



*Project
Capstone*

SALES FORECASTING PROJECT

*Presented by: Nurhidayah
Advised by: Dr Nor Azizah Ali*



CONTENT

01

Problem Statement

02

Data Description

03

Data Preprocessing

04

EDA

05

Model Building

06

Dashboard





PROBLEM STATEMENT

The business problem revolves around improving profitability and operational efficiency through better insights into product sales and profit trends. Two key business insight to be focused on are profitability analysis and sales forecasting.

First, profitability analysis is to identify patterns of profit loss, both geographically and by product category, is critical for the sales department to understand underperforming areas and products. This will help in making strategic decisions regarding product focus, marketing efforts, and regional investments.

Second, sales forecasting is to estimate future sales for January 2022. This will help in planning inventory, marketing, and resource allocation. Accurate forecasting ensures that



DATA DESCRIPTION

Provided with 3 Excel files (Orders_2017.xlsx, Orders_2018.xlsx, Orders_2019.xlsx) and 2 CSV files (Orders_2020.csv, Orders_2021.csv) that contain information about order date, customer, region, product category, sub-category, sales, profit, shipping cost, and return status.



SAMPLE DATA

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	Category	City	Country/Region	Discount	Order ID	Postal Code	Product Name	Profit	Quantity	Region	Sales	Segment	Ship Mode	State	Sub-Category	OrderDate	ShipDate	Days to Ship	Scheduled
2	Furniture	Fort Lauderdale	United States	45	US-2017-108	33311	Bretford CR4500 Series Slim F	-383	5	South	958	Consumer	Standard Class	Florida	Tables	11/10/2017	18/10/2017	6	
3	Office Supplies	Fort Lauderdale	United States	20	US-2017-108	33311	Eldon Fold 'N Roll Cart System	3	2	South	22	Consumer	Standard Class	Florida	Storage	11/10/2017	18/10/2017	6	
4	Office Supplies	Fort Worth	United States	80	US-2017-118	76106	Holmes Replacement Filter fo	-124	5	Central	69	Home Office	Standard Class	Texas	Appliances	22/11/2017	26/11/2017	6	
5	Office Supplies	Fort Worth	United States	80	US-2017-118	76106	Storex DuraTech Recycled Pla	-4	3	Central	3	Home Office	Standard Class	Texas	Binders	22/11/2017	26/11/2017	6	
6	Furniture	Orem	United States	0	CA-2017-106	84057	Bretford CR4500 Series Slim F	240	3	West	1045	Consumer	Standard Class	Utah	Tables	25/9/2017	30/9/2017	6	

1 Row ID,Order ID,Order Date,Ship Date,Ship Mode,Customer Name,Segment,Country/Region,City,State,Postal Code,Region,Category,Sub-Category,Product Name,Sales,Quantity,Discount,Profit,Days to Sh
 2 1,CA-2020-152156,8/11/2020 0:00,11/11/2020 0:00,Second Class,Claire Gute,Consumer,United States,Henderson,Kentucky,42420,South,Furniture,Bookcases,Bush Somerset Collection Bookcase,261.96,2,
 3 2,CA-2020-152156,8/11/2020 0:00,11/11/2020 0:00,Second Class,Claire Gute,Consumer,United States,Henderson,Kentucky,42420,South,Furniture,Chairs,"Hon Deluxe Fabric Upholstered Stacking Chairs
 4 3,CA-2020-138688,12/6/2020 0:00,16/6/2020 0:00,Second Class,Darrin Van Huff,Corporate,United States,Los Angeles,California,90036,West,Office Supplies,Labels,Self-Adhesive Address Labels for
 5 14,CA-2020-161389,5/12/2020 0:00,10/12/2020 0:00,Standard Class,Irene Maddox,Consumer,United States,Seattle,Washington,98103,West,Office Supplies,Binders,Fellowes PB200 Plastic Comb Binding
 6 22,CA-2020-137330,9/12/2020 0:00,13/12/2020 0:00,Standard Class,Ken Black,Corporate,United States,Fremont,Nebraska,68025,Central,Office Supplies,Art,Newell 318,19.46,7,0,5.0596,6
 7 23,CA-2020-137330,9/12/2020 0:00,13/12/2020 0:00,Standard Class,Ken Black,Corporate,United States,Fremont,Nebraska,68025,Central,Office Supplies,Appliances,"Acco Six-Outlet Power Strip, 4' C
 8 26,CA-2020-121755,16/1/2020 0:00,20/1/2020 0:00,Second Class,Eric Hoffmann,Consumer,United States,Los Angeles,California,90049,West,Office Supplies,Binders,Wilson Jones Active Use Binders,11
 9 27,CA-2020-121755,16/1/2020 0:00,20/1/2020 0:00,Second Class,Eric Hoffmann,Consumer,United States,Los Angeles,California,90049,West,Technology,Accessories,Imation 8GB Mini TravelDrive USB 2.
 10 36,CA-2020-117590,8/12/2020 0:00,10/12/2020 0:00,First Class,Gene Hale,Corporate,United States,Richardson,Texas,75080,Central,Technology,Phones,GE 30524EE4,1097.544,7,0.2,123.4737,1
 11 37,CA-2020-117590,8/12/2020 0:00,10/12/2020 0:00,First Class,Gene Hale,Corporate,United States,Richardson,Texas,75080,Central,Furniture,Furnishings,"Electrix Architect's Clamp-On Swing Arm L
 12 43,CA-2020-101343,17/7/2020 0:00,22/7/2020 0:00,Standard Class,Ruben Ausman,Corporate,United States,Los Angeles,California,90049,West,Office Supplies,Storage,"Eldon Base for stackable storag

DATA PREPROCESSING



DATA PREPROCESSING

Import Dataset

```
1 #import libraries
2 import gdown
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import chardet
```

```
1 # download data from drive
2 shared_url = 'https://drive.google.com/file/d/
1y35PeUIW-B02RDqY0zp-o9g7ddWdnKXZ/view?usp=sharing'
3 file_url    = f'https://drive.google.com/uc?id={shared_url.split("/")[-2]}'
4 output = 'ORDERS_2021.csv'
5 gdown.download(file_url, output, quiet=False)
6
7 shared_url = 'https://drive.google.com/file/d/
1dRj1YrNvHknNIvT6eWNhe9w0_9GU6YXN/view?usp=sharing'
8 file_url    = f'https://drive.google.com/uc?id={shared_url.split("/")[-2]}'
9 output = 'ORDERS_2020.csv'
10 gdown.download(file_url, output, quiet=False)
11
12 shared_url = 'https://docs.google.com/spreadsheets/d/
1yn7J3i5gVfYEZEDXCVbOT8s-ITyLiBzN/edit?usp=sharing&ouid=115467567108565679378&
rtpof=true&sd=true'
13 file_url    = f'https://drive.google.com/uc?id={shared_url.split("/")[-2]}'
14 output = 'ORDER_2019.xlsx'
15 gdown.download(file_url, output, quiet=False)
16
17 shared_url = 'https://docs.google.com/spreadsheets/d/
19AjIVfQWU6i8Qlif_EUkt3mi0BNvlB01/edit?usp=sharing&ouid=115467567108565679378&
rtpof=true&sd=true'
18 file_url    = f'https://drive.google.com/uc?id={shared_url.split("/")[-2]}'
19 output = 'ORDER_2018.xlsx'
20 gdown.download(file_url, output, quiet=False)
21
```

DATA PREPROCESSING

Explore Dataset

```
1 # add ORDER_2018 data into dataframe
2 df_order2018 = pd.read_excel('ORDER_2018.xlsx')
3 df_order2018.head()
```

	Category	City	Country/Region	Discount	Order ID	Postal Code	Product Name	Profit	Quantity	Region	Sales	Segment	Ship Mode	State	Sub-Category	OrderDate	ShipDate	Days to Ship Scheduled
0	Furniture	Henderson	United States	0	CA-152156	42420.0	Bush Somerset Collection Bookcase	42	2	South	262	Consumer	Second Class	Kentucky	Bookcases	2018-11-08	2018-11-11	3.0
1	Furniture	Henderson	United States	0	CA-152156	42420.0	Hon Deluxe Fabric Upholstered Stacking Chairs,...	220	3	South	732	Consumer	Second Class	Kentucky	Chairs	2018-11-08	2018-11-11	3.0
2	Office Supplies	Los Angeles	United States	0	CA-138688	90036.0	Self-Adhesive Address Labels for Typewriters b...	7	2	West	15	Corporate	Second Class	California	Labels	2018-06-12	2018-06-16	3.0
3	Office Supplies	Seattle	United States	20	CA-161389	98103.0	Fellowes PB200 Plastic Comb Binding Machine	133	3	West	408	Consumer	Standard Class	Washington	Binders	2018-12-05	2018-12-10	6.0

```
1 # check entries and column
2 df_order2021.shape, df_order2020.shape, df_order2019.shape, df_order2018.shape, df_order2017.shape
```

```
((3312, 20), (2587, 20), (4278, 18), (3179, 18), (2522, 18))
```

DATA PREPROCESSING

Data Filtration

```
1 # Drop extra columns
2 df_order2021.drop(['Row ID', 'Customer Name'], axis=1, inplace=True)
3 df_order2020.drop(['Row ID', 'Customer Name'], axis=1, inplace=True)
4
5 df_order2021.shape, df_order2020.shape
```

DATA PREPROCESSING

Data Cleaning

```
1 # rename columns in df_order2019, df_order2018, and df_order2017
2 dfList = [df_order2017,df_order2018,df_order2019]
3 for df in dfList:
4     df.rename(columns={'OrderDate': 'Order Date', 'ShipDate':'Ship Date'},
5               inplace=True)
5     print(df.columns)
```

```
1 # Convert datatype object to datetime for columns in df_order2021 and
2 # df_order2020
3 dfList = [df_order2021,df_order2020]
4 for df in dfList:
5     df['Order Date'] = pd.to_datetime(df['Order Date'], format="%d/%m/%Y %H:%M")
6     df['Ship Date'] = pd.to_datetime(df['Ship Date'], format="%d/%m/%Y %H:%M")
7     print(df['Order Date'].dtypes)
8     print(df['Ship Date'].dtypes)
```

DATA PREPROCESSING

Export

```
1 #save into excel  
2 order_return_df.to_excel('order_return_df_3.xlsx', index=False)
```

EXPLORATORY DATA ANALYSIS

```
1 # Descriptive Statistics for Numerical Columns
2 order_return_df.describe()
```

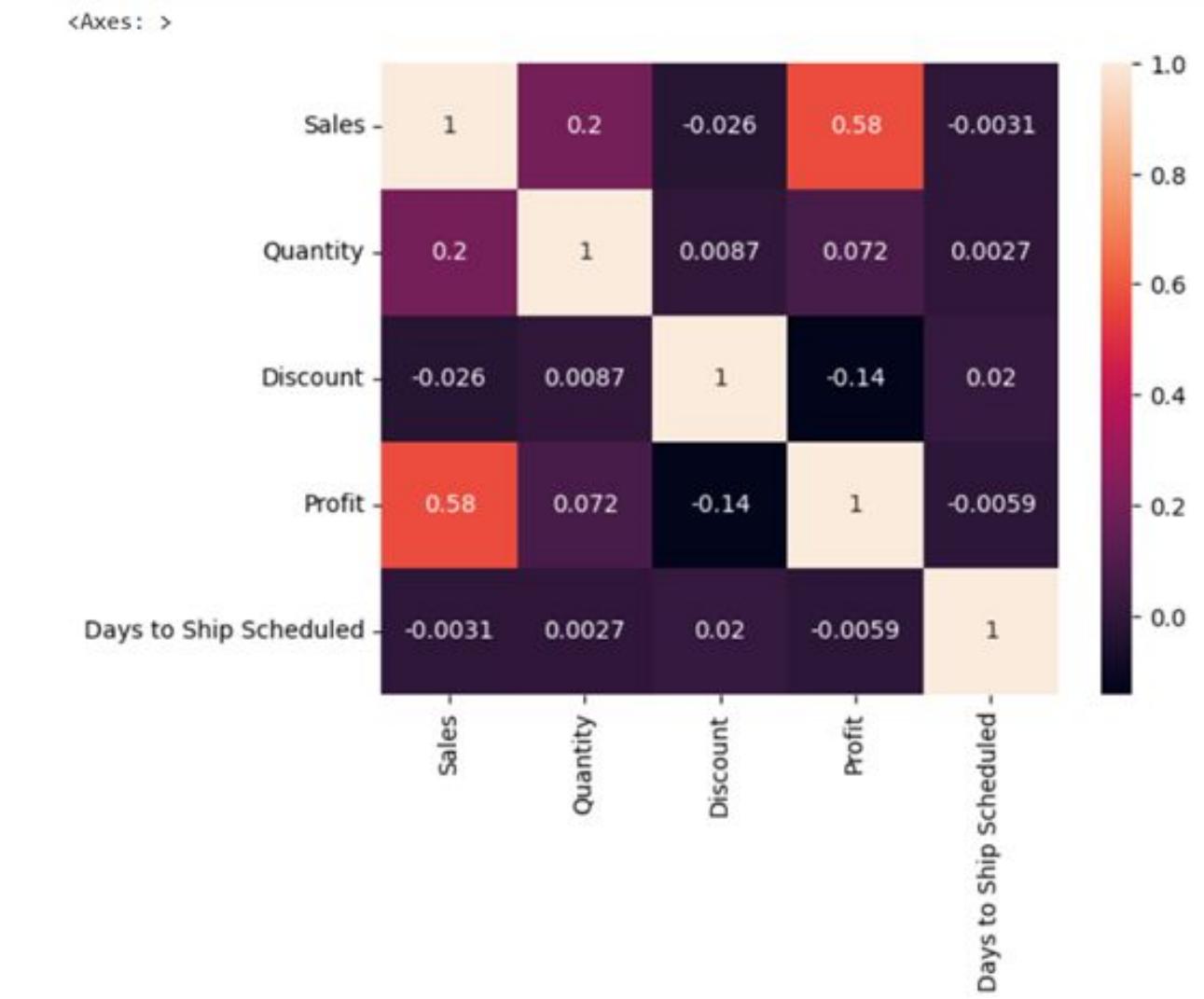
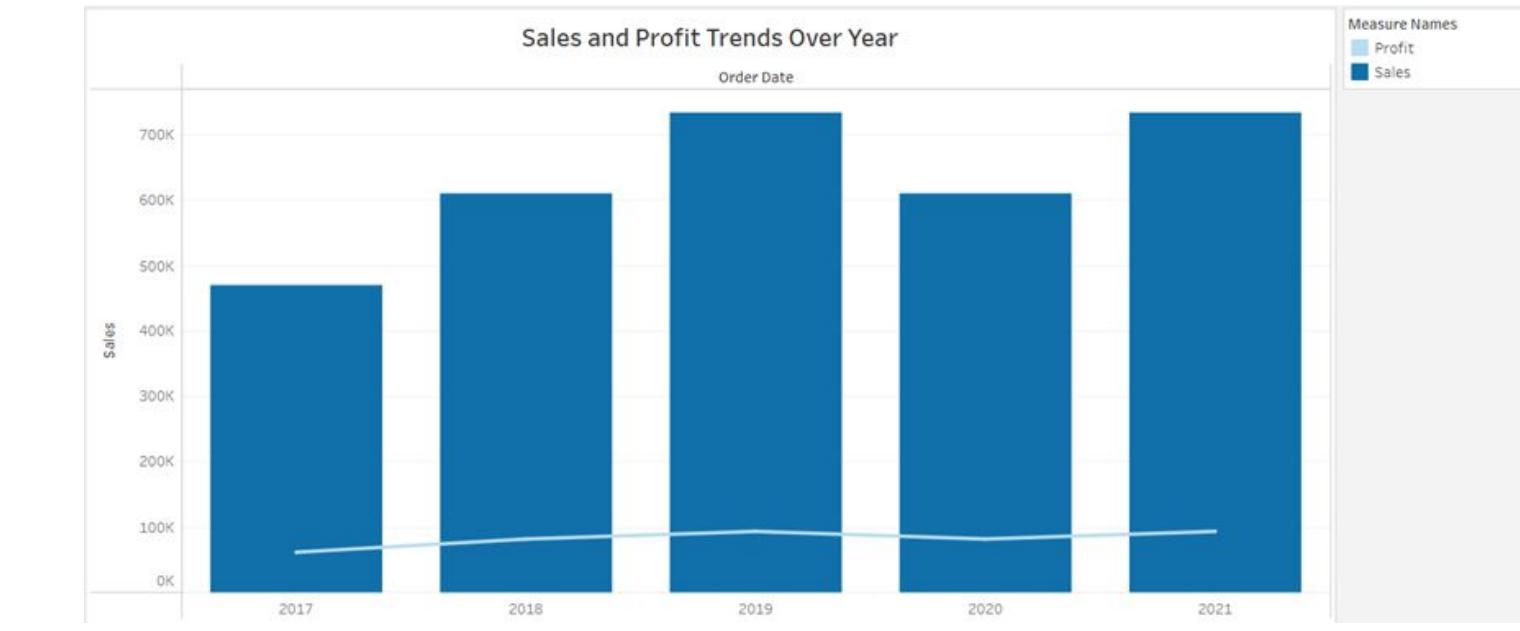
	Order Date	Ship Date	Sales	Quantity	Discount	Profit	Days to Ship	Scheduled
count	13900	13900	13900.000000	13900.000000	13900.000000	13900.000000	13900.000000	13900.000000
mean	2019-10-11 14:48:20.719424512	2019-10-15 13:19:39.971222784	227.015817	3.784532	9.027449	29.644708	4.491649	
min	2017-01-02 00:00:00	2017-01-04 00:00:00	0.000000	1.000000	0.000000	-6600.000000	1.000000	
25%	2018-09-04 00:00:00	2018-09-07 00:00:00	17.380000	2.000000	0.000000	2.000000	3.000000	
50%	2019-10-19 00:00:00	2019-10-22 00:00:00	54.686000	3.000000	0.200000	8.750700	6.000000	
75%	2020-12-13 00:00:00	2020-12-18 00:00:00	207.036000	5.000000	20.000000	29.000000	6.000000	
max	2021-12-30 00:00:00	2022-01-05 00:00:00	17500.000000	14.000000	80.000000	8400.000000	6.000000	
std		Nan	602.037902	2.217501	17.430890	249.769539	1.948143	

```
[40] 1 # Value counts for categorical variables
2 order_return_df['Region'].value_counts()
```

```
count
Region
West    4442
East    4018
Central 3238
South   2202
```

dtype: int64

EXPLORATORY DATA ANALYSIS



EXPLORATORY DATA ANALYSIS

MODEL BUILDING

- Phase where predictive models are constructed, trained, and assessed for their performance.
- Metrics such as accuracy, precision, recall, and F1-score are employed to measure model performance, and techniques like cross-validation help estimate generalization.
- The goal is to build models that provide actionable insights or predictions, ensuring their effectiveness and reliability for real-world applications.

MODEL BUILDING

ARIMA

ARIMA stands for AutoRegressive Integrated Moving Average, a statistical time series model commonly used for forecasting. It's popular in cases where data shows patterns over time and can be decomposed into components such as trend, seasonality, and noise. ARIMA models are especially effective for series that are non-stationary but can be made stationary through differencing.

```
1  from statsmodels.tsa.arima.model import ARIMA  
2  
3  # fit model  
4  model_arima = ARIMA (daily_sales_final['Sales'], order = (1,1,0))  
5  model_fit_arima = model_arima.fit()  
6  
7  # summary of fit model  
8  print(model_fit_arima.summary())
```

```
[▶] 1  # split into train and test sets  
2  X = daily_sales_final['Sales']  
3  size = int(len(X) * 0.80)  
4  train, test = X[0:size], X[size:len(X)]  
5  history = [x for x in train]  
6  
7  
8  # walk-forward validation  
9  prediction_arima = test.copy()  
10 for t in range(len(test)):  
11     model = ARIMA (history, order = (1,1,0))  
12     model_fit = model.fit()  
13     output = model_fit.forecast()  
14     yhat = output[0]  
15     prediction_arima.iloc[t] = yhat # Use iloc to assign by position  
16     obs = test.iloc[t] # Use iloc to access by position  
17     history.append(obs)  
18     print('predicted=%f, expected=%f' % (yhat, obs))
```

MODEL BUILDING

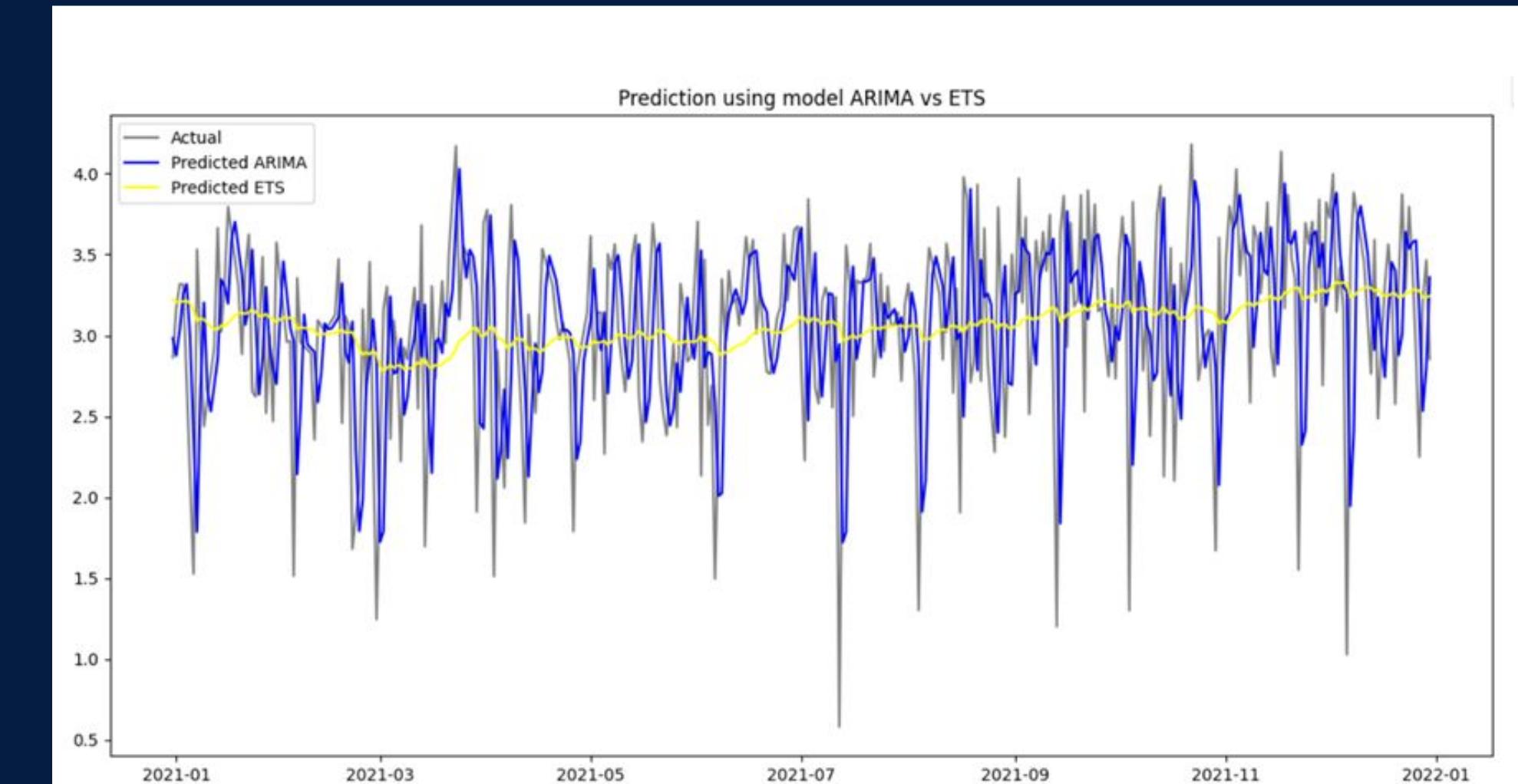
ETS

The ETS model is a part of Exponential Smoothing methods, which assign exponentially decreasing weights to past observations. The ETS framework is particularly good for modeling time series that exhibit trends and seasonality. It automatically identifies the nature of the trend and seasonality components to produce a forecast.

```
1 # fit an ETS model and plot residual errors
2 from statsmodels.tsa.holtwinters import ExponentialSmoothing
3
4 # fit model
5 model_es = ExponentialSmoothing(daily_sales_final['Sales'],
6                                 trend='add',
7                                 seasonal='add',
8                                 seasonal_periods=365)
9 model_fit_es = model_es.fit()
10 # summary of fit model
11 print(model_fit_es.summary())
```

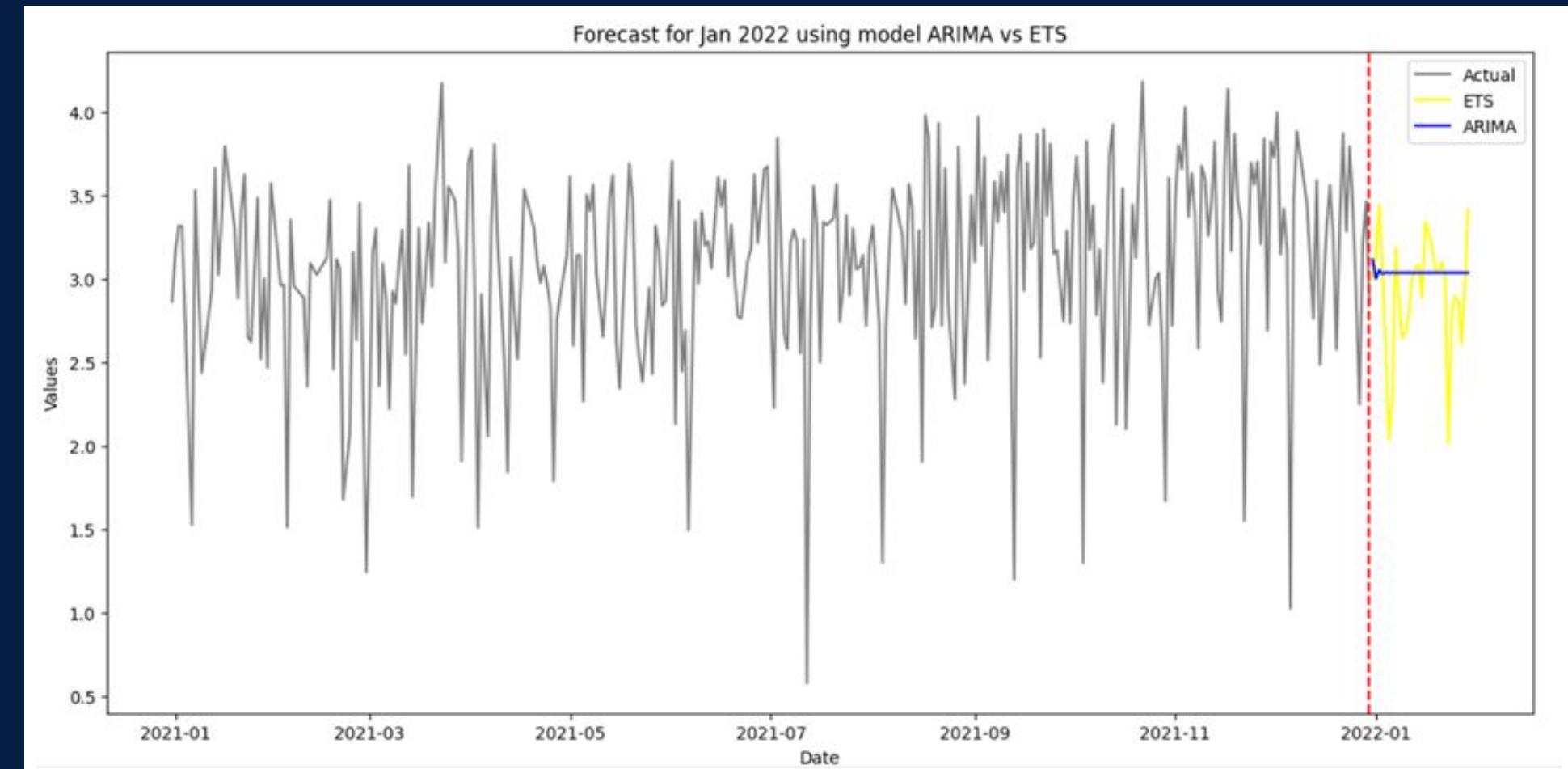
```
1 # split into train and test sets
2 X = daily_sales_final['Sales']
3 size = int(len(X) * 0.8)
4 train, test = X[0:size], X[size:len(X)]
5 history = [x for x in train]
6
7 # walk-forward validation
8 prediction_es = test.copy()
9 for t in range(len(test)):
10     model = ExponentialSmoothing(history)
11     model_fit = model.fit()
12     output = model_fit.forecast()
13     yhat = output[0]
14     prediction_es.iloc[t] = yhat
15     obs = test.iloc[t]
16     history.append(obs)
17     print('predicted=%f, expected=%f' % (yhat, obs))
```

MODEL BUILDING EVALUATION



Model	RMSE
0 ARIMA	0.682820
1 ETS	0.569452

MODEL BUILDING EVALUATION

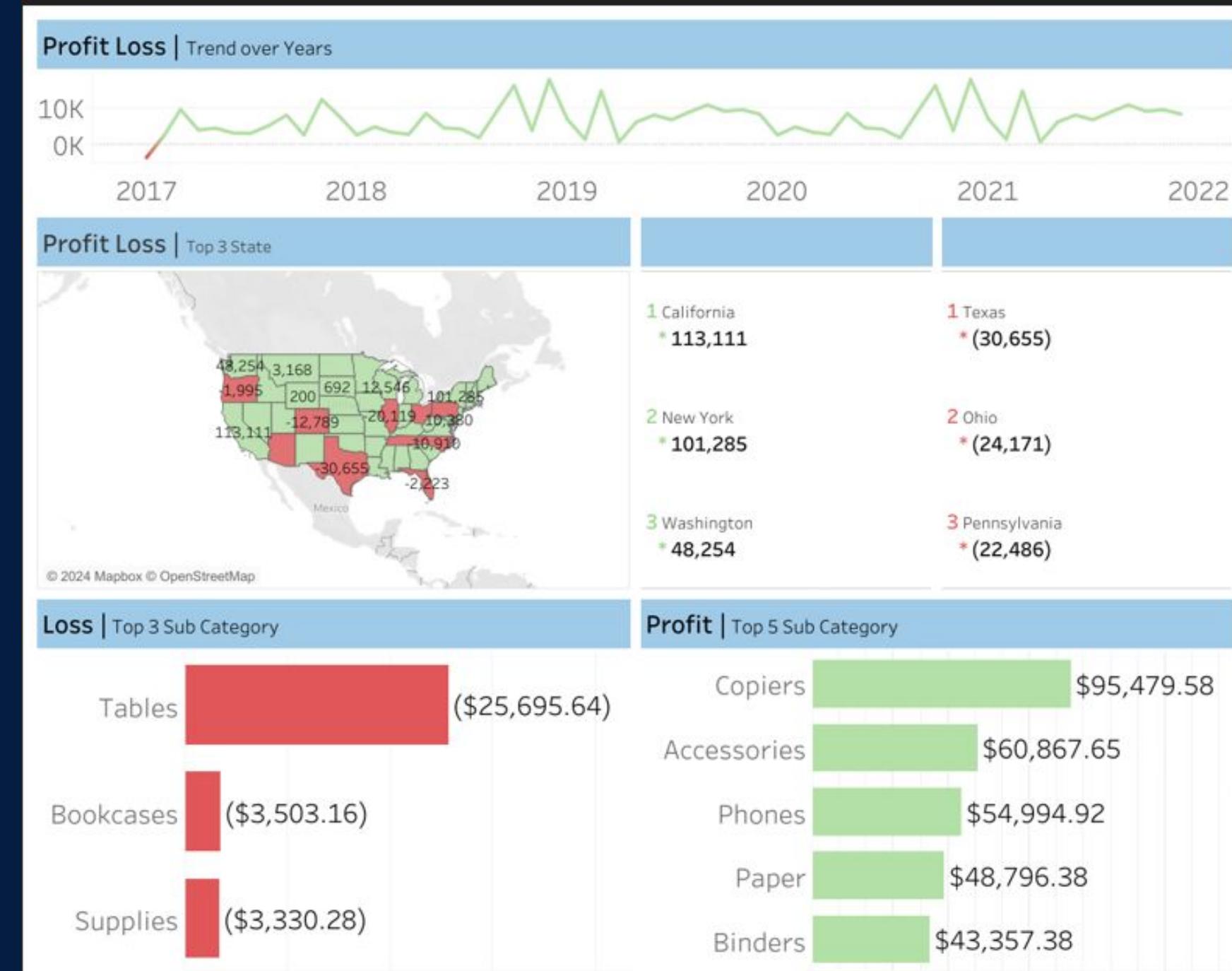


However, forecasting accuracy is critical. We can see that ARIMA forecast getting flat after a few days while ETS showing consistency through the end of month. Thus, ETS might be preferred over ARIMA based on above forecast for Jan 2022 and slightly lower RMSE value. We can conclude that the ETS model performs slightly better than the ARIMA model in terms of accuracy for this dataset

DASHBOARD



DASHBOARD





Thank You!

