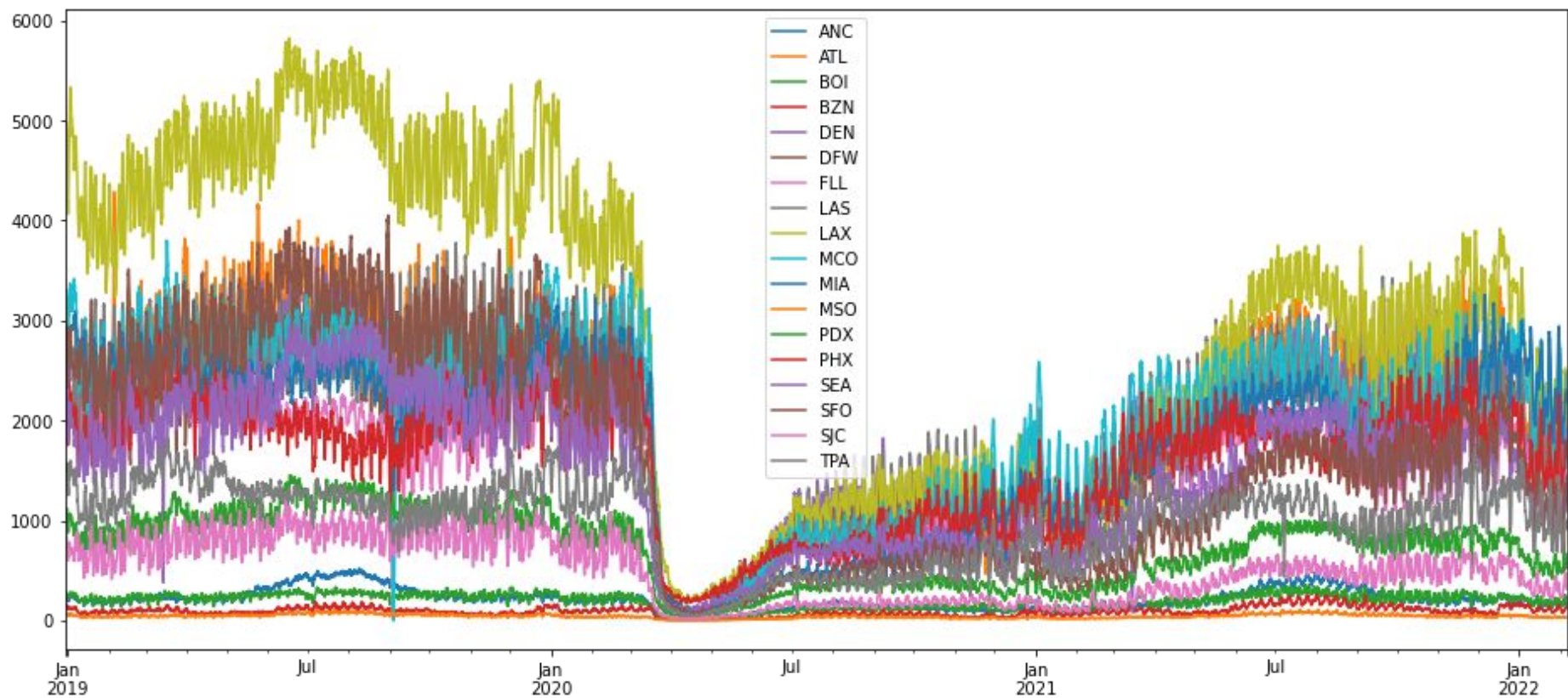# Predict TSA Throughput

A time series analysis

# Problem Identification and Business Context

# Predictive model - TSA Throughput

Purpose:

- Allow airports and airlines make better resource allocation and staffing decisions

- TSA efficiency and customer service

- Recovering from Covid-19 travel restrictions beginning in March 2020

## Stakeholders:

- Airports
- Airlines
- Airport businesses

## Challenges:

- Extreme weather
- War
- Pandemics

# Data Collection

Data gathering: Github repository (source)

- Individual Csv files for each airport

- Columns for each gate

- Values: TSA throughput (number of people going through security)

# Data Collection

Building final Dataframe:

- Aggregate all gates within an airport

- A column for each airport

- Datetime Index (hourly)

# Data Exploration and Cleaning

# Data

**18 columns**: US International Airports

**27216 rows**: Hourly Data between

　　　　　December 30th 2018 – February 5th 2022

**Values**: Number of people going through security

# Airports

ANC - Anchorage

ATL - Atlanta

BOI - Boise

BZN - Bozeman

DEN - Denver

DFW - Dallas Fort Worth

FLL - Fort Lauderdale

LAS - Las Vegas

LAX - Los Angeles

MCO - Orlando

MIA - Miami

MSO - Missoula Montana

PDX - Portland

SEA - Seattle

SFO - San Francisco

SJC - San Jose

TPA - Tampa

# Boxplots



Boxplot of TSA Throughput for 18 US Airports

# Histograms

# Histograms



SFO

# Histograms



MCO

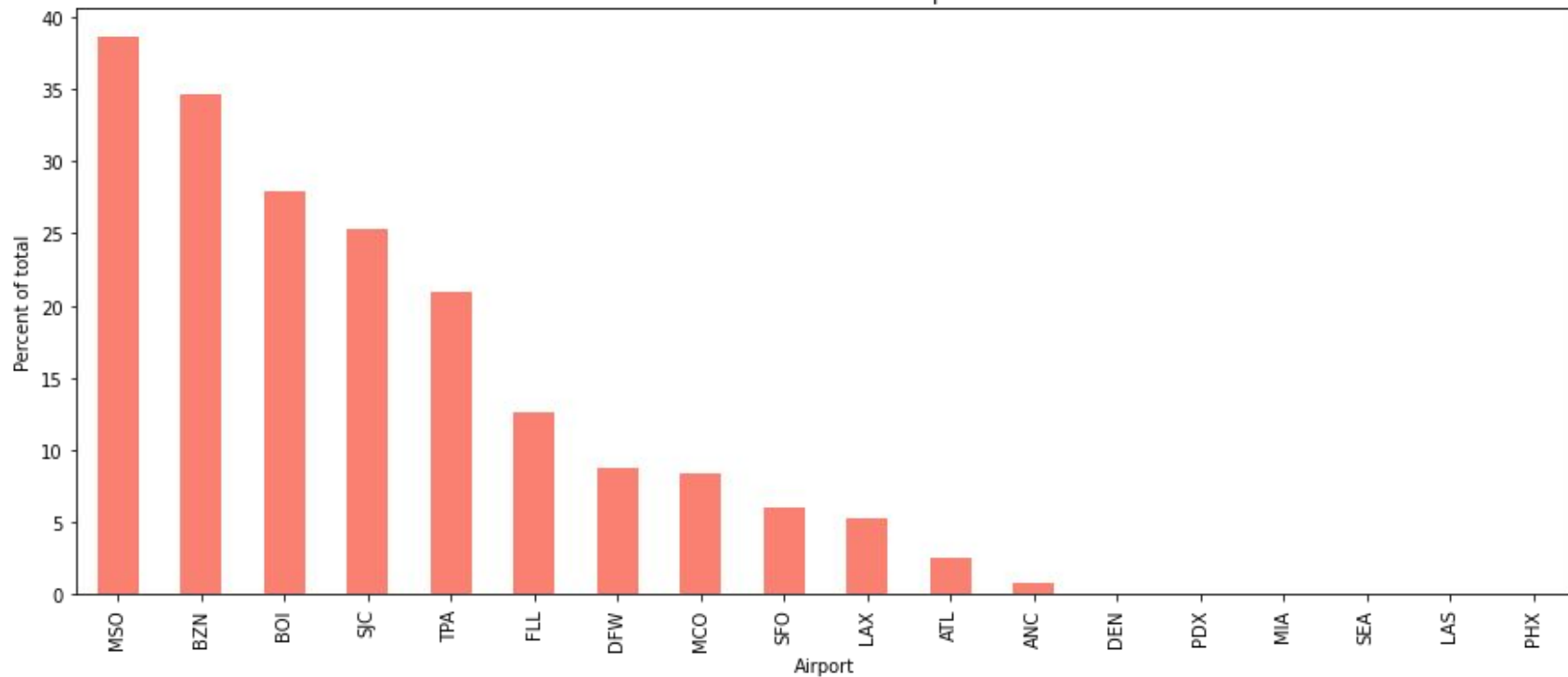# Histograms

# Density Plots



Density plots for each airport

# Data Distribution

- First peak for relatively low counts

- Large airports have a second peak at higher counts

- Extremely busy times with relatively large throughput are much more rare

# Null values



Percent null for each airport

Null value distribution per hour for each airport

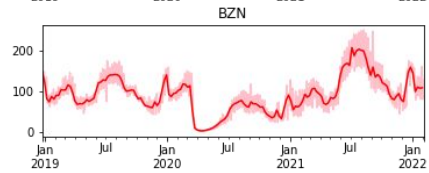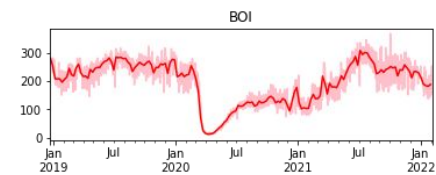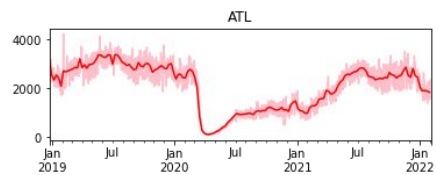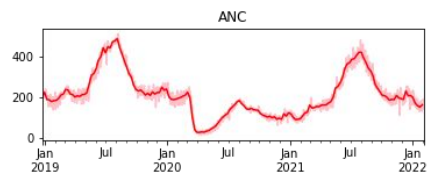TSA Throughput Timeseries (2-day slice)

# Null value treatment

- Small airports have a higher ratio of values missing than large airports

  (Outgoing flight times for small vs large airports)

- Missing values are concentrated in late evening to early morning hours

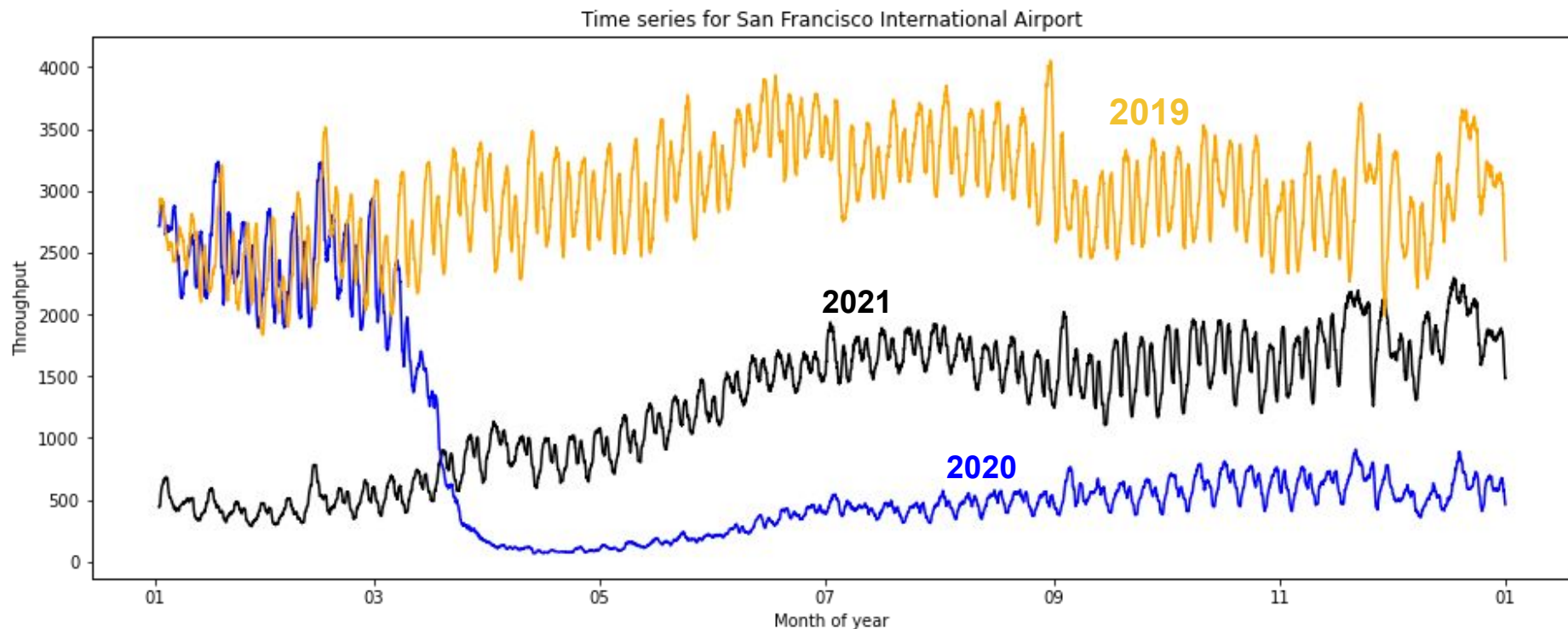  (Least common times for outgoing flights for large airports)

**Impute null values with ZERO**
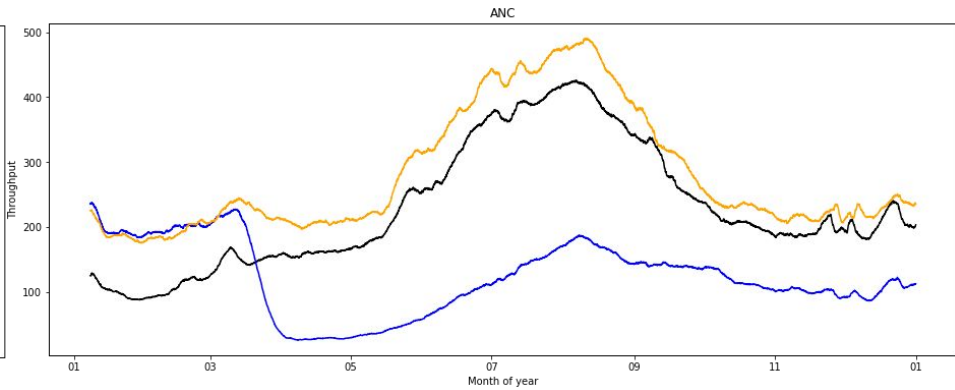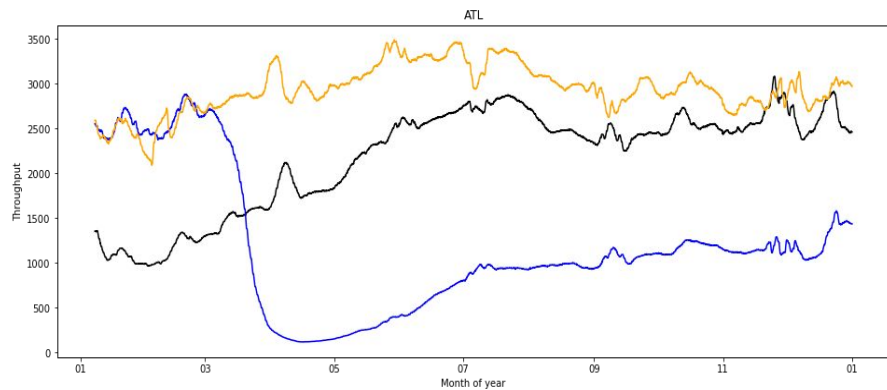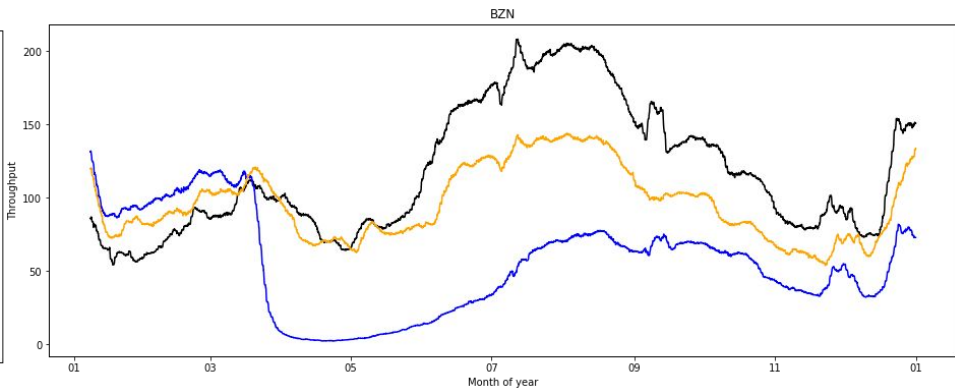
# Time Series Plot - All Airports



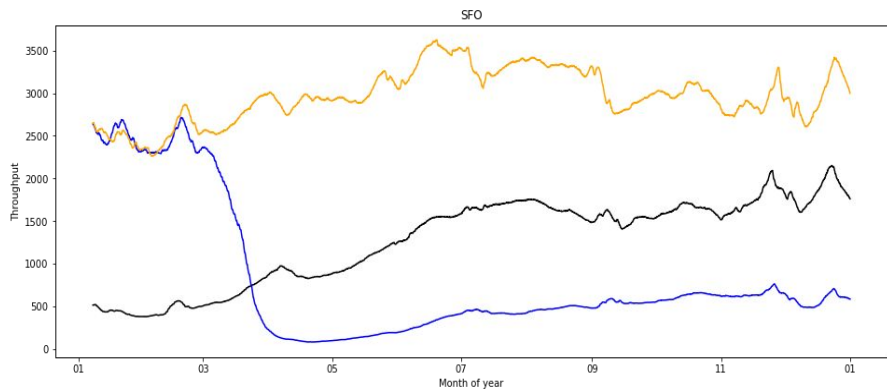TSA Throughput Timeseries

# Yearly Trends - SFO



Time series for San Francisco International Airport

# Yearly Trends

# Preprocessing and Training

Autocorrelation

# Seasonality - SFO

Strong seasonality with period of 24 (hours)

# Augmented Dicky-Fuller Test

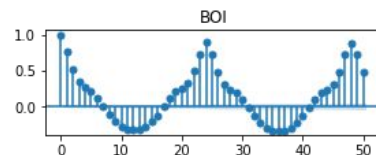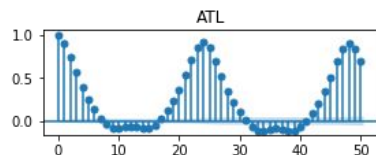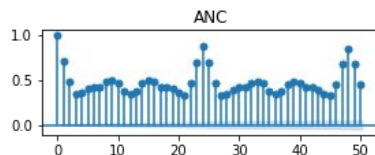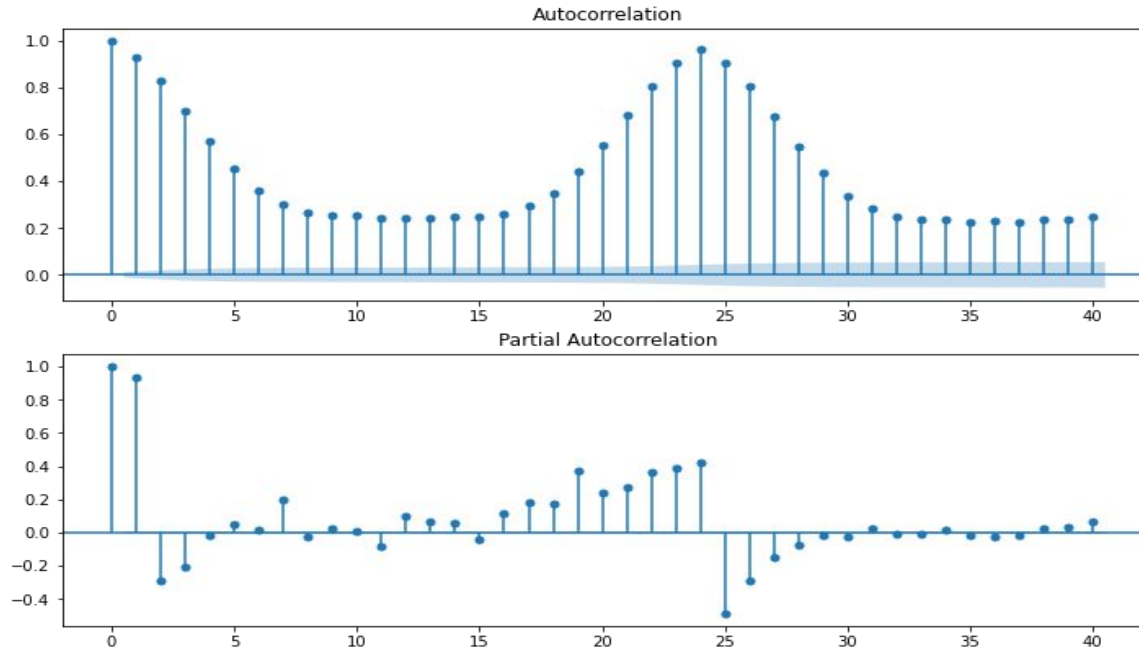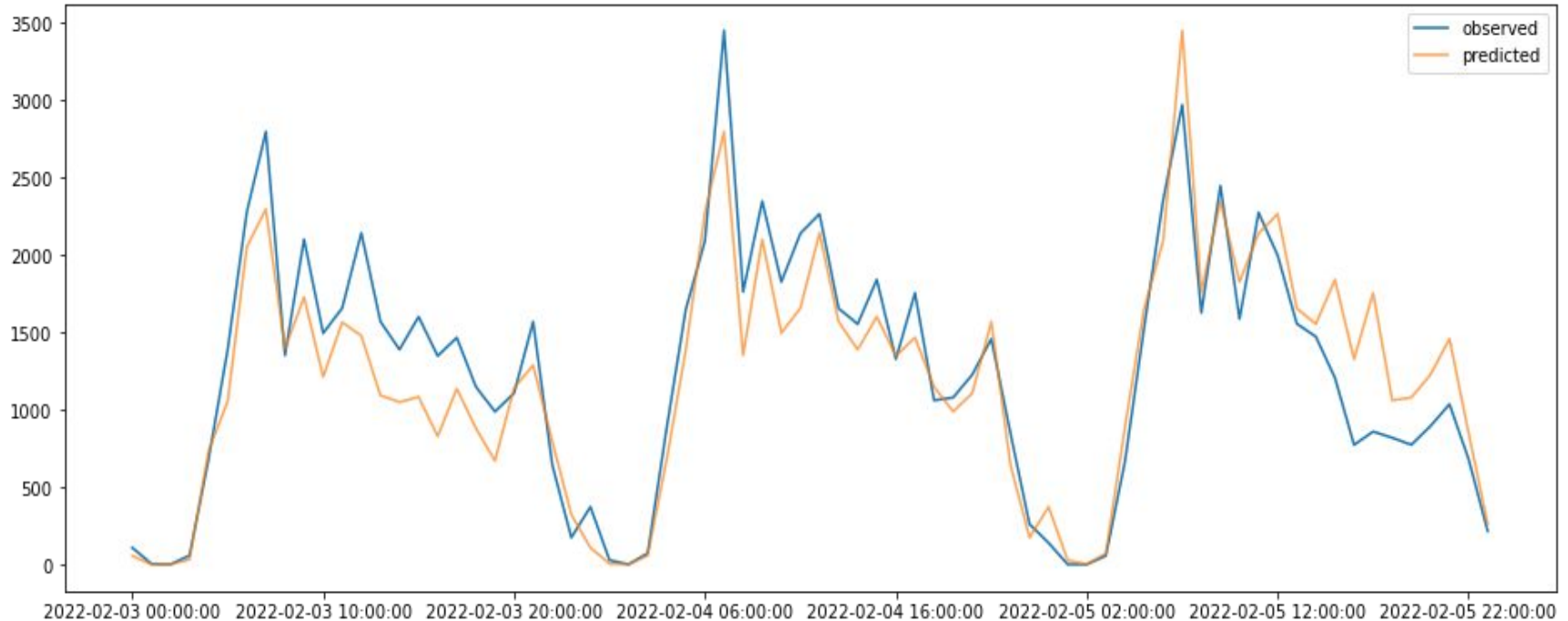- Null hypothesis - time series is non-stationary
- Only tests for trend
- Reject null if p-value is small (less than 5%)

# Baseline Model - Yesterday's values

Mean Error = 276 people

# Modeling

# ARIMA Model

Grid search

$(p, d, q)$ X $(P, D, Q, S)$

# Evaluate Model

# Prediction

# Future work: