

The background of the slide features a digital display with numbers in orange and white. An orange horizontal bar is positioned above the title.

# Forecasting

## WHO Flu Data

---

Adnan Muhammed

# Agenda

---

Introduction to Forecasting

---

Forecasting Methods

---

Industry Introduction and Use Cases

---

Methods

---

Code Walkthrough

---

Challenges

---

Conclusions

---

# Intro to Forecasting

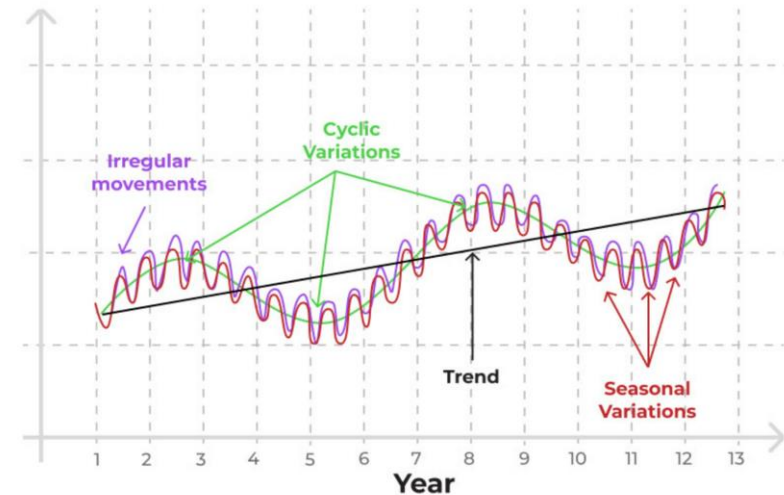
Making scientific predictions for future outcomes based on historical time data

Building models through historical analysis and using them to drive future strategic decision-making

Future outcome completely unavailable, must use evidence-based priors

Applications include: Stock Price forecasting, Weather forecasting, Healthcare predictions, Business Planning

**TIME SERIES DATA COMPONENTS**







MANIFOLD MARKETS



# Forecasting Methods

- **Online forecasting aggregation platforms:**
  - *Manifold Markets*
  - *Metaculus*
  - *Others*
- **Common models for forecasting:**
  - *Tree-based algorithms (Decision Trees, XGBoost)*
  - *Linear Regression*
  - *ARIMA*
  - *Exponential Smoothing Model*
  - *Neural Networks (LSTM, MLP, RNN, CNN)*
- **Univariate versus multivariate forecasting**
  - *Univariate Time Series Forecasting – only using previous values of time series*
  - *Multivariate Time Series Forecasting – using predictors other than the series (a.k.a exogenous variables)*

# Industry Introduction and Use Cases

## Forecasting in Healthcare

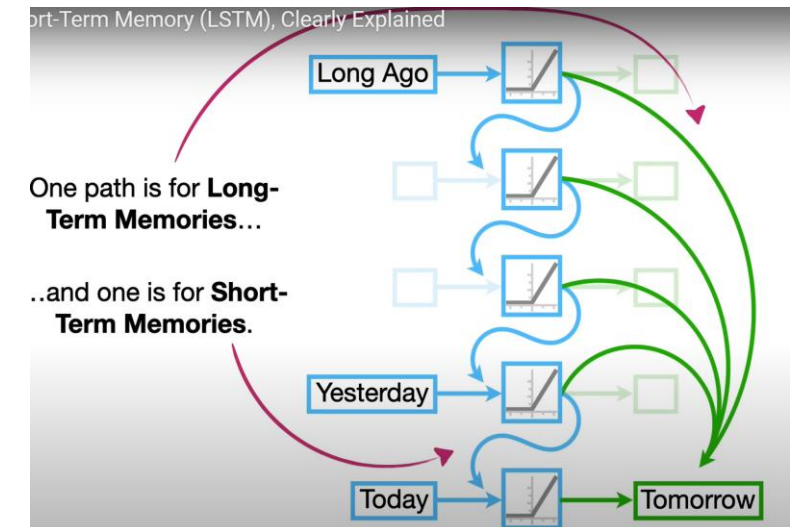
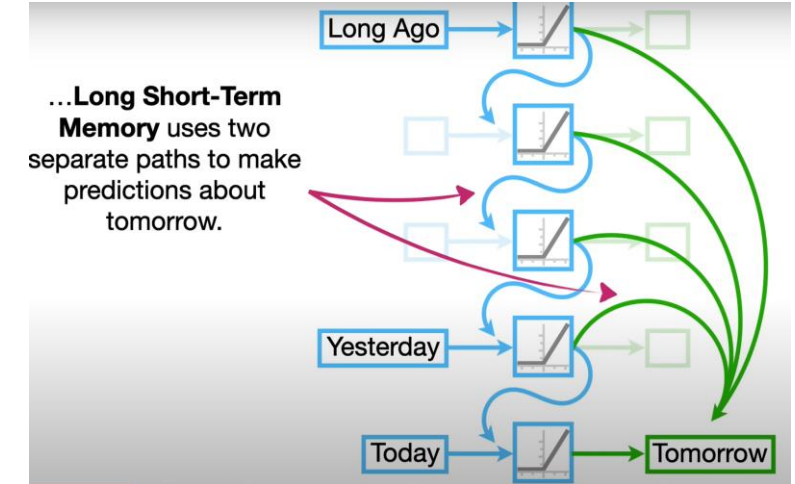
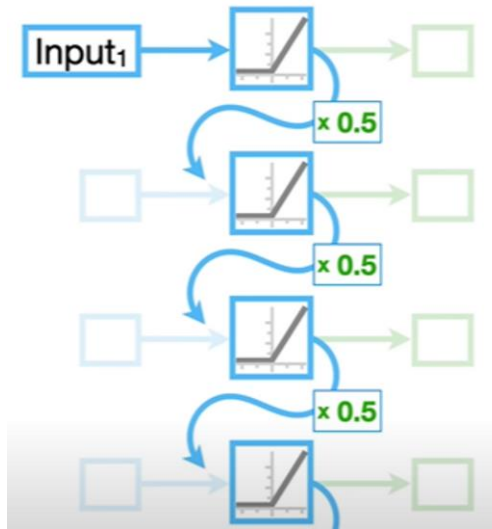
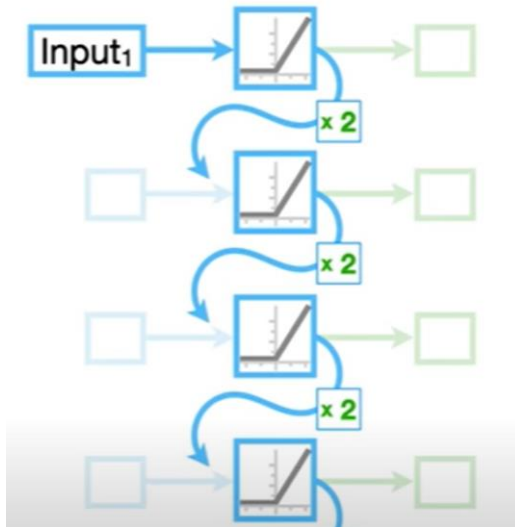
- Lots of publicly available data
- Predicting number of COVID deaths, flu cases, hospital beds, etc.
- Used for seasonal illnesses
- Monitoring and predicting disease spread
- Informing decisions about resource allocation (vaccines, traveling nurses, etc.)

## Models in Healthcare

- ARIMA / SARIMA
- LSTM
- Linear Regression
- Multilayer perceptron
- Vector autoregression

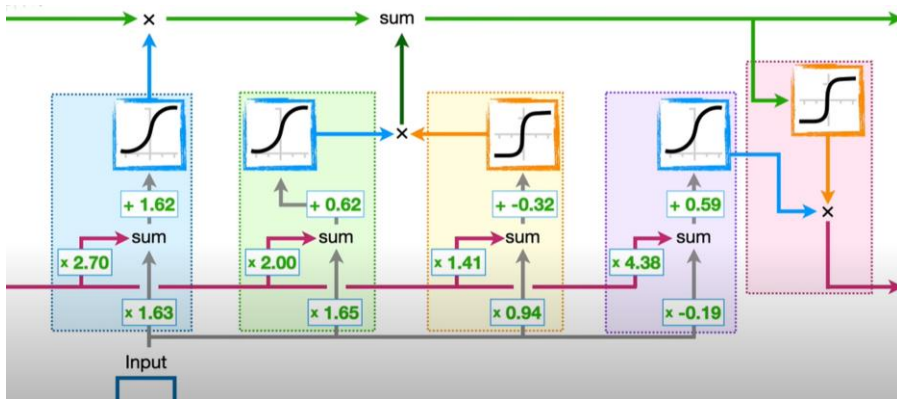
# Forecasting Methods: Neural Networks

- RNN: using the same feedback loop connection, so the gradient can explode or vanish. It uses different paths for long and short-term memories.
- LSTM: a type of RNN, which is designed to avoid exploding or vanishing

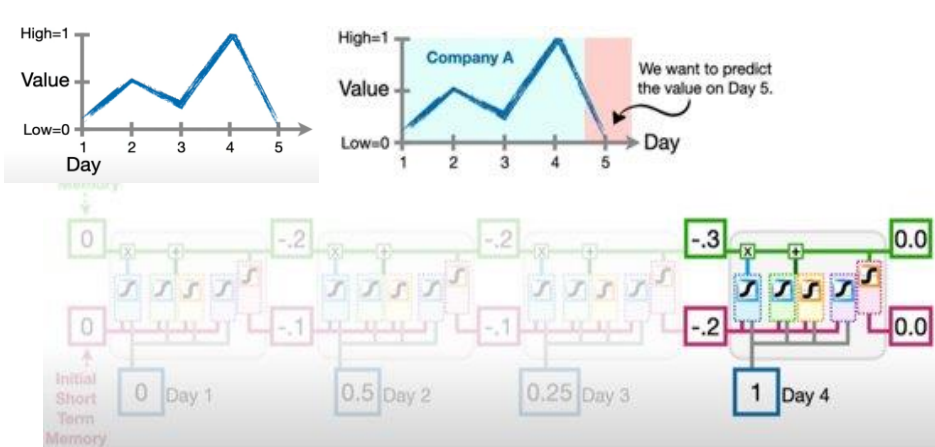


# Forecasting Methods: LSTM

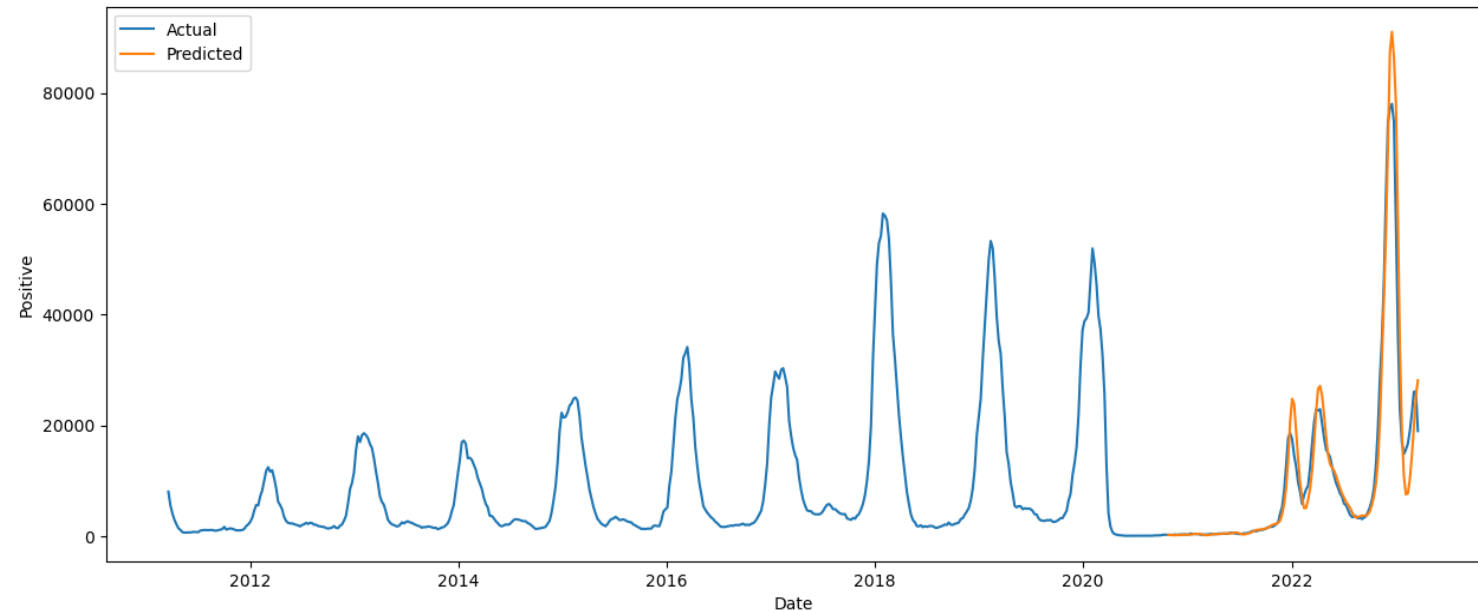
## How the LSTM unit works



- Our work using LSTM for influenza forecasting
- Good forecasting
- Complicated and time consuming (with the help of GPT)



Example: Forecasting stock price at Day5



# Forecasting Methods: ARIMA

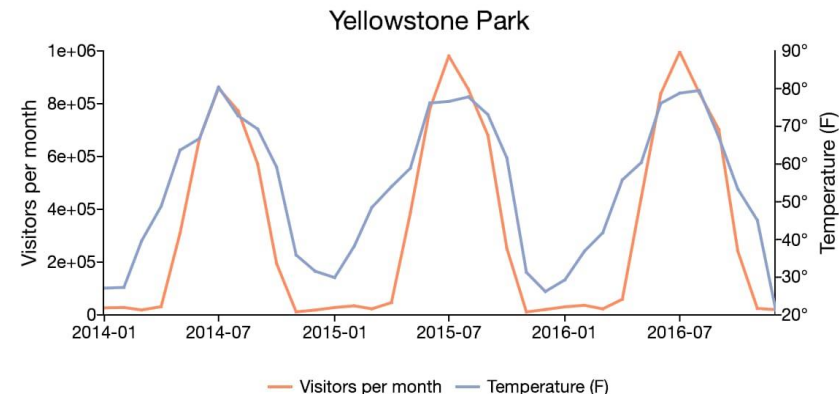
---

## Autoregressive Integrated Moving Average Model

- a form of regression analysis that gauges the strength of one dependent variable relative to other changing variables.
- ARIMA models use differencing to convert a non-stationary time series into a stationary one, and then predict future values from historical data.

### Variations:

- SARIMA
  - Seasonal ARIMA
- SARIMAX
  - SARIMA using exogenous variable(s)





# ARIMA

## Parameters:

- Autoregression (p) - a model that shows a changing variable that regresses on its own lagged, or prior, values.
- Integrated (d) - the differencing of raw observations to allow the time series to become stationary (i.e., data values are replaced by the difference between the data values and the previous values).
- Moving average (q) - dependency between an observation and a residual error from a moving average model applied to lagged observations.

## Pros and Cons:

### Pros

- Good for short-term forecasting
- Only needs historical data
- Models non-stationary data
- Easy to understand and interpret

### Cons

- Not built for long-term forecasting
- Poor at predicting turning points
- Computationally expensive
- Parameters are subjective



[Back to Top](#)

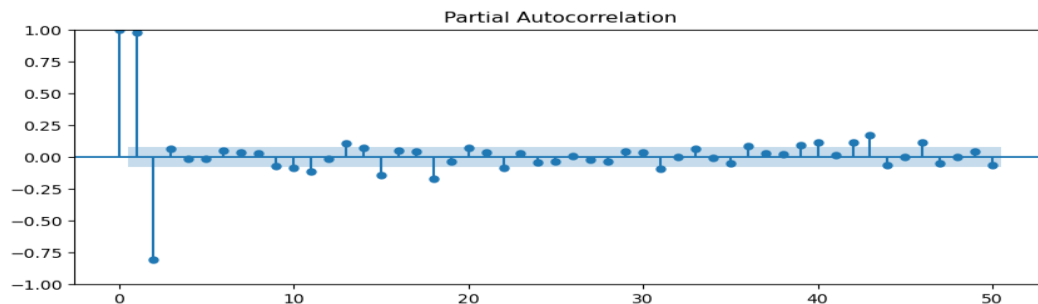
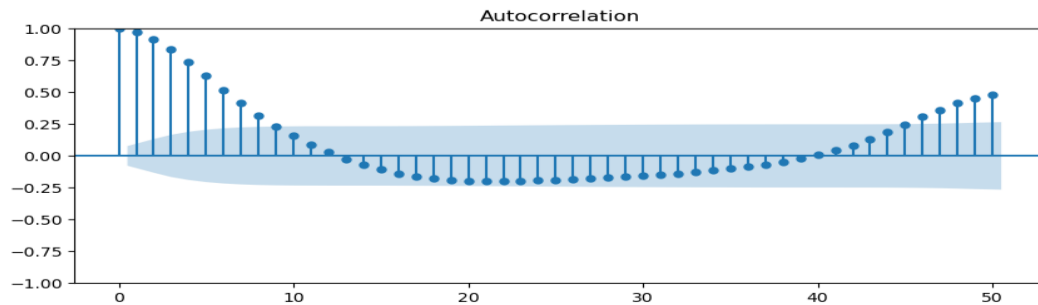
```
In [221]: df = df.drop(columns = ['ISO_WEEK', 'SPEC_RECEIVED_NB', 'SPEC_PROCESSED_NB', 'AH1',  
                                'AH5', 'INF_A', 'ANOTSUBTYPED', 'BVIC', 'BYAM', 'BNOTDETER',  
                                'INF_NEGATIVE', 'ILI_ACTIVITY', 'WHOREGION', 'ITZ', 'COUNTRY',  
                                'COUNTRY_CODE', 'ORIGIN_SOURCE', 'ISO_YEAR']) # Dropping all
```

```
In [222]: df.info()
```

# Code Walkthrough

Augmented Dickey-Fuller Statistic: -7.814637  
p-value: 0.000000  
critical values at different levels:  
1%: -3.441  
5%: -2.866  
10%: -2.569

	Positive
Date	
2010-12-05 18:30:00+00:00	3699.0
2010-12-12 18:30:00+00:00	5755.0
2010-12-19 18:30:00+00:00	9384.0
2010-12-26 18:30:00+00:00	11594.0
2011-01-02 18:30:00+00:00	15185.0



## ARIMA

```
model = ARIMA(df['Positive'], order=(1,1,1))  
results = model.fit()  
forecast = results.predict(start='2018-01-21 18:30:00+00:00', end='2023-03-19 18:30:00+00:00')  
# plt.plot(df['Positive'])
```

## ARIMA Seasonal Ordered

```
model = ARIMA(df['Positive'], order=(1,1,1), seasonal_order=(1,1,1,52))  
results = model.fit()  
forecast = results.predict(start='2018-01-21 18:30:00+00:00', end='2023-03-19 18:30:00+00:00')  
# plt.plot(df['Positive'])
```

A Basic Sarimax model is now run to see if it fits better.

```
# Sarima Model  
model = sm.tsa.statespace.SARIMAX(df, order=(1,1,1), seasonal_order=(1,1,1,52))
```

A Model is Tuned using ranges 1 to 5 for the order.

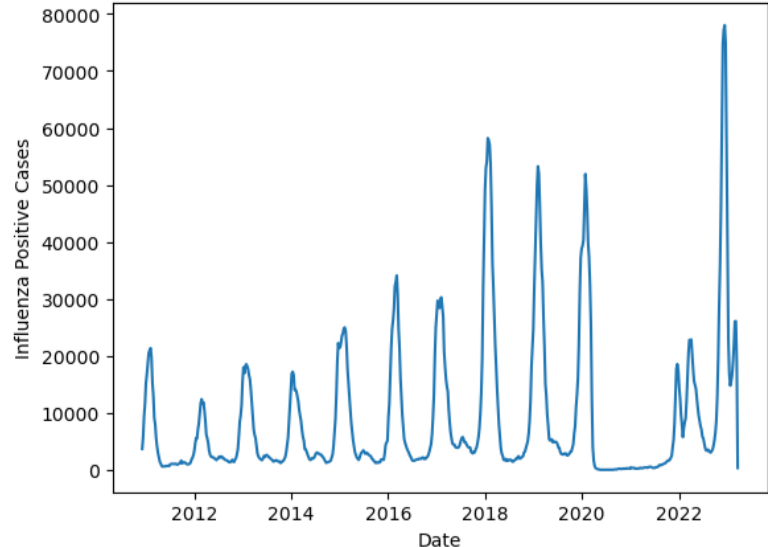
```
# Model Tuning using AutoArima  
tunedmodel = pm.auto_arima(df, start_p=1, start_q=1, max_p=5, max_q=5, m=52, seasonal=True,  
                           trace=True, error_action='ignore', suppress_warnings=True, stepwise=True)
```

Performing stepwise search to minimize aic

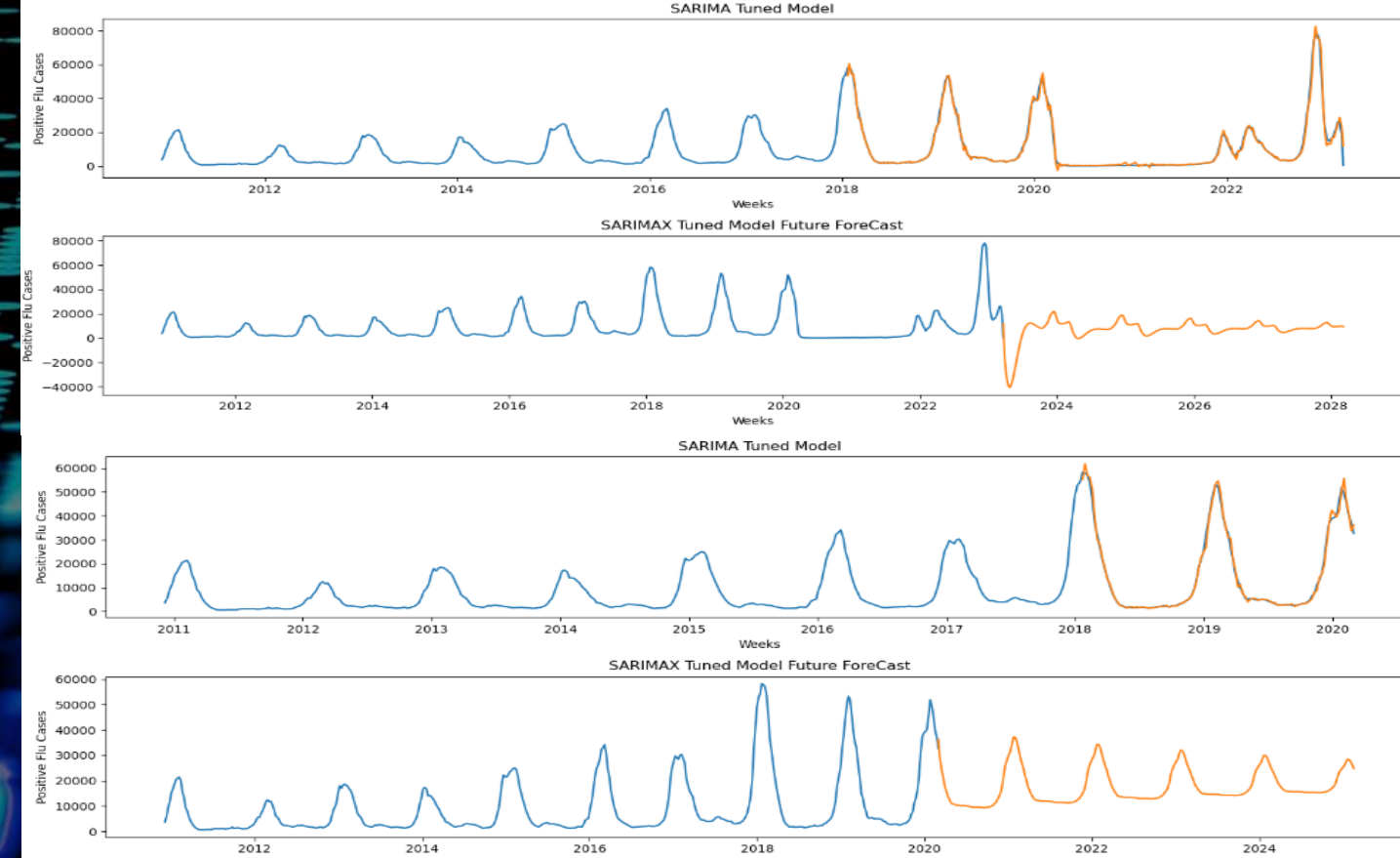
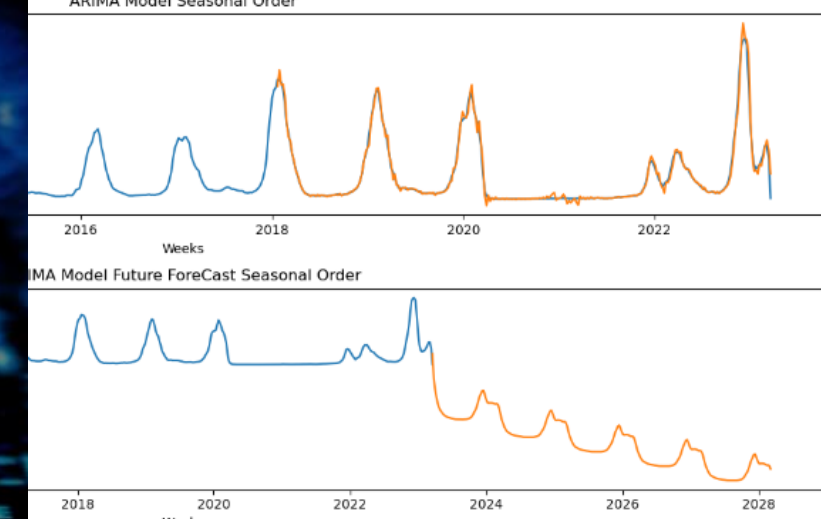
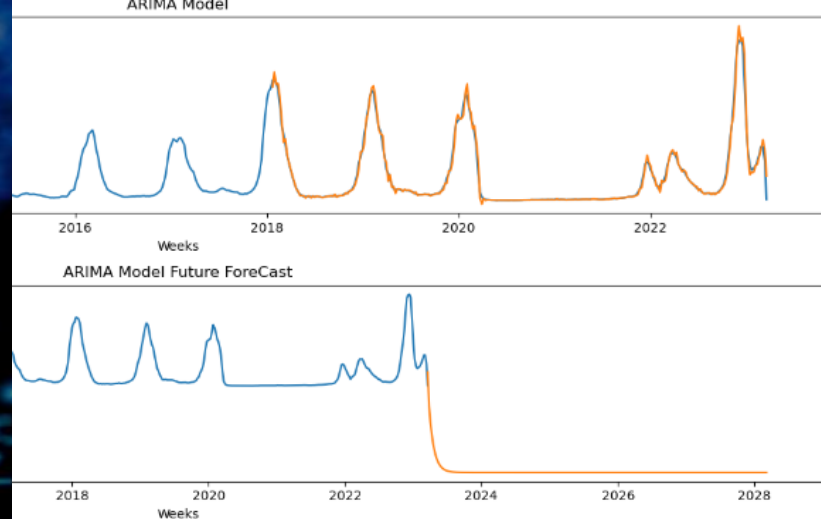
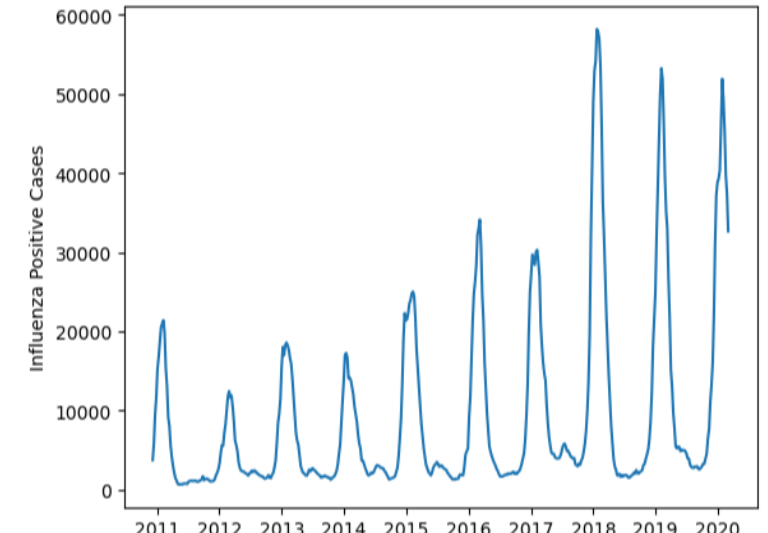
```
ARIMA(1,1,1)(1,0,1)[52] intercept : AIC=11222.623, Time=1.88 sec  
ARIMA(0,1,0)(0,0,0)[52] intercept : AIC=11985.915, Time=0.04 sec  
ARIMA(1,1,0)(1,0,0)[52] intercept : AIC=11267.792, Time=1.07 sec  
ARIMA(0,1,1)(0,0,1)[52] intercept : AIC=11476.258, Time=2.69 sec  
ARIMA(0,1,0)(0,0,0)[52] intercept : AIC=11983.918, Time=0.04 sec
```

# fit a SARIMAX model using the optimal order parameters

```
model = sm.tsa.statespace.SARIMAX(df, order=tunedmodel.order, seasonal_order=tunedmodel.seasonal_order)  
results = model.fit()  
forecast = results.predict(start='2018-01-21 18:30:00+00:00', end='2023-03-19 18:30:00+00:00')
```



# Graph Walkthrough





# Forecasting Metrics

- Mean Squared Error (MSE)
  - penalizes large errors or outliers more compared to the small errors
- Root Mean Squared Error (RMSE)
  - RMSE value is in the same unit as the forecasted value
  - easier to understand compared to MSE
- Mean Absolute Percent Error (MAPE)
  - measures the average of forecast errors in percentages
  - works best with data without zeros and extreme values
- Root Mean Squared Log Error (RMSLE)
  - not scale-dependent and is useful across a range of scales
  - It is not affected by large outliers.
  - It considers only the relative error between the actual value and the predicted value.

Metrics used to evaluate the final model

Mean Squared Error (MSE)	2564043.054
Root Mean Squared Error (RMSE)	1601.262956
Mean Absolute Percentage Error (MAPE)	7.837775809
Root Mean Squared Log Error (RMSLE)	0.101783719

Other metrics used in forecasting

- Huber Loss
- Log Cosh Loss
- Quantile Loss

# Challenges and Conclusions

## Challenges

- Publicly available data messy or limited
- Hyperparameter tuning
- Disruption in seasonal patterns from COVID pandemic
- Different metrics for forecasting models (no precision, recall, F1)
- Google Collab problems
- LSTM failed lots of times, very slow
- Choosing between ARIMA and SARIMA
- Controlling Dataset

## Conclusions

- Forecasting is a powerful tool
- Time series models are unique
- MSE lower for SARIMA than ARIMA
- Dataset not necessarily accurate due to Covid Disruptions

## Next Steps

- Improve Model using Exogenous variables.
- Learn to read ACF & PACF Plots better to tune orders.