

CHAPTER 1

**TIME SERIES FORECASTING OF MONTHLY ACTIVE
USERS WITH LONG SHORT-TERM MEMORY NETWORKS**

Matthew Louis Rosendin
University of California, Berkeley
Department of Industrial Engineering and Operations Research

April 13, 2018

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 4 |
| 1.1 | Goal | 4 |
| 1.2 | Motivation | 4 |
| 1.3 | Time Series Forecasting | 5 |
| 1.4 | Autoregressive Models | 5 |
| 1.5 | Recurrent Neural Networks | 5 |
| 2 | Object-Oriented Abstraction | 6 |
| 2.1 | Decomposition | 7 |
| 2.2 | Encapsulation | 7 |
| 3 | Hyperparameter Optimization | 7 |
| 3.1 | Cross-validation | 7 |
| 3.2 | Grid Search | 8 |
| 3.3 | Results | 8 |
| 4 | Systems Engineering | 8 |
| 4.1 | Server Infrastructure | 9 |
| 4.2 | Distributed Task Queue | 9 |
| 5 | User Interface | 10 |
| 5.1 | Forecast Visualization | 10 |
| 5.2 | Documentation | 10 |
| 6 | Next Steps | 11 |
| 6.1 | Websockets | 11 |
| 6.2 | Data Pipeline Automation (Apache Hadoop) | 11 |
| 6.3 | Dimensionality Reduction by Feature Selection | 11 |
| 6.4 | Improving Hyperparameter Optimization Efficiency | 11 |
| 6.5 | Confidence Interval Visualization | 12 |
| 6.6 | Distributed Machine Learning Clusters | 12 |
| 7 | Conclusion | 12 |

List of Figures

| | | |
|---|---|----|
| 1 | Diagram of a one-unit Long Short-Term Memory (LSTM) network [4] | 6 |
| 2 | Example output: 5.79% MAPE | 9 |
| 3 | Host machine system information | 10 |
| 4 | Dashboard user interface | 11 |
| 5 | Product 1 trend: 6.24% MAPE | 12 |
| 6 | Product 1 forecast: 6.24% MAPE | 13 |

1 Introduction

Our project, in partnership with Adobe Systems Inc., is an interactive dashboard that forecasts monthly active users for any Adobe software product. We've accomplished forecasts with excellent accuracy as far as 5 months into the future. The growth model's forecast can be used for a number of reasons, including as an input to forecasting revenue or as a measure of product health. While forecasting monthly active users is nothing new, our application of artificial neural networks (ANNs) seems to be a novel approach that has yielded promising results.

In this chapter of the paper I analyze and discuss the implications of our model's performance and describe my work in deploying our model. I begin by explaining the concrete goals of this work, then outlining the motivation, before finally diving into a technical introduction to our model. Following the basic conceptual discussion around our model, I start the discussion about how I refined the "raw" model code into a flexible, scalable, and configurable object. The translation of the model into a manageable machine learning software system led to additional improvements to the model, such as hyperparameter optimization. A byproduct of the optimization was an exhaustive experimentation of model performance metrics. The underpinning of the preceding work is the user interface. I discuss my work in front-end and systems engineering that resulted in the final deliverable: a forecasting dashboard powered by machine learning.

1.1 Goal

Our team's foremost goal is to create a practical and interactive product growth model for Adobe senior management. My individual goal was to design and create a user interface to interact with the model. Throughout my contributions, I've added additional value around testing the performance of our model and gathering relevant metrics.

1.2 Motivation

The motivation for creating a dashboard is straightforward. The dashboard enables the user to train the machine learning model resulting in performance metrics and the monthly active user forecast. The motivation for the section on hyperparameter optimization was simply to improve the results of our model. From the resulting tests, we've found that our model's forecasting capability had improved. Part of the discussion in this paper is reserved for analyzing the source of model improvement.

1.3 Time Series Forecasting

In order to explain why we chose our particular model (a long short-term memory network) I will explain the characteristics of the problem we are solving through a brief taxonomy. The problem is best framed as, "how do we forecast monthly active usage multiple weeks into the future?". The most salient characteristic of this problem is that the data is a long sequence (i.e., data is chronologically ordered) of values, known as a time series. What we are trying to achieve is a forecast of time series data, so the solution to the problem is known as time series forecasting.

1.4 Autoregressive Models

The simplest model might be an autoregressive (AR) model in which the forecast values are regressed on prior data:

$$X_t = c + \sum_{i=1}^p \phi_i X_{t-i} + \epsilon_t \quad (1)$$

where c is a constant, ϵ_t is an error term, and ϕ_1, \dots, ϕ_p are the parameters to be estimated by linear regression. To forecast a value for the present time t , we specify a parameter p , known as the "order" or "lag". With $p = 3$ the corresponding model would look like this:

$$X_t = c + \phi_1 X_{t-1} + \phi_2 X_{t-2} + \phi_3 X_{t-3} + \epsilon_t \quad (2)$$

We would call this an autoregressive model with a lag of 3 time steps where time steps are defined as the number of sequential data points. The autoregressive model can be notated as $AR(3)$ to indicate its order with $p = 3$.

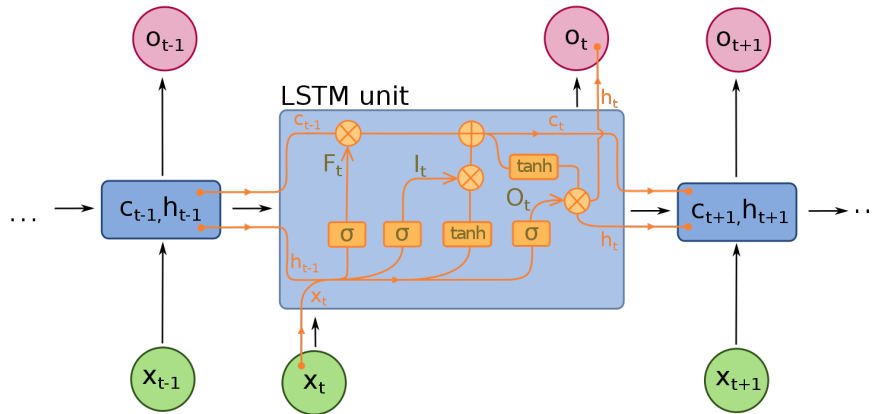
In our work we use the concept of "lag" from the autoregressive model in our supervised learning model. Our dataset is multivariate, meaning there are more than two observations for each time step. Since we would like to forecast multiple weeks into the future, our model is multi-stepped.

1.5 Recurrent Neural Networks

There are many good reasons to justify using an artificial neural network. Much of the recent literature on time series forecasting is focusing on the advantage of ANNs versus other methods, including autoregressive interactive moving average (ARIMA) models. For instance, Ahmed et al. and Zhang show that ANNs are shown to be superior for generalized time series problems [2].

Although complex models such as neural networks can be hard to interpret, our project team values performance over interpretability. Furthermore, we have designed features similar to those in financial time series models. For financial stock data, Adebaye et al. show that ANNs are superior to ARIMA models [1]. Lastly, ANNs are capable of modeling non-linearities which can improve the forecast.

Figure 1: Diagram of a one-unit Long Short-Term Memory (LSTM) network [4]



Recurrent neural networks (RNNs) are a class of artificial neural network where connections between units form a directed graph along a sequence. In other words, RNNs retain state at one time to the next, using the previous state's output for the current estimate. This characteristic enables multi-stepped time series forecasting. A particular architecture of RNNs known as Long Short-Term Memory (LSTM) allows the model to recognize and retain short-term patterns for long periods of time. This architecture is used best where data from an earlier state needs to be recalled at a later state. Examples of LSTM networks in industry include Google Voice [3].

2 Object-Oriented Abstraction

The working LSTM model was created in a Jupyter Notebook. While serving as a minimum proof of concept, the model should not require the end-user engage directly with code. Furthermore, the following criteria were set: 1) the lifecycle of the model (from training, testing, to final predictions) needed to be automated and callable with one command; 2) the Jupyter Notebook code (which executes similarly to a script) needed to be organized into an object-oriented structure to store data within the object instance; 3) the model needed to be capable of loading

datasets of arbitrary products. Criterion (2) is more a matter of preference and style to the author, as other programming paradigms could be considered. The rest of this section briefly discusses the process involved in restructuring the model code to meet the criteria.

2.1 Decomposition

The goal of decomposition is to break a complex system into parts that are easier to understand, program, and maintain. The process of decomposing the Jupyter Notebook code involved defining separate modules for loading datasets (`data.py`), feature engineering (`features.py`), training/testing/predicting (`core.py`), and accessing utilities (`utils.py`). Furthermore, the code was organized into functions, and eventually into methods (discussed below).

2.2 Encapsulation

After breaking the code into manageable chunks, I created a Python class wrapper to store internal state and to allow the end-user to specify runtime options. The entire model can be run by using the wrapper's simple application programming interface (API).

3 Hyperparameter Optimization

Hyperparameter optimization deals with the learning parameters of a model rather than weights for features. While I made a significant effort in this area, the optimization was only run on one machine. Increasing compute resources to perform a more extensive search may yield further interesting results.

3.1 Cross-validation

In time series problems, cross-validation is performed chronologically with splits at fixed time intervals. In my trials, I ran cross-validation with 2, 3, and 10 splits. I wanted to understand the change in our key metric, MAPE (mean absolute percent error), as the number of splits S increases. According to the `scikit-learn` documentation, the training set Tr for the i th split has a size

$$Tr_i = \left\lfloor \frac{i \times N}{(S+1)} \right\rfloor + N \bmod (S+1) \quad (3)$$

where N is the number of samples, S is the number of splits, and $\lfloor expression \rfloor$ denotes the floor function. Meanwhile, the test set Te has the following size:

$$Te_i = \left\lfloor \frac{N}{(S+1)} \right\rfloor \quad (4)$$

For time series splits with $S > 1$, the sum of the first train/test set is strictly less than N and the last train/test set is equal to N :

$$Tr_i + Te_i < N \quad \forall i \in \{1, \dots, S-1\} \quad (5)$$

$$Tr_S + Te_S = N \quad (6)$$

These equations describe the train/test split pairs. It is important to note that the chronological indexing of the data is preserved for both input data and additional "held-out" data. This splitting strategy helps validate the results of the exhaustive hyperparameter search that is described in the next section.

3.2 Grid Search

This search algorithm is simple. It exhaustively generates candidate learning parameters and selects the best candidate based on the scoring function (MAPE). Each candidate undergoes the time series splitting strategy discussed in Section 3.1 to help validate the candidate's score. Further discussion about alternative algorithms is in Section 6.4.

3.3 Results

An example of the output this process is shown in Figure 2.

4 Systems Engineering

I worked to productionize our model by building a back-end that is capable of extracting and preparing data, training and validating the model, and producing forecasts. The back-end performs the business logic computation without end-user interaction. The back-end was written entirely in Python and its component parts are discussed in this section.

Figure 2: Example output: 5.79% MAPE

```
Parameters: {  
  'batch_size': 150,  
  'epochs': 200,  
  'optimizer': 'adam',  
  'units': 100  
}  
Test loss: [128.97596047141334, 128.97596047141334]  
Testing set RMSE: 24318.87  
Testing set MAPE: 5.79%  
Testing set sMAPE: 0.06%
```

4.1 Server Infrastructure

All software written for this project was developed and tested on local machines. System specifications of the host machine producing the results in this report are shown in Figure 3. The system runs Python 3.6.4 and we've tested on local machines (macOS) and a Docker virtual host container running Ubuntu 16.04 (see Section 6.6 for further discussion about Docker). A simple web server framework, Flask (<http://flask.pocoo.org/>), is employed to interact with the front-end (see in Section 5). A basic API has been developed to enable asynchronous communication with the front-end client device.

4.2 Distributed Task Queue

Since the machine learning processes are long-running processes, an open-source distributed task queue called Celery (<http://www.celeryproject.org/>) is used to manage the scheduling, status, and output. When the user uses the front-end to make a new forecast, an API endpoint queues the task. In our case, there was only one client utilizing the queue at a time. However, the additional benefit of using a task queue is to coordinate background jobs. When the end-user begins a new forecast task, they receive a response from the server immediately. The code on the client polls another API endpoint every second to check the status of the forecast task. When the status appears successful, the server responds with the necessary data to render the visualizations on the client device, as well as summary statistics.

Figure 3: Host machine system information

5 User Interface

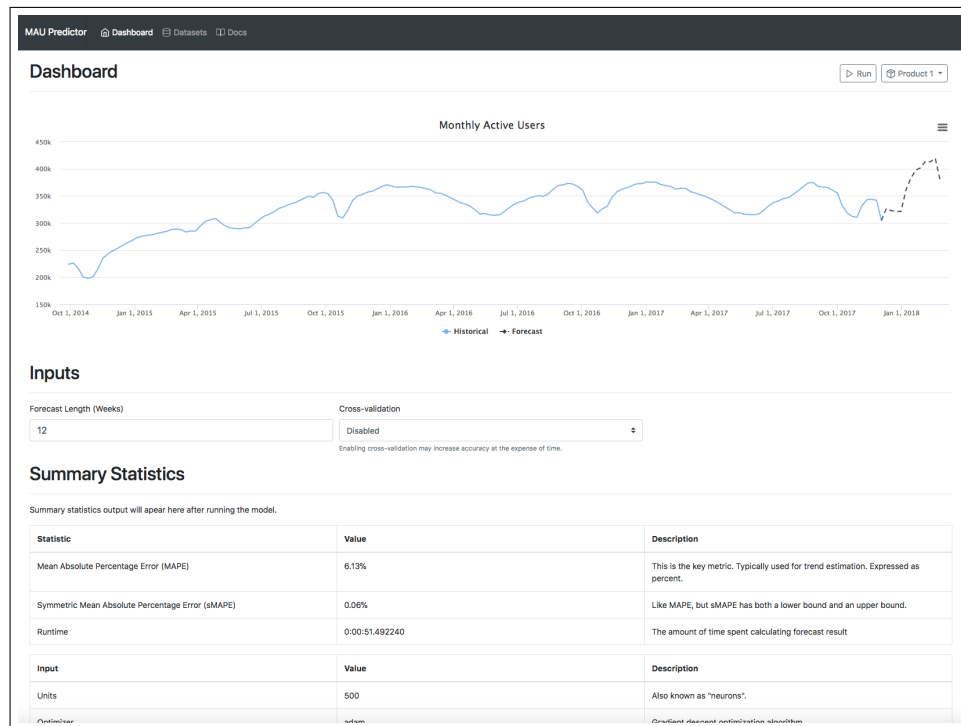
5.1 Forecast Visualization

Figure 5 shows the historical monthly active user trend of Product 1 with the 12 week forecast in the dashed line.

Figure 6 shows the forecast for Product 1 over a 12 week period (December 3rd - February 25th 2018).

5.2 Documentation

Wrote documentation (akin to a usage guide)

Figure 4: Dashboard user interface

6 Next Steps

6.1 Websockets

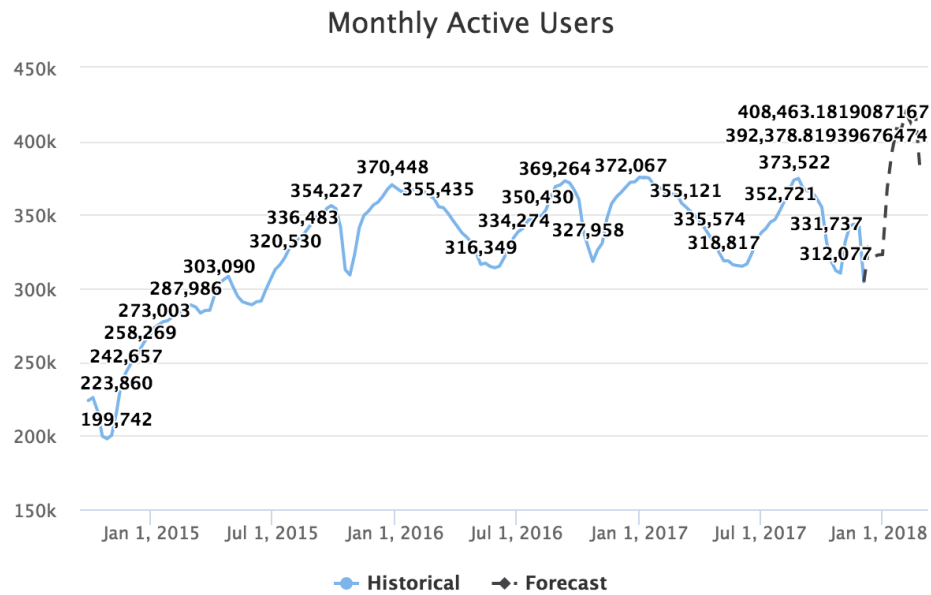
Improve job monitoring; rather than polling every second with HTTP calls, open a websocket channel.

6.2 Data Pipeline Automation (Apache Hadoop)

6.3 Dimensionality Reduction by Feature Selection

6.4 Improving Hyperparameter Optimization Efficiency

The brute-force search algorithm used to optimize hyperparameters can probably be replaced with another. The `GridSearchCV` class from `scikit-learn` implements the algorithm that was used for our project, but other hyperparameter optimizers are available in the `model_selection` module of `scikit-learn`. The other options include a randomized search where a subset of

Figure 5: Product 1 trend: 6.24% MAPE

the parameter grid is sampled. This requires passing in distributions of the hyperparameters themselves, which could be an interesting opportunity to avoid an exhaustive search.

6.5 Confidence Interval Visualization

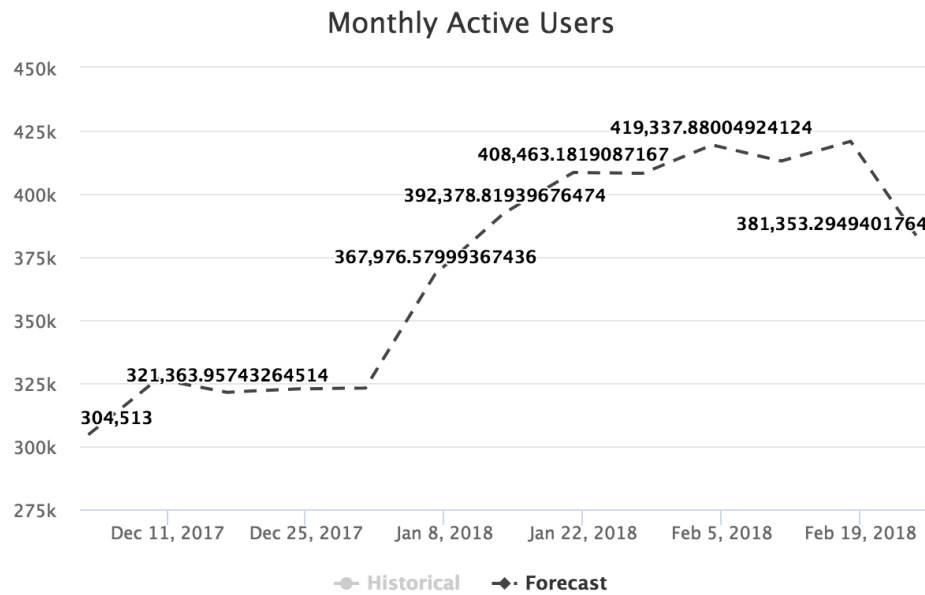
A major upgrade would be a visualization of forecast confidence intervals on the current graph.

6.6 Distributed Machine Learning Clusters

Docker containers can be scaled on nodes in a distributed system.

7 Conclusion

In conclusion, ...

Figure 6: Product 1 forecast: 6.24% MAPE

References

- [1] Adebiyi, A. A., Adewumi, A. O., & Ayo, C. K. (2014). Comparison of ARIMA and Artificial Neural Network Models for Stock Price Prediction. *Journal of Applied Mathematics*, 2014.
- [2] Ahmed, N. K., Atiya, A. F., Gayar, N. E., & El-Shishiny, H. (2010). An Empirical Comparison of Machine Learning Models for Time Series Forecasting. *Econometric Reviews*, 29(5-6).
- [3] Beaufays. (2015, August 11). The neural networks behind Google Voice transcription [Blog post]. Retrieved from <https://research.googleblog.com/2015/08/the-neural-networks-behind-google-voice.html>.
- [4] Deloche, François. (2017). Diagram of a one-unit Long Short-Term Memory (LSTM) network. Retrieved from <https://commons.wikimedia.org/w/index.php?curid=60149410>.