

Web Traffic Forecasting

Time Series Analysis

Shree Karimkolikuzhiyil; Programming; Statistics, M.S. (Distance); shreejesh@tamu.edu

Jingcheng Xia; Computations; Computer Science, B.S. (On-Campus); sixtyfour64@tamu.edu

Jackson Smith; Analysis; Statistics, M.S. (Distance); jackson.t.smith@tamu.edu

Samuel Burge; Writing; Statistics, M.S. (Distance); samuelburge@tamu.edu

Max Kutschinski; Theory; Statistics, M.S. (Distance); mwk556@tamu.edu

Introduction and Motivation

The objective of this analysis is to forecast daily unique visitors to an academic website over a 30-day horizon. Predicting website traffic allows IT departments to manage project throughput and prioritize maintenance and enhancements to website functionality and effectively allocate web server resources. Web traffic is also a key indicator of customer growth and expansion, as well as sustaining recurring customers and ingrained growth. The details provided by web traffic throughput reports contain many metrics, including page loads, returning visitors, and unique visits, each of which conveys a different picture and set of information for an organization. As well, having a picture of expected throughput and confirming (or denying) expectations with reality allows a business to understand unexpected growth and/or unexpected decay in business development.

The data contains five years of daily time series data of user visits. There are four features in the data set, which include daily counts for the number of page loads, first-time visitors, returning visitors, and unique visitors.¹ An initial plot of the data shows strong seasonality and volatility, but doesn't appear to have any discernible trend or cyclical behavior. An explanation for this could be due to the nature of the website. Students would likely be the largest share of users for a website of this nature, and the seasonality seems associated with the academic calendar typically seen at academic institutions.

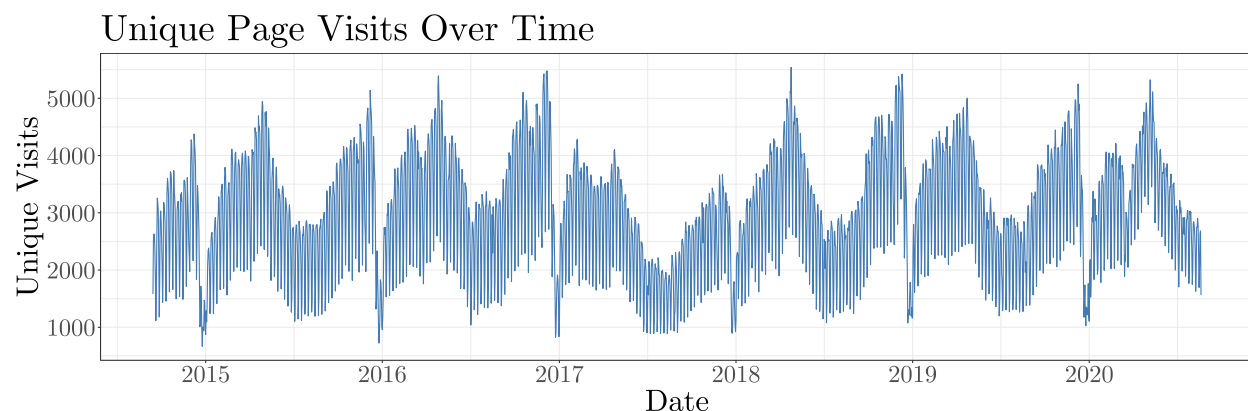


Figure 1

¹A visit is defined as a stream of hits on one or more pages on the site on a given day by the same user within a 6-hour window, identified by the IP address of the specific device. Returning visitors are identified through allowed cookies on a user's device, and the total number of returning and first-time visitors is, by definition, the number of unique visitors.

Modeling

SARIMA

Stationarity is a common assumption underlying many time series procedures. As such, it is important to assess the level of stationarity prior to modeling and make the appropriate adjustments if necessary.

Shumway and Stoffer [2019] describe a stationary time series as one whose properties do not depend on the time at which the series is observed. More specifically,

- (i) *the mean value function $\mu_t = E(x_t)$ is constant and does not depend on time t*
- (ii) *the autocovariance function $\gamma(s, t) = \text{cov}(x_s, x_t) = E[(x_s - \mu_s)(x_t - \mu_t)]$ depends on times s and t only through their lagged difference.*

The strong seasonality that is apparent in Figure 1 is indicative of non-stationarity, since seasonality will affect the value of the time series at different times. Seasonality is defined as a recurring pattern at a fixed and known frequency based on the time of the year, week, or day.

Figures 2 and 3 aim to identify the types of seasonality present in the data. Figure 2 plots a subset of the first several weeks and indicates that there exists weekly seasonality, whereas Figure 3 uses locally weighted scatterplot smoothers (Lowess) to emphasize the inherent yearly seasonality.

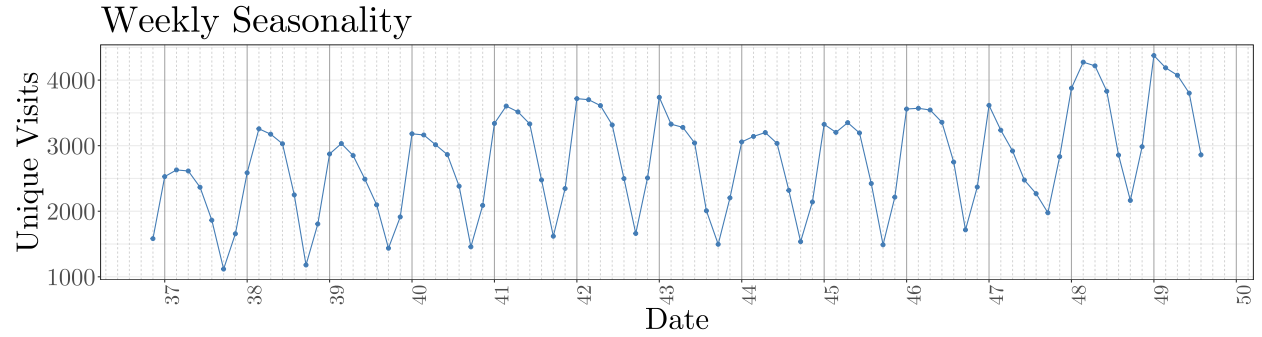


Figure 2: Sample of weekly page visits

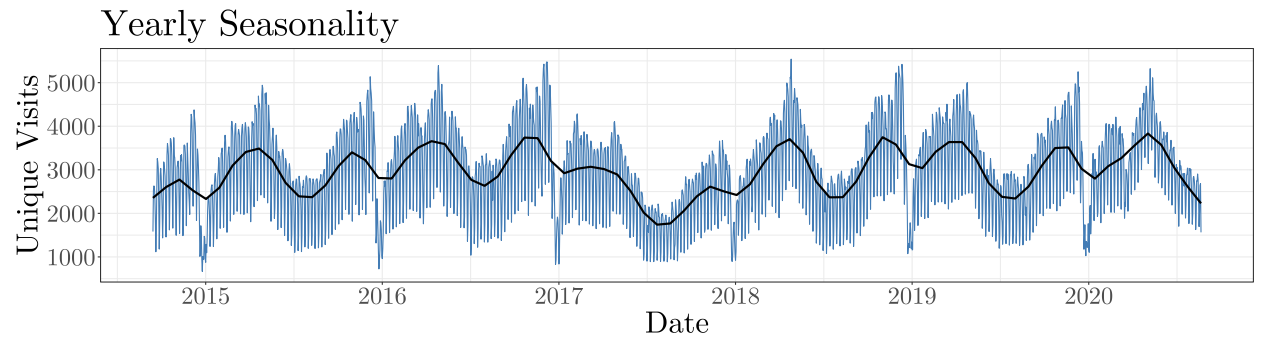


Figure 3: Smoothing via Lowess

A popular approach in addressing non-stationarity due to seasonality is to eliminate these effects via seasonal differencing. The seasonal difference of a time series is the series of changes from one season to the next, which is defined as follows:

$$\nabla x_t = x_t - x_{t-m} \quad (1)$$

Hence, yearly and weekly seasonality will be handled by computing the lag 365 and lag 7 seasonal difference, respectively.

$$\nabla x_t^* = (x_t - x_{t-365}) - x_{t-7} \quad (2)$$

Time plots of our series at different levels of differencing are displayed in Figure 4. The differenced series appears to be stationary with constant mean and variance.

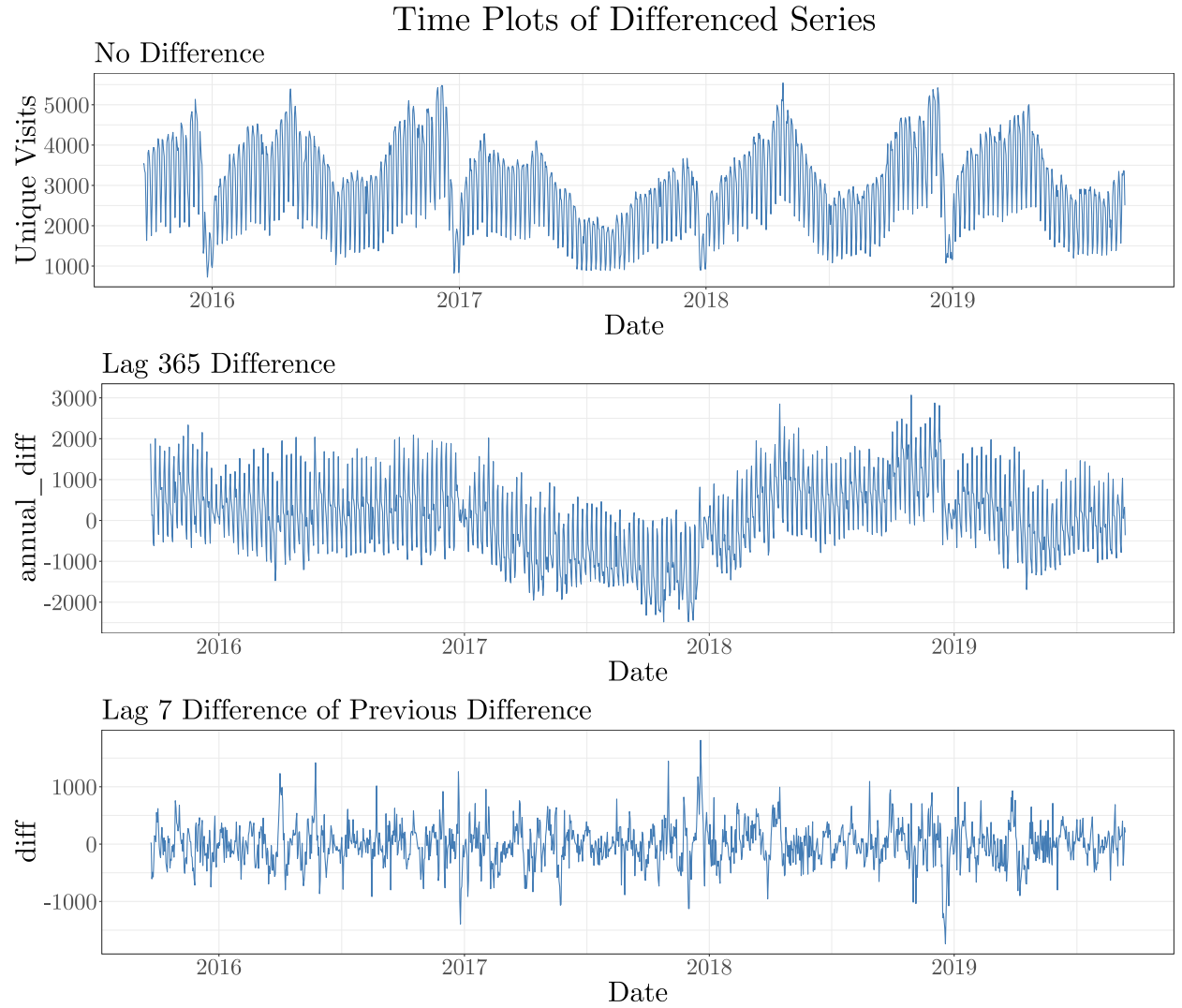


Figure 4: Time plots of differenced series

The ACF and PACF of the differenced series ∇x_t^* are displayed in Figure 5. Neither the ACF nor the PACF seems to cut off after a certain lag, which would be indicative of an AR or MA process. Rather, both of them appear to tail off over time.

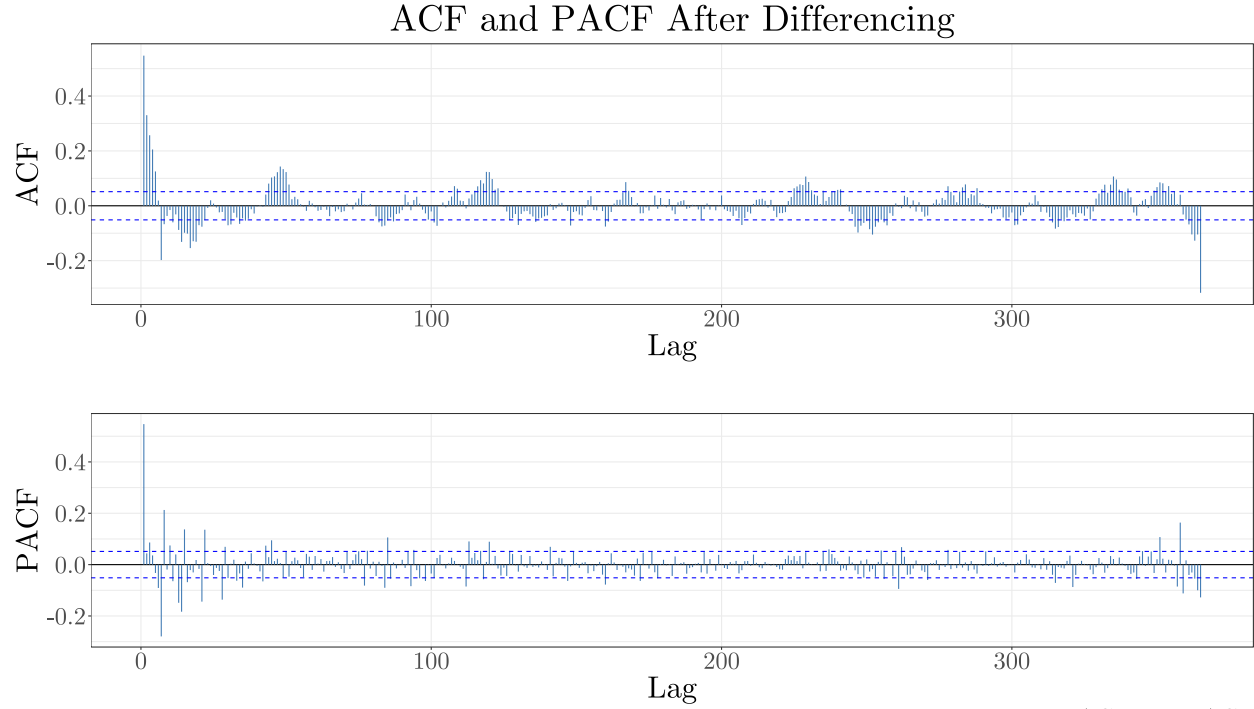


Figure 5: ACF and PACF

Both ACF and PACF show a slow decay making it difficult to determine specific orders for the family of ARMA(p,q) models defined in Eq(3) as:

$$x_t = \alpha + \phi_1 x_{t-1} + \cdots + \phi_p x_{t-p} + w_t + \theta_1 w_{t-1} + \cdots + \theta_q w_{t-q} \quad (3)$$

where $\phi_p \neq 0, \theta_q \neq 0, \sigma_w^2 > 0$, and the model is causal and invertible.

We opted to fit a range of ARMA(p,q) models with small orders, with the final model selected based on the accuracy of the model forecasts test or hold-out set since our primary goal for this analysis is forecasting (not necessarily inferences or hypothesis testing). The model selection criterion for the fitted ARMA models is shown in the table below.

```
## # A tibble: 4 x 6
##   .model    sigma2 log_lik    AIC    AICc    BIC
##   <chr>      <dbl>  <dbl>  <dbl>  <dbl>  <dbl>
## 1 ARMA(1,2) 66504. -10138. 20287. 20287. 20319.
## 2 ARMA(1,1) 66465. -10137. 20287. 20287. 20324.
## 3 ARMA(2,1) 67068. -10144. 20299. 20299. 20331.
## 4 ARMA(2,2) 90509. -10359. 20727. 20727. 20754.
```

The parameters of the ARMA(1,2) model are estimated via conditional least squares and are displayed in Eq(4).

$$\hat{x}_t = 0.44x_{t-1} + 0.13x_{t-2} + w_t + 0.05\omega_{t-1} - 0.09\omega_{t-2} \quad (4)$$

Figure 6 displays of plot of the residuals of the fitted ARMA(1,2) model. Initially, the residuals seem to behave like white noise, being centered around zero with a constant variance. However, further analysis of autocorrelation plot and formal testing using the Box-Ljung test indicate that the innovations are not uncorrelated (i.e., they are not white noise). The autocorrelation in the innovations does appear small for most lags, given how similar the various models are this is likely the best fit we can obtain from the ARMA(p,q) modelling procedure.

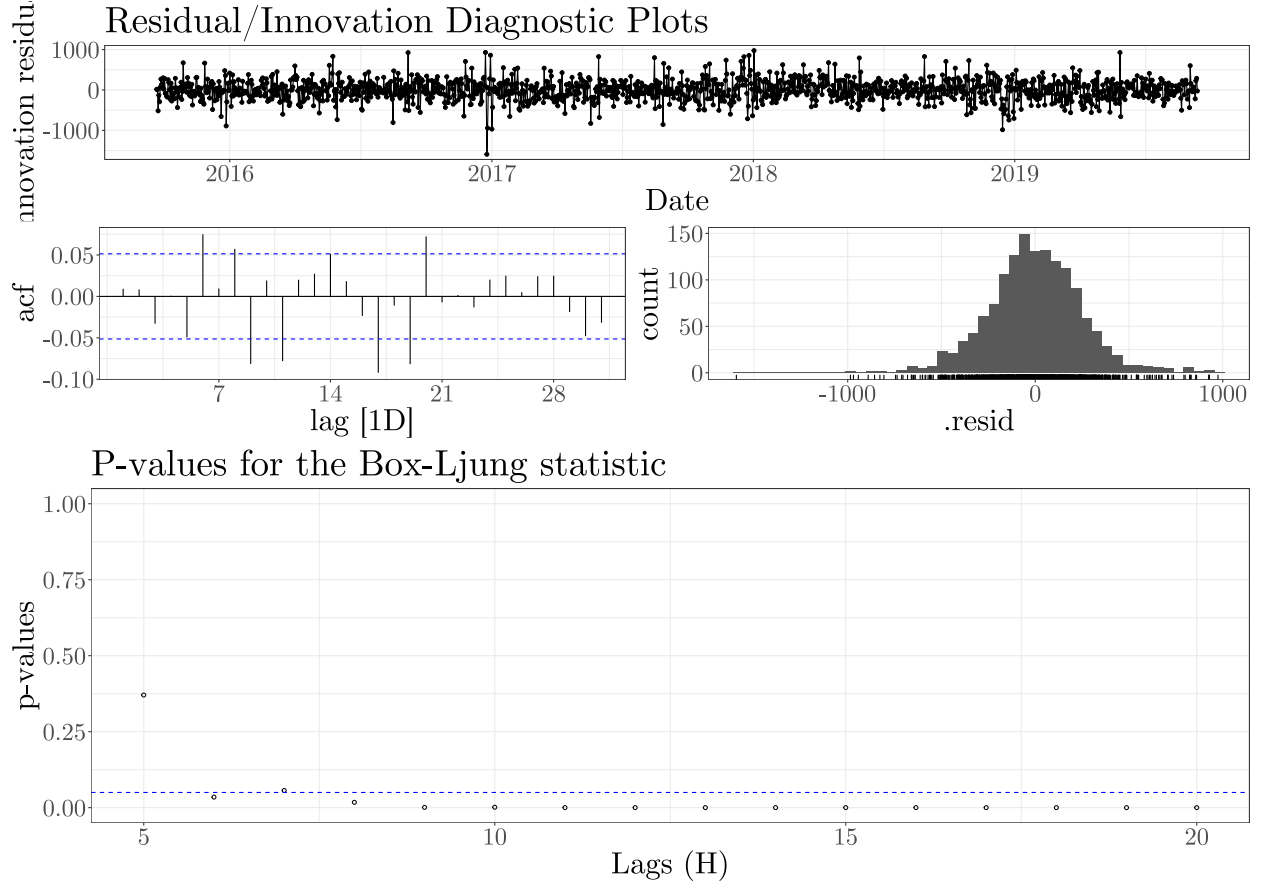


Table 1 contains the model's polynomial roots. Since they appear to be different from each other by a reasonable margin, we can conclude that there is no parameter redundancy in the model.

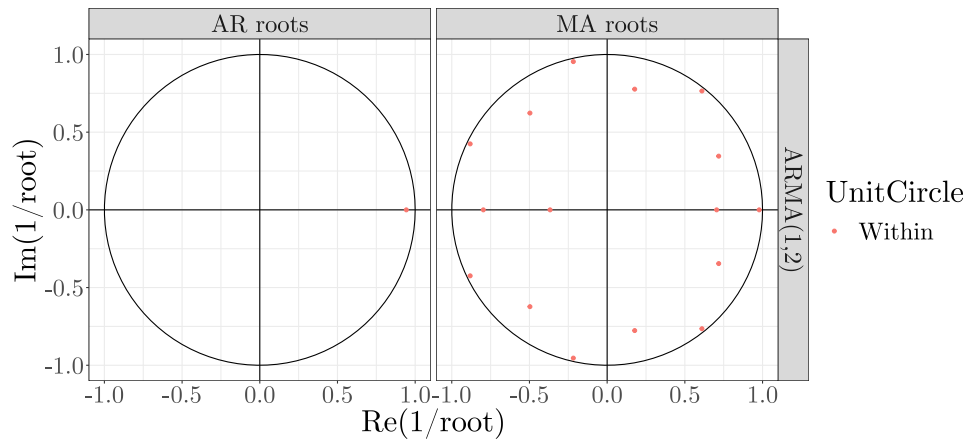
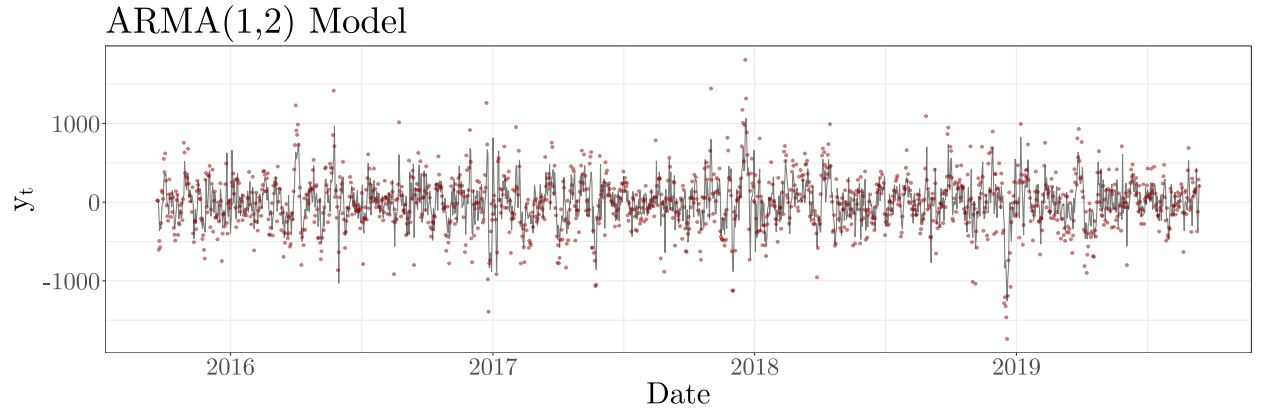


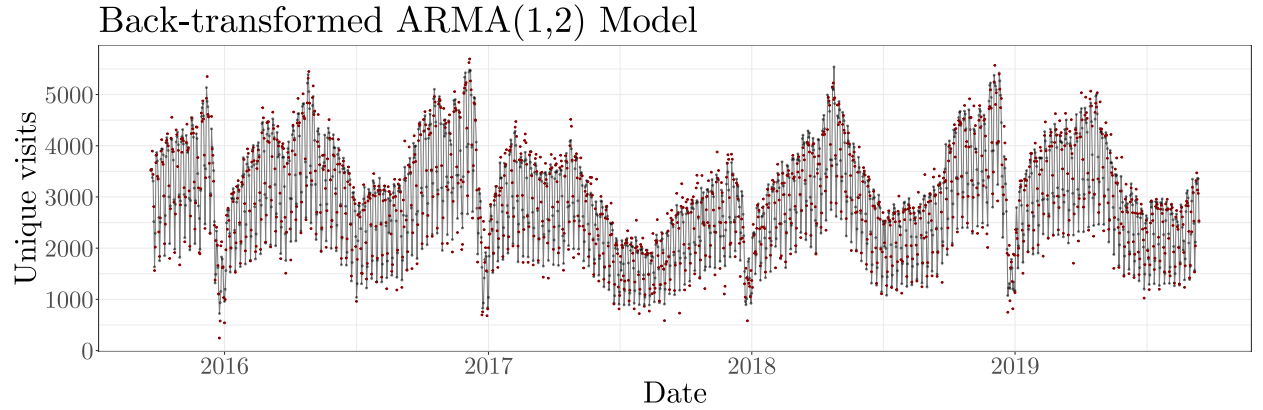
Table 1: Roots of polynomials

AR	MA
1.06056372	0.6372060+0.7990310i, -0.9207893+0.4434288i, -0.2274162-0.9963754i, 0.6372060-0.7990310i,
	-0.2274162+0.9963754i, -0.9207893-0.4434288i, -0.7824690-0.9811850i,
	1.0219991+0.0000000i, 0.2792599+1.2235177i, -1.2549828-0.0000000i, 0.2792599-1.2235177i,
	1.1307004+0.5445166i, -0.7824690+0.9811850i, 1.1307004-0.5445166i, 1.4191005-0.0000000i,
	-2.7228424-0.0000000i

The fitted values of the ARMA(1,2) model are plotted against the actual values of the seasonally differenced series ∇x_t^* (Figure 7) and the.



Fitted model shown by black line and the actual values of the training set are shown in red.



Fitted model shown by black line and the actual values of the training set are shown in red.

Prophet

Taylor and Letham [2018]

O'Hara-Wild [2020]

Results

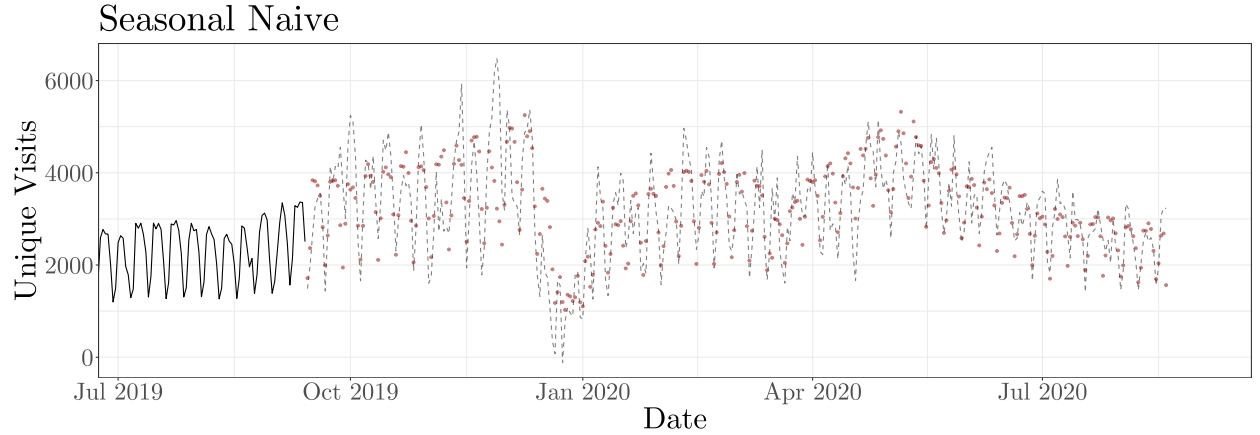
All the above models were trained on a training set, and the predictive accuracy was then evaluated on a test set. The training set consists of page visits starting from 2015-09-21 until 2019-09-13, and test set contains the data from 2019-09-14 to 2020-08-19, making up the last 341 observations of the data. Our primary evaluation metrics for model comparison are the root mean squared error (RMSE) and mean absolute error (MAE), which both have the advantage of being measured on the same scale as the data (i.e., the number of website visits). Using these are our primary accuracy measures gives us more interpretable results. For clarification, the RMSE and MAE are defined as

$$RMSE = \sqrt{\left(\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2\right)}$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

In addition to these accuracy measures, a common-sense baseline serves as a sanity check, and is often used as a benchmark for more advanced time series models. Given daily data with yearly seasonality, a common-sense baseline is to predict the number of unique visits at time t to be equal to the number of unique visits at $t-365$. In other words, a random walk model making a constant prediction with yearly seasonality, which is known as a seasonal naive model (Eq. X)

$$\hat{x}_t = x_{t-365}$$



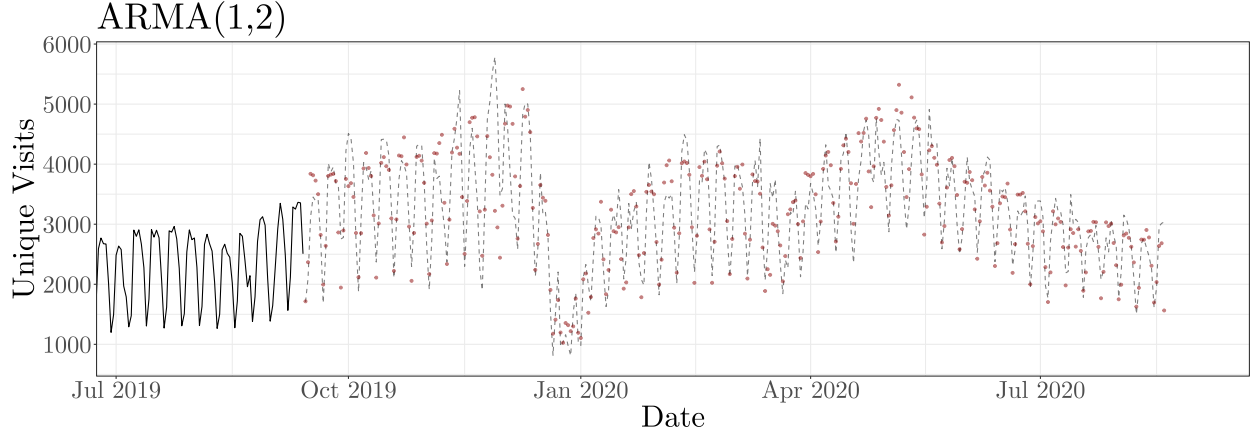
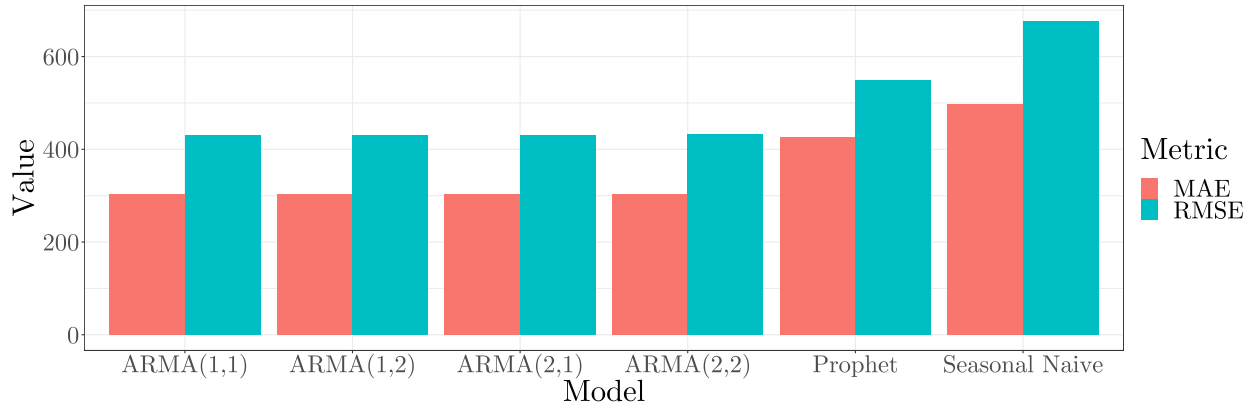


Table 2: Accuracy measures from the test set performance of fitted models.

Model	ME	RMSE	MAE	MPE	MAPE	MASE	RMSSE
ARMA(1,2)	-1.25	430.61	302.12	99.33	100.46	0.71	0.77
ARMA(1,1)	-1.33	430.73	302.23	99.38	100.61	0.71	0.77
ARMA(2,1)	-1.57	430.96	302.51	99.61	100.84	0.71	0.77
ARMA(2,2)	-3.93	432.13	302.79	99.45	99.45	0.71	0.77
Prophet	272.59	548.81	426.81	7.86	15.49	1.52	1.27
Seasonal Naive	9.59	676.70	498.06	53.16	536.76	1.17	1.21



Selected accuracy measures for the various fitted models.

Discuss results here

Conclusion

References

- Mitchell O'Hara-Wild. *fable.prophet: Prophet Modelling Interface for 'fable'*, 2020. URL <https://fable.tidyverts.org>. R package version 0.1.0.
- 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>.