

Machine Learning Final Assignment

Constantin Pusch - 20304226

Contents

Part One	2
Introduction	2
Data Preprocessing	2
Model Used	2
Prophet Technical Aspects	3
Results	4
Impact on Usage during the pandemic period	5
Impact on Usage after the Pandemic Period	6
Model Performance	7
Part Two	8
Question I	8
Question II	8
Question III	8
Question IV	8
Appendix	9
Cross Validation First Model	9
Cross Validation Second Model	10
Code Used	11

Part One

Introduction

The goal of this assignment is to assess the impact of the covid pandemic on Dublin's city-bike scheme. The first major decision was to define the date range of the pandemic. The dates I chose for the pandemic were 29-02-2020 as the start of the pandemic and 21-01-2022 as the end. The first date correlates to the first covid case that was detected in Ireland and the latter is when the government removed all restrictions and no new ones implemented at a later date. One thing to note with these dates is that the Dublin bike dataset only has data until 01/01/2022 so there is no post pandemic data and, as discussed later, it is difficult to forecast the post pandemic impact.

Data Preprocessing

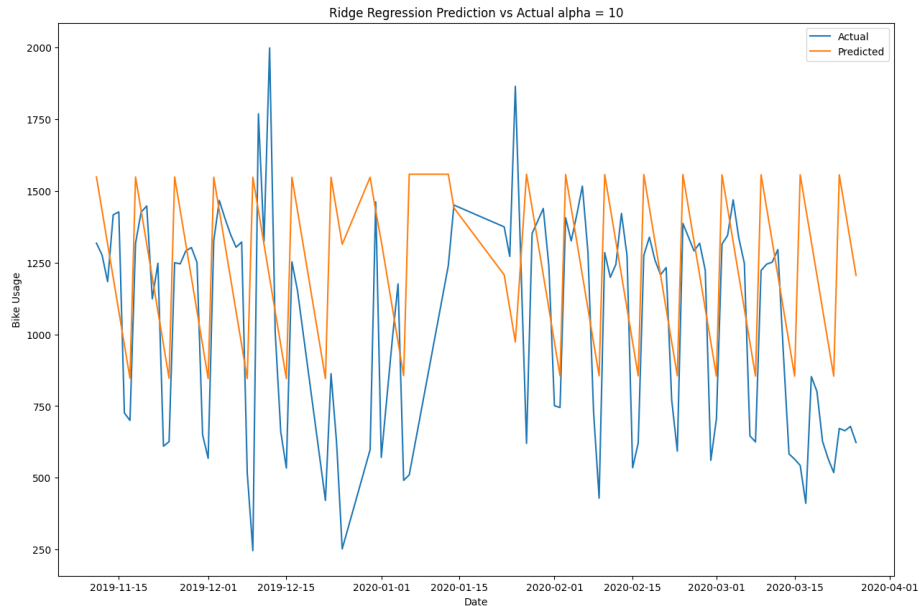
As mentioned in the problem statement, it is crucial to determine a method of measuring bike usage. To measure usage, I decided that each time a bike is removed from a station it counts as one use. This figure is calculated by taking the inverse difference of the Available Bikes per day. Since the bike has to be returned eventually, at an unspecified time, I did not count a return as a use.¹ I decided to keep this a simple measurement since there are multiple unknown variables that could affect bike usage. To prove this point I will give two examples. When I was processing the data, I noticed that the total number of bike stands would fluctuate. This could have an impact on the total number of bikes in the system, thus potentially increasing usage, or have no impact at all. Another point is that bikes are removed from the system for maintenance. It is unclear whether this is recorded in the data or not. Overall, my point is that there are many unknowns and therefore the bike usage figure is simple and will give a general overview not an exact usage figure.

To use the data in the model I combined all the datasets into one and aggregated all the five-minute timestamps of recorded data for all stands into a daily usage figure. Once I compiled this combined data set, I calculated the usage figure for each day as described above. I believe that a daily usage figure is enough granularity for making a long-term forecast.

Model Used

Originally, I intended to use a Ridge regression to forecast the effect that the pandemic had on the bike scheme. After optimizing the alpha value, I was getting a model with not ideal results that did not seem to accurately predict far enough out into the future. I had added additional parameters to attempt to capture the impact of seasonality, such as day of the week and month of the year. These hyperparameters increased the accuracy of the model but I still felt that it was overgeneralizing too much. In order to overcome these shortcomings, I decided to do some research to find a better model to use.

¹ As a quick note: throughout the report I will interchangeably use 1 'usage' as a rider even though these do not exactly correlate.



This graph shows the ridge regression model predicting using the training data and not even forecasting. This gave me little hope for accurate forecasting and caused me to search for a better model to use.

The machine learning model I chose was developed by Facebook and is called Prophet.² Prophet is specifically designed for time series forecasting which gives it multiple advantages that allow for more accurate long-term predictions. The main advantages that made sense for this project were the following: Prophet accelerates at incorporating seasonality into its forecasts and is designed to work with daily data. This is perfect for the Dublin bikes data set. Prophet is also good at handling longer time horizons which is ideal when forecasting out bike usage for this scenario.

Prophet Technical Aspects

Prophet is based on an additive model where linear/non-linear trends are fit with yearly, weekly, and daily seasonality. The model can automatically detect points in the time-series data where trends change and then it can select the appropriate linear or logistic growth model.

Predictions are generated by splitting the time series into three main components: trend, seasonality, and holidays. For this scenario I have not had the model account for holidays. The trend aspect models non-periodic changes in the data using a piecewise linear or logistic growth curve. This allows for the model to predict different types of growth trends in the data. As mentioned in the lectures, the seasonality captures regular patterns that appear in the data and the model can account for daily, weekly, and yearly seasonality. The daily hyperparameter is used for intraday data so therefore I only selected yearly and weekly hyperparameters. Prophet uses an additive model where these separate components are combined to make the final forecast.

² <https://facebook.github.io/prophet/>

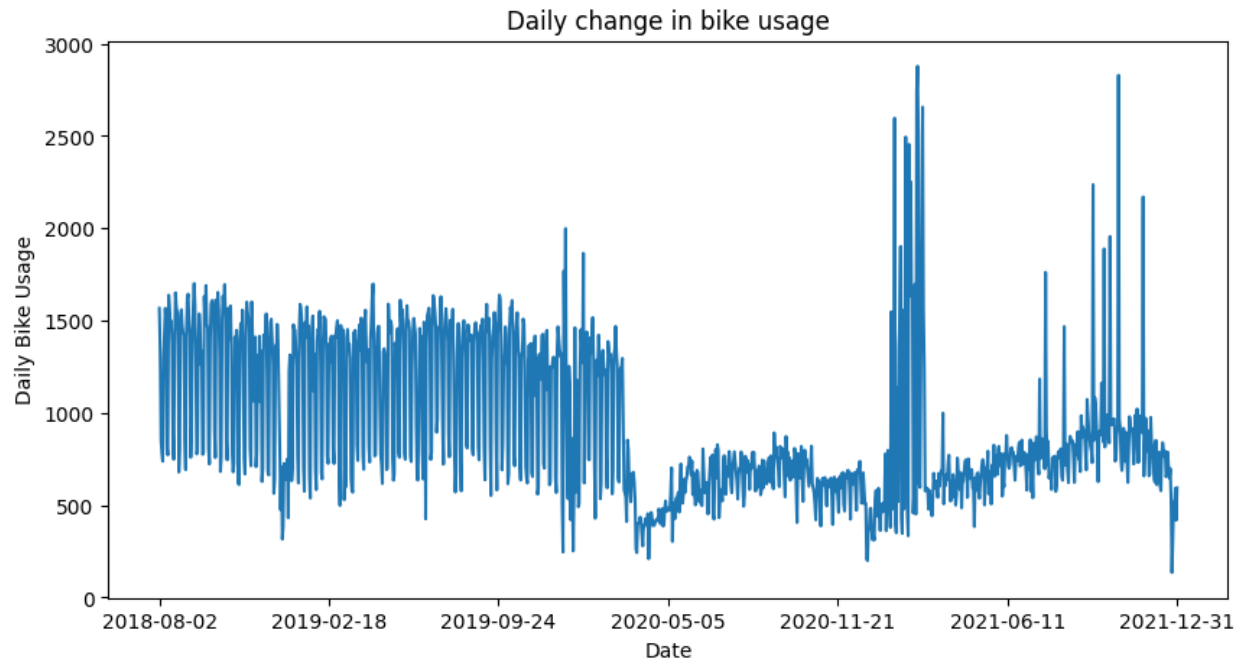
When creating the prophet model the cost function is automatically selected as it performs this step automatically. It typically varies between evaluating the results using mean absolute error or mean squared error.³

Results

Below is the raw data that was used to perform the predictions. As mentioned above the dates that I have defined as 'the pandemic' are from 29-02-2020 to 21-01-2022. Due to the few extreme outliers in usage (more explained in the chart caption) I have calculated the descriptive statistics using the median and not mean. Using the mean in this case gives a much more accurate representation of actual usage. The date range of recorded data is 02-08-2018 to 12-31-2021. The first day and last day of the actual dataset are omitted due to there not being a full day of recording on these days.

These are the descriptive statistics results:

Median over the whole period: 766.0, Pre Covid-Median: 1336.0, During Covid-Median: 675.0



This plot shows the original data that I have gotten using the bike usage statistic mentioned from above. There are some weird outliers during the covid period but since I could not pinpoint down why they were so high I decided to keep them in the data set. This definitely had an impact on predictions, but it was not trend breaking so I kept these points in since they well could be real usage. I do know that over this time period there were some updates being made to the system so this could explain these massive jumps.⁴

The graph and statistics clearly show that as soon as lockdown initiated there was a severe reduction in the usage of Dublin city bikes that had not recovered to pre-pandemic levels before the end of the

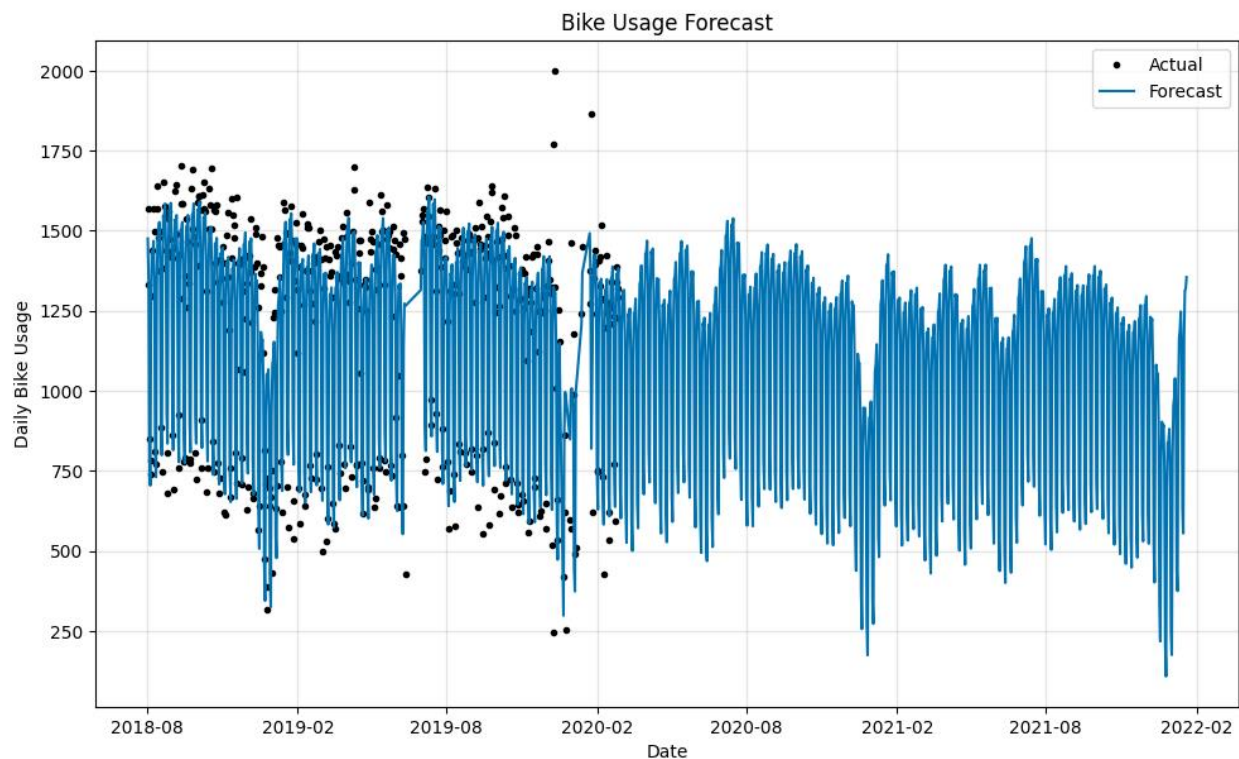
³ Taylor SJ, Letham B. 2017. Forecasting at scale. PeerJ Preprints 5:e3190v2
<https://doi.org/10.7287/peerj.preprints.3190v2>

⁴ <https://www.facebook.com/DublinCityCouncil/posts/please-note-that-due-to-a-system-upgrade-dublinbikes-will-not-be-available-from-/10159750050654625/>

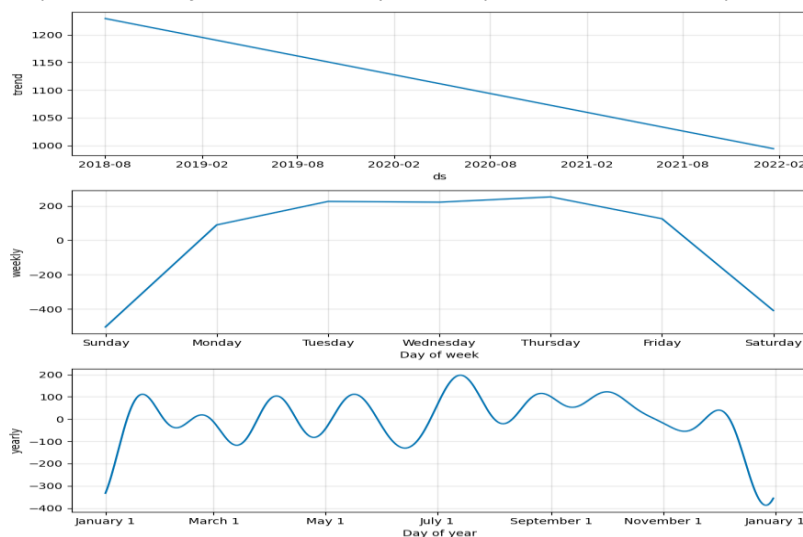
pandemic. The results show that the median daily bike usage dropped by 670 riders from before the pandemic to during. *This is a roughly 50% drop!*

Impact on Usage during the pandemic period

In order to investigate the impact of the pandemic prophet was trained on all datapoints up until the first recorded case of covid which was 29-02-2020. As mentioned, this is the day that the first covid case was detected and before then I found no real news that would correlate to a reduction in city bike usage as everyone was still working and traveling. Using Prophet, I forecasted 692 days until the day that all restrictions were lifted. The first graph represents the forecast using the prophet model and the following three graphs show the overall trend of the pre pandemic data.



This graph represents the forecast made by the Prophet model over the pandemic period.



As can be seen in the three graphs there was clearly a reducing trend of bike usage already before the pandemic. The weekly and yearly seasonality can be seen in the second and third graphs respectively. Clearly most users are using the bikes to commute as can be seen in the weekly seasonality as well as the strong cut off in the Christmas period. The Christmas cut off could be attributed to weather, but this dip really only occurs in December where temperatures are similar to November and January/February.

Based off the model the median daily usage is as following:

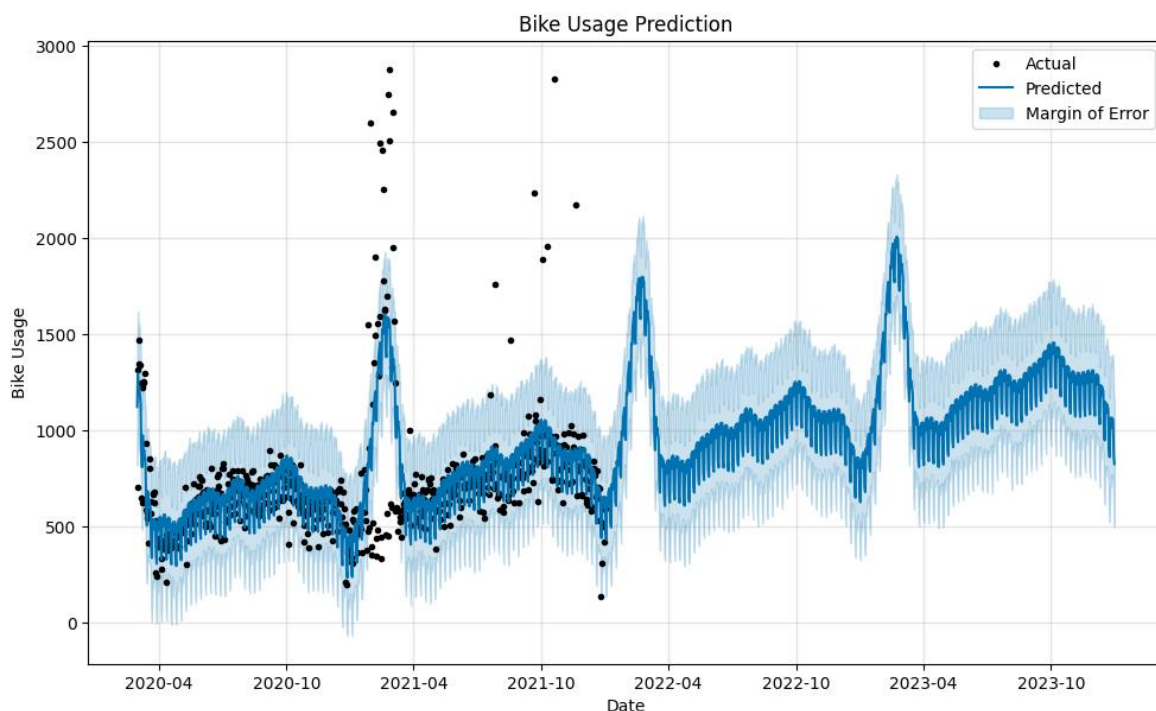
Forecasted median during covid: 1190.54, Pre covid forecast median: 1321.47, Median over the whole forecast period: 1236.84

As mentioned in the caption of the graph. There had been a reducing trend in overall bike usage even before the pandemic. This trend is exemplified in the forecast. Using the pre covid median of the recorded data there is a median daily usage of 1336, but the overall forecasted median is 1237. This is a drop of around 100 daily users and represents a ridership drop of about 7.8%. This is a significant drop that Dublin city council should investigate. When using the mean, as shown in the graph, there is a 200 rider drop from August 2018 to February 2022

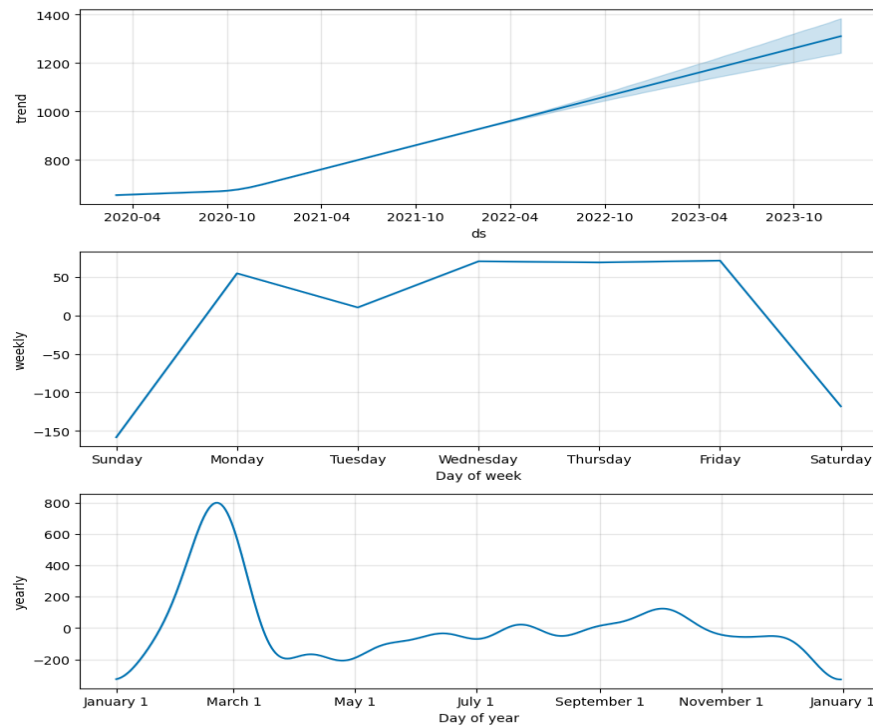
As for the pandemic, despite the predicted declining ridership numbers, it had a massive impact on ridership. Had the pandemic not happened Dublin bikes would have had an additional 515 additional daily riders.

Impact on Usage after the Pandemic Period

As in the previous section the dates of the pandemic remain the same: 29-02-2020 until 21-01-2022. Since there is no data that can be used to model past the end of the pandemic only the trends towards the end of the pandemic where restrictions were being eased can be observed. Using the usage data starting from 29-02-2020 until the end of the dataset, 31-12-2021, Prophet was used to forecast 730 days into the future until 01-01-2024. As before, the first graph represents the forecast until 01-01-2024 and the next three graphs represent the overall trend and weekly/yearly seasonality.



This graph represents the forecast of bike usage after the pandemic. As mentioned earlier there are some extreme outliers that clearly have an effect on the forecast, but I did not want to remove them.



As above these three graphs represent the trend of the forecast. Starting in October of 2020 there is a linear growth of usage. This couples with the fact that there were still people commuting during the pandemic which shows a positive post pandemic recovery for Dublin bike. I have included the yearly graph, but it is somewhat useless due to the skewing of the data around the January/February window.

Based on the forecast these are the median ridership numbers:

Forecast median: 895.78, post covid forecast median: 1090.11, October 2023 forecast median: 1333.91,

Obviously throughout the entire forecast period the median ridership will be below pre-pandemic levels. The reason I chose to show the median ridership for October 2023 is that it can be compared to the levels in October of 2019. The median ridership for October 2019 was 1318. This shows that over the forecasted month of October 2023 there was basically equal ridership as pre pandemic levels. This gives a roughly two-year time frame for ridership to reach normal levels once again. The reason I chose October is that there is a slight peak in the recorded data around then and this is also shown in the forecast post covid.

Model Performance

Cross validation was used to evaluate the performance of the model. The initial training size for both models was 365 days, period was 30 days, and the horizon was 30 days as well. This means that cross-validation was performed on the Prophet model by training the model on the first 365 days of data, then making and evaluating predictions for the next 30 days and repeating this process every 30 days. Looking at the mean absolute error, both models had an error of around 110 which I would consider quite accurate and am happy with. The full table of the cross-validation results is in the appendix.

Part Two

Question I

An ROC (receiver operating characteristic) curve is a visual tool that can be used to evaluate the performance of binary classifiers. What it does is plot the true positive rate, or sensitivity, against the false positive rate and multiple thresholds that can be set. You can use the ROC curve to compare the curve of a baseline to one of the classifier you are testing and then can visually assess how much better the classifier is compared to baseline. The area under the ROC curve can provide a metric for summarizing the performance. AN ROC curve would be used over a accuracy metric because it is less affected by class imbalance and since one can see the classifiers performance across different threshold levels

Question II

As per its name a linear regression should be used in a situation where there is a linear correlation between the data. One example would be when attempting to predict the value of an exponentially growing population. In this case linear regression would fail. In this case using a logistic regression would make sense as it matches the data better. Another example where linear regression can fail is when there are extreme outliers. If one is, for example, attempting to predict the values of homes there might be some mansions or houses in a desirable location that have a higher price. These outliers could disproportionately influence the model and lead to inaccuracies. In this case one could remove the outliers to improve the accuracy of the model.

Question III

In SVMs a kernel is a function used to transform non-linearly separable data into higher dimensional space where it becomes linearly separable where it is then possible to classify them effectively. A few types of kernels for SVMs are linear, polynomial, and radial basis. For CNNs a kernel is a small matrix that is used to perform convolution operations on the input data like an image which is then used to extract features from the input. Different types of CNN kernels can extract different types of features from the input. SVM kernels are used for transforming and classifying data and CNN kernels for feature extraction usually in deep learning.

Question IV

The idea behind resampling the dataset multiple times is to be able to utilize the entire dataset for training and testing which allows for evaluating the model's general performance better. For example, If I had a dataset of 100 points and I use a 5-fold cross-validation I split this dataset into 5 parts. For each iteration of the cross-validation, one of these parts is used as the test data and the rest for training. These multiple cycles can then be used to derive the model's generalized accuracy. The cross-validation also can mitigate the risk of bias that can occur if the dataset were to be only split once. Cross-validation can be especially useful when one has a small dataset as one can use the dataset more efficiently then. Somewhere it would not be appropriate to use cross-validation would be in time series data. This is because the temporal order of the datapoints is important and randomly splitting up the dataset would disrupt the order.

Appendix

Cross Validation First Model

horizon	mse	rmse	mae	mape	mdape	smape	coverage
3 days	44841.01	211.757	132.5995	0.158792	0.074711	0.135259	0.849206
4 days	34875.54	186.7499	118.0059	0.153618	0.067082	0.128794	0.849206
5 days	58182.5	241.2105	144.9008	0.151196	0.060531	0.158661	0.796296
6 days	43923.25	209.5787	121.253	0.118355	0.051188	0.139644	0.857143
7 days	64985	254.9216	160.2823	0.170943	0.057458	0.176969	0.769841
8 days	46712.65	216.1311	142.1531	0.153473	0.079269	0.148118	0.833333
9 days	92972.39	304.9137	193.1506	0.362415	0.112661	0.213808	0.801587
10 days	71809.54	267.973	151.6371	0.299234	0.086416	0.160366	0.888889
11 days	69707.76	264.0223	143.1068	0.302993	0.068751	0.16144	0.888889
12 days	48777.57	220.8565	124.4771	0.09551	0.043677	0.101377	0.888889
13 days	38195.67	195.4371	115.6514	0.094635	0.075186	0.096996	0.944444
14 days	33083.33	181.8882	110.898	0.092517	0.072911	0.095776	0.953704
15 days	9355.891	96.72585	82.60332	0.082539	0.053273	0.083387	0.944444
16 days	10155.4	100.774	83.09267	0.081242	0.058966	0.084021	0.888889
17 days	10562.68	102.7749	83.36925	0.085823	0.059801	0.087371	0.896825
18 days	8692.716	93.23474	75.8216	0.078913	0.059894	0.079309	0.936508
19 days	8615.094	92.81753	79.9671	0.080254	0.0704	0.079587	0.981481
20 days	7911.352	88.94578	77.13832	0.070472	0.069876	0.070159	1
21 days	11583.45	107.6264	92.16983	0.083812	0.07831	0.080994	0.944444
22 days	12424.64	111.4659	92.85964	0.093065	0.069876	0.091769	0.944444
23 days	10581.56	102.8667	82.46753	0.088335	0.057887	0.086406	0.944444
24 days	14254.09	119.3905	80.29199	0.092752	0.050373	0.085392	0.944444
25 days	45351.7	212.9594	111.984	0.247417	0.049585	0.135695	0.888889
26 days	111476.8	333.8815	170.2392	0.270859	0.055253	0.175595	0.849206
27 days	131648.4	362.8338	199.9767	0.317738	0.068167	0.205704	0.833333
28 days	128112.6	357.9282	207.8277	0.235238	0.079738	0.191065	0.825397
29 days	66902.66	258.6555	151.2347	0.183021	0.067874	0.14157	0.87963
30 days	31576.31	177.6972	110.8556	0.131502	0.065732	0.106959	0.898148

Cross Validation Second Model

horizon	mse	rmse	mae	mape	mdape	smape	coverage
3 days	107323.2	327.6022	187.1131	0.250744	0.125872	0.208447	0.868966
4 days	75614.8	274.9815	153.1554	0.25355	0.116923	0.183692	0.9
5 days	75240.76	274.3005	159.4186	0.265149	0.115081	0.198846	0.9
6 days	79247.36	281.5091	165.7829	0.280746	0.110658	0.213919	0.9
7 days	60752.93	246.4811	148.6664	0.242801	0.085404	0.195282	0.9
8 days	110529	332.459	181.5503	0.24566	0.085404	0.218665	0.865517
9 days	103277.4	321.3679	177.2623	0.231752	0.102462	0.209383	0.862069
10 days	99497.28	315.4319	178.1417	0.227755	0.120624	0.211277	0.862069
11 days	36248.68	190.3909	132.0787	0.193304	0.112343	0.169402	0.896552
12 days	33932.79	184.2086	129.5213	0.183518	0.120101	0.173094	0.865517
13 days	35692.07	188.9235	130.2537	0.189082	0.112343	0.185502	0.865517
14 days	43342.91	208.1896	136.2209	0.18071	0.120101	0.187453	0.868966
15 days	40675.26	201.6811	145.7107	0.188102	0.166516	0.19716	0.9
16 days	35775.77	189.1448	145.3847	0.184657	0.15546	0.190041	0.934483
17 days	20988.02	144.8724	124.7357	0.170741	0.153162	0.172357	0.968966
18 days	279461.4	528.6411	251.4044	0.1904	0.111311	0.223512	0.896552
19 days	283040.9	532.0159	249.354	0.184718	0.076854	0.227055	0.862069
20 days	284575.5	533.4562	251.0542	0.189535	0.076854	0.239485	0.862069
21 days	26665.22	163.2949	115.6537	0.151711	0.094355	0.175638	0.931034
22 days	52361.4	228.8261	144.1287	0.166713	0.124839	0.191619	0.931034
23 days	51735.77	227.455	149.2709	0.17673	0.148265	0.190963	0.934483
24 days	49258.13	221.9417	145.9852	0.277826	0.155263	0.210385	0.934483
25 days	21151.05	145.434	110.2301	0.267737	0.114466	0.185104	0.965517
26 days	21335.24	146.0659	106.0778	0.258047	0.111279	0.183036	0.968966
27 days	16340.27	127.8291	92.4109	0.15762	0.0844	0.154502	1
28 days	19466.75	139.5233	107.173	0.164857	0.102969	0.167724	1
29 days	22162.2	148.8697	117.7591	0.187933	0.106471	0.184269	0.965517
30 days	19719.08	140.4246	112.3507	0.178483	0.110105	0.167508	0.965517

Code Used