

Machine Learning(MO444)

Final Series Forecasting

Felipe Galvão
RA:116790
felipelemes@outlook.com

William Tustumi
RA:120281
whatust@gmail.com

Leo Yuuki Omi Omi
RA 138684
leoyuuki@gmail.com

João Pedro Ramos Lopes
RA 139546
email address

Abstract—

I. INTRODUCTION

This project will focus on the problem of Web Traffic Time Series Forecasting hosted on Kaggle [2]. The problem focuses on the problem of forecasting the future values of multiple time series. Sequential or temporal observations emerge in many key real-world problems, ranging from biological data, financial markets, weather forecasting, to audio and video processing. The field of time series encapsulates many different problems, ranging from analysis and inference to classification and forecast. On our project we will forecast future web traffic for approximately Wikipedia pages, with a dataset provided by Kaggle [3].

II. DATASET ANALYSIS

A. Dataset

The dataset provided is comprised of 803 days of visits measurements from 145063 combination of wikipedia page and access agent. An page name and header example can be seen on table I. The entries with NaN are correspondent to the days when the wikipedia page did not existed yet, thus those values where substituted by 0 in order to treat this edge case.

	2015 07/01	2015 07/02	2015 07/03	2015 07/04	2015 07/05
2NE1_zh.wikipedia.org_...	18 key	11	5	13	14
2PM_zh.wikipedia.org_a...	11	14	15	18	11
3C_zh.wikipedia.org_al...	1	0	1	1	0
4minute_zh.wikipedia.or...	35	13	10	94	4
52_Hz_I_Love_You_zh...	NaN	NaN	NaN	NaN	NaN
5566_zh.wikipedia.org_...	12	7	4	5	20
91Days_zh.wikipedia.org...	NaN	NaN	NaN	NaN	NaN
A'N'D_zh.wikipedia.org..	118	26	30	24	29

Table I
DATA EXAMPLE FROM KAGGLE SERIES FORECAST

The first step was performing an analysis of the data to understand which patterns could leverage our learning methods. Plotting graphs of the traffic from a few pages we observed that different pages show remarkably distinguish trends in traffic. This suggests that a good model needs to explicitly differentiate pages. For specific entries we note a strong year-to-year autocorrelation. We also observed weaker week-to-week and month-to-month trends. One clear example of this case was the Thanksgiving holiday page that showed an yearly

spike around a Thursday on late November, which is the day of the holiday.

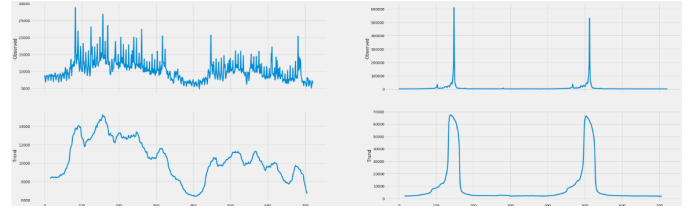


Figure 1. Observed Traffic and Traffic Trend from two different pages

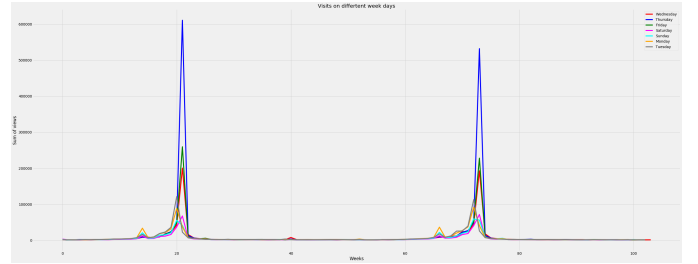


Figure 2. Observed Traffic separated by weekday for the thanksgiving page.

Additionally, we have to deal with spikes in traffic not following a regular pattern. We note a short-term dependency following those spikes, which means the days immediately before the prediction are important to consider.

III. BACKGROUND

A. Time Series

B. Multi-Layer Fully Connected

C. Recurrent Neural Network

A Recurrent Neural Network (RNN) is very similar to a feed-forward network, except that the connections also links backwards on the network. So the output is defined similarly to equation 1. Note that the result from last step contributes to the next and so on, therefore the i th step depends on all previous steps, which takes a form of a memory [1].

$$Y_{(t)} = \phi(X_{(t)}\dot{W}_x + Y_{(t-1)}\dot{W}_y + b) \quad (1)$$

The training method analog to the feed-forward networks, the first step is a forward pass that unroll the network time

loop, the values for all Y'_s in the time series are also calculated in the process. The gradient value is calculated using all the results in the cost function, for example if the time series generates Y_1, \dots, Y_n and the cost function uses only the last three values of Y , then the gradient is calculated using only Y_{n-2}, Y_{n-1}, Y_n .

Recurrent neural network tends to suffer from exploding and vanishing gradients more than a convolution or fully connected network, since the number of steps in time it unrolls also adds to the depth of the network. To solve this problem techniques like batch normalization, non-saturating activation functions and gradient clipping are exploited.

Another problem with the simple approach shown previously happens when the network has to remember inputs far early in the time steps, most of the information is lost by transversing the network. To overcome this problem a more sophisticated cell was developed, Long Short Term Memory (LSTM). The full explanation of how it works is given on subsection III-D.

D. Long Short Term Memory

The main improvement LSTM cell wants to achieve is learn what to store in the long term memory, what to discard and what to use from the input. The LSTM cell keeps two state matrix that holds information from previous time sequences, one for long an the other for short term memory, those states are than used as input for the next iteration. The figure 3 show a LSTM diagram, the $c_{(t-1)}$ represents the long term memory input and the $h_{(t-1)}$ represent the short term memory, note that the short term output is also the $Y_{(t)}$.

The LSTM cell uses one fully connected layer with hyperbolic tangent activation function as the main flow of input and three other fully connected with logistic activation function as controllers. Each controller is attached to a gate, basically an element-wise multiplication, since the logistic function output range is 0 to 1 the gate is either close when 0 or open when 1.

Each gate receives a name based on their function and is controlled by one of the matrix controllers:

- 1) *Forget Gate*: controls what information from the input is retained on the long term state, it receives $f_{(t)}$ (logistic function controllers) and $c_{(t-1)}$ (long term state) as input.
- 2) *Input Gate*: controls which information on the main flow will be added to the long term memory and the output memory, receives $g_{(t)}$ (main input flow) and $i_{(t)}$ (logistic controller).
- 3) *Output Gate*: controls the information that goes to the output and short term memory, receives the addition result from *Input Gate* and *Forget Gate* after an hyperbolic tangent activation function and $o_{(t-1)}$ (logistic function controller).

IV. EXPERIMENT RESULTS & DISCUSSION

We first separated or data in train, validation and test datasets, comprising of respectively 70%, 20% and 10% of the total data. We broke our data in the time axis, so the validation and test data are referents to the later days instead of different pages. After we separated the dataset we end up

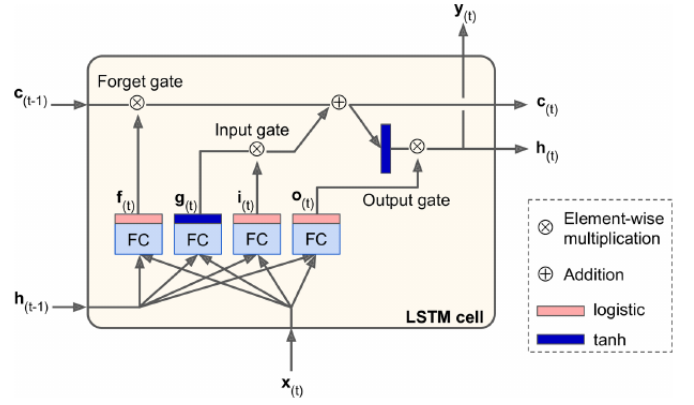


Figure 3. LSTM Diagram from [1]

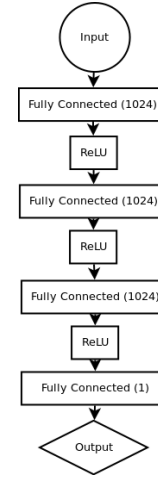


Figure 4. Diagram of model MLP's architecture

with, 559 samples for the training set, 159 samples for the validation set and 79 samples for the test set.

The code of the experiments were written in python with the module for neural network keras using the tensorflow as back-end. The models were ran in a intel core i7-6700k, with 16GB of memory and Geforce-980 Ti.

Due to memory constraints we decided to limit the number of pages on the test data, from 145063 to 50000.

A. MLP

The MLP model receives an input composed of an unique identifier of the wikipage accessed, the number of days passed since the first measurement was collected, and the number of visits from the last 7 days.

The MLP is built with one input layer of size 1024, two hidden layers of size 1024 and a output layer of size 1, a diagram of the model is shown on figure 4. The optimizer chosen was Adam with learning rate of 10^{-7} during 20 epochs and using batch size of 1000. The learning rate does not decrease over time. The loss graph can be seen on figure 5.

It is easy to see on figure 6 that the model achieved a pretty good accuracy on the maintaining close to the behavior of the

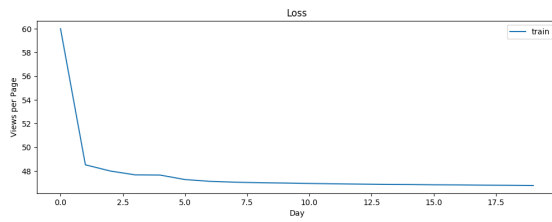


Figure 5. Loss graph from model MLP

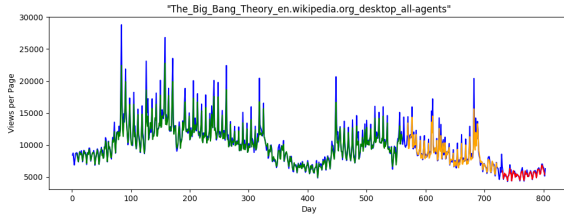


Figure 6. Plot of the values of visits per day compared to the prediction of MLP of the page The Big Bang Theory (en) of all desktops

curve, however the model undershoot the largest spikes, which probably leads to the biggest percentage of lost score.

The graph for thanksgiving is a very odd one, the only spikes happens on thanksgiving day, which is once a year and does not occur neither on validation data or test data, leaving the hardest prediction to appear only on training data.

B. LSTM-stateless

The LSTM model only uses as input, the page id, the number of days pass the first measurement and the last count of visits of that page.

The LSTM model uses 4 layers of LSTM cells of size 256, a normalization batch, two fully connected layers of size 1024 with ReLU as activation function and an output layer. parallel to the LSTM layer there is a fully connected layer of size 1024 to try improving the use of recent data, a batch normalization is also used after the fully connected layer too, figure 8 shows the model architecture.

The batch normalization layers are used to contain the exploding gradients of deeper neural networks, as well as the ReLU activation function. The Adam optimizer was chosen using learning rate of 10^{-5} , whitout learning rate decay, the model was trained during 5 epochs with batch size of 1000.

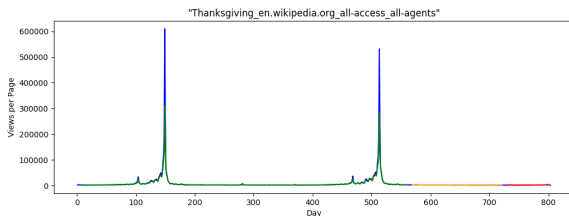


Figure 7. Plot of the values of visits per day compared to the prediction of MLP of the page Thanksgiving (en) of all desktops

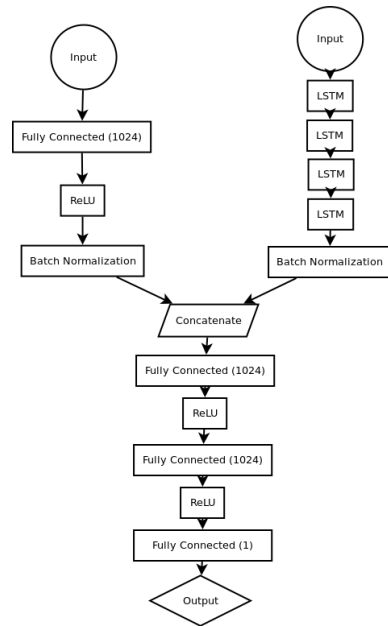


Figure 8. Diagram of model MLP's architecture

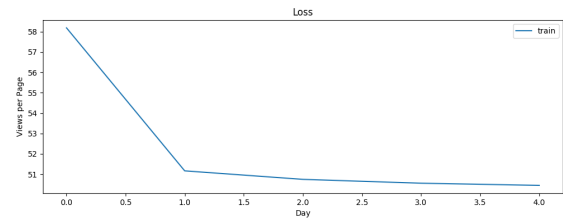


Figure 9. Loss graph from model LSTM(stateless)

The loss curve from the model training is displayed on figure 9.

the figure 10 shows that the LSTM model also follows the general format of the visit graph like the MLP model, however it undershoots not only the high spikes but, also some lower and more plateau segments. which increases the loss value of the model.

On the figure 11 it is possible to see that the LSTM model has a lot of difficult to follow the hight spikes from some pages.

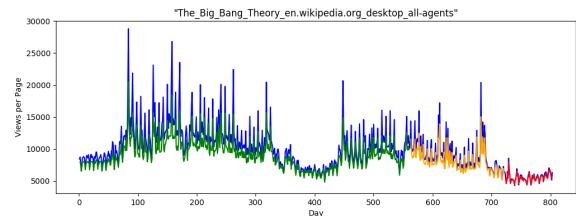


Figure 10. Plot of the values of visits per day, of the page The Big Bang Theory (en) of all desktops

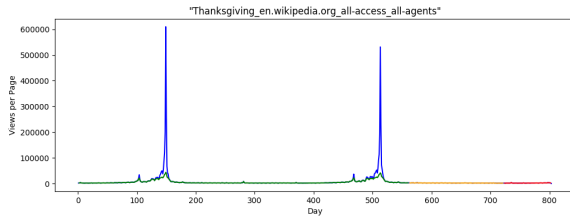


Figure 11. Plot of the values of visits per day, of the page Thanksgiving (en) of all desktops

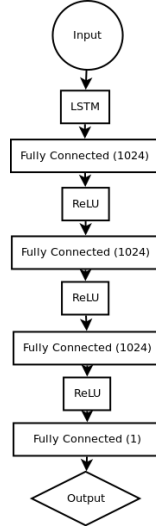


Figure 12. Diagram of model MLP's architecture

C. LSTM-stateful

This model passes the state from one batch to another allowing it to learn correlations between series of events. It receive as input the same data as the LSTM stateless model. It has one LSTM layer of size 256, followed by two fully connected hidden layer fo size 1024, and an output layer, figure 12 shows a diagram of the model. The optimizer used was the Adam with learning rate of 10^{-4} . The model was trained during 50 epoch with batch size of 50000 so the time steps would align their position on the batch index. The loss curve from the model training is displayed on figure 13.

From the loss graph on figure 13 the model seems to learn something, however from the visits prediction on figure 14 and 15, it is easy to see that the model simply cannot fit to any

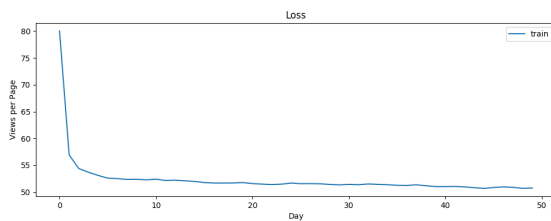


Figure 13. Loss graph from model LSTM(stateful)

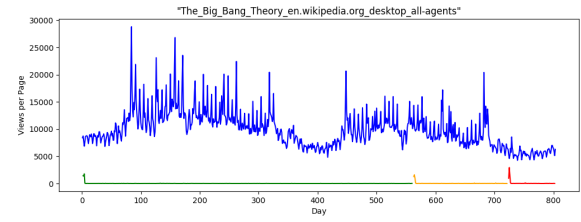


Figure 14. Plot of the values of visits per day, of the page The Big Bang Theory (en) of all desktops

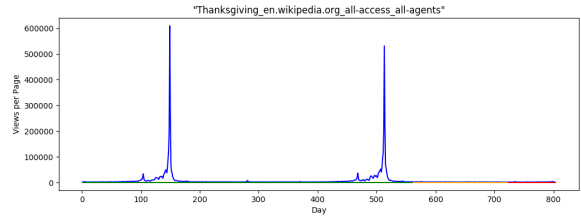


Figure 15. Plot of the values of visits per day, of the page Thanksgiving (en) of all desktops

curve of the prediction. This probably happens because the time step we are using is to small, which does not allow the model to unroll enough steps, and learn the relation between output of one state and the input of another state.

Table II show the loss value for all models on the training, validation and test dataset. The best performing model was the first one. All models performed better on validation and test dataset, this probably is a result of easier to predict data on those sets. As an example the wikipedia page from Thanksgiving has its bots spikes in the training data set while validation and test is basically stable, another example can be seen in figure 16 where there is a big spike on the training data that never happens again.

	Model_1	model_2(stateless)	model_3(stateful)
Train	46.76071	50.50748	59.08432
Validation	34.91314	38.97160	46.33987
Test	35.16409	39.14184	44.61450

Table II
TABLE WITH LOSS VALUE OF ALL MODELS ON TRAINING, VALIDATION AND TEST DATASET.

V. CONCLUSION & FUTURE WORK

Our best model was the MLP which used values from 7 days previous the one intended to be predicted, however this model has a obvious limitation, because it cannot learn predictions and relation between visits longer than the 7 days

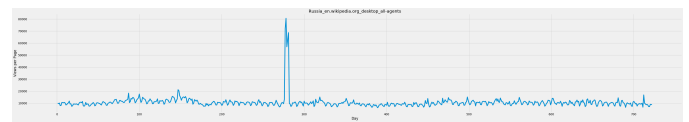


Figure 16. Plot of visits count for page Russia (en) on all desktops

it receives as input. In other words no long time sequence is learned using this model only small (7-days), page and day visit relation, since both data are also inputs. The stateless LSTM functions closely to an LSTM and end-up achieving a very similar performance to the MLP model, specially with the fully connected bypassing the LSTM layers.

The stateful LSTM takes a lot longer to train and since we could only test one architecture, it is unclear if a model based on this component can achieve significant better results, since our only try was unable to achieve anything. This was a good base work to future experiments on the wikipedia page visit forecast problem.

Possible ways to improve our models are:

- Use embedding to cluster pages with similar behaviors (all models);
- Add dropout layers to improve prediction (all models);
- Use visit measurements from days before 7 days prior, by calculating mean or media of a window of days (MLP model);
- Change hyper parameters on the layer size; (LSTM stateful model);
- Use a different ways to feed inputs to the model, by increasing the time-steps in each backpropagation rollback (LSTM stateful model);
- increase the number of LSTM cells (LSTM stateful model);

REFERENCES

- [1] Aurélien Géron. Hands-on machine learning with scikit-learn and tensorflow: concepts, tools, and techniques to build intelligent systems, 2017.
- [2] Kaggle. Web traffic time series forecasting. <https://www.kaggle.com/c/web-traffic-time-series-forecasting>.
- [3] Kaggle. Web traffic time series forecasting data.