# Washington Metropolitan Area Transit Authority Analysis

# Washington Metropolitan Area Transit Authority Analysis

**Problem statement:-**

The Washington Metropolitan Area Transit Authority (Metro) was created by an interstate compact in 1967 to plan, develop, build, finance, and operate a balanced regional transportation system in the national capital area. Metro began building its rail system in 1969, acquired four regional bus systems in 1973, and began operating the first phase of Metrorail in 1976. Today, Metrorail serves 91 stations and has 117 miles of track. Metrobus serves the nation's capital 24 hours a day, seven days a week with 1,500 buses. Metrorail and Metrobus serve a population of approximately 4 million within a 1,500-square mile jurisdiction. Metro began its paratransit service, MetroAccess, in 1994; it provides about 2.3 million trips per year.

The Washington Metropolitan Area Transit Authority wants to know some insights from its old data recorded. We have data from 2012 to 2016 recorded with seven features and 26785 records of problems faced in different tracks due to different problems. Each problem has a different time durations to clear it.

**Opportunity:-**

We have several records of data available to perform some analysis and few questions can be answered such as
- In which direction most of the metro moves in?
- What are the available tracks, which track is the busiest one?
- In which location we had more problems?
- What kind of problems does WMATA metro face?
- Which problem does occur most of the time?
- How much delay does it take for each problem?
- What is the amount of time does most of the problems get resolved?

**Tools used:-**

1. Python 3.6
2. Numpy
3. Pandas
4. Matplotlib
5. Seaborn
6. Time Series Forecasting

**Dataset:-**

We have a dataset from WMATA [here](#). This data set contains several records with features like ID, Time, location, direction, color, problem and delay. We have another file with scraped metro disruptions of size 1.4MB.

**Process:-**

Use the tools that we have like Numpy, Pandas, Matplotlib, Seaborn and Python to get insights as asked above from the data. Let's answer each and every point one by one.

Let's explore each and every column that can contribute to our analysis.

**Direction:-**

We have a feature called Direction that Metro moves towards.



From the graph, you can see that we have most of the metros that go in the direction of Shady grove faces more problems.

**Track Color:-**

This color feature tells us about the track that the metro travels. We have several tracks in the Washington metro station like blue, green, silver, orange, red, yellow, Greed.

**Track color vs Number of Probelms**



You can see from the graph that we have most of the problems on the Red track.

**Location:-**

This feature tells us that the problem occurring locations. We have locations such as Greenbelt , Vienna, New Carrollton, Largo Town Center, Huntington, Branch Avenue, West Falls Church, Fort Totten, Wiehle-Reston East, Franconia-Springfield, Rhode Island Avenue, King Street, Silver Spring, Metro Center, Farragut North, Stadium-Armory, Gallery Place, College Park, Takoma, Rosslyn and so on….

Let's explore this feature and try to find out the location that faces more problems.

## percentage of problems comes in locations



From the graphs, you can see that we have more problems in locations Greenbelt with Vienna in second place.

**Problem:-**

So we are going to explore what kind of problems do Washington metro faces and what is the most occurring problem. There are several problems like a brake problem, operational problem, door problem, equipment problem, track problem, police activity, and emergency problem.

Let's plot the graph and see which problem caused a number of times.

## problem vs number of times



From the graph, you can see that we have more problems because they did not operate. The most caused second problem is a brake problem, which is dangerous to the lives of passengers travelling in Wmata if not identified earlier.

**Delay in minutes:-**

We have a feature called delay in our dataset that is caused due to different problems. Each problem has a different time duration delay. Let's see the amount of delay that occurred most times due to problems.

**Delay in time vs Number of times**

From the graph, it is observed that for most of the problems the maximum time delay is just 6 minutes.

**Month:-**

We have Monthly Average of Delay in minutes from the years 2012 to 2016. From the graph below you can see that we have 26 Minutes delay in the direction of Metro_center in February month. Also, we can see that we have 24 Minutes delay in the direction of the stadium Armory in the Month of March.

We have Average Monthly delay in minutes from 4 minutes to 24 minutes. Least Monthly average delay occurred in September month in the direction of Francois to Springfield and Also in the direction of Takoma.

## Dealy average by direction and month

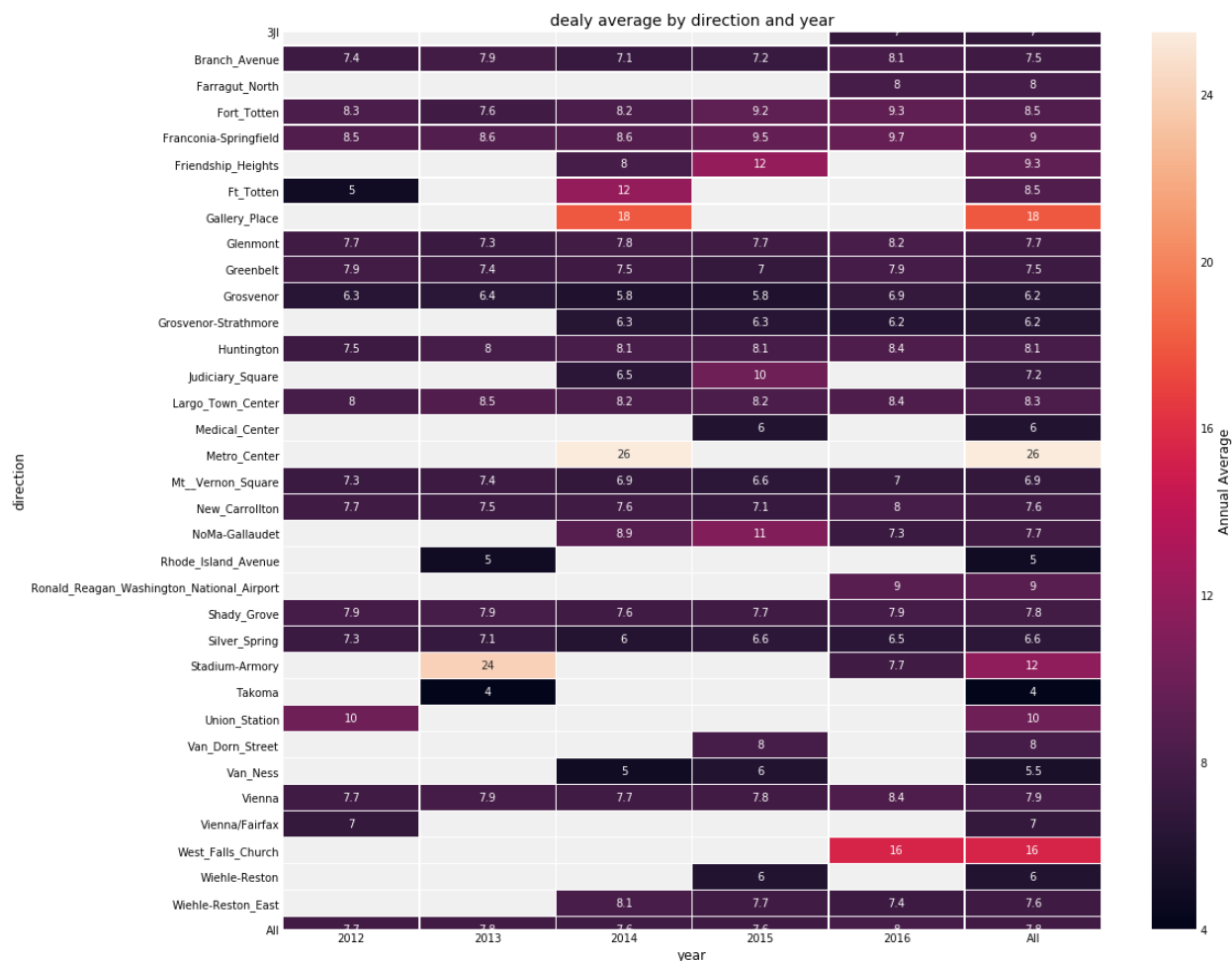| direction | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | All |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3l | | | | | | | | 7 | | | | | 7 |
| Branch_Avenue | 8.4 | 7.7 | 7.6 | 7.5 | 7.9 | 7.1 | 7.6 | 7 | 7.7 | 7.3 | 7.9 | 7 | 7.5 |
| Farragut_North | | 9 | | 7 | | | | | | | | | 8 |
| Fort_Totten | 8.7 | 9 | 10 | 8.1 | 7.9 | 7.8 | 8.7 | 9.5 | 7.9 | 10 | 6.6 | 8.3 | 8.5 |
| Franconia-Springfield | 8.5 | 9.1 | 11 | 8.8 | 9 | 8.5 | 8.8 | 9 | 9.2 | 9.2 | 9 | 8.9 | 9 |
| Friendship_Heights | | | 12 | | | | | 12 | 4 | | | | 9.3 |
| Ft_Totten | | | | | | | | 5 | | | | 12 | 8.5 |
| Gallery_Place | | | | | 18 | | | | | | | | 18 |
| Glenmont | 8 | 8 | 7.7 | 8 | 7.4 | 7.7 | 7.3 | 8 | 8.1 | 7.6 | 7.7 | 7.5 | 7.7 |
| Greenbelt | 7.7 | 7.7 | 7.1 | 7.8 | 7.2 | 7.4 | 7.5 | 7.5 | 7.3 | 7.6 | 7.7 | 7.5 | 7.5 |
| Grosvenor | 6 | 5.5 | 5.7 | 4.9 | 6.6 | 6.7 | 5.9 | 7.1 | 6.4 | 6 | 5.9 | 6.3 | 6.2 |
| Grosvenor-Strathmore | 5.7 | 5.6 | 5.1 | 7.5 | 8 | 7.1 | 4.9 | 7.2 | 6 | 5.3 | 4 | 6.8 | 6.2 |
| Huntington | 8.2 | 8.1 | 8.5 | 8.6 | 8.6 | 7 | 8.8 | 8.4 | 8.1 | 8.2 | 6.8 | 8 | 8.1 |
| Judiciary_Square | 10 | | | 6.5 | | | | | | | | | 7.2 |
| Largo_Town_Center | 8 | 8.9 | 8.7 | 8.2 | 8 | 8.2 | 8.2 | 8.3 | 8.5 | 8.3 | 8.5 | 7.9 | 8.3 |
| Medical_Center | | | | | | 6 | | | | | | | 6 |
| Metro_Center | | 26 | | | | | | | | | | | 26 |
| Mt__Vernon_Square | 7 | 7.3 | 6.9 | 6.5 | 6.5 | 7.7 | 6.6 | 7.2 | 6.5 | 6.7 | 6.6 | 7.4 | 6.9 |
| New_Carrollton | 8 | 7.9 | 7.2 | 7.6 | 7.9 | 7.5 | 6.9 | 7.6 | 7.6 | 7.9 | 7.7 | 7.2 | 7.6 |
| NoMa-Gallaudet | | | | 12 | 9.8 | 9.5 | 10 | 11 | | 8.6 | 6.5 | | 7.7 |
| Rhode_Island_Avenue | | | | | | | | | | 5 | | | 5 |
| Ronald_Reagan_Washington_National_Airport | | | | | | | 9 | | | | | | 9 |
| Shady_Grove | 8.4 | 8.4 | 7.5 | 7.5 | 7.9 | 7.4 | 7.2 | 7.9 | 8.4 | 7.5 | 8.4 | 7.4 | 7.8 |
| Silver_Spring | 6.6 | 7 | 6.2 | 7.5 | 6.6 | 6.5 | 6.9 | 6.4 | 6.2 | 6.4 | 6.7 | 6.6 | 6.6 |
| Stadium-Armory | 6 | 11 | 24 | | | | | | | | | | 12 |
| Takoma | | | 4 | | | | | | | | | | 4 |
| Union_Station | | | | | | | | 10 | | | | | 10 |
| Van_Dorn_Street | | 8 | | | | | | | | | | | 8 |
| Van_Ness | 6 | | | | 5 | | | | | | | | 5.5 |
| Vienna | 8.4 | 7.8 | 7.2 | 7.5 | 8 | 7.6 | 7.7 | 7.6 | 8.6 | 8.8 | 8.1 | 7.7 | 7.9 |
| Vienna/Fairfax | | | | | 7 | | | | | | | | 7 |
| West_Falls_Church | | | | | | | | | 25 | 6 | | | 16 |
| Wiehle-Reston | | | | | | | | | | | 6 | | 6 |
| Wiehle-Reston_East | 7.5 | 7 | 7.3 | 6.5 | 8.2 | 7.9 | 7.9 | 7.2 | 8.7 | 8.4 | 7.4 | 7.7 | 7.6 |
| All | 8.1 | 8.2 | 7.7 | 7.7 | 7.8 | 7.6 | 7.5 | 7.7 | 8 | 7.8 | 7.8 | 7.6 | 7.8 |

**Month**

**Annual Average**

## Year:-

We have most of the problems that occurred in the year 2015 compared to others in the time duration of 2012 to 2016.

Down below we have Yearly Average Delay in minutes to each and every direction we have. We have the Highest yearly Average delay in the year of 2014 in the direction of Metro center.

dealy average by direction and year

| direction | 2012 | 2013 | 2014 | 2015 | 2016 | All |
|---|---|---|---|---|---|---|
| 3JI | | | | | | |
| Branch_Avenue | 7.4 | 7.9 | 7.1 | 7.2 | 8.1 | 7.5 |
| Farragut_North | | | | | 8 | 8 |
| Fort_Totten | 8.3 | 7.6 | 8.2 | 9.2 | 9.3 | 8.5 |
| Franconia-Springfield | 8.5 | 8.6 | 8.6 | 9.5 | 9.7 | 9 |
| Friendship_Heights | | | 8 | 12 | | 9.3 |
| Ft_Totten | 5 | | 12 | | | 8.5 |
| Gallery_Place | | | 18 | | | 18 |
| Glenmont | 7.7 | 7.3 | 7.8 | 7.7 | 8.2 | 7.7 |
| Greenbelt | 7.9 | 7.4 | 7.5 | 7 | 7.9 | 7.5 |
| Grosvenor | 6.3 | 6.4 | 5.8 | 5.8 | 6.9 | 6.2 |
| Grosvenor-Strathmore | | | 6.3 | 6.3 | 6.2 | 6.2 |
| Huntington | 7.5 | 8 | 8.1 | 8.1 | 8.4 | 8.1 |
| Judiciary_Square | | | 6.5 | 10 | | 7.2 |
| Largo_Town_Center | 8 | 8.5 | 8.2 | 8.2 | 8.4 | 8.3 |
| Medical_Center | | | | 6 | | 6 |
| Metro_Center | | | 26 | | | 26 |
| Mt__Vernon_Square | 7.3 | 7.4 | 6.9 | 6.6 | 7 | 6.9 |
| New_Carrollton | 7.7 | 7.5 | 7.6 | 7.1 | 8 | 7.6 |
| NoMa-Gallaudet | | | 8.9 | 11 | 7.3 | 7.7 |
| Rhode_Island_Avenue | | 5 | | | | 5 |
| Ronald_Reagan_Washington_National_Airport | | | | | 9 | 9 |
| Shady_Grove | 7.9 | 7.9 | 7.6 | 7.7 | 7.9 | 7.8 |
| Silver_Spring | 7.3 | 7.1 | 6 | 6.6 | 6.5 | 6.6 |
| Stadium-Armory | | 24 | | | 7.7 | 12 |
| Takoma | | 4 | | | | 4 |
| Union_Station | 10 | | | | | 10 |
| Van_Dorn_Street | | | | 8 | | 8 |
| Van_Ness | | | 5 | 6 | | 5.5 |
| Vienna | 7.7 | 7.9 | 7.7 | 7.8 | 8.4 | 7.9 |
| Vienna/Fairfax | 7 | | | | | 7 |
| West_Falls_Church | | | | | 16 | 16 |
| Wiehle-Reston | | | | 6 | | 6 |
| Wiehle-Reston_East | | | 8.1 | 7.7 | 7.4 | 7.6 |
| All | 7.7 | 7.8 | 7.6 | 7.6 | 8 | 7.8 |

year

**Final observations:-**

- most of the problems are coming to the trains that travel to the direction of shady grove
- we have most of the problems on the Redline.
- Most of the problems are coming at the location Greenbelt.
- at the location, Vienna comes more problems after Greenbelt.
- most metros get the problem of did not operate.
- The second-most problem is a Brake problem.

this "Brake problem" is dangerous and risk to people's lives.

- most of the trains get delayed for 6 min with the 4min second place.
- we have most of the problems in the month of 7,8 i.e July, august.
- we had 62highest problems in 30- august- 2016.
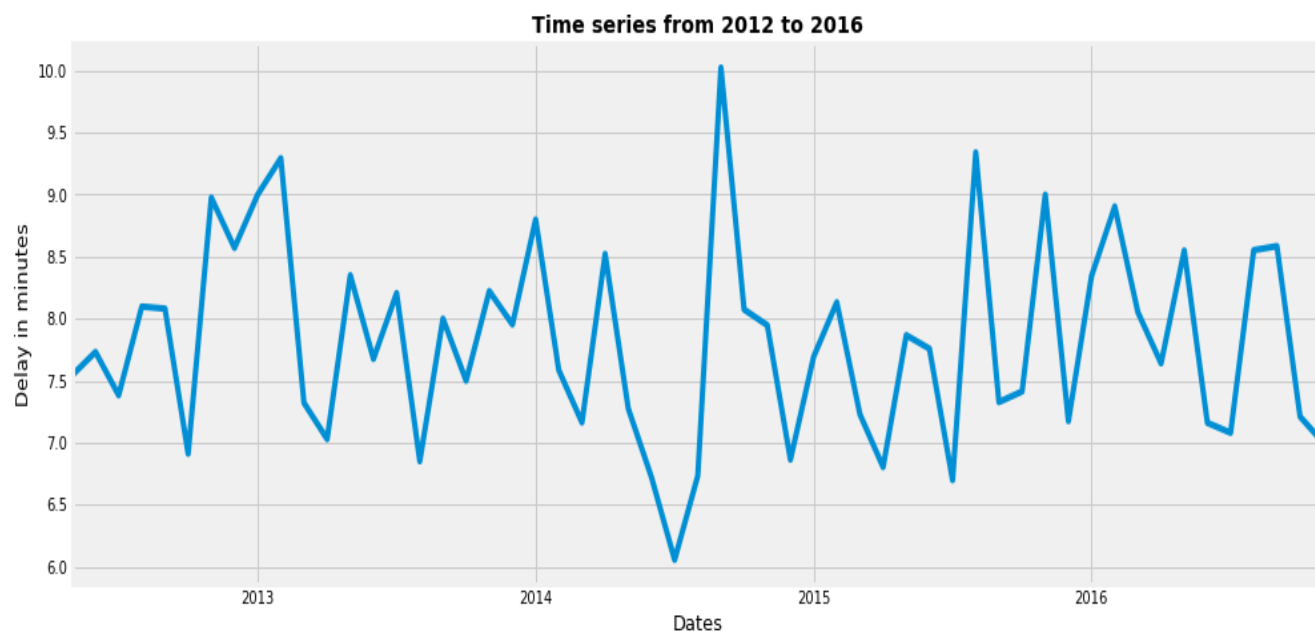- we had the highest problems in the year 2015.

- we got Most of the problems in the office hours which said to be the busiest hours.
  - 7 and 8 am in the morning.
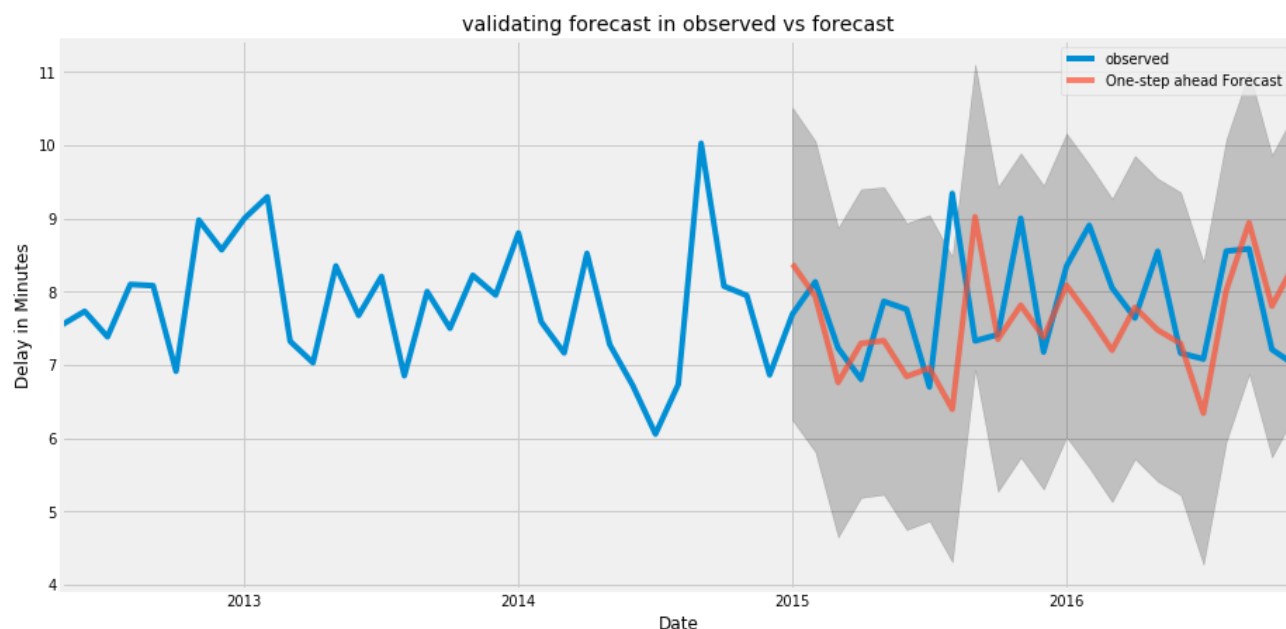  - 4 and 5 pm in the evening.

**Time series forecasting:-**

Now we have all the data cleaned for time series forecasting. As we have several directions, take one direction and predict the average delay that is going to occur for every month for 2017 to 2025. We have the data available from 2012 to 2016.

Lets plot the data we have from 2012 to 2016 with its average delay per month in minutes.
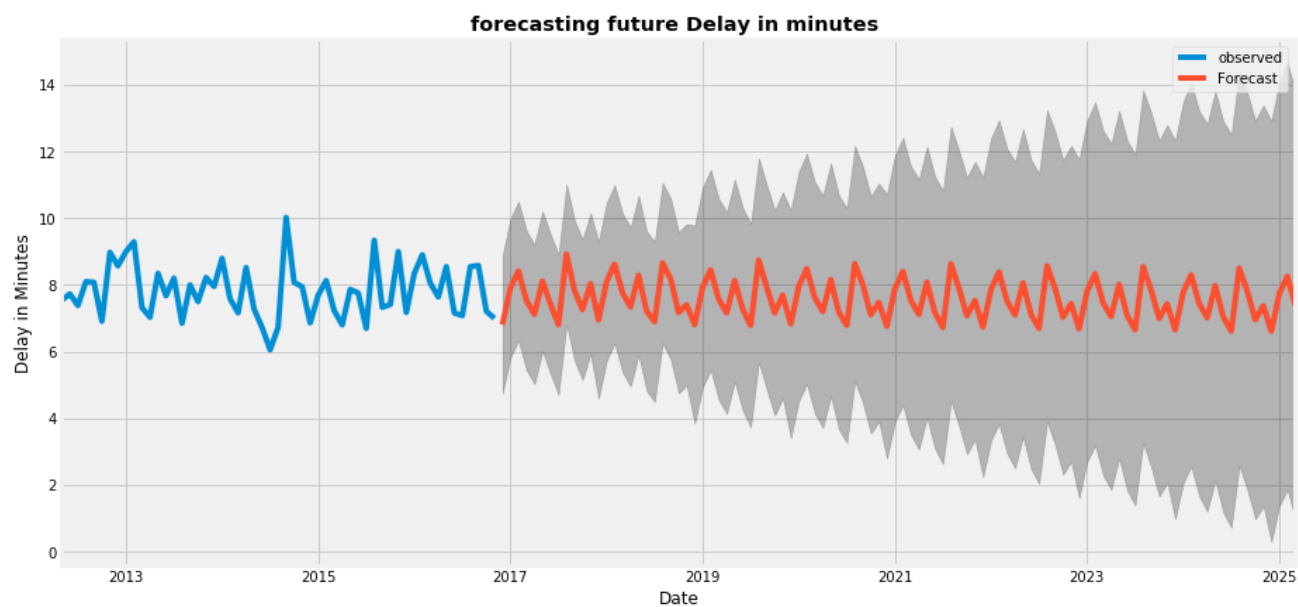


This is how it looks once you plot the data. You can see that we have an increase in delay from the second part of the year.

Let's use ARIMA - Auto-Regressive Integrated Moving Average to validate the forecast and see the observed vs forecasted.

validating forecast in observed vs forecast

Our error metric Mean Absolute Percentage Error gives only 9.91 which means our model is predicting to 91% accuracy here.

Now let's forecast the future average delay per month in minutes from 2017 to 2025.



forecasting future Delay in minutes

Check out the predicted average delay in minutes per month here below.

| Month | Predicted Average Delay |
|---|---|
| 2019-02-01 | 8.450399 |
| 2019-03-01 | 7.577566 |
| 2019-04-01 | 7.160183 |
| 2019-05-01 | 8.133552 |
| 2019-06-01 | 7.274166 |
| 2019-07-01 | 6.788310 |
| 2019-08-01 | 8.745031 |
| 2019-09-01 | 7.930906 |
| 2019-10-01 | 7.160700 |