

Seoul Bike Demand Analysis

William K Davis III

Max Kutschinski

Pei-Yin Yang

2022-08-02

Abstract

City-wide bike sharing programs have gained in popularity among major cities in recent years, providing an attractive alternative to cars and other forms of transportation. A key factor in the program’s success is the efficient allocation of rental bikes across the city. This study aims to predict bike demand based on time differences by making hourly forecasts over a 24-hour horizon, enabling policymakers to get a better picture of the next day’s bike demand. We developed time series models, as well deep learning and tree-based models, and compared them using time series cross-validation with MAE as the evaluation criterion. The best performing model is based on the LSTM architecture, which yielded a test prediction Mean Absolute Error less than half that of the naïve baseline forecast. This study addresses the time-dependent aspect of bike demand, but further studies are needed to investigate spatial variations.

Introduction

While the concept of renting transportation has been around for decades in the form of rental car companies, such as those often found at airports, it has also been adapted to the modern “sharing economy” in the form of bikeshare programs. These programs, often offered by local governments, allow users to rent bicycles for individual (point-to-point) trips, such as commuting to and from work. Cycling can be a faster mode of transportation than walking and might even be faster than driving or taking a taxi in the most congested of cities. Finally, bicycles offer a lower-pollution alternative to driving, which can be appealing to cities struggling to contain emissions.

A key challenge faced by administrators of bikesharing programs is the efficient allocation of available bicycles. Bikes must be available in the places that people need them and at the time they are needed in order for the program to be effective. In order to efficiently allocate bikes, Schuijbroek et al. (2017) suggest administrators can periodically transport bikes from areas of low demand to areas of high demand. Demand in a bikeshare system is a function of both time and location. In this paper we focus on the time component of the demand for bicycles at each hour of the day, ignoring the spatial component of demand. We will use the Seoul Bike Sharing Demand data from UCI (2020). Our focus will be limited to predicting demand at each hour of the day, ignoring inferential aspects of the analysis. While much of the recent research using this dataset has focused on machine learning methodologies, we will take a multifaceted approach that incorporates advances in time series modeling in addition to the popular machine learning methodologies.

This paper begins with a review of the relevant literature, including studies of the Seoul data specifically. Next, we present the results of the exploratory data analysis, including a description of the dataset and relevant profiles of the features. The fourth section describes the modeling techniques to be applied to the data. The fifth section describes the evaluation techniques to be used for measuring model performance and selecting the best model. The penultimate section presents the results of the analysis and a discussion in the context of the problem to be solved. Finally, the conclusion provides suggestions for action and application of the results, while also highlighting potential areas of future research.

Literature

The popularity and accessibility of the Seoul bike dataset has resulted in its use for numerous studies. A majority of these studies have focused on the use of various machine learning algorithms. E and Cho (2020) found that a CUBIST model, which combines tree- and regression-based methods into a series of rules, performed best on the Seoul data when measured by R^2 and RMSE on the testing dataset. Gao and Chen (2022) found that another tree-based method, random forest (RF), performed best on a similar bikeshare dataset when measured by R^2 and RMSE. Both studies further showed that weather-related variables, such as temperatures and precipitation, were among the most important for predicting demand. Gao and Chen’s results highlight the importance of selecting a relevant evaluation metric and explanatory variables. When socioeconomic variables were included in the model, the RF outperformed the support vector machine (SVM) when measured by both RMSE and MAE. However, when the socioeconomic variables were excluded from the model, the RF outperformed the SVM when measured by RMSE, but the SVM performed better when measured by MAE. This indicates that without the socioeconomic variables included, the SVM was prone to a few errors that were quite large in magnitude, while the RF was more prone to smaller but more frequent errors. This reinforces the importance of using multiple metrics when evaluating predictions.

Considerably fewer researchers have made use of traditional time series methodologies when predicting demand of a similar nature to the bikeshare data. Both E and Cho (2020) and Gao and Chen (2022) make use of temporal variables such as hour, day of the week, and holidays. In each case they were found to be of moderate or high importance. Gao and Chen (2022) applied linear regression using temporal variables such as a weekend indicator as a predictor, but this model greatly underperformed the machine learning methods. Further, there is no discussion of any attention paid to stationarity and autoregression, which are common in time series data but may also result in violations of the standard assumptions of linear regression.

Our analysis looks to build on this work by using newer machine learning methods such as long short-term memory and recurrent neural networks, as well as modern time series techniques that leverage the underlying structure of the data.

Exploratory Analysis

Table 1: Variable definitions

Variable	name	Type	Definition
Hour	Hour	datetime	year-month-day hour:minute:second
Rented Bike count	BikeCount	numeric	Count of bikes rented at each hour
Temperature	Temperature	numeric	Temperature in Celsius
Humidity	Humidity	numeric	% humidity
Windspeed	WindSpeed	numeric	meters/second
Visibility	Visibility	numeric	in 10m
Dew point temperature	Dewpoint	numeric	Celsius
Solar radiation	SolarRadtion	numeric	MJ/m2
Rainfall	Rainfall	numeric	mm
Snowfall	Snowfall	numeric	cm
Seasons	Seasons	categorical	Winter, Spring, Summer, Autumn
Holiday	Holiday	categorical	Holiday/No holiday
Functional Day	FunctionalDay	categorical	NoFunc(Non Functional Hours), Fun(Functional hours)
Workday	Workday	categorical	A workday is a weekday that is not a holiday.

The dataset consists of 8,760 hourly observations of 12 variables from 2017-12-01 00:00:00 to 2018-11-30 23:00:00. There are 295 observations where *BikeCount*=0 due to the bikeshare system not functioning. There

are no other periods where $BikeCount=0$. Table 1 contains information on the variables in the dataset, including the

Bike Count

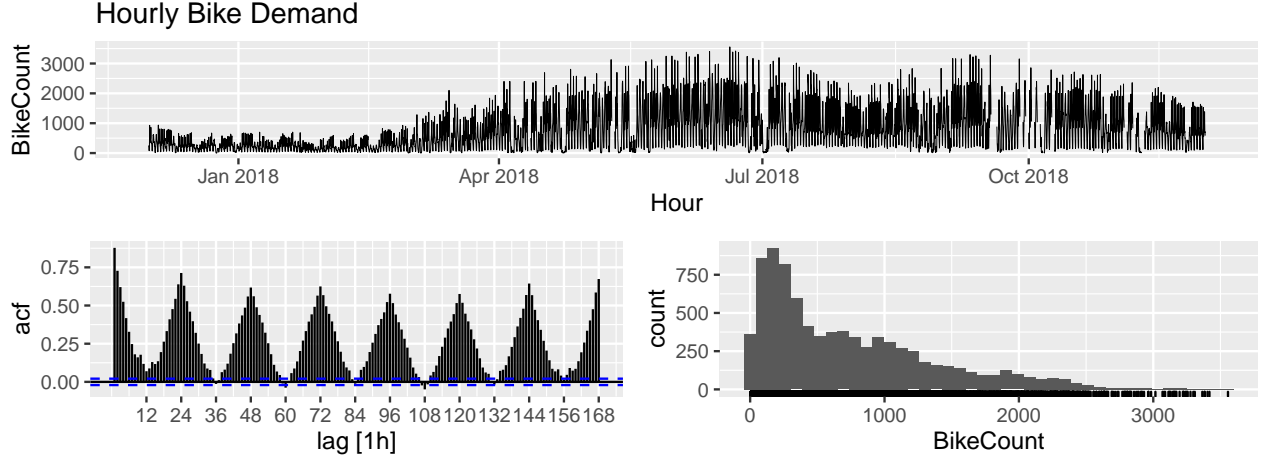


Figure 1: Hourly bike demand

The data shows increasing demand and variability during the summer months. The bike demand data are counts, meaning they are technically discrete. Based on the histogram and the count nature of the data, it appears that a poisson distribution would be most appropriate for the data. If we were to model the bike demand on a continuous scale, a log-normal distribution might be appropriate.

Figure 1 highlights the incredible variation in demand over time. The variance in demand appears to increase during the summer months and then decrease again in the autumn. This heteroskedasticity violates the constant variance assumption of ordinary least squares regression and will have to be corrected if regression-based methods are to be used. We can also see that the mean of the series appears to increase during the summer months before decreasing again in the autumn.

Kwiatkowski et al. (1992)’s test yields $p = 0.01$, meaning there is strong evidence in favor of the presence of a unit root and the data is likely non-stationary. Ljung and Box (1978)’s test was conducted up to 24 lags which resulted in a test statistics of 39743 with $p = 0$, indicating strong evidence that there is serial correlation in the hourly bike count. The acf plot in Figure 1 shows that autocorrelation is significant at most lags out to 168 hours, which represents the same hour of the same day in the previous week. The strong serial correlation makes this dataset a good candidate for time series techniques.

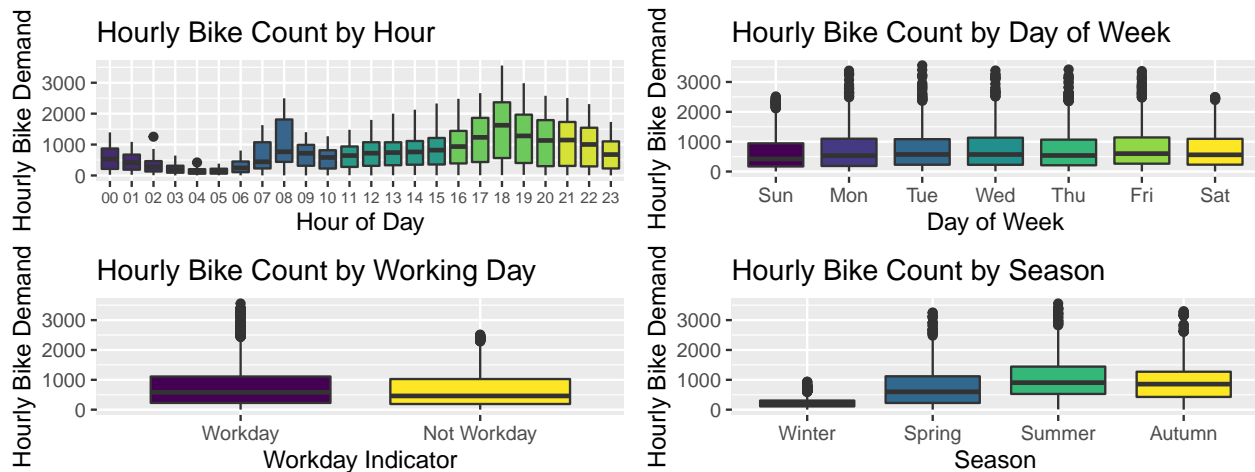


Figure 2: Hourly bike demand by time period

Figure 2 highlights various seasonal and temporal patterns in the data. There is certainly an effect based on the hour of day, with demand increasing sharply in the morning, presumably during the morning commute, before decreasing into the lunchtime hour. From there, demand rises steadily through the evening commute before peaking around dinner time and then falling through the nighttime hours. Variability also appears greater during the evening hours, which is consistent with the idea heteroscedasticity; the variance increases with the level of the series. Demand appears slightly higher during workdays (where workdays are weekdays that are not holidays) and weekdays, though the difference does not appear particularly large. There are a number of large positive outliers during the weekdays. Finally, demand appears largest during the summer months and smallest during the winter, as biking would be a less desirable option in the cold.

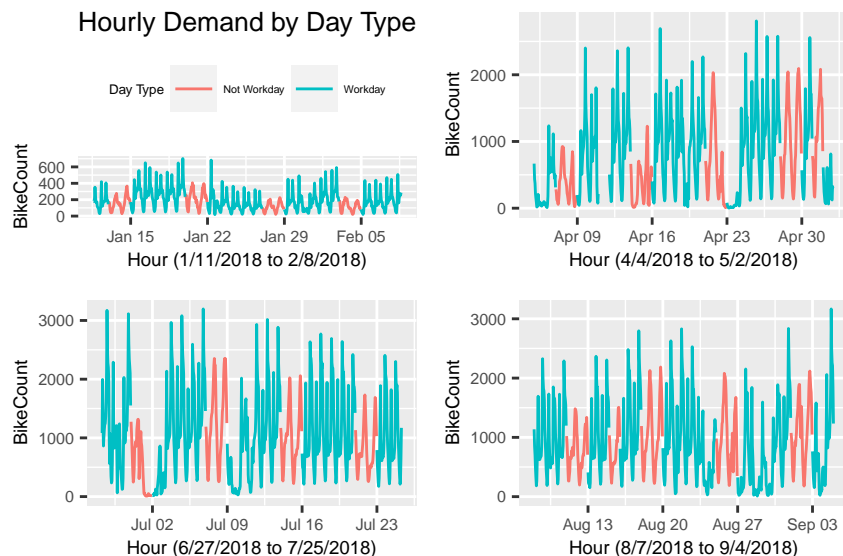


Figure 3: Hourly bike demand by type of day, highlighting differing seasonality between workdays and non-workdays

Figure 3 provides insight into the seasonality of the bike demand. Each panel is a 4-week sample of hourly bike demand, with non-workdays (weekends and holidays) highlighted in red. While these days appear to follow a consistent hourly pattern much like the workdays, we can see that there is a difference between the seasonality of workdays and the seasonality of non-workdays. Most notably, the troughs in demand appear

consistent between the types of days while the peaks are consistently higher for workdays as compared to non-workdays. This changing seasonality based on type of day will need to be captured in our model.

Continuing with the impact of seasonality and temperature on demand, we will next explore covariates included with the dataset.

Covariates

This dataset includes a number of covariates to aid in modeling bike demand. All of these covariates are listed in Table 1, along with their definitions. The covariates fall into one of two broad categories: weather and social. Weather covariates include variables such as temperature, humidity, precipitation, and others. Social variables capture the impact of calendar-based human behavior, such as holidays and weekends.

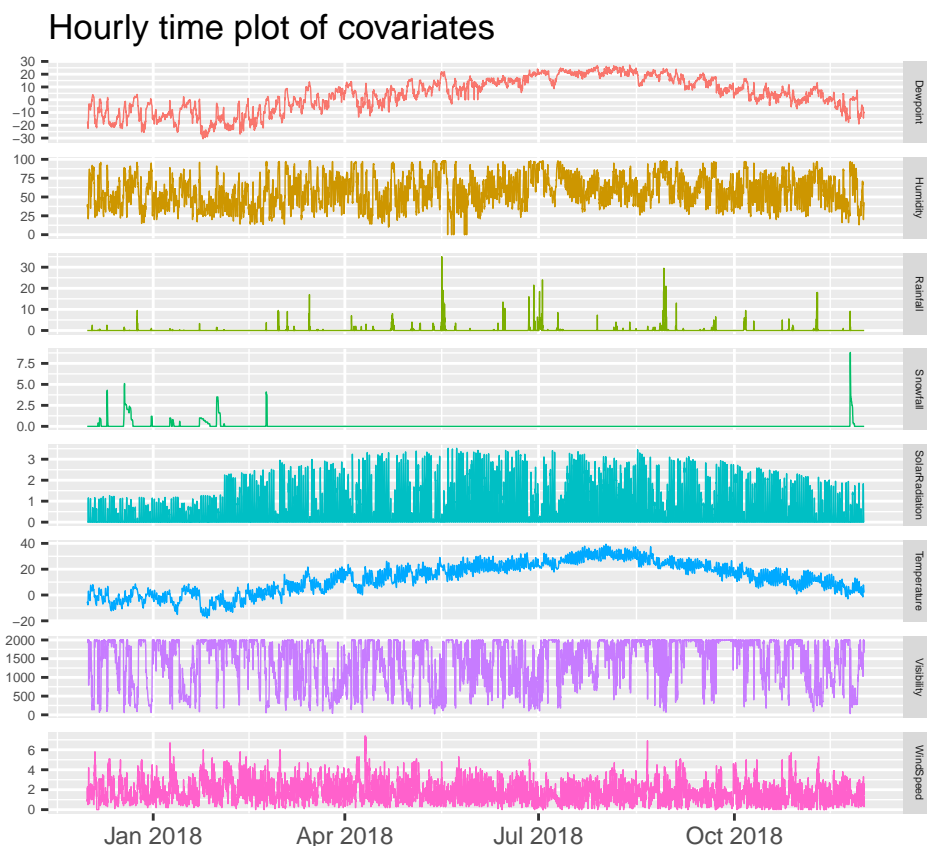


Figure 4: Hourly time plots of covariates

Time plots of the covariates (Figure 4) show a dynamic set of variables. Most of the covariates appear to be typical time series data with varying degrees of trend, seasonality, cyclicity, autocorrelation, and heteroskedasticity. Dewpoint and Solar Radiation follow a predictable pattern that mirrors temperature throughout the year, with an upward trend peaking in the summer months and downward trend that hits a trough in the winter months. Humidity, visibility, and Wind Speed appear to have less of a trend throughout the year. Precipitation (rainfall and snowfall) appears to have a much more random pattern throughout the year, with a large number of periods having no precipitation. Depending on the predictive performance of the raw continuous precipitation features, it may prove more performant to convert them to binary variables that simply indicate if precipitation occurred during that period.

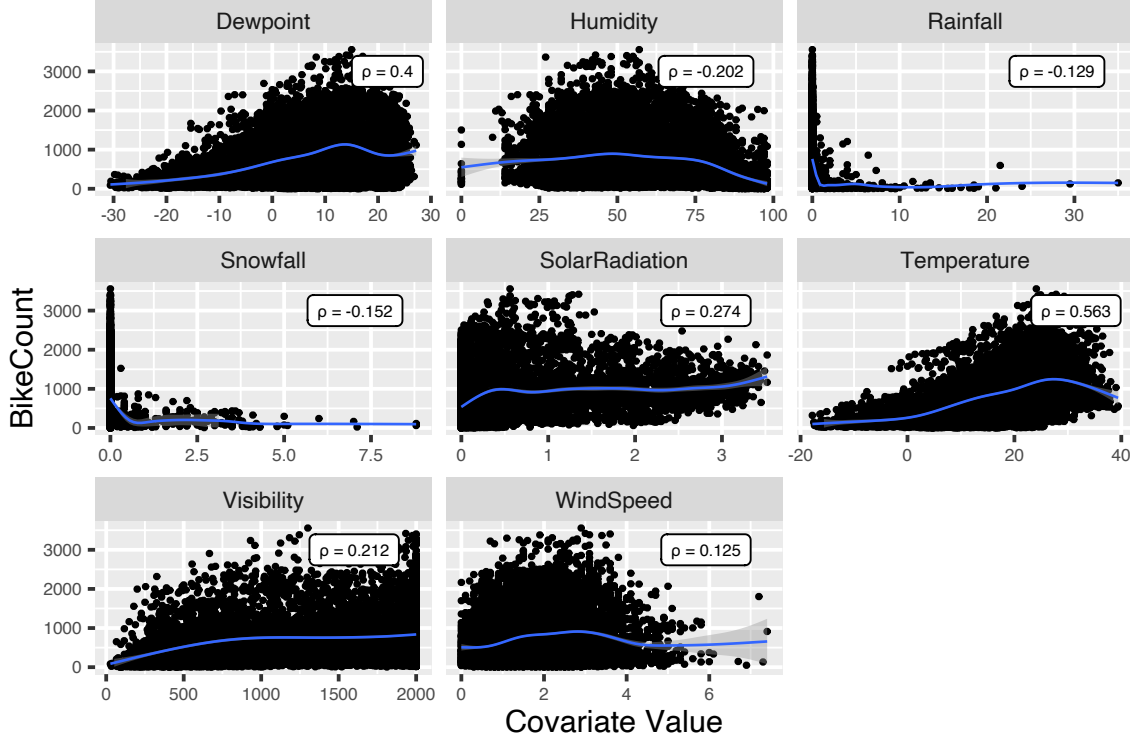


Figure 5: Scatterplots and smoothed GAM fits of Bike Count vs. each covariate

Figure 5 highlights a number of important relationships in the data. Bike count appears moderately linearly correlated with both temperature ($\rho = 0.563$) and dewpoint ($\rho = 0.400$), though a higher-order interaction might better represent this relationship. However, there is significant multicollinearity between temperature and dewpoint ($\rho = 0.913$), so additional analysis will be required to isolate the effect of each variable. Finally, we can see that a number of the covariates, such as wind speed, visibility, solar radiation, rainfall, and snowfall all appear to follow non-normal distributions. As such, their relationship to bike count might be best represented with a higher-order polynomial or other non-linear relationship.

Modeling

The modeling methods we have chosen for predicting bike demand fall into roughly two categories: parametric time series modeling and non-parametric machine learning. The parametric time series models include traditional time series methods like ARIMA and modern methods including the regression-based Prophet and dynamic linear model FASSTER. The non-parametric machine learning methods include random forests, boosted trees, and Long Short-Term Memory neural networks.

The stated goal of this analysis is to predict bike demand, not to conduct statistical inference (on the factors that impact bike demand). Therefore, the non-parametric methods can be used freely despite their typical lack of inferential ability. Further, we do not have to constrain ourselves to the typical assumptions of parametric methods that are required for inference, namely independent observations and normally distributed residuals with mean zero and constant variance. The assumption of independent observations is obviously violated in time series data and most often appears in the form of autocorrelated errors. Therefore, the ability to disregard the assumptions required for inference in favor of a focus on prediction will greatly expand the number of viable methods for modeling the data.

Time Series Modeling

Figure 1 highlights the heteroskedasticity in the bike demand. This can be corrected using a Box and Cox (1964)'s transformation. If we let x_t represent the value of raw Bike Count data and y_t represent the transformed data, we have

$$y_t = \frac{x_t^\lambda - 1}{\lambda} \quad (1)$$

where $\lambda = 0.1478452$ was selected using Guerrero (1993)'s method.

SARIMAX

Traditional linear regression models can be adapted to handle the autocovariance structure of time series by assuming that the errors follow a seasonal ARIMA (SARIMA) process instead of the traditional $\epsilon \sim \text{iid } N(0, \sigma^2)$. This results in the modified regression equation

$$y_t = \beta_0 + \sum_{j=1}^k \beta_j z_{jt} + x_t \quad (2)$$

where:

- y_t is the response variable at time t
- β_0 is the traditional intercept
- z_{1t}, \dots, z_{kt} are the k exogenous regressors observed at time t .
- β_1, \dots, β_k are the regression coefficients.
- x_t are the regression errors, which are assumed to follow an SARIMA process as in (3).

(2) is often referred to as SARIMAX for Seasonal ARIMA with eXogenous regressors.

$$\Phi_P(B^S)\phi(B)\nabla_S^D\nabla^d x_t = \Theta_Q(B^S)\theta(B)w_t \quad (3)$$

where:

- $\Phi_P(B^S)$ are the P seasonal autoregressive components
- $\phi(B)$ are the p autoregressive components
- ∇_S^D are the D seasonal differences
- ∇^d are the d differences
- $\Theta_Q(B^S)$ are the Q seasonal moving average components
- $\theta(B)$ are the q moving average components
- $w_t \sim \text{iid } N(0, \sigma_w^2)$ is the traditional Gaussian white noise

as defined by Shumway and Stoffer (2019). SARIMA models of this form are often written $ARIMA(p, d, q) \times (P, D, Q)_S$. A KPSS test indicates that the data is non-stationary and requires $D = 1$ seasonal difference. Autocorrelation and partial autocorrelation plots can be used to determine the ARIMA order.

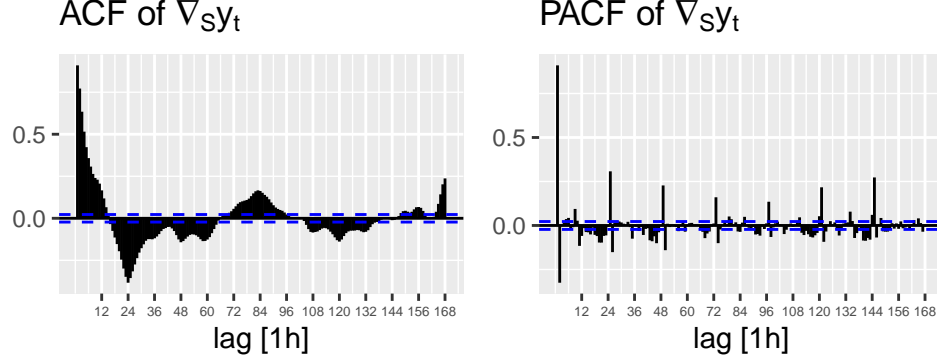


Figure 6: Correlograms for ∇sy_t

Figure 6 indicates a seasonal autoregressive model with $P = 6$ and $p \geq 12$. This is a large number of seasonal autoregressive terms and will almost certainly result in characteristic roots inside the unit circle, causing the model to be unstable. A quick search of the parameter space using the R function `ARIMA` O’Hara-Wild et al. (2021)’s `fable` package to automatically select values for p, d, q, Q results in no stable $ARIMA(p, d, q) \times (6, 1, Q)_{24}$ models being found.

When the data exhibit higher frequency or multiple types of seasonality (such as daily and weekly, in our case), an alternative solution for modeling seasonality suggested by Hyndman and Athanasopoulos (2021) can be to use fourier terms. Introducing the fourier terms will result in a model of the form

$$y_t = \beta_0 + \sum_{j=1}^k \beta_j z_{jt} + s_d(t, m) + s_w(t, n) + x_t \quad (4)$$

$$s_d(t, m) = \sum_{i=1}^m \left[\alpha_i \sin\left(\frac{2\pi it}{24}\right) + \beta_i \cos\left(\frac{2\pi it}{24}\right) \right] \quad (5)$$

$$s_w(t, n) = \sum_{l=1}^n \left[\gamma_l \sin\left(\frac{2\pi lt}{168}\right) + \delta_l \cos\left(\frac{2\pi lt}{168}\right) \right] \quad (6)$$

where:

- $\beta_0 + \sum_{j=1}^k \beta_j z_{jt}$ are the intercept, independent variables and their coefficients, as defined in (2)
- $s_d(t, m)$ are the m daily seasonality fourier terms
- $s_w(t, n)$ are the n weekly seasonality fourier terms
- x_t are the model errors, which are assumed to follow a $ARIMA(p, d, q)$ process (note that in contrast to (2) this is a non-seasonal ARIMA model; the seasonality is omitted from the ARIMA process because it is captured by the fourier terms).

The R function `fable::ARIMA` is used to fit this model, which, as (<https://stats.stackexchange.com/users/159/rob-hyndman>) notes, handles missing values automatically using a Kalman filter. As a result of the missing values, the point (parameter) estimates will likely be biased. The use of covariates may capture some of the pattern in the missing data, if there is any, which would help to reduce the bias. However, the final point forecasts will likely still have some bias.

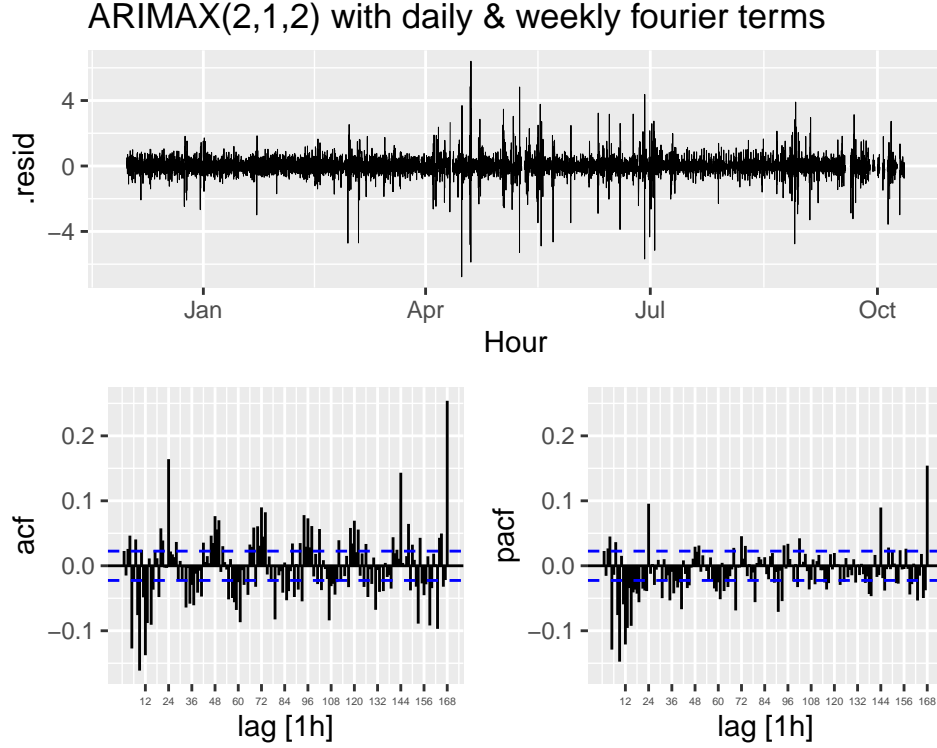


Figure 7: Residuals for an ARIMAX(2,1,2) model with daily and weekly fourier terms

The fourier model resulted in the selection of an $ARIMA(2, 1, 2)$ process for the regression residuals. The residual ACF and PACF plots from the fourier model (Figure 7) still show significant autocorrelation. Autocorrelation in the residuals violates the Gauss-Markov assumptions, meaning the model should not be used for inference. The goal of this analysis is to generate accurate (point) forecasts, which are still valid even when the errors exhibit autocorrelation.

Next we will fit a linear regression with SARIMA errors to the data where all parameters p, d, q, P, D, Q are chosen automatically using `fable::ARIMA`, which relies on the algorithm developed by Hyndman and Khandakar (2008). This resulted in the selection of an $SARIMA(1, 0, 1) \times (5, 1, 0)_{24}$ for x_t .

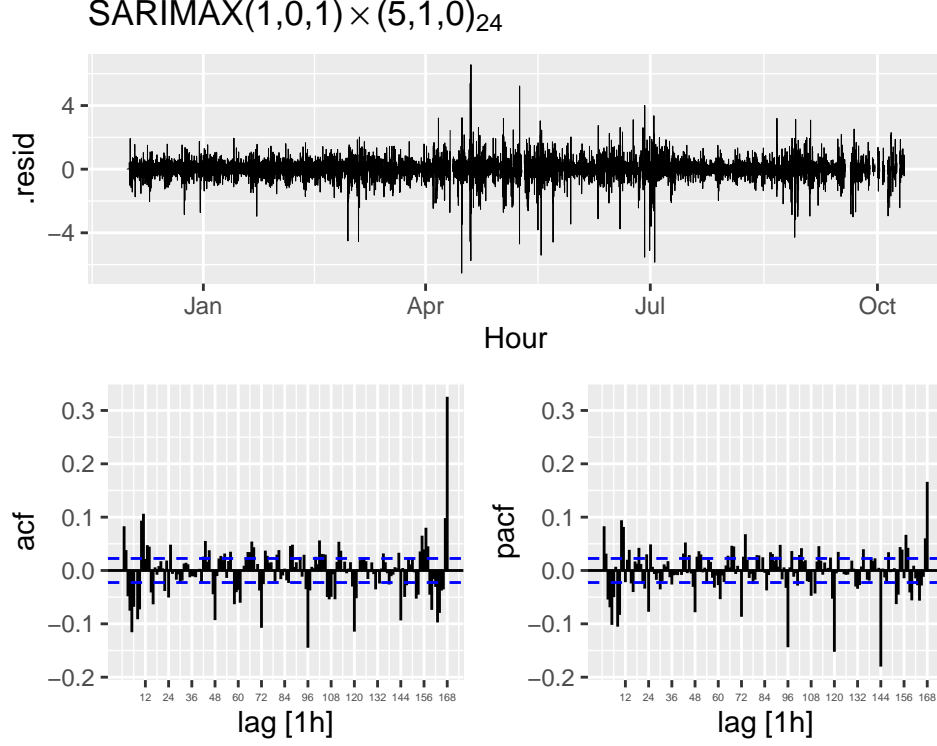


Figure 8: Residuals for an $SARIMAX(1, 0, 1) \times (5, 1, 0)_{24}$

Figure 8 shows that, while the autocorrelation (and partial autocorrelation) for a $SARIMA(1, 0, 1) \times (5, 1, 0)_{24}$ model has been reduced over the $ARIMAX(2, 1, 2)$ with fourier terms, there is still repeating statistically significant autocorrelation.

Finally, we will test a SARIMAX model that makes use of two derived variables. First, as was seen in Figure 5, there appears to be a higher-order relationship between temperature and bike count. Therefore, we will add a $Temperature^2$ term to the model. Additionally, as can be seen in 4, Rainfall and Snowfall are both continuous variables that have a large number of zero values. We believe it is a reasonable assumption that the presence or absence of Rainfall or Snowfall will have the largest impact on the bike count, and the quantity of rain or snow is only marginal in effect compared with the presence or absence rain or snow. Further, we believe that any precipitation, regardless of type, will impact bike count, and the type (and quantity) of precipitation is of much less importance. Therefore, we will create a derived indicator variable **Precipitation_Flag** to identify the presence of any precipitation, defined as follows:

$$\text{Precipitation Flag} = \begin{cases} 1, & \text{Rainfall} > 0 \text{ or } \text{Snowfall} > 0 \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

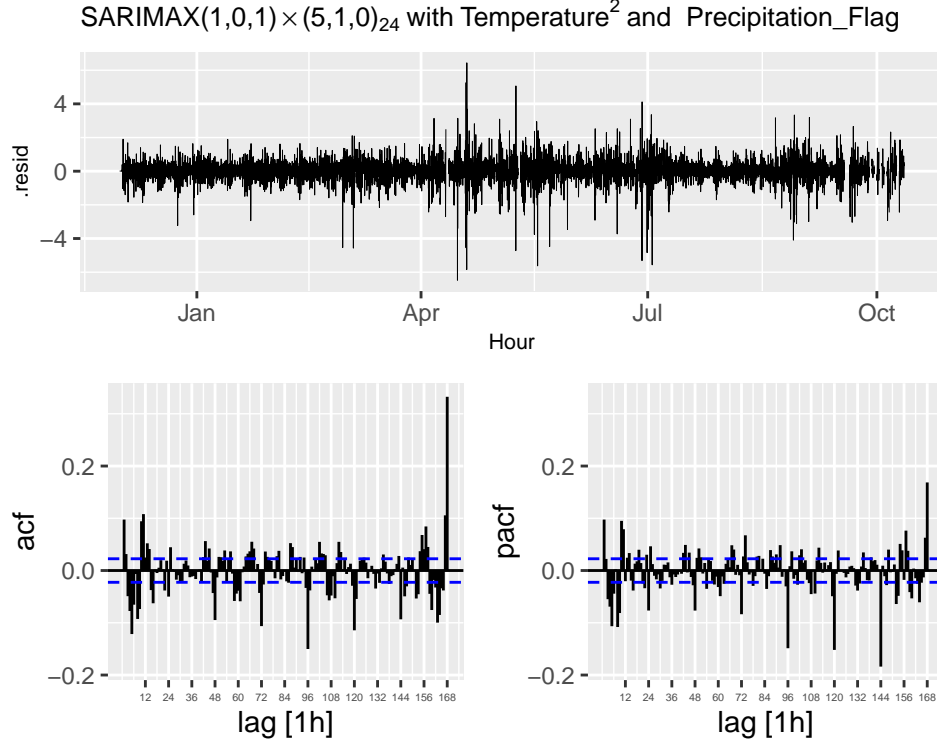


Figure 9: Residuals for an $SARIMAX(1, 0, 1) \times (5, 1, 0)_{24}$ with $Temperature^2$ and Precipitation Flag

The addition of $Temperature^2$ and Precipitation Flag have reduced AIC from 14833 to 14789, AICc from 14833 to 14789, and BIC from 14965 to 14927, when compared with the model from Figure 8. Therefore the addition of these two variables has helped improve the model fit. Unfortunately, as 9 shows there is still statistically significant autocorrelation in the residuals. It is possible that the seasonal patterns in the data are too complex or too high frequency to be captured by an ARIMA, ARIMAX, or SARIMAX model. As a result of the residual variance structure, prediction intervals cannot be generated with these models, though point estimates are still valid.

In an effort to overcome the current residual structure we will explore more modern time series techniques.

Prophet

Taylor and Letham (2018)'s Prophet uses a generalized additive model (GAM) to capture the different features of the time series. Our model will take the form

$$y(t) = g(t) + s_d(t, m) + s_w(t, n) + \sum_{j=1}^k \beta_j z_{jt} + \epsilon_t \quad (8)$$

where:

- $g(t)$ is a piecewise constant function to represent the trend in the series
- $s_d(t, m)$ as in (5)
- $s_w(t, n)$ as in (6)
- $\sum_{j=1}^k \beta_j z_{jt}$ is as in 2
- ϵ_t is the model error

One major benefit of (8) is that it is much faster to fit than a model with ARIMA terms. The downside is that it does not explicitly capture (non-seasonal) $AR(p)$ and $MA(q)$ terms. As Taylor and Letham (2018)

note, because this is a curve-fitting exercise rather than a generative model, missing values do not need to be manually imputed. Instead, the curves $g(t)$, $s_d(t, m)$, and $s_w(t, n)$ are fit to the raw data and then used to impute the missing values. Further, as with the ARIMA model, the presence of the covariates helps to control for any missingness that might be related to one of the covariates. However, the missing value interpolation will still result in biased parameter estimates for the model.

We will use `prophet` from O'Hara-Wild (2020)'s package `fable.prophet` to train the model in R. The `R` function sets a number of desirable default parameter values as recommended by Taylor and Letham (2018). In order to properly tune the model to our dataset we will use TSCV (see Model Selection and Test Error) to select m and n . The TSCV will consist of 63 folds/time-slices. We will search over the parameter space $m \in \{1, 2, 5, 10, 20, 50\}$ and $n \in \{3, 5, 10, 20, 50\}$, giving us $6 \times 5 = 30$ possible parameter combinations.

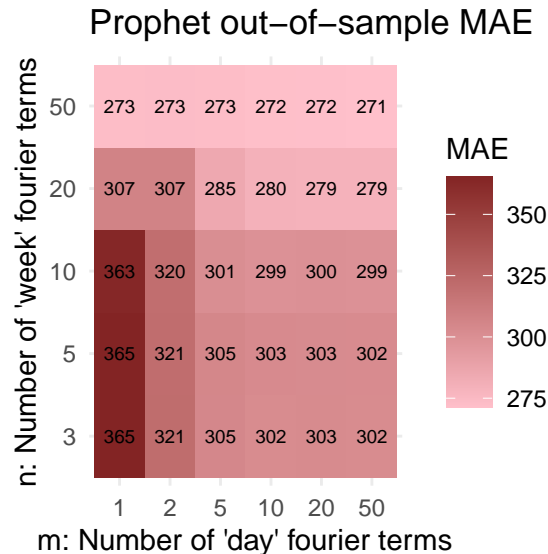


Figure 10: Cross-validated MAE for Prophet models with various numbers of (daily and weekly) fourier terms

Figure 10 identifies $m = 50$ and $n = 50$ as the optimal number of fourier terms (they minimize MAE). We fit this model to the entire training dataset, including the two derived variables discussed in the SARIMAX section, and analyze the residuals.

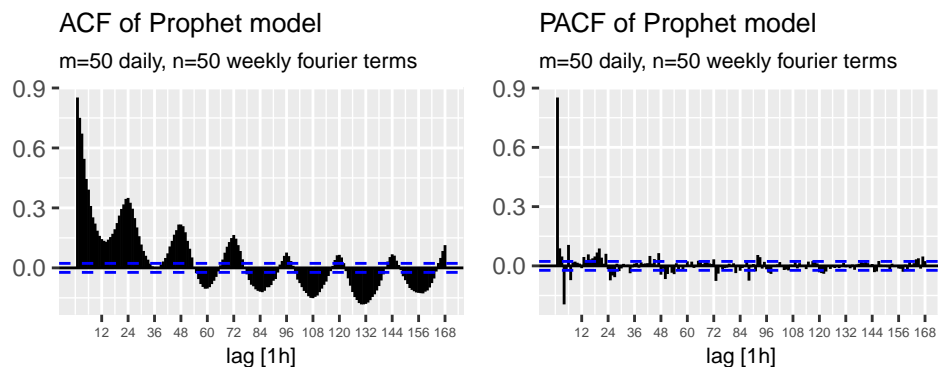


Figure 11: Residual correlograms for a Prophet model with $m = 50$ (daily) and $n = 50$ (weekly) fourier terms

As can be seen from the PACF of the residuals in Figure 11, there is still autocorrelation not captured by the model. Therefore, the model point estimates will be valid but prediction intervals cannot be generated due to the error structure. This is to be expected, as Taylor and Letham (2018) notes that the design of

Prophet as a curve-fitting exercise, rather than a generative model (such as ARIMA) that captures the temporal dependence of the data, can be expected to give up some inferential capabilities (in exchange for computational and other benefits). Given that our goal is forecasting/prediction rather than inference, this should not pose a problem for the analysis.

fasster

One interesting feature of the residual PACFs from both the SARIMAX model (Figure 8) and the Prophet model (Figure 11) is that the seasonal “spikes” are parabolic in shape. That is, the PACF starts high at lag 24, then decreases until reaching a min around lag 72 or lag 96 before rising again to another peak at lag 168 (which is the one week lag). This shape might be why the AR and fourier terms are not capturing all of the autocorrelation in the data: the strength of the correlation appears to be a non-linear function of time. One reason for this shape in the PACF plot could be the presence of switching seasonality. More specifically, it’s possible that the seasonality effects (not just the series itself) is different on working days and non-working days. A hand-wavy analysis would say that, given there are seven days in a week, any single day is, on “average”, $7/2 = 3.5$ days away from any other day. For weekdays specifically, they are 3.5 days away from a weekend (again, on average). Therefore, the trough in the PACF curve occurring between 3 and 4 days might be due to that lag most frequently correlating weekdays with weekend days, resulting in lower correlation than when weekdays are compared with other weekdays, which would occur most frequently at lags < 3 and > 4 . This is certainly not a technical explanation of the pattern, but it is plausible. The previously employed techniques (SARIMAX and Prophet) don’t allow for changing seasonality.

O’Hara-Wild and Hyndman (2022)’s fasster methodology allows for the seasonal (and trend) component(s) of the model to be switched based on the “state” of a given observation, which is defined via a discrete variable passed to the model. In our case, this variable would be *Workday*, which identifies the observation as a workday or non-workday (weekend or holiday). Different seasonality coefficients can then be fit based on whether a bike count is observed on a workday or non-workday. This is an improvement over the previous models, which simply include the *Workday* variable as a regression coefficient, reducing its effect to a shift in the level of the series. Like some of the previous methods, fasster uses fourier terms to capture seasonality. We will use 10 and 5 sin-cos for the daily and weekly seasonality, respectively, based on the recommendation of Hyndman and Athanasopoulos (2021). Finally, as with the other time-series methods, fasster uses Kalman filtering and smoothing to interpolate missing values, which will result in biased parameter estimates, though the bias will be less than that from a model fit by excluding the missing values.

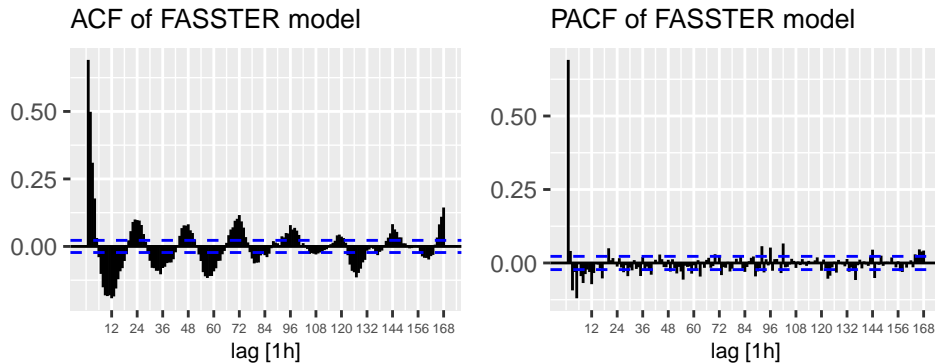


Figure 12: Residual correlograms from the fasster model.

As with all of the other time-series models, fasster is not able to capture the autogression in the series, resulting in significantly autocorrelated errors. Therefore, point forecasts from the fasster model are valid, but inference may not be conducted with the model due to the impact of the autocorrelation on the covariance structure of the model.

Machine Learning

Tree Methods

Random Forests Random forests is an ensemble method for regression and classification tasks that is built on decision trees. It extends on the idea of bootstrap aggregation, or bagging, which is used to reduce the variance of a statistical learning method via averaging. In the context of regression trees, this is done by constructing B unpruned trees from B bootstrap samples and averaging the predictions \hat{f} as displayed in Eq. (9).

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x) \quad (9)$$

One notable drawback of bagging is that the individual trees can be very correlated depending on how strong the predictors are. Random forests address this issue by randomly selecting a subset of the features at each split and thus de-correlating the trees.

Random Forests (RF) is implemented using the ranger package for increased computational speed. Fernández-Delgado et al. (2014) notes that RF is an algorithm that is known to provide good results in the default settings. The arguably most influential hyperparameter is *mtry*, the number of randomly drawn features that are available at each split. In the regression case, $p/3$ is the default setting. When executing ranger via caret it automatically performs a grid search of *mtry* its entire parameter space. By default, the algorithm evaluates 3 points in the parameter space (smallest and largest possible *mtry*, as well as their mean) with 25 bootstrap iterations as an evaluation strategy and chooses the value with the lowest MSE.

The hyperparameters of our model are evaluated using TSCV, yielding an optimal *mtry* value of 53. Probst et al. (2019) suggests that such a high value is usually indicative of a high number of relevant predictors.

XGBoost XGBoost is an advanced implementation of Gradient Boosting, which is an ensemble method for regression and classification tasks that combines multiple weak learners into a stronger learner. In the context of decision trees, a weak learner is defined as a tree with a small number of terminal nodes. The trees are grown sequentially and they are fit on the residuals of the current fit as opposed to the outcome Y . This has the effect of capturing signal that is not yet accounted for by the current set of trees. In addition, each weak learner is shrunk down by some shrinkage factor before it is used, making boosting a “slow” learning approach.

The first step is to initialize the model $F(x)$ with a constant value γ , which can be obtained by minimizing it with respect to a loss function L , as displayed in the following optimization problem:

$$F_0(x) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, \gamma) \quad (10)$$

After specifying the number of base learners M , the following steps are repeated for each base learner from $m=1$ to $m=M$:

First, the pseudo-residuals r_{im} are calculated for each i th training example.

$$r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad (11)$$

Then a base learner $h_m(x)$ is fit to the pseudo-residuals using the modified training set $\{(x_i, r_{im})_{i=1}^n\}$.

Lastly, the model is updated as follows:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x) \quad (12)$$

$$\text{where } \gamma_m = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$$

XGBoost is a more regularized form of Gradient Boosting which uses L1 and L2 regularization to improve model generalization capabilities. It also allows for parallel processing, and has a built-in routine for handling missing values via its sparsity-aware split finding algorithm.

There are a variety of booster parameters in XGBoost that can be optimized via TSCV. Rather than tuning all parameters simultaneously, Banerjee found it often helps to make small changes incrementally. The general idea is to start with a high learning rate and a small number of base learners, then tune other parameters, and finally decrease the learning rate while proportionally increasing the number of trees. The other adjusted parameters are the maximum depth of a tree *max_depth*, the minimum required loss reduction *gamma*, the fraction of columns to be subsampled *colsample_bytree*, the minimum sum of weights of all observations required in a child *min_child_weight*, and the fraction of observations to be randomly sampled per tree *subsample*.

Long Short-Term Memory (LSTM) Recurrent Neural Network

LSTM is a variety of Recurrent Neural Network (RNN). According to Hochreiter and Schmidhuber (1997), the benefit of the Long Short-Term Memory (LSTM) network over other recurrent networks comes from “constant error back propagation”, an improved method of back-propagating the error.

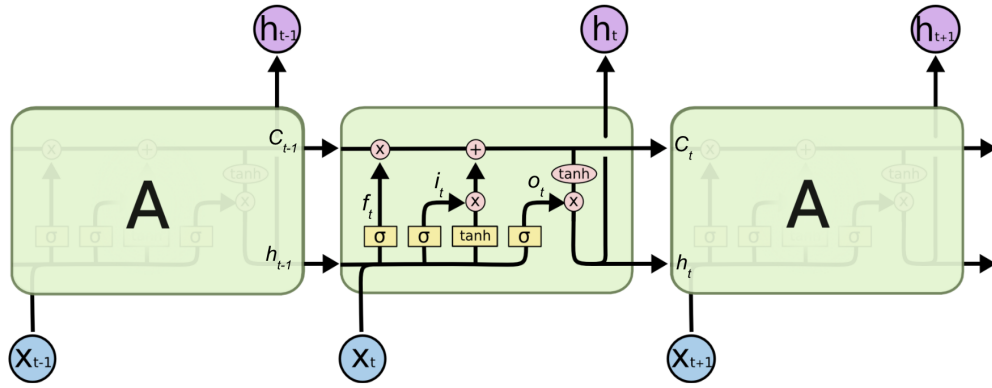


Figure 13: The repeating module in an LSTM contains four interaction layers. Images borrowed from @colah’s blog.

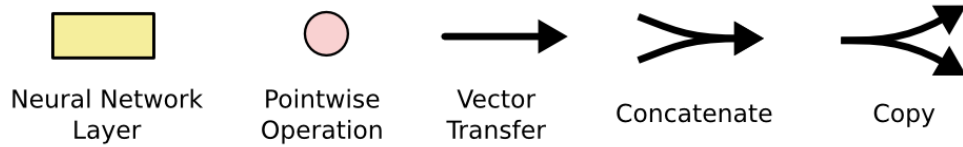


Figure 14: The repeating module in an LSTM contains four interaction layers. Images borrowed from @colah’s blog.

The key to LSTMs is the cell state, the horizontal line running through the top of the graph. An LSTM features three gates: input, forget & output. Gates are a way to let information through optionally. They comprise a Sigmoid neural net layer and a pointwise multiplication operation.

Let X_t be the input and h_t be the output at time t . At time t , the equation of gates, input and output of the LSTM cell are as follows:

$$i_t = \sigma(W_i[x_t] + R_i[h_{t-1}] + b_i) \quad (13)$$

$$f_t = \sigma(W_f[x_t] + R_f[h_{t-1}] + b_f) \quad (14)$$

$$o_t = \sigma(W_o[x_t] + R_o[h_{t-1}] + b_o) \quad (15)$$

$$g_t = \tanh(W_x[x_t] + R_x[h_{t-1}] + b_x) \quad (16)$$

$$h_t = o_t \times \tanh(c_t) \quad (17)$$

where i_t , f_t , o_t denotes input gate, forget gate and output gate respectively. The g_t and h_t refer to input and output of the LSTM cell respectively, and h_{t-1} refers to the output of the previous LSTM cell. Sigmoid (σ) and \tanh are the activation functions to map the non-linearity.

c_t represents the Constant Error Carousel (CEC), the magic of the LSTM as it prevents vanishing gradients, it defines the memory cell to store the past state. It is denoted as follows:

$$c_t = c_{t-1} \times \text{forget gate} + \text{new input} \times \text{input gate} \quad (18)$$

In the case of regular RNNs during backpropagation, the derivative of an activation function, such as a logistic, will be less than one. Therefore over time, the repeated multiplication of that value against the weights will lead to a vanishing gradient.

In the case of an LSTM, we only multiply cell state \times forget gate, which acts as both the weights and the activation function for the cell state. As long as forget gate = 1, the information from the previous cell state passes through unchanged. This is why LSTM can deal with more intricate problems than the RNN, by keeping a constant flow of error throughout the backpropagation from cell to cell.

Data processing Brownlee (2020) suggest it is a good idea for machine learning algorithms that fit a model that uses a weighted sum of input variables, such as linear regression, logistic regression, and artificial neural networks (deep learning) to normalize the data. Here, by using the `scikit-learn` object `MinMaxScaler`, the normalization scales each input variable separately to the range $[0, 1]$, which is the range for floating-point values where we have the most precision. A value is normalized as follows:

$$y = \frac{x - \min\{x\}}{\max\{x\} - \min\{x\}} \quad (19)$$

Model Training `keras`, a Python library developed by Chollet et al. (2015), is used for model building. Each input data was a list of lagged hours of bike rental count and the output data for that particular input was the bike demand for the next hour. One-hour, 12-hour, and 24-hour lagged timesteps are used for the input data. The processed dataset was partitioned into train and test set where about 86% of the total data was used for model training and the remaining 14% was used for testing.

The following hyperparameters must be tuned:

- **nodes(neurons)**: units accepting a vector of real-valued inputs and producing a single real-valued output.
- **batch size**: how many obs. are used at each step, converge faster with a larger value.
- **epoch**: an iteration over all training obs.
- **dropout**: regularization method to prevent over-fitting
- **learning rate**: define how quickly the network updates its parameters

A summary of the LSTM model setup is below.

- **Prediction Target**: Bike rental demand forecasting in the city of Seoul for the next 24 hours.
- **Input Variables**
 - Observed daily bike rental count and weather condition data during 12/01/2017 - 11/30/2018O

- Lagged values of all variables for the previous 1 hour, 12 hours, and 24 hours
- **Training Parameters**
 - Learning Rate: 0.001
 - Number of nodes: 50
 - Number of epochs: 60
 - Batch size: 24
 - Dropout: 0.2

Forecast performance is assessed using Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and the learning curve. A good fit is identified by a training and validation loss that decreases to a point of stability with a minimal gap between the two final loss values.

Lag 1 Hour Model Model Architecture:

- 1) LSTM with 64 neurons in the first visible layer
- 2) LSTM with 32 neurons in the second visible layer
- 3) Dropout 50%
- 4) 1 neuron in the output layer for predicting Bike Rental Demand.
- 5) The input shape will be 1 time step with 14 features.
- 6) Used the Mean Absolute Error (MAE) loss function and the efficient Adam version of stochastic gradient descent.
- 7) The model will be fit for 100 training epochs with a batch size of 50.

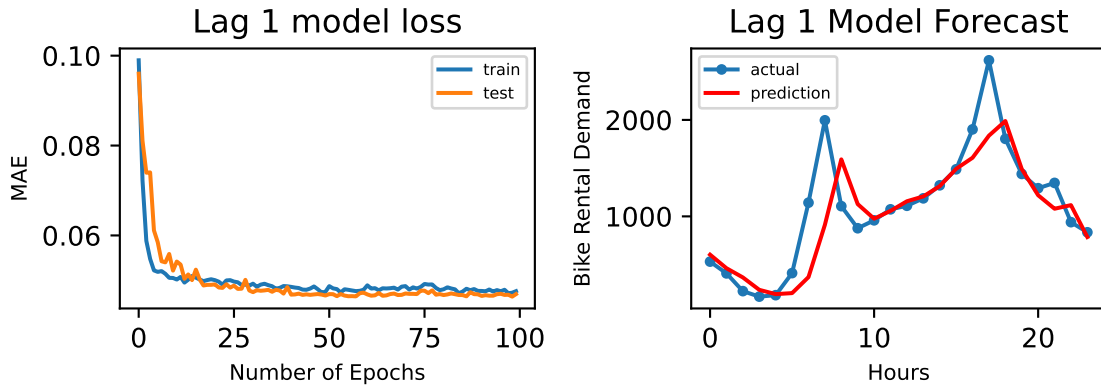


Figure 15: LSTM Lag 1 Model Loss and Forecasted Bike Rental Demand of one day

Although in Figure 15 we can see both errors converge fairly fast, the proposed network is not a good fit since there is a gap between the training error and test error.

Lag 12 Hours Model Model Architecture:

- 1) LSTM with 64 neurons in the first visible layer
- 2) LSTM with 32 neurons in the second visible layer
- 3) Dropout 50%
- 4) 1 neuron in the output layer for predicting Bike Rental Demand.
- 5) The input shape will be 12 time step with 14 features.
- 6) Used the Mean Absolute Error (MAE) loss function and the efficient Adam version of stochastic gradient descent.
- 7) The model will be fit for 100 training epochs with a batch size of 50.

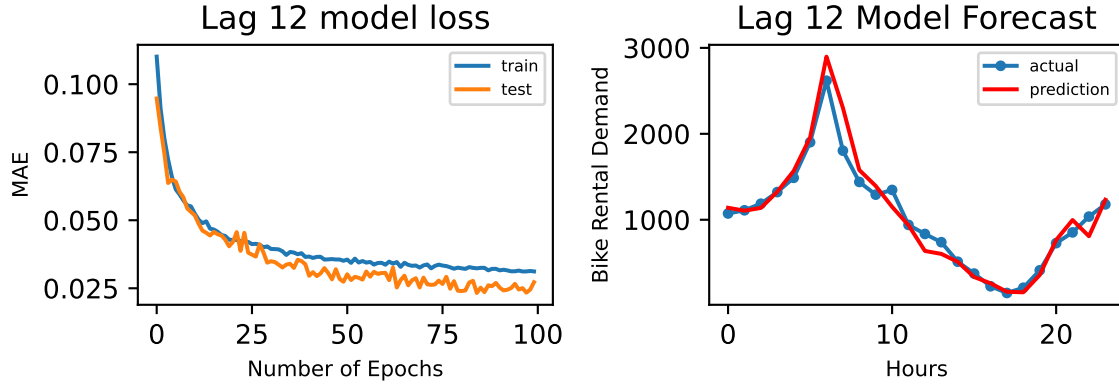


Figure 16: LSTM Lag 12 Model Loss and Forecasted Bike Rental Demand of one day

The Loss vs Epoch curve is shown in Figure 16 in which the progress of mode while training is represented. Both training and testing loss decreases in a smooth fashion fairly quickly.

Lag 24 Hours Model Model Architecture:

- 1) LSTM with 64 neurons in the first visible layer
- 2) LSTM with 32 neurons in the second visible layer
- 3) Dropout 50%
- 4) 1 neuron in the output layer for predicting Bike Rental Demand.
- 5) The input shape will be 24 time step with 14 features.
- 6) Used the Mean Absolute Error (MAE) loss function and the efficient Adam version of stochastic gradient descent.
- 7) The model will be fit for 100 training epochs with a batch size of 50.

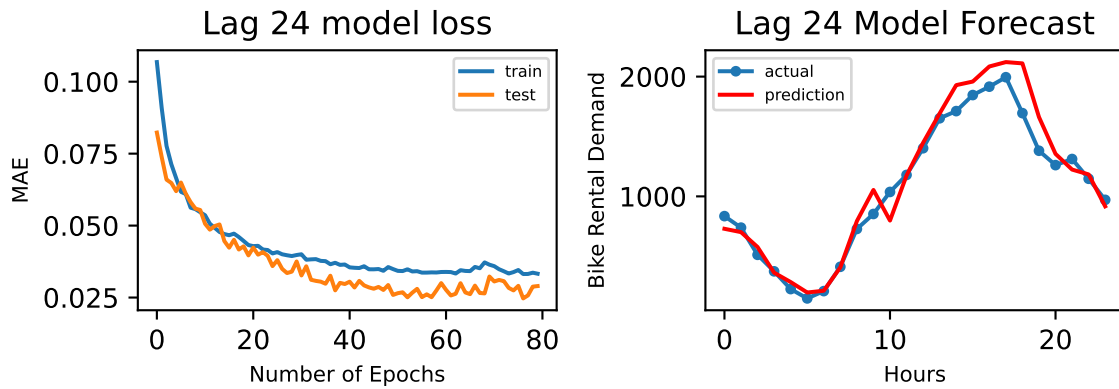


Figure 17: LSTM Lag 24 Model Loss and Forecasted Bike Rental Demand of one day

Evaluation

Error

Given our stated use for these models is forecasting (prediction), when discussing model evaluation we must first define the notion of forecast (prediction) error rate. Hyndman and Koehler (2006) define a number of

different metrics to use for forecast error, each with their own benefits and drawbacks. We will use Mean Absolute Error (MAE) to select the best model from among the list of candidate models. Unlike percent errors, which have the general form $100 \times \frac{\hat{e}}{y_t}$, MAE is defined when $y_t = 0$. Further, MAE is on the same scale as the original dataset (number of bikes per hour), which gives it nice properties of interpretability. We will also report Root Mean Squared Error (RMSE), as it is a common error metric. According to Hyndman and Athanasopoulos (2021), using MAE will result in forecasts of the median while RMSE results in forecasts of the mean. Because our data appear skewed, the median is a more appropriate measure of center. Therefore, MAE is the best error metric for selecting the best forecasting model and estimating test error.

For the purposes of discussing model error rate, let y_t be the observed bike count in period t , \hat{y}_t be the forecast bike count in period t . Then,

$$\hat{e}_t = y_t - \hat{y}_t \quad (20)$$

$$\text{Mean Absolute Error} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m |e_{ij}| \quad (21)$$

$$\text{Root Mean Squared Error} = \sqrt{\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m e_{ij}^2} \quad (22)$$

Where $m = 24$ is the number of periods for which bike count is forecast (the horizon) and $n = 50$ is the number of iterations (folds) in the time series cross-validation.

Model Selection

Models will be selected and evaluated in the context of their use for forecasting rather than statistical inference. As such, we will use time series cross-validation (TSCV). Hyndman and Athanasopoulos (2021) define TSCV as selecting a point or series of points from the dataset as test sets, then selecting all prior points as the training set.

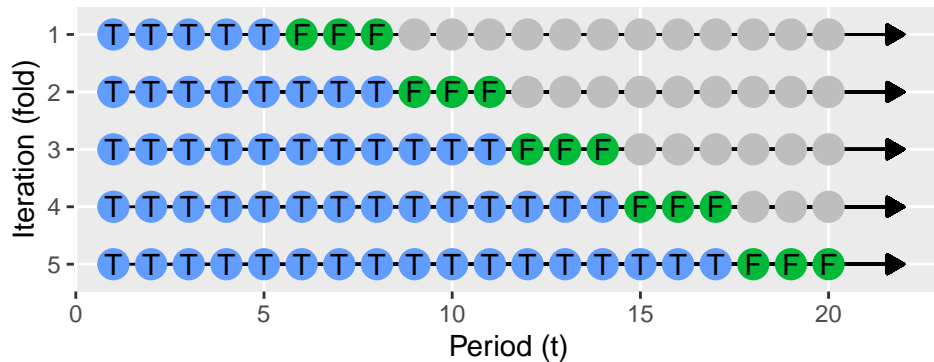


Figure 18: Visualization of time series cross-validation example

Using the example in Figure 18, the first iteration trains on the first 5 observations (blue) and generates forecasts on the next 3 observations (green). In the second iteration the model trains on the first $5 + 3 = 8$ observations and forecasts on the next 3 observations. This continues until the final iteration, where the model trains on all but the last 3 observations and then generates forecasts for the final 3 observations. In this example 5 is the initialization value (the number of training observations in the first iteration), 3 is the step size (the number of observations that are added to the training set each iteration), and 3 is also the horizon (the number of periods for which we generate a forecast). Each of the forecasts is then compared with actual (observed) values to evaluate the forecast.

Our evaluation will use the last 50 days of data as the test folds, with each test fold having 24 observations. The first iteration will train on the 7560 observations from 2017-12-01 00:00:00 to 2018-10-11 23:00:00 and test (forecast) on the 24 observations from 2018-10-12 00:00:00 to 2018-10-12 23:00:00. The second iteration will train on the 7584 observations from 2017-12-01 00:00:00 to 2018-10-12 23:00:00 and test (forecast) on the 24 observations from 2018-10-13 00:00:00 to 2018-10-13 23:00:00. This continues until the final iteration where it trains on the 8736 observations from 2017-12-01 00:00:00 to 2018-11-29 23:00:00 and tests (forecasts) on the 24 observations from 2018-11-30 00:00:00 to 2018-11-30 23:00:00. This means our total sample size for calculating the error metrics for each model will be $24 \times 50 = 1200$. Table 1 depicts the time periods used and sample sizes for a sample of the folds.

Table 2: Train and test set (fold) details for the time series cross-validation process used to calculate estimated test error (MAE and RMSE).

fold	Train Start	Train End	Train n	Test Start	Test End	Test n
1	2017-12-01 00:00	2018-10-11 23:00	7560	2018-10-12 00:00	2018-10-12 23:00	24
2	2017-12-01 00:00	2018-10-12 23:00	7584	2018-10-13 00:00	2018-10-13 23:00	24
3	2017-12-01 00:00	2018-10-13 23:00	7608	2018-10-14 00:00	2018-10-14 23:00	24
...
24	2017-12-01 00:00	2018-11-03 23:00	8112	2018-11-04 00:00	2018-11-04 23:00	24
25	2017-12-01 00:00	2018-11-04 23:00	8136	2018-11-05 00:00	2018-11-05 23:00	24
...
50	2017-12-01 00:00	2018-11-29 23:00	8736	2018-11-30 00:00	2018-11-30 23:00	24

Results

A common-sense, non-machine learning baseline serves as a sanity check and is often used to establish a baseline of comparison for more advanced machine learning models.

In this case, bike demand can be assumed to be season with a daily (24 hours) period. Given hourly data, a common-sense baseline is to predict the bike demand at time t to be equal to the bike demand at time $t - 24$, the same hour the previous day. In the time series paradigm this is known as a Seasonal NAIVE (SNAIVE) model.

$$\hat{x}_t = x_{t-24} \quad (23)$$

Table 3: Test Error by Model Type

Model	MAE	RMSE
Baseline		
Baseline	190	330
Time Series Methods		
SARIMAX (1,0,1) (5,1,0)	201	300
fasster	188	286
Prophet	169	242
Tree Methods		
Random Forest	165	232
XGBoost	129	172
LSTM Models		
LSTM Lag 1	167	271
LSTM Lag 12	97	149
LSTM Lag 24	103	148

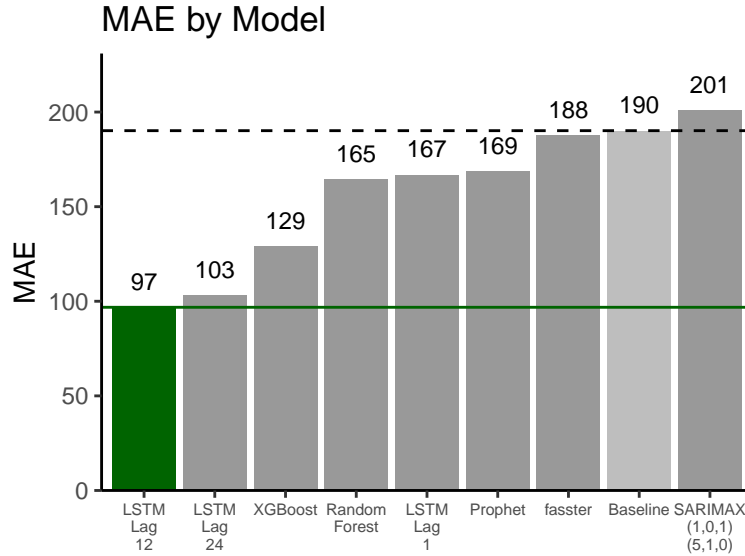


Figure 19: Mean Absolute Error by Model

Table 3 and Figure 19 report the MAE and RMSE for each model. Perhaps unsurprisingly, the machine learning methods performed better than the time series methods, on average. The LSTM Lag 12 model performed the best (smallest MAE), with $MAE = 97$ bikes demanded per hour. This was followed by LSTM Lag 24 and XGBoost. The baseline SNAIVE model outperformed both ARIMA-based models, which further reinforces the difficulty that ARIMA methods have with such high-frequency data.

Conclusion

The focus of this analysis has been on the ability to accurately forecast hourly bike demand for the next 24 hours using data from Seoul, South Korea’s bike sharing program. This analysis shows that a LSTM Lag 12 generates the most accurate forecast of demand, as measured by mean absolute error, with $MAE = 97$ bikes per hour. Assuming that this error rate is within the operational constraints of the bikeshare management organization (i.e. they can still make bike allocation decisions even with the given error rate), this model

could be implemented immediately to help ensure an adequate number of bikes are available over the next 24 hours.

As with most machine learning models, future research into this application could focus on the inclusion of additional predictor variables and more recent bike usage data. One key deficit of the models presented in this analysis are that they do not account for spatial variation in demand. That is, what is the geographic distribution of bike demand at each hour. Understanding not just the total number of bikes demanded but also *where* in the city the bikes are needed would allow the administrators of the bikeshare to reallocate available bikes from areas of low demand to areas of high demand. Adding a geospatial dimension to the demand would certainly increase the complexity of the prediction because it would require either numerous individual models (one for each geography) or a single model with a multivariate response that is a vector of the same length as the number of “locations” within the city where demand is captured. Although this would be a more complex problem, it would also be of significant more value to the bikeshare administrators, making it a strong candidate for subsequent research.

References

- Prashant Banerjee. A guide on xgboost hyperparameters tuning. URL <https://www.kaggle.com/code/prashant111/a-guide-on-xgboost-hyperparameters-tuning/notebook>.
- G. E. P. Box and D. R. Cox. An analysis of transformations. *Journal of the Royal Statistical Society. Series B (Methodological)*, 26(2):211–252, 1964. ISSN 00359246. URL <http://www.jstor.org/stable/2984418>.
- Jason Brownlee. How to use standardscaler and minmaxscaler transforms in python, 06 2020. URL <https://machinelearningmastery.com/standardscaler-and-minmaxscaler-transforms-in-python/>.
- François Chollet et al. keras, 2015.
- Sathishkumar V E and Yongyun Cho. A rule-based model for seoul bike sharing demand prediction using weather data. *European Journal of Remote Sensing*, 53(sup1):166–183, 2020. doi: 10.1080/22797254.2020.1725789. URL <https://doi.org/10.1080/22797254.2020.1725789>.
- Manuel Fernández-Delgado, Eva Cernadas, Senén Barro, and Dinani Amorim. Do we need hundreds of classifiers to solve real world classification problems? *Journal of Machine Learning Research*, 15(90):3133–3181, 2014. URL <http://jmlr.org/papers/v15/delgado14a.html>.
- Chang Gao and Yong Chen. Using machine learning methods to predict demand for bike sharing. In Jason L. Stienmetz, Berta Ferrer-Rosell, and David Massimo, editors, *Information and Communication Technologies in Tourism 2022*, pages 282–296, Cham, 2022. Springer International Publishing. ISBN 978-3-030-94751-4.
- Victor M. Guerrero. Time-series analysis supported by power transformations. *Journal of Forecasting*, 12(1):37–48, 1993. doi: <https://doi.org/10.1002/for.3980120104>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/for.3980120104>.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9:1735–1780, 11 1997. doi: 10.1162/neco.1997.9.8.1735.
- Rob Hyndman (<https://stats.stackexchange.com/users/159/rob-hyndman>). Training arima model on data with missing values. Cross Validated. URL <https://stats.stackexchange.com/q/267264>. URL:<https://stats.stackexchange.com/q/267264> (version: 2017-03-13).
- Rob Hyndman and G. Athanasopoulos. *Forecasting: Principles and Practice*. OTexts, Melbourne, Australia, 3rd edition, 2021. URL [OTexts.com/fpp3](https://otexts.com/fpp3).
- Rob J. Hyndman and Yeasmin Khandakar. Automatic time series forecasting: The forecast package for r. *Journal of Statistical Software*, 27(3):1–22, 2008. doi: 10.18637/jss.v027.i03. URL <https://www.jstatsoft.org/index.php/jss/article/view/v027i03>.
- Rob J. Hyndman and Anne B. Koehler. Another look at measures of forecast accuracy. *International Journal of Forecasting*, 22(4):679–688, 2006. ISSN 0169-2070. doi: <https://doi.org/10.1016/j.ijforecast.2006.03.001>. URL <https://www.sciencedirect.com/science/article/pii/S0169207006000239>.
- Denis Kwiatkowski, Peter C.B. Phillips, Peter Schmidt, and Yongcheol Shin. Testing the null hypothesis of stationarity against the alternative of a unit root: How sure are we that economic time series have a unit root? *Journal of Econometrics*, 54(1):159–178, 1992. ISSN 0304-4076. doi: [https://doi.org/10.1016/0304-4076\(92\)90104-Y](https://doi.org/10.1016/0304-4076(92)90104-Y). URL <https://www.sciencedirect.com/science/article/pii/030440769290104Y>.
- G. M. Ljung and G. E. P. Box. On a measure of lack of fit in time series models. *Biometrika*, 65(2):297–303, 08 1978. ISSN 0006-3444. doi: 10.1093/biomet/65.2.297. URL <https://doi.org/10.1093/biomet/65.2.297>.
- Mitchell O’Hara-Wild. *fable.prophet: Prophet Modelling Interface for ‘fable’*, 2020. URL <https://fable.tidyverse.org>. R package version 0.1.0.
- Mitchell O’Hara-Wild and Rob Hyndman. *fasster: Fast Additive Switching of Seasonality, Trend and Exogenous Regressors*, 2022. URL <https://github.com/mitchelloharawild/fasster>. R package version 0.1.0.9100.

- Mitchell O’Hara-Wild, Rob Hyndman, and Earo Wang. *fable: Forecasting Models for Tidy Time Series*, 2021. URL <https://fable.tidyverts.org>. R package version 0.3.1.
- Philipp Probst, Marvin N. Wright, and Anne-Laure Boulesteix. Hyperparameters and tuning strategies for random forest. *WIREs Data Mining and Knowledge Discovery*, 9(3):e1301, 2019. doi: <https://doi.org/10.1002/widm.1301>. URL <https://wires.onlinelibrary.wiley.com/doi/abs/10.1002/widm.1301>.
- J. Schuijbroek, R.C. Hampshire, and W.-J. van Hove. Inventory rebalancing and vehicle routing in bike sharing systems. *European Journal of Operational Research*, 257(3):992–1004, 2017. URL <https://EconPapers.repec.org/RePEc:eee:ejores:v:257:y:2017:i:3:p:992-1004>.
- R.H. Shumway and D.S. Stoffer. *Time Series: A Data Analysis Approach Using R*. A Chapman & Hall book. CRC Press, Taylor & Francis Group, 2019. ISBN 9780367221096.
- Sean J. Taylor and Benjamin Letham. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018. doi: 10.1080/00031305.2017.1380080. URL <https://doi.org/10.1080/00031305.2017.1380080>.
- UCI. Seoul bike sharing demand data set, Mar 2020. URL <https://archive.ics.uci.edu/ml/datasets/Seoul+Bike+Sharing+Demand>.