

Springboard--DSC Capstone Project 2

Lease Delinquency Time Series

Garrick Skalski

September 2021

1 INTRODUCTION

Leasing is an alternative method to borrowing for financing business equipment. In terms of the core economics, leases are very similar to loans, but may offer more flexibility, different tax treatment, and lower near-term costs relative to equipment purchased with loans. Similarly to loans, when bank customers are delinquent and then default on their lease payments, bank lenders must record losses and bank profitability is reduced.

The focus of this project is to assess whether changes in lease delinquency rates can be forecast so that bank lenders can more proactively reduce their risk of credit losses. Using time series and regression methods, I analyzed time series on lease delinquency rates together with time series on several economic variables (features) to evaluate whether these economic variables have predictive value for modeling and forecasting lease delinquency. The main business use case is to forecast lease delinquency rate a short time into the future such as one or a few calendar months or quarters.

The key findings are that (i) several of the economic variables are correlated with and predictive of lease delinquency rate, and (ii) a simple linear regression model that makes use of several of the lagged features is the best predictor of lease delinquency rate among the models considered and a simple univariate time series model without features also performs well.

Full details of the analysis with Python code can be found in Jupyter notebooks on GitHub (https://github.com/gskalski267/springboard/blob/main/Capstone_Project_2/notebooks/).

2 APPROACH

The statistical approaches I used closely follow Hyndman and Athanasopoulos 2021 (hereafter, [HA 2021]) and were implemented using Python modules pmdarima, statsmodels, and sklearn.

2.1 Data Acquisition

The target (response) variable is quarterly, seasonally-adjusted lease delinquency rate for all banks from the U.S. Federal Reserve as part of their public data releases. The economic variables used as predictor variables, or features, are seven standard economic series that are publicly available.

The data and sources are:

response variable

'delinq' = Lease Delinquency Rate, all banks; www.federalreserve.gov/datadownload/

predictor variables (= exogenous variables = features)

'ls_rcvbl' = Lease Receivable Balance, all banks; www.federalreserve.gov/datadownload/

'pmi_man' = ISM Purchasing Managers Index (PMI); www.quandl.com/data/ISM/MAN_PMI-PMI-Composite-Index

'cons_sent' = Univ of Michigan Index of Consumer Sentiment; www.sca.isr.umich.edu/tables.html

'close_price' = S&P 500 Price Index; finance.yahoo.com

'stdtght_ci_smll_netpct' = Loan Standard Tightening, net percentage, commercial and industrial loans from small banks; www.federalreserve.gov/datadownload/

'tot_bus_inv' = Total Business Inventories; www.census.gov/economic-indicators/

'ret_sales' = Retail Sales; www.census.gov/economic-indicators/

I standardized all data sets to quarterly series to match the frequency of the lease delinquency rate time series.

2.2 Data Exploration

I explored the data in three main ways: (i) graphs, (ii) time series diagnostics, and (iii) correlation scatterplots. Full details are available on GitHub.

A simple plot of series values versus time helps to understand the basic structure of the data, series for a few of the variables are shown in Figure 1. For lease delinquency rate (variable labeled 'delinq'), there is a clear increase during recessionary economic periods and a clear decrease during expansionary economic periods. The feature variables show similar patterns. For example, PMI, a measure of economic activity, declines near recessionary periods and increases near expansionary periods.

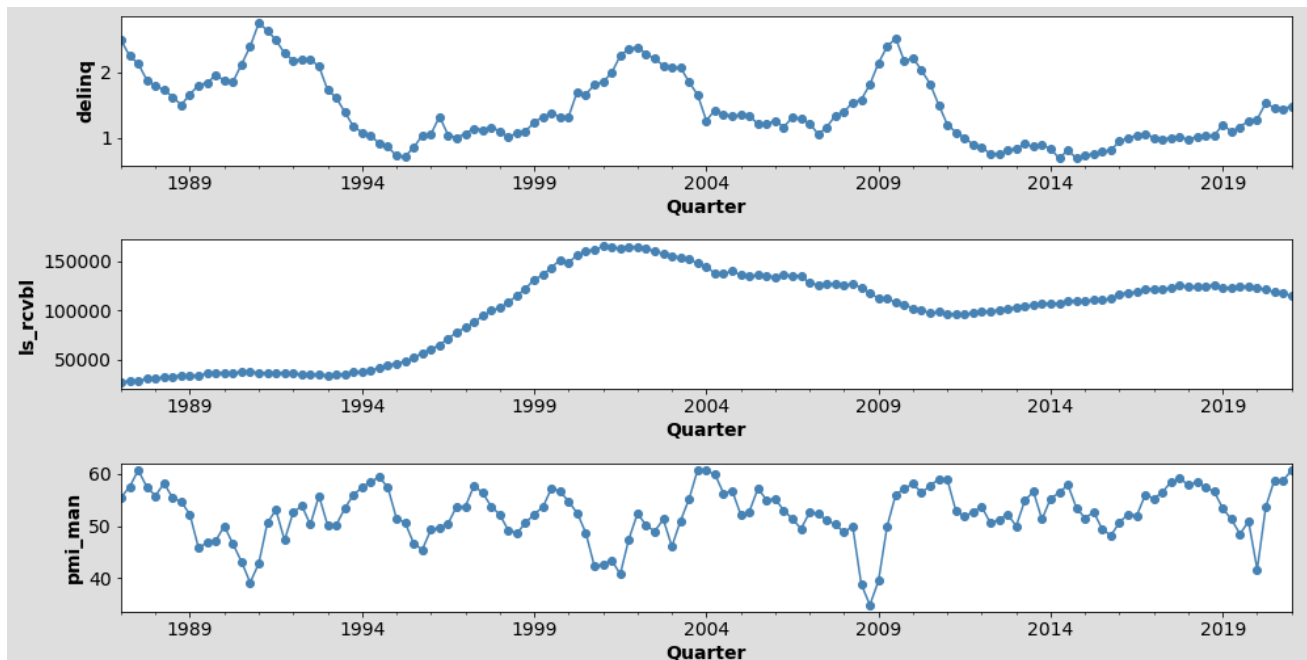


Figure 1

Time series plot of Lease Delinquency Rate ('delinq'), Lease Receivable Balance ('ls_rcvbl'), and ISM Purchasing Managers Index (PMI) ('pmi_man') by calendar quarter.

A standard approach in time series analysis is to transform the data to satisfy the statistical assumptions of time series models such as stationarity [HA 2021]. Accordingly, I natural-log transformed the data and calculated first differences. I then analyzed the transformed data for autocorrelation (via partial autocorrelation) and tested for stationarity using the KPSS test (available in pmdarima). Overall, there is little evidence of strong departure from stationarity except for variable Lease Receivable Balance which shows some trend remaining in the series even after the data were transformed (Figure 2).

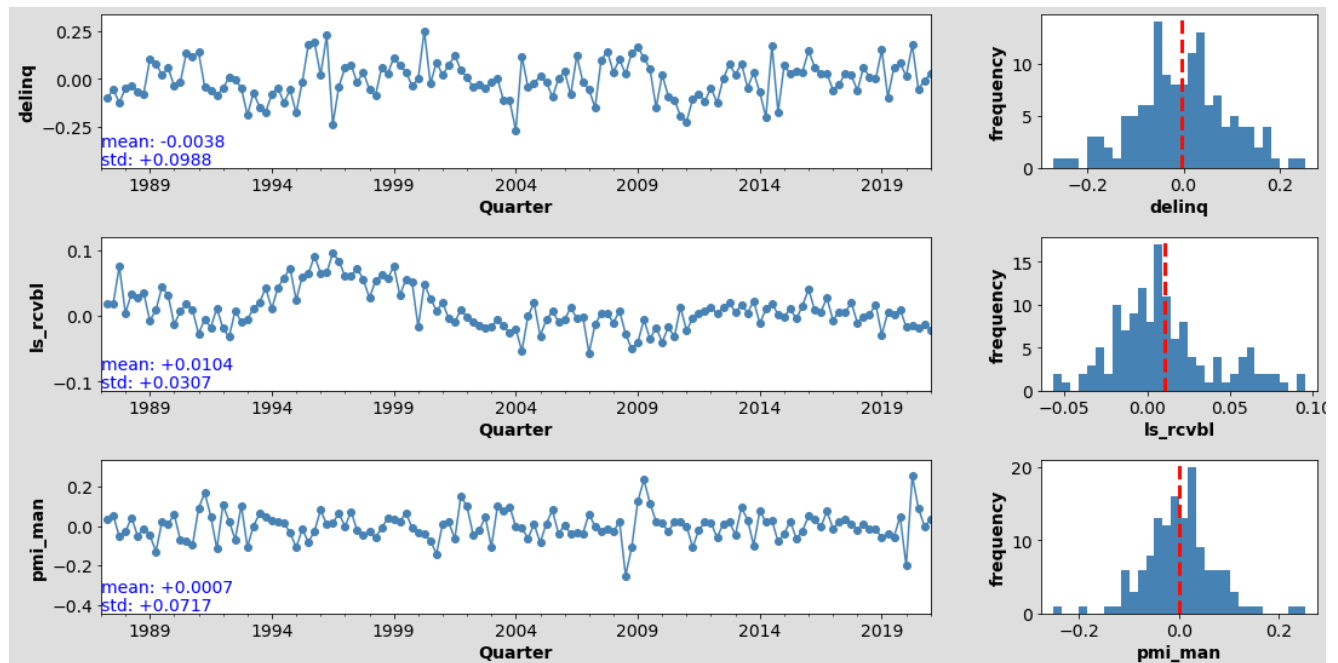


Figure 2

Time series plots and histograms of residuals of natural-log, first differenced transformed variables Lease Delinquency Rate ('delinq'), Lease Receivable Balance ('ls_rcvbl'), and ISM Purchasing Managers Index (PMI) ('pmi_man') by calendar quarter.

Analysis of partial autocorrelation indicates that several variables show some evidence of autocorrelation, with most autocorrelation at shorter lags, but also some cases at longer lags. For example, Lease Delinquency Rate shows autocorrelation at lags 1, 2, and 8 and Lease Receivable Balance shows autocorrelation at lags 1, 2, 3, 4, 9 and 10 (Figure 3).

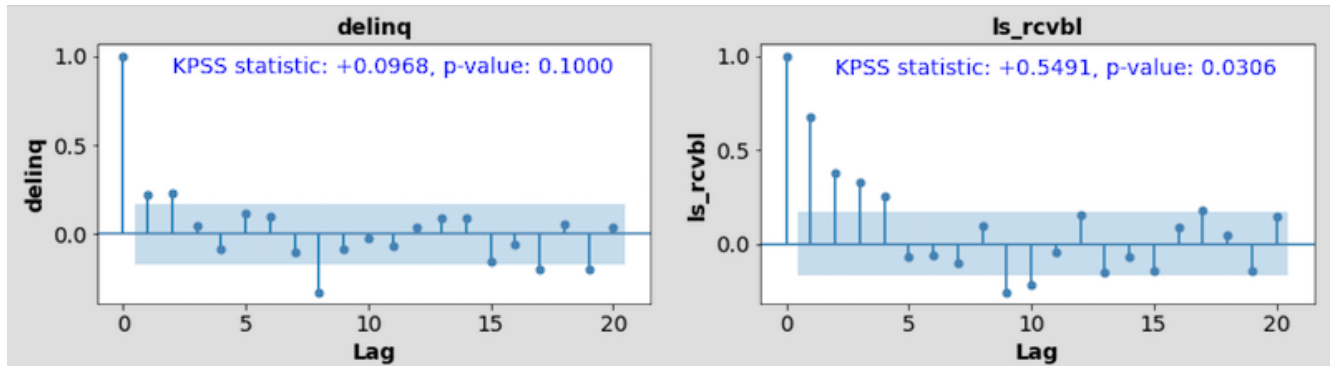


Figure 3

Partial autocorrelation plot of variables Lease Delinquency Rate ('delinq') and Lease Receivable Balance ('ls_rcvbl') by lag.

Finally, using correlations and scatterplots, I analyzed the pairwise relationships among all variables for lags 0 to 10. Several of the feature variables are correlated with lease delinquency rate and each other. In particular, there are correlations between delinquency rate and feature variables with short lags (1, 2, 3) and also for features near lag 8, with correlation coefficients as high as about 0.38 in absolute magnitude. An example of some of the correlation plots is shown in Figure 4.

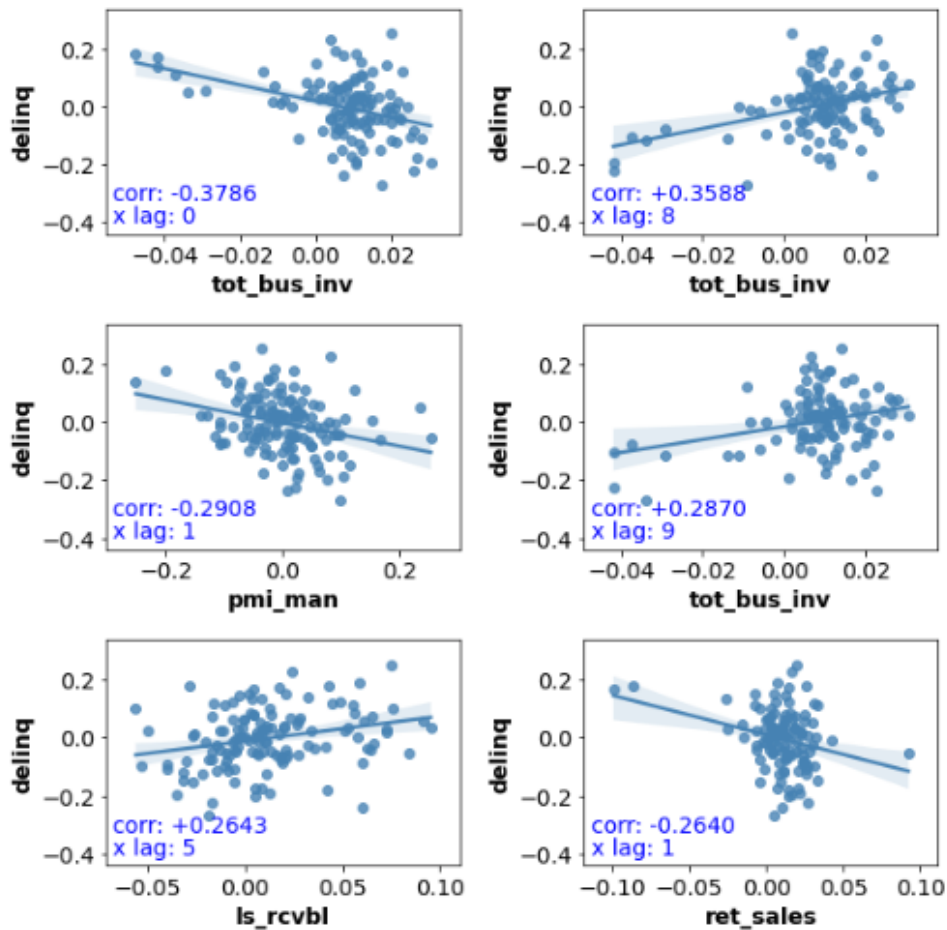


Figure 4

Pairwise scatterplots between lease delinquency rate and several feature variables at different lags.

2.1 Time Series Modeling

Model Fitting Methods

Using the Python time series package pmdarima and drawing on the methods in [HA 2021], I fit a collection of time series models to the lease delinquency data and the seven economic variables data. I defined the following three categories of models:

1. "Baseline ARIMA" Model - an ARIMA(p,1,q) model
2. "Null ARIMA Exog" Model - an ARIMA(0,1,0) model with exogenous variables
3. "Full ARIMA" Model - an ARIMA(p,1,0) with exogenous variables

I fit the models to the training data using cross validation with a rolling forecast window implemented in pmdarima [HA 2021] and evaluated fit using mean squared error (MSE). Because the most immediately useful business case is a short-term forecast, I used the last four quarters in the data ('2020Q1' to '2020Q2') as test data and thus a rolling forecast window of four quarters in cross validation on the training data. Model selection criteria AIC and BIC were also used to help inform model fit [HA 2021]. Because the full parameter space for all of the models is quite large if all combinations of features with various lags are evaluated, I searched a subset of this full parameter space as an initial evaluation of these models. Specifically,

Baseline ARIMA Model

I fit ARIMA(p,1,q) models with lags 0 to 8 for p and q (i.e., 81 different parameter combinations).

Null ARIMA Exog Model

I fit ARIMA(0,1,0) with exogenous variables models using a forward stepwise model selection approach [HA 2021] with eight iterations. In each iteration, I added to the model the feature-lag combination as an exogenous variable that provided the best-fit (smallest MSE) by cross validation. In this way, one feature variable was added to the model for each of the eight iterations, resulting in an eight-feature model.

Full ARIMA Model

for this fuller model, I searched the parameter space based on the results from the prior two models. Specifically, I fit ARIMA(p,1,0) with exogenous variables models in which p was in the set { 0, 1, 2 } and q was set to zero (motivated by results for the best Baseline ARIMA model) and the best four variables from the Null ARIMA Exog model were considered { tot_bus_inv_8, ret_sales_2, pmi_man_1, ls_rcvbl_8 } as exogenous variables (the labeling convention is of the form featurename_lag). I fit all combinations of p and exogenous variables from these sets of parameters.

Finally, I measured model predictions versus test data by calculating absolute percent error (APE) and its mean over observations (MAPE).

Model Fitting Results

There are not large qualitative differences among the three models in terms of visual fit to the data. Figure 1 illustrates the fit of the Baseline ARIMA model to the training and test data along with the residuals. The best fit models by cross validation and BIC are shown.

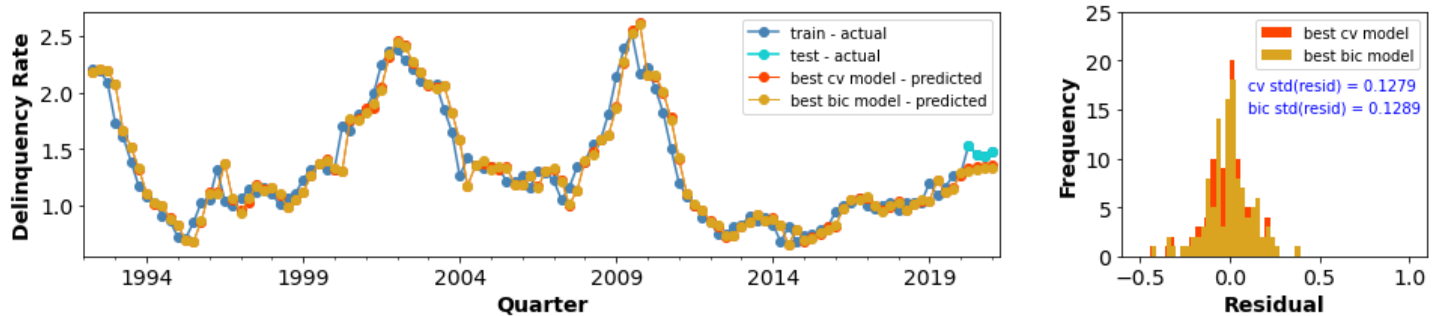


Figure 5

Predictions and residuals from the best-fit Baseline ARIMA models by cross validation (red) and BIC (gold) are shown. Actual data are shown in blue (train) and aqua (test).

The best-fit autocorrelation, moving average and features for each model are:

Baseline ARIMA

ARIMA(2,1,0) (no features included as exogenous variables)

Null ARIMA Exog

ARIMA(0,1,0) { tot_bus_inv_8, ret_sales_2, pmi_man_1, ls_rcvbl_8 }

Full ARIMA

ARIMA(1,1,0) { pmi_man_1, ls_rcvbl_8 }

For the Baseline ARIMA and Null ARIMA Exog models, several of the parameterizations resulted in very similar values for MSE. In these cases, I made a judgment call, also informed by the AIC and BIC statistics, to select as the best model a parameterization that accounted for most of the improvement in fit. Adding further parameters to these models improves fit, but there are diminishing returns to this improvement. Further details on these analyses are on GitHub.

For the best-fit features, the time series results are consistent with the exploratory data analysis. For example, delinquency rate shows negative autocorrelation, PMI at lag 1 has a negative regression coefficient, and Lease Receivable Balance at lag 8 has a positive regression coefficient.

2.4 Regression Modeling

The analysis of the three alternative time series models helped identify features that are the most important for forecasting lease delinquency rate. The time series results suggest that both autocorrelation in lease delinquency rate and features with different lags as exogenous variables are useful in the models.

In this analysis, I extended the time series work by fitting a simple linear regression model for lease delinquency rate that incorporates the best features previously identified. In this way, we can compare the time series model to a basic regression model to help further understand the impact of the different features for predicting lease delinquency rate and assess the performance of the time series models relative to the classic linear regression model as a kind of standard benchmark.

I fit the linear regression model using the best features identified in the time series model as well as the same variables with nearby lags with the idea that if a given lag is important then perhaps nearby lags are important for the same variable. Accordingly, I fit a model with the following 8 predictor variables:

Regression

delinq_1, delinq_2, pmi_man_1, pmi_man_2, pmi_man_3, ls_rcvbl_7, ls_rcvbl_8, and ls_rcvbl_9.

Figure 6 shows the actual and predicted lease delinquency time series with training and test data for this regression model along with a histogram of residuals. Qualitatively, the results are similar to the time series models.

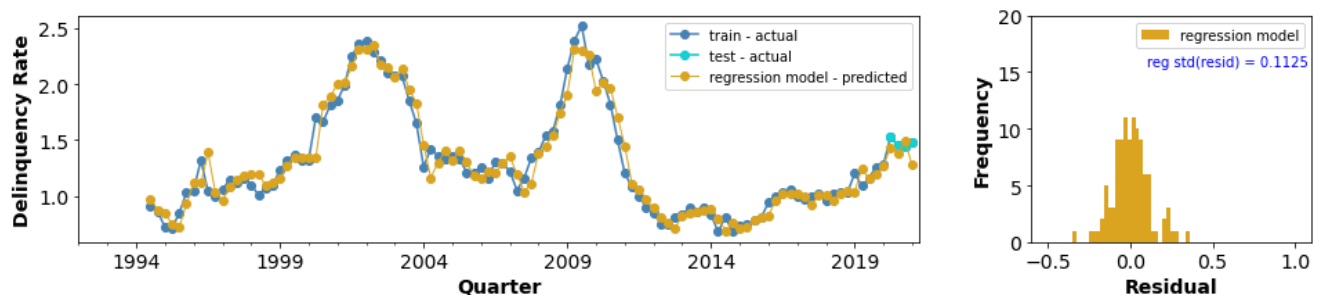


Figure 6

Predictions and residuals from the best-fit Regression model (gold). Actual data are shown in blue (train) and aqua (test).

To better understand how feature variables are impacting predictions, using the Regression model I calculated feature importance using the product of the regression coefficient and the standard deviation of the feature. This calculation measures the impact of a one unit standard deviation change in the feature variable on the response variable, lease delinquency rate. Figure 7 shows a scaled version of this measure of importance for each feature. When sorted by magnitude, the first four features have the largest impact on lease delinquency rate.

Again, as with the time series analysis, the magnitude and direction of the impact of each feature on lease delinquency rate is consistent with the results of the exploratory data analysis.

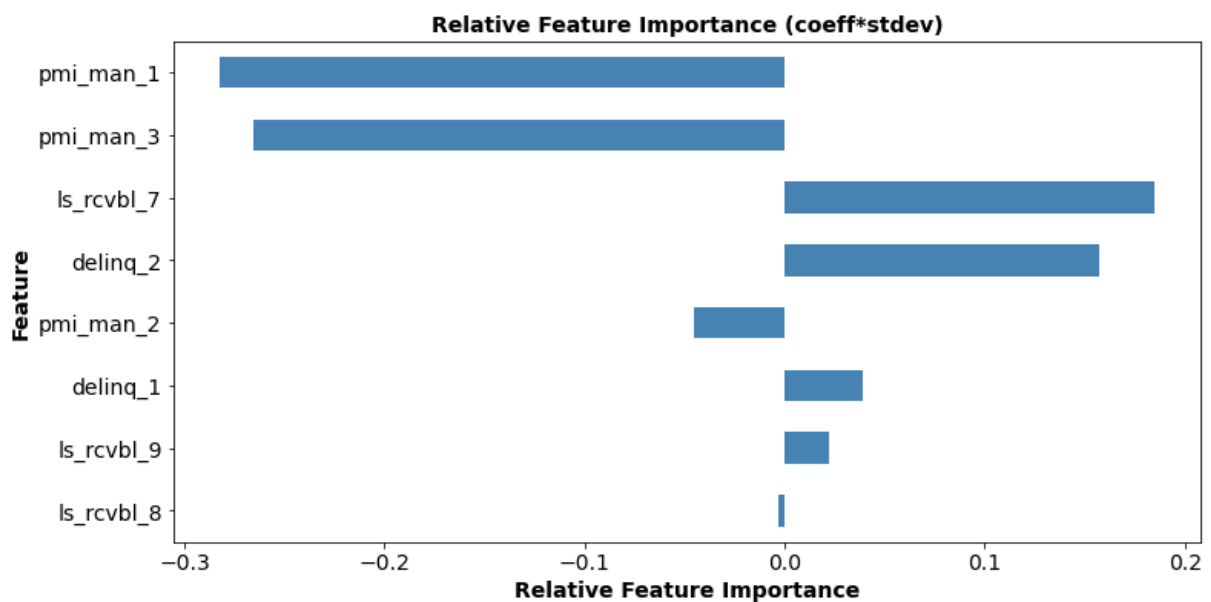


Figure 7
Relative feature importance for the best-fit Regression model.

3 FINDINGS

I fit four models to the data, three ARIMA time series models ('Baseline ARIMA', 'Null ARIMA Exog', 'Full ARIMA') and one regression model ('Regression'). These four models fit the training data similarly well, which can be seen qualitatively in the time series plots and more quantitatively in the distribution of absolute percent error (APE) on the training data (Figure 8). However, the fit in terms of APE differs somewhat more among the models for the test data (Figure 8). Hence, the models varied in their robustness for predicting the new values in the test data.

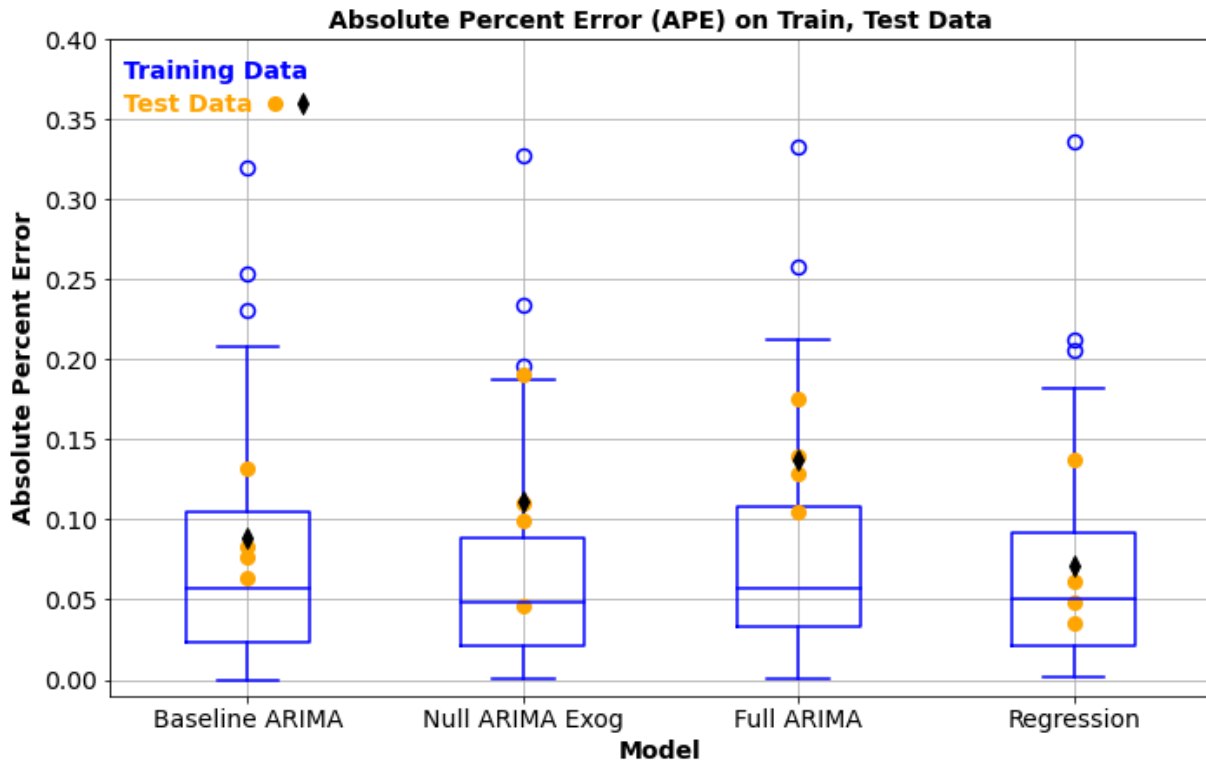


Figure 8

Absolute percent error of the four models on the training data (blue boxplot) and test data (orange circles). Mean absolute percent error (MAPE) on the test data is shown for each model (black diamonds).

For the test data, which is the last four quarters in the time series, 2020Q1 to 2020Q2, the Regression model fits best with a MAPE of 7.1%, followed by the Baseline ARIMA with a MAPE of 8.8% (Figure 8, Table 1). Considering a one-quarter-ahead forecast, which would forecast the first quarter in the test data, 2020Q1, forecast values and absolute percent error are also shown in Table 1. For this single data point, which is the most important business case, the Regression model prediction (1.44%) differs from the actual value (1.53%) by only about 6.2% (note that APE is a relative error), but it's very important to note that this is just a single forecasted data point. Interestingly, all four models predict increases in lease delinquency rate in 2020Q1 relative to 2019Q4 (per time series plots here - Figure 5, 6 - and additional plots on GitHub).

Table 1

For each of the four models: mean absolute percent error (MAPE) over the four quarters of test data and forecast lease delinquency rate and associated absolute percent error for the first quarter in the test data, 2020Q1. Actual delinquency rate in 2020Q1 was 1.53%.

Model	MAPE - Test Data	Forecast - 2020Q1	APE - 2020Q1
Baseline ARIMA	8.8%	1.31%	14.5%
Null ARIMA Exog	11.2%	1.38%	10.0%
Full ARIMA	13.7%	1.37%	10.5%
Regression	7.1%	1.44%	6.2%

4 CONCLUSIONS AND FUTURE WORK

The exploratory data analysis and model fitting make clear that there is structure in these time series data that correlates with time trends in lease delinquency rate. There is autocorrelation structure at short lags near lags 1 and 2 as well as longer lags near lag 8. Further, the feature variables are correlated with lease delinquency rate over these same sets of lags. These are the most important modeling results because they form a foundation for what features can be examined to explain and predict lease delinquency rate.

Given the structure in the data just described, an important next step is to better account for the inclusion of different features at different lags while at the same time not exploding the dimensionality of the model. The results of the Regression model suggest it might be important to include more features in the model simultaneously (versus the stepwise model selection approach considered here). One approach might be to try using parametric functions to weight lagged variables using fewer parameters (versus having a coefficient parameter for each lagged variable). A second approach might involve a dimensionality reduction approach, such as Principal Components Analysis or a similar technique.

5 BUSINESS RECOMMENDATIONS

The results of this analysis suggest the following business recommendations.

1. Use the Regression model for one-quarter-ahead forecasting.
2. Monitor related economic time series as qualitative leading indicators for future trends in Lease Delinquency Rate. Specifically,
 - 2.1. ISM Purchasing Managers Index (PMI) is negatively correlated with Lease Delinquency Rate with a 1 to 3 quarter lag, so decreases in PMI in the current quarter forecast possible increases in delinquency in the next 1 to 3 quarters.
 - 2.2. Lease Receivable Balance, all banks, is positively correlated with Lease Delinquency Rate with a 7 to 9 quarter lag, so increases in Lease Receivable Balance in the current quarter forecast possible increases in delinquency in the next 7 to 9 quarters
3. In future analyses, consider making enhancing the macroeconomic data with internal, company-specific data on lease delinquency rates and associated lease attributes.

6 CONSULTED RESOURCES

Hyndman, R.J., & Athanasopoulos, G. (2021) Forecasting: principles and practice, 3rd edition, OTexts: Melbourne, Australia. [OTexts.com/fpp3](https://otexts.com/fpp3).

PMDARIMA Python package for time series (<https://pypi.org/project/pmdarima/>)

STATSMODELS Python package for statistics (<https://www.statsmodels.org/>)

SKLEARN Python package for machine learning (<https://scikit-learn.org/>)