

Website Traffic Time Series

Author: Aditya Srivastava
Date: Feb 25th 2022

Introduction

Often companies want to analyze the brand value they possess in the digital space or quantify their online user engagement. One of the ways is by monitoring the traffic that their website produces, tracking growth based on new users, analyzing trend of existing users for retention and churn rate, uncovering certain trend patterns and forecasting the traffic that is supposed to follow to be better equipped for the future and take business action with computational decisions accordingly. This project is based on analyzing the current trend and seasonalities as well as predict the future traffic landings on a website from new users (first time visits) to quantify growth using machine learning time series model.

Data Source

The pageviews data is that of the website <http://staffforecasting.com/> of almost 6 years of data from 2014 to 2020. The data is publicly available and can be downloaded from [Kaggle](https://www.kaggle.com/).

Data Dictionary

Column	Description	Data Type
Row	Row number	Integer
Day	Day of the month in text	String
DayOfWeek	Day of the week in numeric	Integer
Date	Date in mm/dd/yyyy format	String
PageLoads	Daily number of pages loaded	String
UniqueVisits	Daily number of visitors from whose IP addresses there haven't been hits on any page in over 6 hours	String
FirstTimeVisits	Number of unique visitors who do not have a cookie identifying them as a previous customer	String
ReturningVisits	Number of unique visitors minus first time visitors	String

Understanding Data Structure

```
In [1]: # Importing libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from statsmodels.formula.api import ols
from datetime import datetime
from statsmodels.tsa.stattools import adfuller
from statsmodels.tsa.stattools import kpss
from statsmodels.tsa.seasonal import seasonal_decompose
import statsmodels.api as sm
from statsmodels.tsa.arima.model import ARIMA
import pandas as pd
import itertools
from prophet import Prophet
from prophet.plot import add_seasonpoints_to_plot, plot_cross_validation_metric
from prophet.diagnostics import plot_validation_performance_metrics
import warnings
import logging

logging.disable(logging.INFO)
warnings.filterwarnings('ignore')
%matplotlib inline
```

```
In [63]: # Loading data
df = pd.read_csv('daily-website-visitors.csv')
df.head()
```

```
Out[63]:
```

Row	Day	DayOfWeek	Date	PageLoads	UniqueVisits	FirstTimeVisits	ReturningVisits
0	1	Sunday	1 9/14/2014	2146	1582	1430	152
1	2	Monday	2 9/15/2014	3621	2528	2297	231
2	3	Tuesday	3 9/16/2014	3698	2630	2352	278
3	4	Wednesday	4 9/17/2014	3667	2614	2327	287
4	5	Thursday	5 9/18/2014	3316	2366	2130	236

```
In [3]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2167 entries, 0 to 2166
Data columns (total 8 columns):
 #   Column              Non-Null Count  Dtype
---  ---
 0   Row                 2167 non-null  int64
 1   Day                 2167 non-null  object
 2   DayOfWeek           2167 non-null  int64
 3   Date                2167 non-null  object
 4   PageLoads           2167 non-null  object
 5   UniqueVisits        2167 non-null  object
 6   FirstTimeVisits     2167 non-null  object
 7   ReturningVisits     2167 non-null  object
dtypes: int64(2), object(6)
memory usage: 115.6+ KB
```

There are 2167 rows and 8 columns where each row represents the number of pageviews and visitors on a daily basis.

Based on the structure of the dataframe, the following code removes redundant and no-value columns, changes data types of page metrics to integers and standardizes the column names.

```
In [64]: # making structural changes for standardization
df.drop(['Row', 'DayOfWeek'], axis = 1, inplace = True)
df = df.replace('.', '', regex = True)
columns = df.columns
for column in columns[1:]:
    df[column] = df[column].astype('int64')
df['date'] = pd.to_datetime(df['Date'], format = '%m/%d/%Y')
df.rename(columns = {
    'Day': 'day',
    'Date': 'date',
    'PageLoads': 'n_page_loads',
    'UniqueVisits': 'n_unique_visits',
    'FirstTimeVisits': 'n_first_visits',
    'ReturningVisits': 'n_return_visits'
}, inplace = True)
```

```
In [51]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2167 entries, 0 to 2166
Data columns (total 6 columns):
 #   Column              Non-Null Count  Dtype
---  ---
 0   day                 2167 non-null  object
 1   date                2167 non-null  datetime64[ns]
 2   n_page_loads        2167 non-null  int64
 3   n_unique_visits     2167 non-null  int64
 4   n_first_visits      2167 non-null  int64
 5   n_return_visits     2167 non-null  int64
dtypes: datetime64(1), int64(4), object(1)
memory usage: 101.7+ KB
```

```
In [61]: df.head()
```

```
Out[61]:
```

	day	date	n_page_loads	n_unique_visits	n_first_visits	n_return_visits
0	Sunday	2014-09-14	2146	1582	1430	152
1	Monday	2014-09-15	3621	2528	2297	231
2	Tuesday	2014-09-16	3698	2630	2352	278
3	Wednesday	2014-09-17	3667	2614	2327	287
4	Thursday	2014-09-18	3316	2366	2130	236

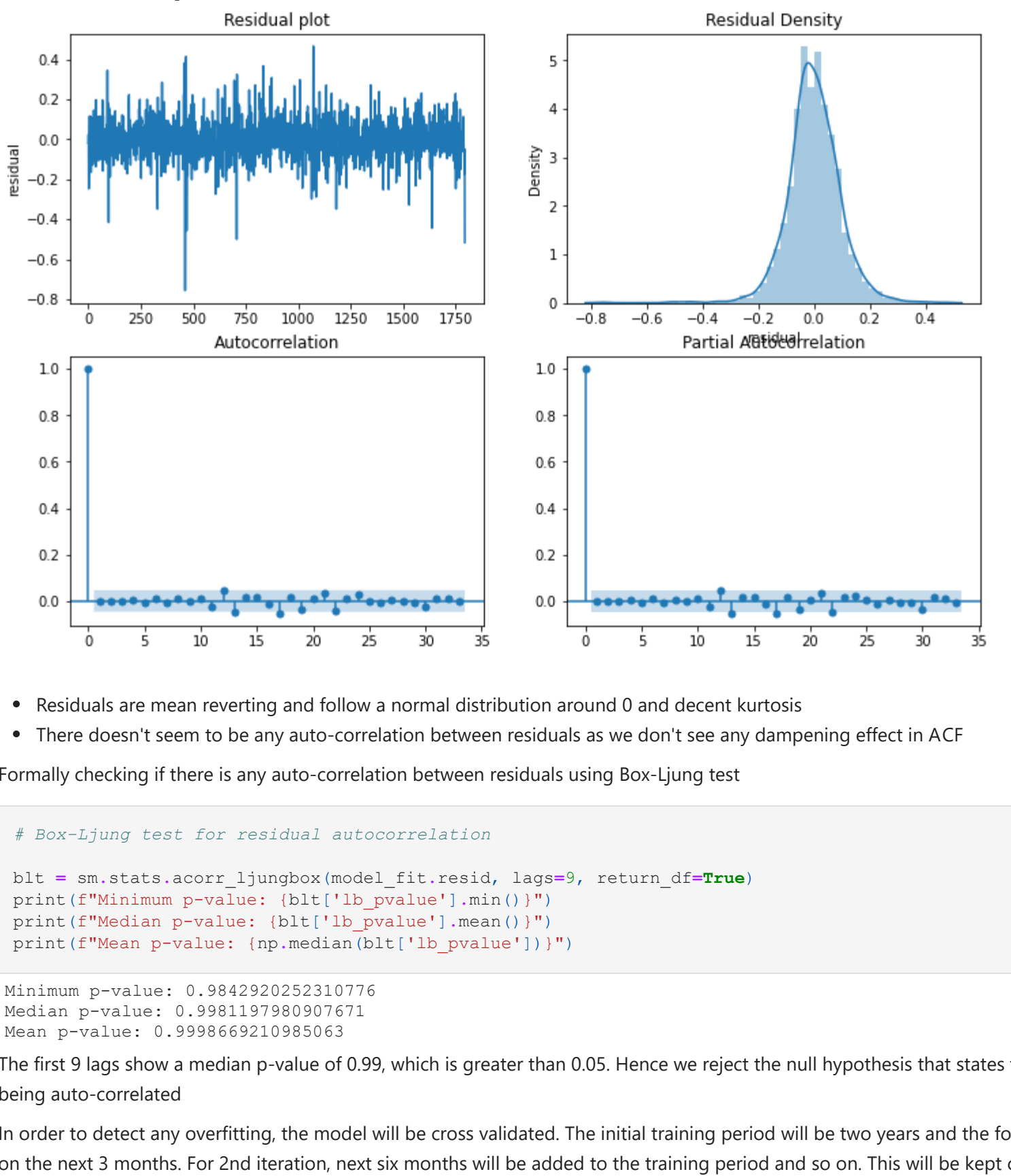
The data is cleaned and ready for further analysis.

```
In [71]: # checking for missing values
df.isna().sum()
```

```
Out[71]:
```

day	date	n_page_loads	n_unique_visits	n_first_visits	n_return_visits
0	0	0	0	0	0
1	1	0	0	0	0
2	2	0	0	0	0
3	3	0	0	0	0
4	4	0	0	0	0
5	5	0	0	0	0
6	6	0	0	0	0
7	7	0	0	0	0
8	8	0	0	0	0
9	9	0	0	0	0
10	10	0	0	0	0
11	11	0	0	0	0
12	12	0	0	0	0
13	13	0	0	0	0
14	14	0	0	0	0
15	15	0	0	0	0
16	16	0	0	0	0
17	17	0	0	0	0
18	18	0	0	0	0
19	19	0	0	0	0
20	20	0	0	0	0
21	21	0	0	0	0
22	22	0	0	0	0
23	23	0	0	0	0
24	24	0	0	0	0
25	25	0	0	0	0
26	26	0	0	0	0
27	27	0	0	0	0
28	28	0	0	0	0
29	29	0	0	0	0
30	30	0	0	0	0
31	31	0	0	0	0
32	32	0	0	0	0
33	33	0	0	0	0
34	34	0	0	0	0
35	35	0	0	0	0
36	36	0	0	0	0
37	37	0	0	0	0
38	38	0	0	0	0
39	39	0	0	0	0
40	40	0	0	0	0
41	41	0	0	0	0
42	42	0	0	0	0
43	43	0	0	0	0
44	44	0	0	0	0
45	45	0	0	0	0
46	46	0	0	0	0
47	47	0	0	0	0
48	48	0	0	0	0
49	49	0	0	0	0
50	50	0	0	0	0
51	51	0	0	0	0
52	52	0	0	0	0
53	53	0	0	0	0
54	54	0	0	0	0
55	55	0	0	0	0
56	56	0	0	0	0
57	57	0	0	0	0
58	58	0	0	0	0
59	59	0	0	0	0
60	60	0	0	0	0
61	61	0	0	0	0
62	62	0	0	0	0
63	63	0	0	0	0
64	64	0	0	0	0
65	65	0	0	0	0
66	66	0	0	0	0
67	67	0	0	0	0
68	68	0	0	0	0
69	69	0	0	0	0
70	70	0	0	0	0
71	71	0	0	0	0
72	72	0	0	0	0
73	73	0	0	0	0
74	74	0	0	0	0
75	75	0	0	0	0
76	76	0	0	0	0
77	77	0	0	0	0
78	78	0	0	0	0
79	79	0	0	0	0
80	80	0	0	0	0
81	81	0	0	0	0
82	82	0	0	0	0
83	83	0	0	0	0
84	84	0	0	0	0
85	85	0	0	0	0
86	86	0	0	0	0
87	87	0	0	0	0
88	88	0	0	0	0
89	89	0	0	0	0
90	90	0	0	0	0
91	91	0	0	0	0
92	92	0	0	0	0
93	93	0	0	0	0
94	94	0	0	0	0
95	95	0	0	0	0
96	96	0	0	0	0
97	97	0	0	0	0
98	98	0	0	0	0
99	99	0	0	0	0
100	100	0	0	0	0
101	101	0	0	0	0
102	102	0	0	0	0
103	103	0	0	0	0
104	104	0	0	0	0
105	105	0	0	0	0
106	106	0	0	0	0
107	107	0	0	0	0
108	108	0	0	0	0
109	109	0	0	0	0
110	110	0	0	0	0
111	111	0	0	0	0
112	112	0	0	0	0
113	113	0	0	0	0
114	114	0	0	0	0
115	115	0	0	0	0
116	116	0	0	0	0
117	117	0	0	0	0
118	118	0	0	0	0
119	119	0	0	0	0
120	120	0	0	0	0
121	121	0	0	0	0
122	122	0	0	0	0
123	123	0	0	0	0
124	124	0	0	0	0
125	125	0	0	0	0
126	126	0	0	0	0
127	127	0	0	0	0
128	128	0	0	0	0
129	129	0	0	0	0
130	130	0	0	0	0
131	131	0	0	0	0
132	132	0	0	0	0
133	133	0	0	0	0
134	134	0	0	0	0
135	135	0	0	0	0
136	136	0	0	0	0
137	137	0	0	0	0
138	138	0	0	0	0
139	139	0	0	0	0
140	140	0	0	0	0
141	141	0	0	0	0
142	142	0	0	0	0
143	143	0	0	0	0
144	144	0	0	0	0
145	145	0	0	0	0
146	146	0	0	0	0
147	147	0	0	0	0
148	148	0	0	0	0
149	149	0	0	0	0
150	150	0	0	0	0
151	151	0	0	0	0
152	152	0	0	0	0
153	153	0	0	0	0
154	154	0	0	0	0
155	155	0	0	0	0
156	156	0	0	0	0
157	157	0	0	0	0
158	158	0	0	0	0
159	159	0	0	0	0
160	160	0	0	0	0
161	161	0	0	0	0
162	162	0	0	0	0
163	163	0	0	0	0
164	164	0	0	0	0
165	165	0	0	0	0
166	166	0	0	0	0
167	167	0	0	0	0
168	168	0	0	0	0
169	169	0	0	0	0
170	170	0	0	0	0
171	171	0	0	0	0
172	172	0	0	0	0
173	173	0	0	0	0
174	174	0	0	0	0
175	175	0	0	0	0
176	176	0	0	0	0
177	177	0	0	0	0
178	178	0	0	0	0
179	179	0	0	0	0
180	180	0	0	0	0
181	181	0	0	0	0
182	182	0	0	0	0
183	183	0	0	0	0
184	184	0	0	0	0
185	185	0	0	0	0
186	186	0	0	0	0
187	187	0	0	0	0
188	188	0	0	0	0
189	189	0	0	0	0
190	190	0	0	0	0
191	191	0	0	0	0
192	192	0	0	0	0
193	193	0	0	0	0
194	194	0	0	0	0
195	195	0	0	0	0
196	196	0	0	0	0
197	197	0	0	0	0
198	198	0	0	0	0
199	199	0	0	0	0
200	200	0	0	0	0
201	201	0	0	0	0
202	202	0	0	0	0
203	203	0	0	0	0
204	204	0	0	0	0
205	205	0	0	0	0
206	206	0	0	0	0
207	207	0	0	0	0
208	208	0	0	0	0
209	209	0	0	0	0
210	210	0	0	0	0
211	211	0	0	0	0
212	212	0	0	0	0
213	213	0	0	0	0
214	214	0	0	0	0
215	215	0	0	0	0
216	216	0	0	0	0
217	217	0	0	0	0
218	218	0	0	0	0
219	219	0	0	0	0
220	220	0	0	0	0
221	221	0	0	0	0
222	222	0	0	0	0
223	223				

RMSE on trained period is: 516.93



- Residuals are mean reverting and follow a normal distribution around zero and decent kurtosis
- There doesn't seem to be any auto-correlation between residuals as we don't see any dampening effect in ACF

Formally checking if there is any auto-correlation between residuals using Box-Ljung test

```
In [76]: # Box-Ljung test for residual autocorrelation

blt = sm.stats.acorr_ljungbox(model_fit.resid, lags=9, return_df=True)
print("Minimum p-value: %f" % blt["lb_pvalue"].min())
print("Median p-value: %f" % blt["lb_pvalue"].median())
print("Mean p-value: %f" % blt["lb_pvalue"].mean())

Minimum p-value: 0.9842920252310776
Median p-value: 0.998197980907671
Mean p-value: 0.999666210985063

The first 9 lags show a median p-value of 0.99, which is greater than 0.05. Hence we reject the null hypothesis that states the residuals being auto-correlated
```

In order to detect any overfitting, the model will be cross validated. The initial training period will be two years and the forecast will happen on the next 3 months. For 2nd iteration, next six months will be added to the training period and so on. This will be kept consistent when modeling using Prophet

```
In [30]: # function to reverse the seasonality difference

def inverse_difference(history, yhat, interval1, interval2):
    return yhat + history[-interval1] + history[-interval2] - history[-(interval1 + interval2)]
```

```
In [70]: # function to forecast

def forecast(n_periods, training_period, interval1, interval2):
    fc, ser, conf = model_fit.forecast(n_periods, alpha=0.05)
    history = [x for x in df_copy["n_first_visits"] if x < training_period]
    for feat in fc:
        inverted = inverse_difference(history, yhat, interval1, interval2)
        history.append(inverted)
    fc_index = np.arange(training_period, training_period + n_periods)
    fc_series = pd.Series(history[training_period: training_period + n_periods], index=fc_index)
    fc_series = np.exp(fc_series) * 1
    return fc_series
```

```
In [71]: # cross-validation

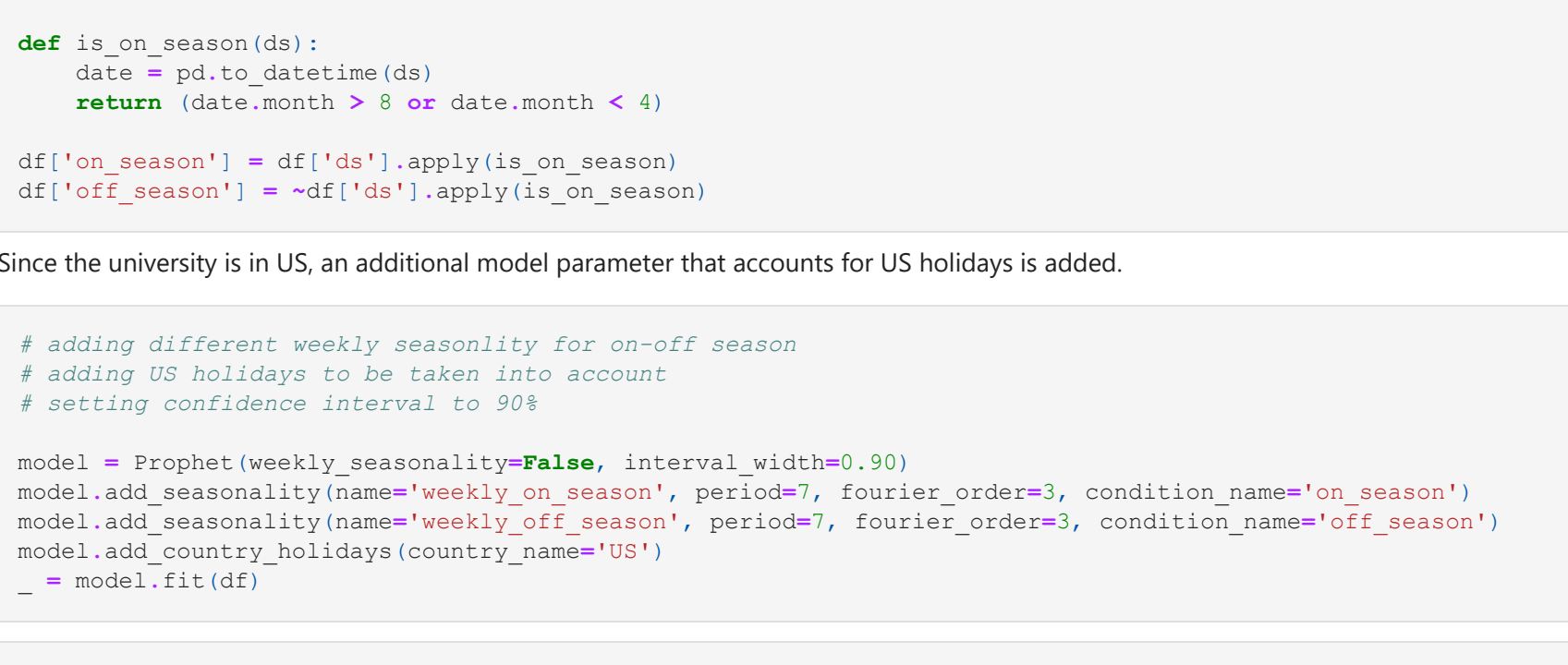
training_period, next_period, n_periods = 730, 180, 90
rmse = []
while training_period + next_period + n_periods <= len(df["n_first_visits"]):
    y_pred = forecast(n_periods, training_period, 7, 365)
    y_train = df["n_first_visits"][training_period: training_period + next_period]
    rmse.append(np.sqrt(np.mean((y_pred - y_train) ** 2)))
    training_period += next_period
    next_period += 60
print("Cross-validation RMSE: %f" % np.mean(rmse) * 2.2)
```

Cross-validation RMSE: 745.73

The difference between trained and cross-validated RMSE is ~230. We will compare this difference with the difference we see in Prophet. Proceeding with the forecast from the model for the next 6 months

```
In [74]: # calculating forecast for the next 6 months

fc_series = forecast(180, len(df_copy["n_first_visits"]), 7, 365)
plt.figure(figsize=(12,5), dpi=100)
plt.plot(df["n_first_visits"], label="training")
plt.plot(fc_series, label="forecast")
plt.title("Forecast")
plt.legend(loc="upper left", fontsize=8)
```



The model does well to capture both a slightly upward trend and seasonality

Facebook Prophet

Prophet by nomenclature understands only 'ds' for date and 'y' for the value to be forecasted. So, the features are renamed.

```
In [77]: # renaming features

df.rename(columns = {
    'date': 'ds',
    'n_first_visits': 'y',
    'y': 'yhat',
    'yhat_lower': 'yhat_lower',
    'yhat_upper': 'yhat_upper'
})
```

From EDA, it was seen that there are seasons where the first time visits are more hence the weekly seasonalities in terms of magnitude might differ as well. In order to capture this, a binary variable is introduced to take care of this

```
In [78]: # defining on-off season

def is_on_season(ds):
    date = pd.to_datetime(ds)
    return (date.month > 5 && date.month < 9)

df["on_season"] = df["ds"].apply(is_on_season)
df["off_season"] = ~df["on_season"]
```

Since the university is in US, an additional model parameter that accounts for US holidays is added.

```
In [79]: # adding different weekly seasonality for on-off season
# adding US holidays to be taken into account
# setting confidence interval to 80%

model = Prophet(weekly_seasonality=False, interval_width=0.90)
model.add_seasonality(name='weekly_on_season', period=7, fourier_order=3, condition_name='on_season')
model.add_seasonality(name='weekly_off_season', period=7, fourier_order=3, condition_name='off_season')
model.add_country_holidays(country_name='US')
model_fit = model.fit(df)
```

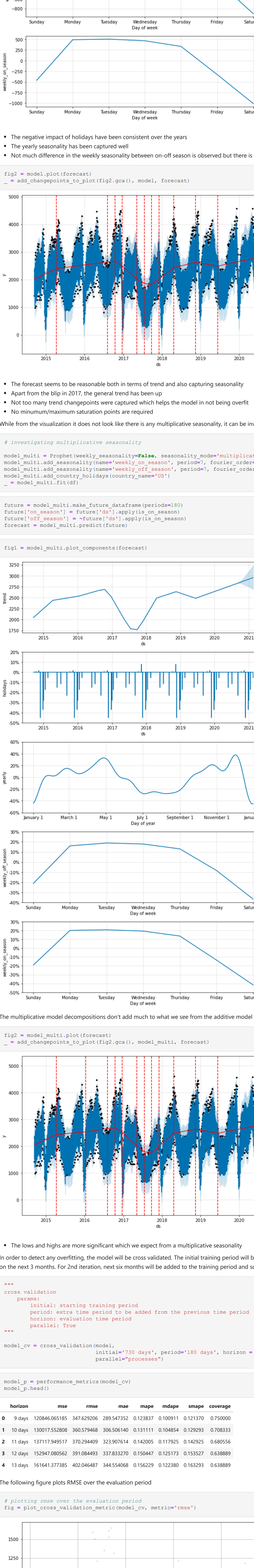
```
In [80]: # predicting for the next 6 months

future = model.make_future_dataframe(periods=180)
future["on_season"] = future["ds"].apply(is_on_season)
future["off_season"] = ~future["on_season"]
forecast = model.predict(future)
forecast["ds", "yhat", "yhat_lower", "yhat_upper"] = forecast[["ds", "yhat", "yhat_lower", "yhat_upper"]].tail()
```

```
Out[80]:
```

ds	yhat	yhat_lower	yhat_upper
2342 2021-02-11	3472.601809	2981.178092	3951.027436
2343 2021-02-12	2833.304079	2338.798560	3336.324910
2344 2021-02-13	2174.100416	1679.242814	2666.299509
2345 2021-02-14	2741.367373	2250.014069	3234.120809
2346 2021-02-15	3584.781715	3085.940610	4089.886075

```
In [81]: fig1 = model_plot_components(forecast)
```



- The negative impact of holidays have been consistent over the years
- The yearly seasonality has been captured well
- Not much difference in the weekly seasonality between on-off season is observed but there is a slight difference in magnitude

```
In [82]: fig2 = model_multi_plot(forecast)
# add changepoints to plot (fig2.gca(), model_multi, forecast)
```



- The forecast seems to be reasonable both in terms of trend and also capturing seasonality
- Apart from the bump in 2017, the general trend has been up
- Not too many trend changepoints were captured which helps the model in not being overfit
- No minimum/maximum saturation points are required

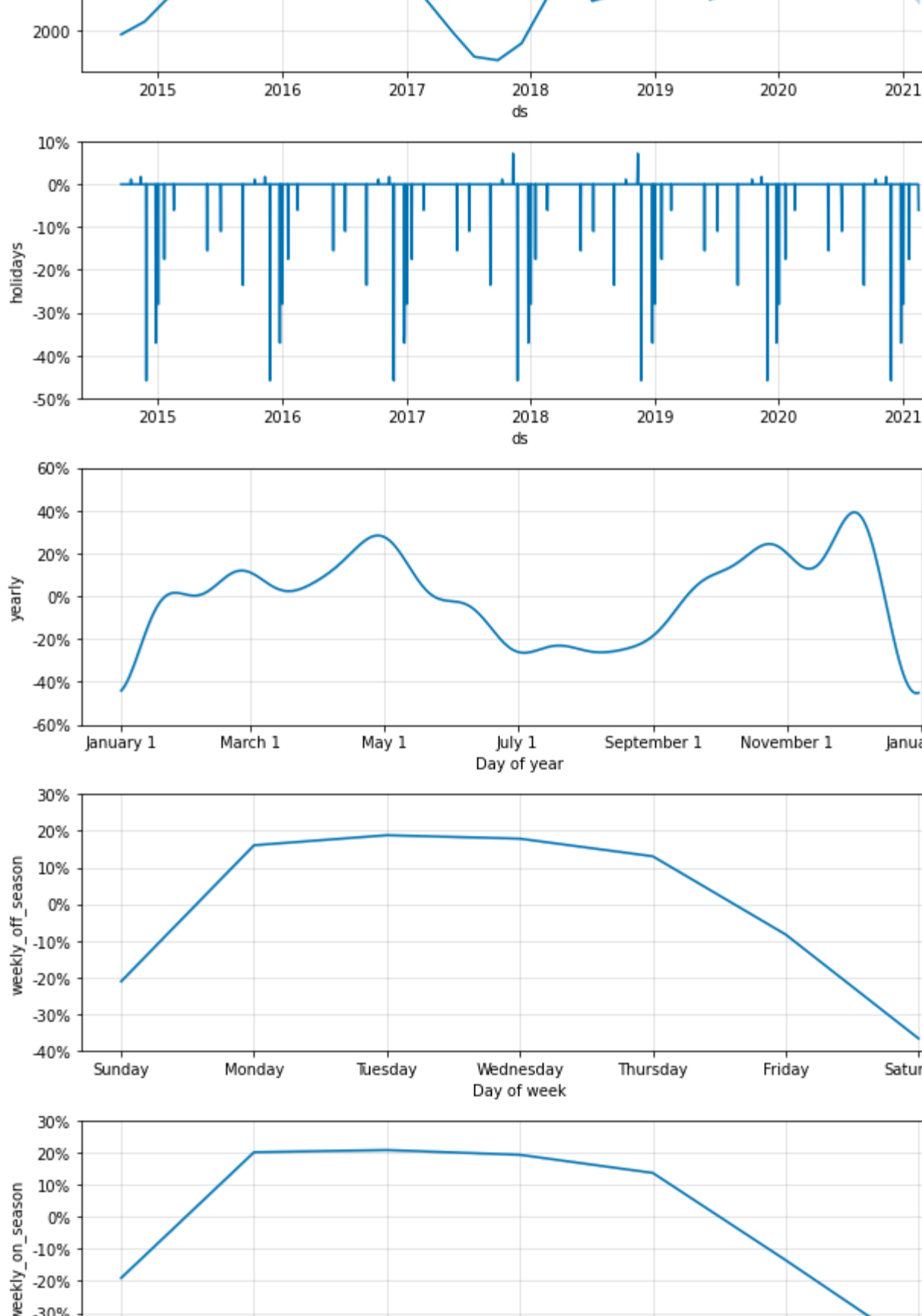
While from the visualization it does not look like there is any multiplicative seasonality, it can be investigated as well.

```
In [83]: # investigating multiplicative seasonality

model_multi = Prophet(weekly_seasonality=False, seasonality_mode='multiplicative', interval_width=0.90)
model_multi.add_seasonality(name='weekly_on_season', period=7, fourier_order=3, condition_name='on_season')
model_multi.add_seasonality(name='weekly_off_season', period=7, fourier_order=3, condition_name='off_season')
model_multi.add_country_holidays(country_name='US')
model_fit_multi = model_multi.fit(df)
```

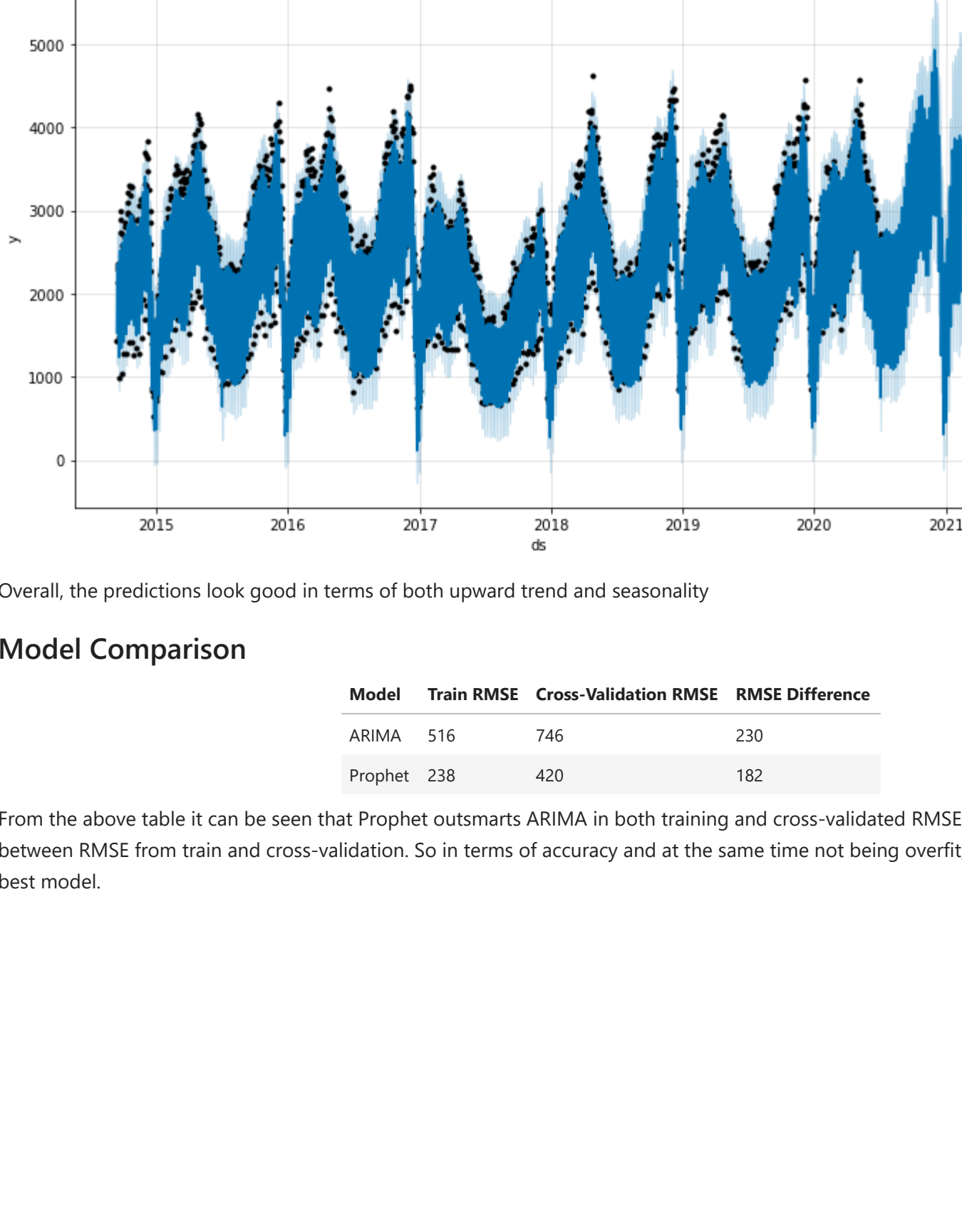
```
In [84]: future = model_multi.make_future_dataframe(periods=180)
future["on_season"] = future["ds"].apply(is_on_season)
future["off_season"] = ~future["on_season"]
forecast = model_fit_multi.predict(future)
```

```
In [85]: fig1 = model_multi_plot_components(forecast)
```



The multiplicative model decompositions don't add much to what we see from the additive model

```
In [86]: fig2 = model_multi_plot(forecast)
# add changepoints to plot (fig2.gca(), model_multi, forecast)
```



The lows and highs are more significant which we expect from a multiplicative seasonality

In order to detect any overfitting, the model will be cross validated. The initial training period will be two years and the forecast will happen on the next 3 months. For 2nd iteration, next six months will be added to the training period and so on.

```
In [44]: # cross validation

params = {
    'initial': starting training period
    'period': extra time period to be added from the previous time period
    'horizon': evaluation time period
    'parallel': True
}

model_cv = cross_validation(model,
    initial='730 days', period='180 days', horizon = '90 days',
    parallel="processes")
```

```
In [45]: model_p = performance_metrics(model_cv)
model_p.head()
```

horizon	mse	rmse	mase	mape	mdape	smape	coverage
0 9 days	120846.065185	347.629206	289.547352	0.128837	0.100911	0.121370	0.750000
1 10 days	130117.552808	360.579468	306.506140	0.131111	0.104854	0.129293	0.708333
2 11 days	137117.949517	370.294409	323.907614	0.142005	0.117925	0.142925	0.680556
3 12 days	152947.080562	391.084493	337.833270	0.150447	0.125173	0.153527	0.638889
4 13 days	161641.373785	402.046487	344.554068	0.156229	0.122380	0.163293	0.638889

The following figure plots RMSE over the evaluation period

```
In [46]: # plotting rmse over the evaluation period
fig = plot_cross_validation_metric(model_cv, metric='rmse')
```



The trend of rolling average of RMSE increases with time which is what is expected as the error term is additive. Given that the RMSE is in the range of 300 - 600 it's doing a decent job.

For model to perform best, the following section does hyperparameter tuning on the complexity of trend and seasonality along with seasonality type with RMSE being the cross-validated evaluation metric.

```
In [47]: # hyperparameter tuning

params = {
    'changepoint_prior_scale': the complexity of trend
    'seasonality_prior_scale': the complexity of seasonality
    'seasonality_mode': additive or multiplicative
}

param_grid = [
    ('changepoint_prior_scale': [0.01, 0.05, 0.1, 0.5],
    'seasonality_prior_scale': [1.0, 5.0, 10.0, 15.0],
    'seasonality_mode': ['additive', 'multiplicative'])
]

all_params = [dict(zip(param_grid.keys(), v)) for v in itertools.product(*param_grid.values())]
rmse = []

for params in all_params:
    model = Prophet(weekly_seasonality=False,
    changepoint_prior_scale=5,
    seasonality_prior_scale=10,
    seasonality_mode='multiplicative',
    interval_width=0.90)
    model.add_seasonality(name='weekly_on_season', period=7, fourier_order=3, condition_name='on_season')
    model.add_seasonality(name='weekly_off_season', period=7, fourier_order=3, condition_name='off_season')
    model.add_country_holidays(country_name='US')
    model_fit = model.fit(df)
    model_cv = cross_validation(model,
    initial='730 days', period='180 days', horizon = '90 days',
    parallel="processes")
    model_p = performance_metrics(model_cv, rolling_window=1)
    rmse.append(model_p['rmse'].values[0])

tuning_results = pd.DataFrame(all_params)
tuning_results['rmse'] = rmse
```

```
In [48]: best_params = all_params[np.argmin(rmse)]
print(best_params)
print("Best RMSE: %f" % (np.min(rmse)))

('changepoint_prior_scale': 0.5, 'seasonality_prior_scale': 10.0, 'seasonality_mode': 'multiplicative')
```

Hyperparameter tuning suggests complexity of trend to be increased and seasonality to be of multiplicative type.

The parameters that gave the best RMSE will be used in the final model for forecasting.

```
In [54]: # forecasting with the best model

model_best = Prophet(weekly_seasonality=False,
    changepoint_prior_scale=5,
    seasonality_prior_scale=10,
    seasonality_mode='multiplicative',
    interval_width=0.90)
model_best.add_seasonality(name='weekly_on_season', period=7, fourier_order=3, condition_name='on_season')
model_best.add_seasonality(name='weekly_off_season', period=7, fourier_order=3, condition_name='off_season')
model_best.add_country_holidays(country_name='US')
model_fit_best = model_best.fit(df)
```

```
In [59]: future = model_best.make_future_dataframe(periods=180)
future["on_season"] = future["ds"].apply(is_on_season)
future["off_season"] = ~future["on_season"]
forecast = model_fit_best.predict(future)
```

```
In [60]: # in-sample RMSE

y_train = df["y"].to_numpy()
y_pred_train = forecast[["yhat"]].to_numpy()[len(y_train):]
print("RMSE on trained period is: %f" % np.sqrt(np.mean((y_pred_train.squeeze() - y_train) ** 2)))

RMSE on trained period is: 238.08
```

```
In [61]: fig1 = model_best_plot_components(forecast)
```



Overall, the predictions look good in terms of both upward trend and seasonality

Model Comparison

Model	Train RMSE	Cross-Validation RMSE	RMSE Difference
ARIMA	516	746	230
Prophet	238	420	182

From the above table it can be seen that Prophet outperforms ARIMA in both training and cross-validated RMSE and also in the difference between RMSE from train and cross-validation. So in terms of accuracy and at the same time not being overfit, Prophet comes out as the best model.