

DeepFlow

IP flow data prediction using deep learning techniques

Dr. Sebastian Fischer - seb1988@stanford.edu

Predicting network traffic is of high relevance for Internet Service Providers, because it enables them to prepare for upcoming peak network utilizations better. This project illustrates how recurrent neural networks, specifically Long Short-Term Memory models, outperform classical autoregressive models in forecasting IP flow data.

PREDICTING

ARIMA stands for "autoregressive integrated moving average" and is a statistical model to forecast time series data. I used the Python library PMDARIMA to run the `auto_arima` function that finds the optimal hyperparameters. This model serves as a baseline.

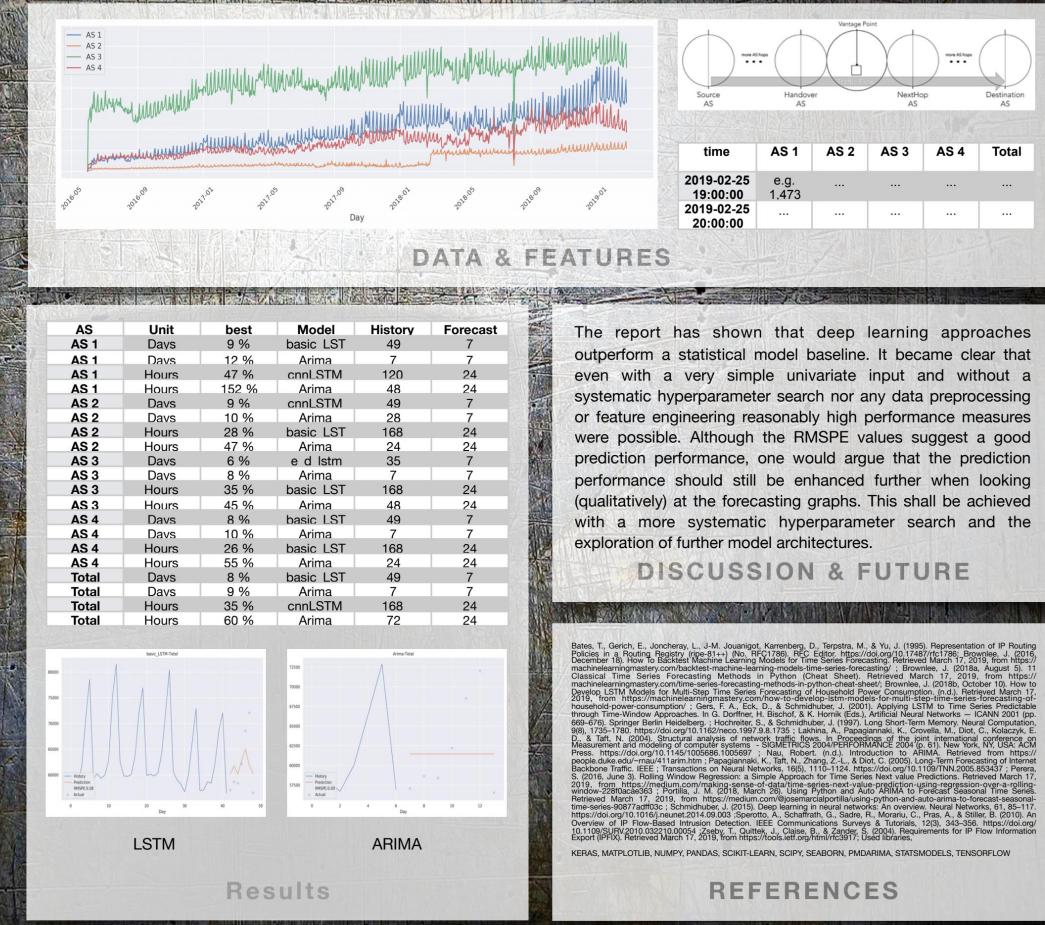
ARIMA (Baseline)

Layer (type)	Output Shape	Param#
<hr/>		
lstm_239 (LSTM)	(None, 200)	161600
dense_367 (Dense)	(None, 100)	20100
dense_368 (Dense)	(None, 7)	707
Basic LSTM		

Layer (type)	Output Shape	Param#
<hr/>		
lstm_211 (LSTM)	(None, 200)	161600
repeat_vector_84 (RepeatVect)	(None, 7, 200)	0
lstm_212 (LSTM)	(None, 7, 200)	320800
time_distributed_167 (TimeDi)	(None, 7, 100)	20100
time_distributed_168 (TimeDi)	(None, 7, 1)	101
Encoder-Decoder LSTM		

Layer (type)	Output Shape	Param#
<hr/>		
conv1d_93 (Conv1D)	(None, 22, 64)	256
conv1d_94 (Conv1D)	(None, 20, 64)	12352
max_pooling1d_47 (MaxPooling)	(None, 10, 64)	0
flatten_47 (Flatten)	(None, 640)	0
repeat_vector_96 (RepeatVect)	(None, 24, 640)	0
lstm_231 (LSTM)	(None, 24, 200)	672800
time_distributed_191 (TimeDi)	(None, 24, 100)	20100
time_distributed_192 (TimeDi)	(None, 24, 1)	101
CNN-Encoder-Decoder LSTM		

MODELS



The report has shown that deep learning approaches outperform a statistical model baseline. It became clear that even with a very simple univariate input and without a systematic hyperparameter search nor any data preprocessing or feature engineering reasonably high performance measures were possible. Although the RMSPE values suggest a good prediction performance, one would argue that the prediction performance should still be enhanced further when looking (qualitatively) at the forecasting graphs. This shall be achieved with a more systematic hyperparameter search and the exploration of further model architectures.

DISCUSSION & FUTURE

Bates, T., Gerlich, E., Jonckheer, L., J-M. Jouary, Kamerling, D., Terpstra, M., & Yu, J. (1998). Representation of IP Routing Policies in a Routing Registry. Ieee8-91-1 (No. RFC1788). IFC Editor. <https://doi.org/10.17477/rfc1788>. Brownlee, J. (2018, August 5). <https://machinelearningmastery.com/backtest-machine-learning-models-time-series-forecasting/>. Brownlee, J. (2018, August 5). <https://machinelearningmastery.com/classical-time-series-forecasting-methods-in-python-cheat-sheet/>. Retrieved March 17, 2019, from <https://machinelearningmastery.com/classical-time-series-forecasting-methods-in-python-cheat-sheet/>. Deville, J. (2018). DeepLSTM Models for Multi-Step Time Series Forecasting of Household Power Consumption. (n.d.). Retrieved March 17, 2019, from <https://www.semanticscience.org/paper/DeepLSTM%20Models%20for%20Multi-Step%20Time%20Series%20Forecasting%20of%20Household%20Power%20Consumption.pdf>. Gerlach, R., & Schmidhuber, J. (2001). Applying LSTM to Time Series Predictable Time Series. In *Advances in Neural Information Processing Systems 14: Advances in Neural Information Processing Systems 14: Volume 1* (pp. 669–676). Springer Berlin Heidelberg. Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>. Lahiri, A., Pandit, S., & Crotwell, M. (2010). *Electric Power Measurement and Modeling of Computer Systems*. Springer US. Li, J., & Liu, Y. (2017). <https://arxiv.org/abs/1708.04097>. Nau, Robert. (n.d.). Introduction to ARIMA. Retrieved from <https://people.duke.edu/~mau411/arma.htm>. Papageorgiou, K., Tait, N., Zhang, Z.-L., & Dot, C. (2005). Long-Term Forecasting of Internet Backbone Traffic. *IEEE Communications Letters*, 9(12), 1133–1135. <https://doi.org/10.1109/LCOMM.2005.1629072>. Portilla, J. M. (2018, March 26). Using Python and Auto ARIMA to Forecast Seasonal Time Series. Retrieved from <https://towardsdatascience.com/using-python-and-auto-arima-to-forecast-seasonal-time-series-22905a0e633>. Portilla, J. M. (2018, March 26). Using Python and Auto ARIMA to Forecast Seasonal Time Series. Retrieved from <https://towardsdatascience.com/using-python-and-auto-arima-to-forecast-seasonal-time-series-22905a0e633>. Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks*, 61, 85–117. <https://doi.org/10.1016/j.neunet.2014.09.003>. Seneviratne, A., Seneviratne, G., Seneviratne, C., Seneviratne, A., & Seneviratne, S. (2010). An Overview of IP Flow-Based Intrusion Detection. *IEEE Communications Surveys & Tutorials*, 12(3), 329–359. <https://doi.org/10.1109/SURV.2010.03210.00054>. Stenzel, T., Quitté, J., Clasen, B., & Zander, S. (2004). Requirements for IP Flow Information Exchange. *Network Protocols*, 16(1), 1–17. <https://doi.org/10.1007/s11043-004-0001-0>.

KERAS, MATPLOTLIB, NUMPY, PANDAS, SCIKIT-LEARN, SCIPY, SEABORN, PMDARIMA, STATSMODELS, TENSORFLOW