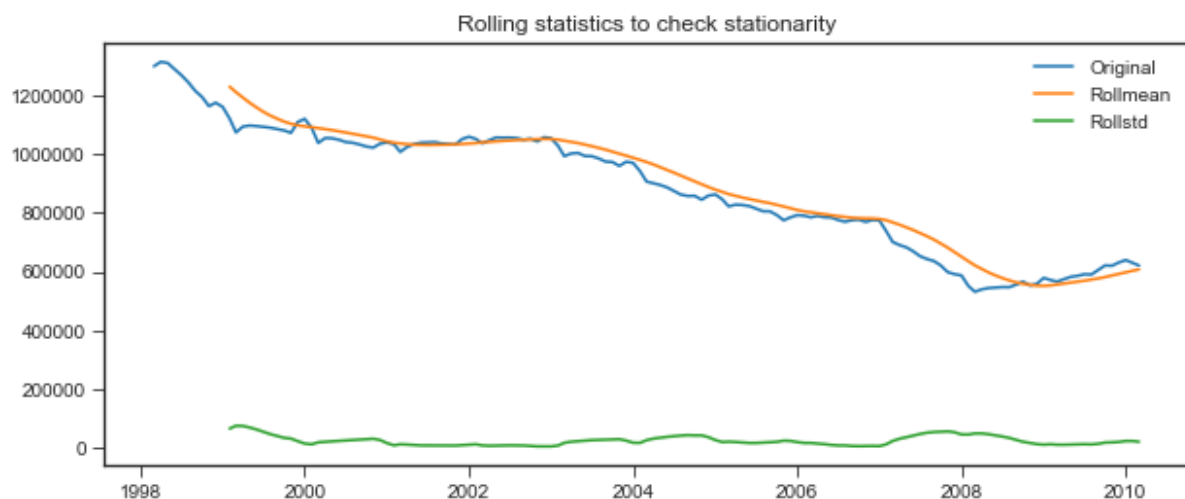


Figure 2(a). Rolling statistics plot of original time series

```

Results of Dickey-Fuller Test:
Test Statistic          -0.573767
p-value                 0.876800
#Lags Used              13.000000
Number of Observations Used 131.000000
Critical Value (1%)     -3.481282
Critical Value (5%)     -2.883868
Critical Value (10%)    -2.578677
dtype: float64
Testing test statistic (-0.5737672740915366) < 1% critical value (99% confidence) (-3.481281802271349)?: False
Testing test statistic (-0.5737672740915366) < 5% critical value (95% confidence) (-2.883867891664528)?: False
Testing test statistic (-0.5737672740915366) < 10% critical value (90% confidence) (-2.5786771965503177)?: False

```

Figure 2(b). Dickey Fuller test results of original time series

The illustration in Figure 2(a) clearly shows varying trend (negative slope of rolling mean denoted in ‘orange’). Also, Dickey Fuller test results shown in Figure 2(b) shows that the test statistic is greater than the critical values. There is a larger p value. Thus, we can say that the data is non-stationary.

3. DATA PREPROCESSING

The non – stationary data needs to be transformed to make it stationary before applying ARIMA model.

The following illustrates transformation steps in order to render original time series stationary.

1. Converting original time series to log scale. Let ‘*ts*’ denote original time series data and ‘*ts_log*’ denote log transform of original time series data (‘*ts*’)

$$ts_log = \log(ts) \quad (1)$$

2. Taking average over the past 3 months (moving average with window size=3) of *ts_log* and subtracting it from the *ts_log*, denoted as ‘*ts_ma*’. Let *t* represent current time instance, then:

$$ts_ma = ts_log(t) - \frac{1}{3} \sum_t (ts_log(t), ts_log(t - 1), ts_log(t - 2)) \quad (2)$$

3. Taking 1st order difference of *ts_ma* denoted as ‘*ts_diff*’

$$ts_diff = diff(ts_ma) \quad (3)$$

The stationarity of the transformed data can further be visualized and confirmed using the rolling statistics plot and Dickey Fuller test as illustrated in Figure 3(a) and 3(b) respectively.

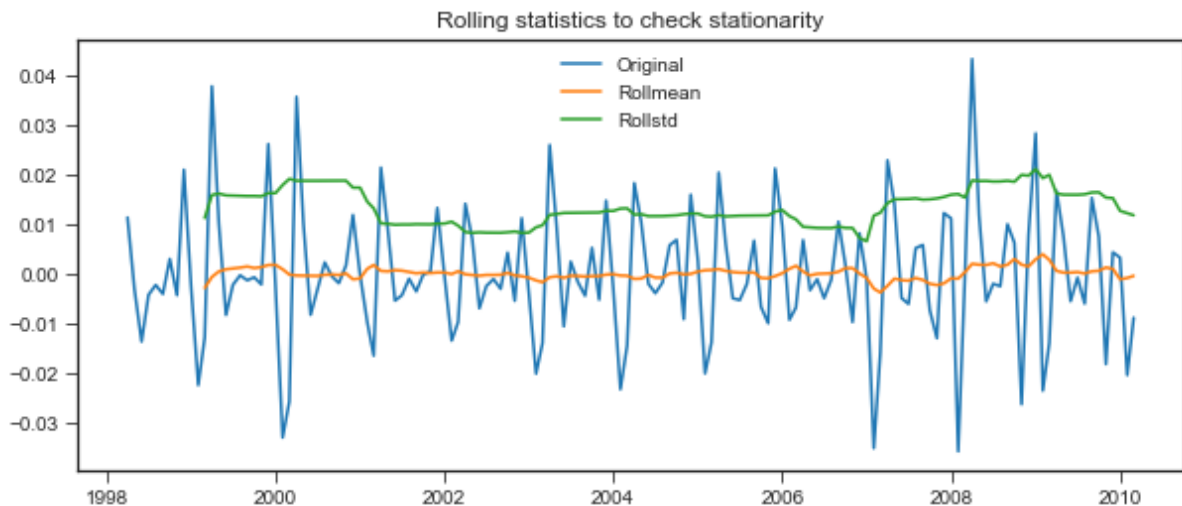


Figure 3(a). Rolling statistics plot transformed time series.

```
Results of Dickey-Fuller Test:
Test Statistic          -3.921429
p-value                 0.001883
#Lags Used              14.000000
Number of Observations Used 129.000000
Critical Value (1%)     -3.482088
Critical Value (5%)     -2.884219
Critical Value (10%)    -2.578864
dtype: float64
Testing test statistic (-3.921428634741788) < 1% critical value (99% confidence) (-3.482087964046026)?: True
Testing test statistic (-3.921428634741788) < 5% critical value (95% confidence) (-2.8842185101614626)?: True
Testing test statistic (-3.921428634741788) < 10% critical value (90% confidence) (-2.578864381347275)?: True
```

Figure 3(b). Dickey Fuller test results of transformed time series.

The illustration in Figure 3(a) shows constant mean (denoted in 'orange'). Also, the test statistic according to Dickey Fuller test is less than 1% critical value which implies that we can say with 99% confidence that the data is stationary. Also, we can observe smaller p-value.

Now we can fit the ARIMA model to the transformed data as it is stationary.

4. ARIMA MODEL

In order to fit the ARIMA model to the transformed data we need to determine optimal parameters i.e. (p, q, d) values.

We have seen that the data is rendered stationary by taking 1st order difference. Therefore, the optimal d value shall be 1.

Further we need to determine p and q values. One of the ways to find p and q values is by plotting partial autocorrelation and auto correlation plots respectively as illustrated below.

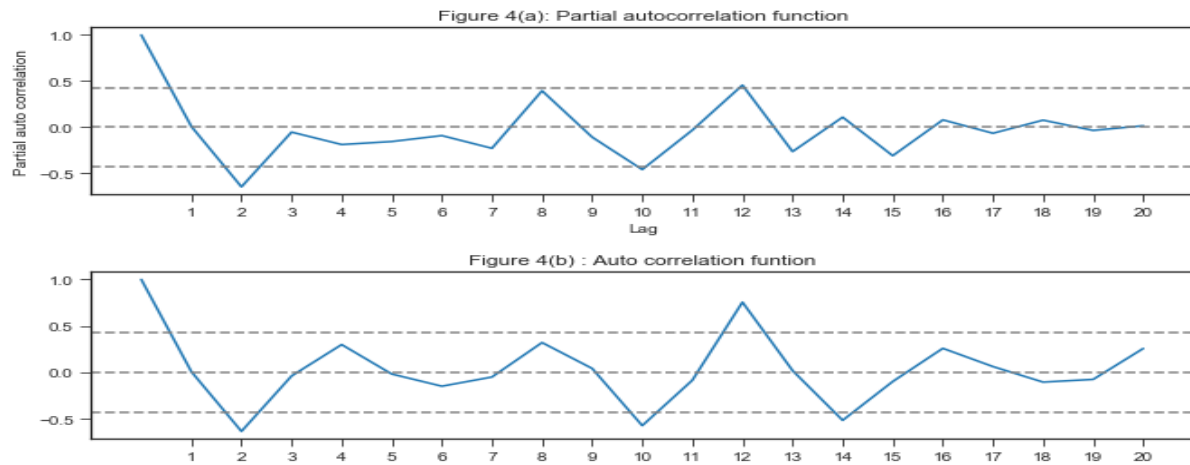


Figure 4. (a) Partial auto correlation plot; (b) Auto correlation plot

From Figure 4(a) we see that the graph drops to zero for the 1st time at around 1, hence we can take $p = 1$

Similarly, from Figure 4(b) we see that the graph drops to zero for the 1st time at around 1, hence we can take $q=1$.

Based on above parameters, the ARIMA model is fitted to the time series data by taking $p = 1, d = 1, q = 1$. Figure 5 illustrates ARIMA model fitted to the transformed data

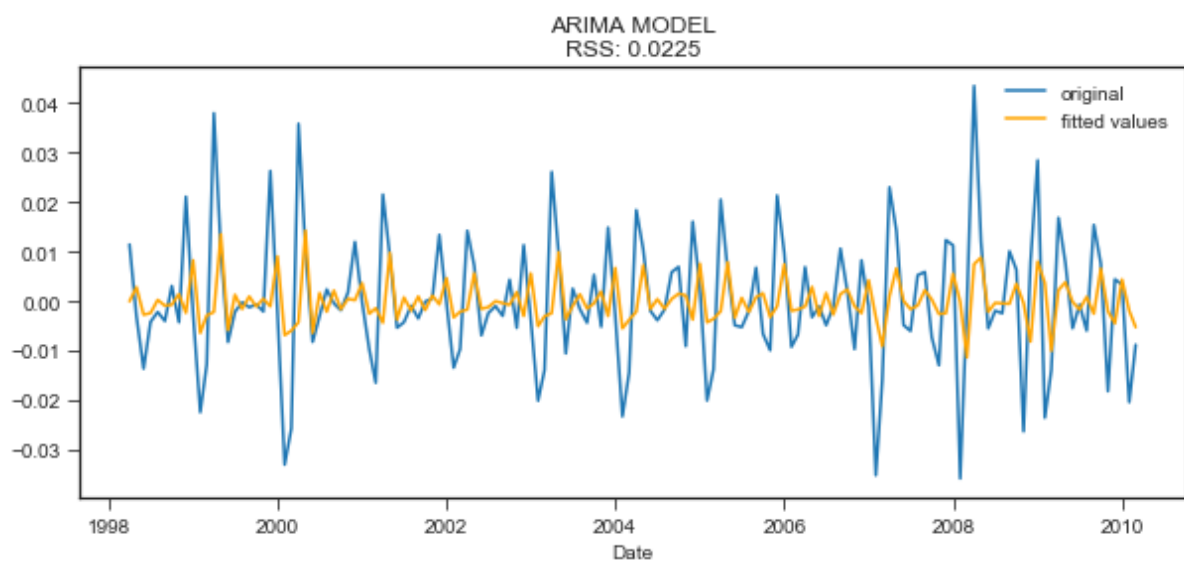


Figure 5. ARIMA model fitted to the transformed data.

From Figure 5 we can observe that the model has fairly fitted the data well indicated by the low Residual sum of squares (RSS) value of 0.0225. The model has also captured the seasonal variations present in the input data fairly well.

The input data for the ARIMA model is a transformed data from which the trend present in the original time series was removed to make it stationary. However, it is necessary to bring back the fitted values to original form in order to incorporate the trend that was removed and make reliable predictions.

This can be achieved by the following steps:

1. First get the fitted values and store it as series. We will notice that the first month will be missing because lag of 1(shift) was taken to carry out 1st order difference.
2. Inverse the differentiation: find the cumulative sum of fitted values and add the base value (usually the first value in the stationary series for which the model is applied) to each term in the time series
3. Add the trend (moving average with window 3) that was initially removed.
4. Take the exponent of the series (anti-log) and square it which will be the predicted values of the time series forecast model as illustrated in Figure 6

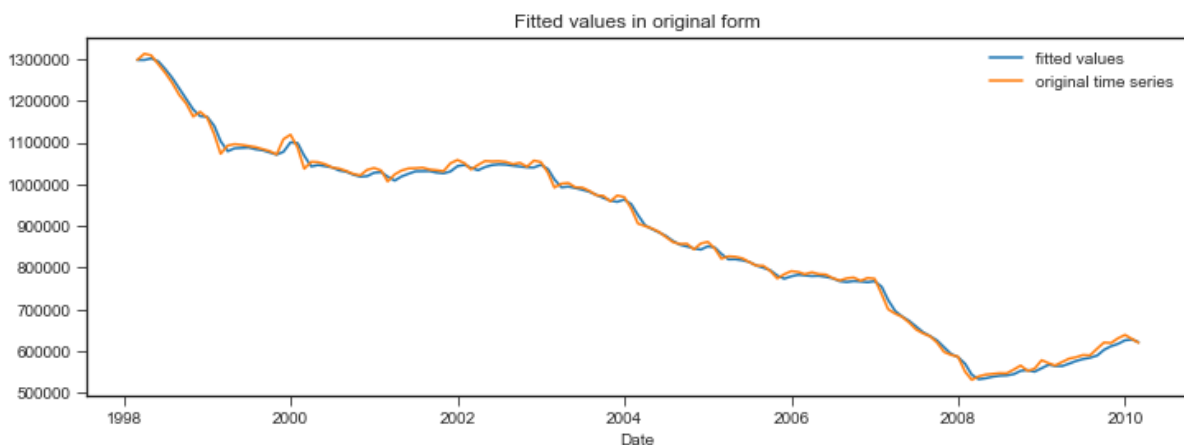


Figure 6. ARIMA model fitted values to original form.

From Figure 6 we can observe a good approximation between actual and fitted values indicating a good fit.

5. FORECAST FOR 3 YEARS

Figure 7 illustrates 3 years forecast after fitting the model.

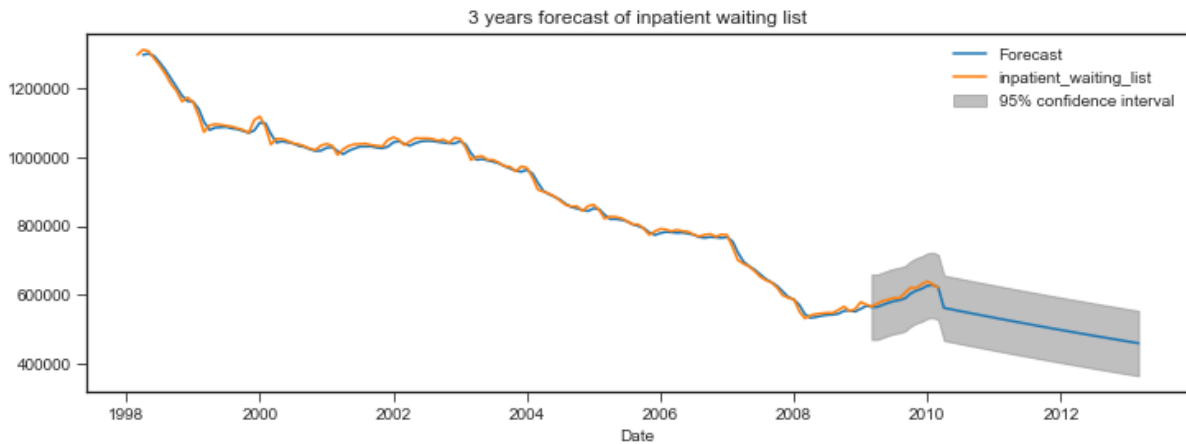


Figure 7. 3 years forecast.

From Figure 7 we can see the trend and seasonal variations fairly approximates the true values. Also, shaded region indicates 95% prediction interval starting from previous year of forecast, which clearly shows that true values lie within the range. Hence, we can say the model has fitted the data well and we have fairly a good forecast

6. CONCLUSION

I have tried to implement ARIMA predictive model for a non-stationary time series. The parameters for the model was taken as $p=1$, $d=1$ and $q=1$ based on rolling statistics plots, Dickey Fuller test, auto correlation and partial auto correlation plots.

The model was able to fit the data fairly well with the selected parameters capturing the trend and seasonal variations. I was also able to get the 3 years forecast with a good approximation to true values.