Walmart

M5 Forecast - Accuracy

Big Data Analysis Exam

Date: 10th June 2020 Maria, Mads, Andy and Emil

All group members have contributed equally to the project.

M5 Forecasting - Accuracy

Objective: Predict 28 days of item sales into the Future for 3049 items

Data:

- Historical data for past 1941 days.
- Calendar
 - Public holidays
 - SNAP days
 - etc.
- Prices
 - Week to week prices for all items.

Motivation: 50.000\$



M5 Forecasting - Accuracy



Chosen Models

- Gradient Boosted Trees (LightGBM)
- Long Short-Term Memory (Keras)
- Playtime Graphs!!



Tree Based - feature engineering (LightGBM)



Tree Based - use predictions in features



example item sale only long-time features 12 true predicted 10 8 sold items 6 2 0 1915 1920 1925 1930 1935 1940 Day

simple rolling mean of last 28 days give: item level: MAE = 1.17 store level = MAE = 721



Long Short-Term Memory (LSTM)

simple rolling mean of last 28 days give: item level: MAE = 1.17 store level = MAE = 721





64% of the sales data is 0 (sparse time series).

Problems and Outlook

1. Items sales seem random.

The total sale is ok. Perhaps, unsupervised learning can cluster the time series.

2. Models are underestimating top-selling items.

The data is unbalanced, more weight on the top-selling items.



Graph Structure

1 graph = 1 day

Only total sales prediction

Good for structuring complex data

Consists of:

- Node Feature Matrix
- Edge Feature Matrix
- Adjacency matrix



Graph Neural Network



MAE = 753

Conditioned GNN

Prediction from 1 day **Conditioned GNN prediction** True 6500 -Prediction Last 28 day prediction 6000 plos 5500-5000 - 4500 -Conditioning added 4000 -Only using prices and sales 3500-1920 1925 19'35 19'15 1930 1940 Day

MAE = 556

Conditioned Time Series GNN



MAE = 417

Evaluation of Graph RNN



2 days / 2 graphs:

Node features matrix: $2*3050 \times 6$ Edge feature matrix: $2*12200 \times 2$ Adjacency matrix: $2*3050 \times 2*3050$ Conditioning: $1*2 \times 9$ 28 days / 28 graphs:

Node features matrix: 85.400 x 6 Edge feature matrix: 341.600 x 2 Adjacency matrix: 85.400 x 85.400 Conditioning: 28 x 9

Does require a lot of RAM

We have tried:

- Import and make data for every iteration --> Slow
- Do not use all data
- Bigger computer / GPU (Kaggle)
- Works for a sinus curve

We have not tried:

- Import data in batches
- Use pytorch's dataloader
- Dimensionality reduction
- Feature minimizing (SHAP)

Graph Based GRU

Only use small amount of data, and only use data from 1 day to predict the next.



Use all the data and use data for 28 days at the time. Load data every iteration.

Very slow, especially because it has to use CPU when loading data for every iteration.



Use all the data and use data for 28 days at the time. Contain everything in the memory.

Cannot be contained in memory. Even on Kaggle (16Gb). Benefits had been large because of GPU.

RuntimeError: CUDA out of memory. Tried to allocate 570.00 MiB (GPU 0; 8.00 GiB total capacity; 6.12 GiB already allocated; 132.25 MiB free; 14.61 MiB cached)

Conclusion

Walmart Save money. Live better.

We are not gonna win...

- But second place?

With great flexibility comes great preprocessing and optimization

Method	Mean Absolute Error
LightGBM	255
LightGBM, Prior prediction as feature	194
LSTM	189
LSTM	381
GNN	753
CGNN	556
3-Day CGNN	417

Appendix

- Network architectures
 - Tree Based methods
 - Edge Graph Neural Network
 - Graph Based GRU
 - James Avery's ESN example from website adopted to Walmart Data
- ESN Results
- Table of expected results of the models

Tree Based - only known features

simple rolling mean of last 28 days give: item level: MAE = 1.17 store level = MAE = 721



example item sale including short-time features 12 true predicted 10 8 sold items 2 0 1915 1920 1925 1930 1935 1940 Day



LightGBM models - technical details

- Drop the first 4 years of data (~1.5 mio. data points left)
- Divide data set into
 - training (all days except last 56)
 - validation (next 28 days)
 - test (last 28 days)
- For HP optimization:
 - 20 combinations using random search
 - validation set used to evaluate

LightGBM learning curves



Including short-time features

LightGBM - feature importances Historical sales features at least 28 days back

- Rolling mean over 28 days shifted 28 days is by far the most dominant feature
- features such as item price (sell_price), item id (item_id) and day in the week (wday) used to tune the average guess

	lag28r28						224002245
	sell_price	8238200					
	lag63r28	7744114					
	lag28r7	6649718					
	item_id	5791264					
	wday	4637822					
	day_int	3588290					
	lag56r28	3076746					
	month	2132179					
ures	lag28	2099931					
	lag35r28	1322224					
	mday	929411					
eati	lag63r7	851649					
ш	lag56r7	773626					
	lag35	759929					
	lag35r7	562334					
	dept_id	550017					
	lag56	477856					
	lag63	462229					
ev	ent_type_1	322978					
eve	nt_name_1	135407					
	year	60939					
	cat_id	15082					
	snap_CA	12192					
	(0.0	0.5	1.0 Feature i	1.5 mportance	2.0	1e8

Feature importance

LightGBM - feature importances' Add historical sales data 1 and 7 days ago

Features

 Rolling means are still by far the most dominant features though with the introduction of less shifted features means that these become more important. Here for rolling means shifted (lag) 1 day (lag1r7 and lag1r28).

lag1r7									158177	966
lag1r28	-						118950354			
wday	5139236									
lag1	3541742									
lag7r28	1951006									
lag28r28	1620313									
sell_price	1453865									
lag7r7	1386034									
lag63r28	1264202									
lag35r28	909967									
mday	890863									
lag56r28	859476									
event name 1	779788									
	763802									
item id	756767									
lag28r7	648410									
lag35	476744									
lag63r7	451904									
lag35r7	450786									
lag56r7	436390									
lag7	-432888									
lag28	411359									
lag63	364252									
lag56	345778									
month	321713									
cat id	296413									
dent id	107816									
event type 1	38769									
snap_CA	8249									
(0.0	0.2	0.4	0.6	0.8	1.0	1.2	1.4	1.6	

Feature importance

LSTM - Technical Details

Dividing the data set:

- Training (all days except last 56)
- Validation (next 28 days)
- Test (last 28 days)

Hyperparameter opt:

- The hyperparameter space of LSTM neural networks is huge, and they are generally to train. Thus, only a few architectures were manually tested.
- We also tested different approaches to ingest time-series into the LSTM network. This includes training on the last 720, 360, 180, 56, 28 days and predicting on all 28 days at once. We also tested a sliding window protocol, using a sliding window of length 28 and 56 to predict only one day in advance.

LSTM - Learning Total Sales

LSTM networks were quite good for directly predicting the total sales of the store with a MAE of 189 Items.



Edge Graph Neural Network



Graph Recurrent Neural Network



Hidden Variables: 1500 Connections in. Sparse Matrix: 100 Spectral Radius : 1.5

Echo State Networks

Predicts well on training data

Does not generalize to validation data



Example2













Tested ML-Algorithms

	Tree-based	LSTM	ESN	Graph-based edgeGNN	Graph-based GRU
Implementation	Easy	Difficult	Easy	Very difficult	Very difficult
Performance for few days	Good	Very good	Good	Medium	Very good
Performance for many days	Medium	Good	Good	Medium	Very good
Memory requirements for few days	Medium	Large	Medium	Small	Small
Memory requirements for many days	Large (grows linearly)	Very Large (grows linearly)	Large	Incredible Large (grows > power2)	Incredible Large (grows > power2)
Training Time	Fast	Slow	Fast	Medium	Slow