

Udacity - Machine Learning Engineer Nanodegree

Dhanaji More

06/23/2017

Capstone Project

Predict Future Sales

Table of Contents

INTRODUCTION.....	3
PROJECT OVERVIEW	3
PROBLEM STATEMENT	3
METRICS.....	3
DATA ANALYSIS.....	4
DATA EXPLORATION.....	4
EXPLORATORY VISUALIZATION.....	5
POPULAR ITEM CATEGORIES AT POPULAR STORE HEATMAP	8
BENCHMARK	9
METHODOLOGY	9
DATA PREPROCESSING.....	10
ALGORITHMS AND TECHNIQUES	12
IMPLEMENTATION	13
REFINEMENT	13
RESULTS	14
MODEL EVALUATION AND VALIDATION	14
JUSTIFICATION	14
CONCLUSION.....	14
REFLECTION.....	14
IMPROVEMENT	15
REFERENCES	15

Introduction

Accurately forecasting product sales and building a sales plan is a common problem all enterprises face. Accurate forecasting helps companies.

- avoid inventory buildup problems
- avoid unforeseen cash flow problems and
- effectively manage production, staff and financing needs.
- maximize profit.

Forecasting is an essential activity for managing a business of any size. It is a month-by-month forecast of the level of product or service sales a company expects to achieve.

Project Overview

In this Kaggle competition, we are challenged to predict sales of a product in a store during the next month in sequence. It has a challenging time-series dataset consisting of daily sales data, kindly provided by one of the largest Russian software firms - 1C Company.

We are provided with daily historical sales data. Our objective is to forecast the total amount of products sold in shops for items in the 35th month in the test set.

<https://www.kaggle.com/c/competitive-data-science-predict-future-sales#description>

We are provided 34 months of daily sales data from 60 stores that have sold over ~2.9 million items. The lists of shops and products changes slightly every month. We have to create a robust model that can handle such situations is part of the challenge.

Problem Statement

The problem is a linear regression forecasting at its core. For each shop_id and item_id pair in the test set we are required to predict the total item sales count.

Metrics

The accuracy of the model will be evaluated based on root mean squared error (RMSE).

$$RMSE_{errors} = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}}$$

\hat{y}_i – is the predicted value of item sales count.

y_i – is the actual value of item sales count.

n – total number of observations.

Data Analysis

Data Exploration

Dataset consists of four relational tables as standard csv files:

Data Filename	Description
sales_train.csv	Training set. Daily historical data from Jan 2013 to Oct 2015.
items.csv	Supplemental information about the items/products.
item_categories.csv	Supplemental information about the item categories.
shops.csv	Supplemental information about the shops.

date	date_block_num	shop_id	item_id	item_price	item_cnt_day
0 02.01.2013	0	59	22154	999.00	1.0
1 03.01.2013	0	25	2552	899.00	1.0
2 05.01.2013	0	25	2552	899.00	-1.0
3 06.01.2013	0	25	2554	1709.05	1.0
4 15.01.2013	0	25	2555	1099.00	1.0

sales_train

shop_name	shop_id
0 !Якутск Орджоникидзе, 56 фран	0
1 !Якутск ТЦ "Центральный" фран	1
2 Адыгея ТЦ "Мега"	2
3 Балашиха ТРК "Октябрь-Киномир"	3
4 Волжский ТЦ "Волга Молл"	4

shops

item_name	item_id	item_category_id
0 ! ВО ВЛАСТИ НАВАЖДЕНИЯ (ПЛАСТ.) D	0	40
1 !ABBY FineReader 12 Professional Edition Full...	1	76
2 ***В ЛУЧАХ СЛАВЫ (UNV) D	2	40
3 ***ГОЛУБАЯ ВОЛНА (Univ) D	3	40
4 ***КОРОБКА (СТЕКЛО) D	4	40

items

item_category_name	item_category_id
0 PC - Гарнитуры/Наушники	0
1 Аксессуары - PS2	1
2 Аксессуары - PS3	2
3 Аксессуары - PS4	3
4 Аксессуары - PSP	4

item_categories

Dataset observation and analysis:

a. Dataset stats

- Shops - 60
- Item Categories - 84
- Items - 22,170
- Median item price - 399
- Min item price - (-1.0)
- Max item price - 307,980
- Total Rows - 2,935,849

- b. Data has no missing values.
- c. Price has negative values
- d. item_sales_count field in sales_train has negative values indicating this maybe returns.

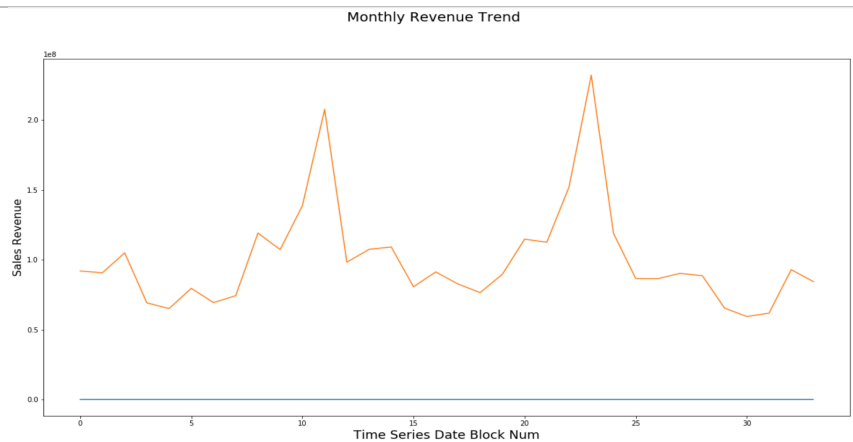
Some interesting information from test set analysis:

- a. Not all shop_id in training set are used in test set. Test set excludes following shops (but not vice versa): [0, 1, 8, 9, 11, 13, 17, 20, 23, 27, 29, 30, 32, 33, 40, 43, 51, 54]
- b. Not all item in train set are in test set and vice versa
- c. In test set, a fixed set of items (5100) are used for each shop_id, and each item only appears one per each shop. This possibly means that items are picked from a generator, which will result in lots of 0 for item count.

Exploratory Visualization

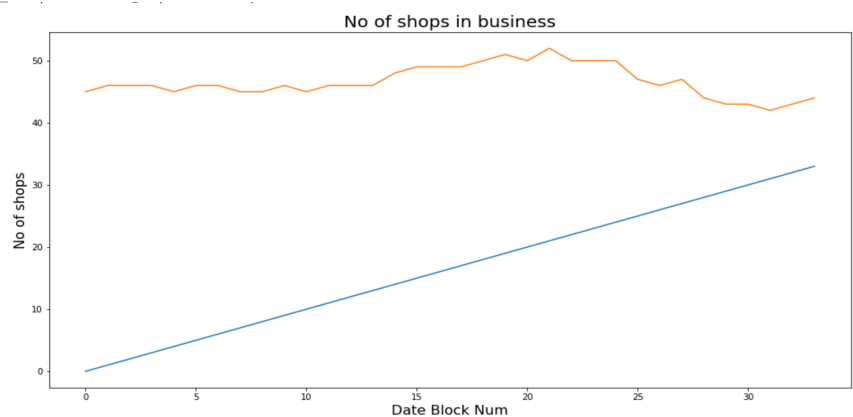
Monthly revenue trend indicates.

- 1. Company has obvious seasonality in their sales cycle during Nov-Dec.
- 2. Sales over the years increased in first 2 years but during its third year has started to slow down.



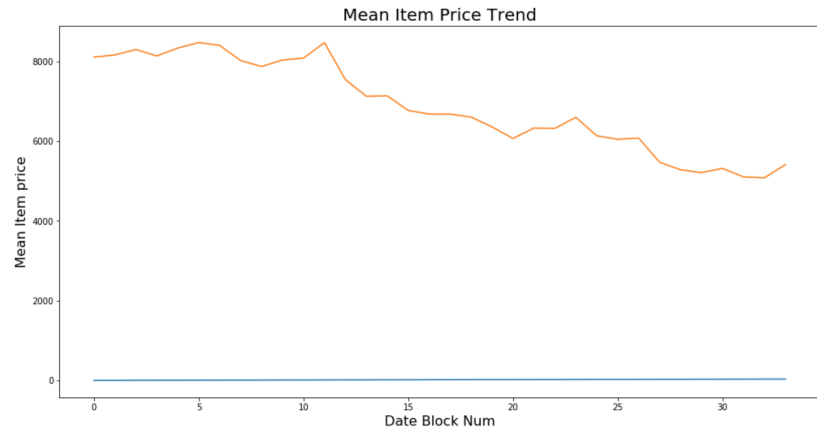
Number of stores in business

- 1. It has stayed relatively stable over the years.



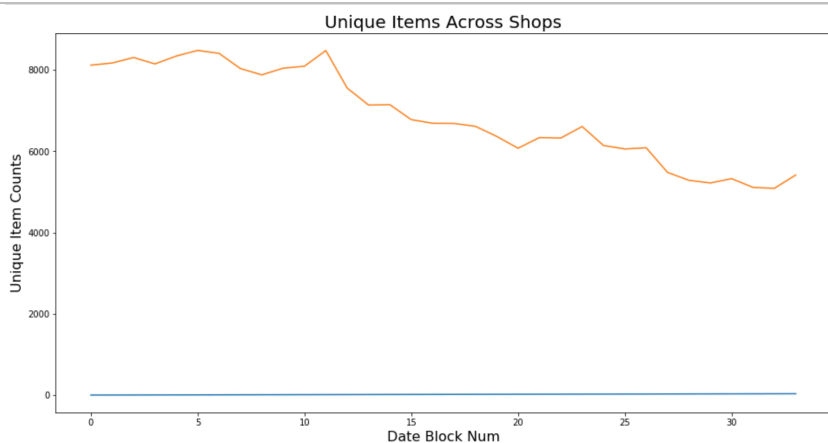
Mean item price trend

- Declining mean price indicates `company is either reducing prices to retain customers or to move inventory or company may be selling less premium products.



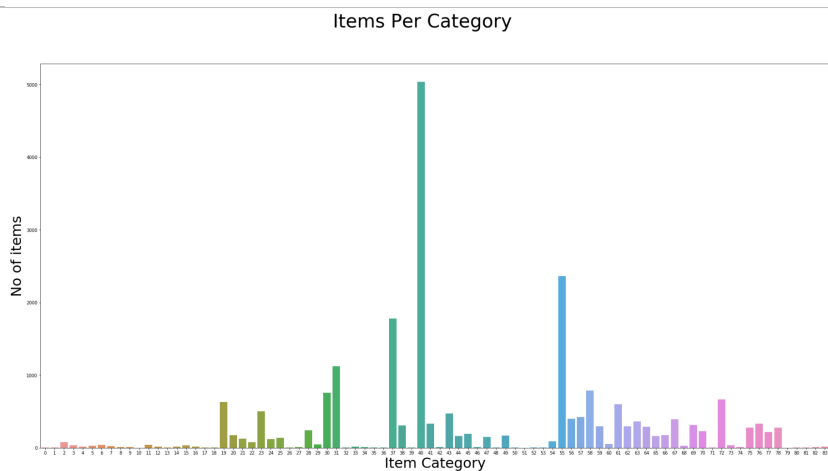
Unique Items sold across shops

- Declining may indicate customers are losing interests and company is unable to introduce new item categories and items.



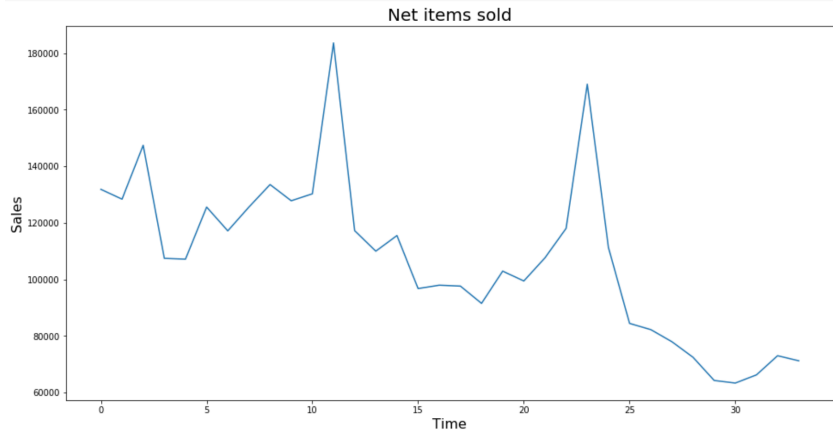
Items per Category shows

- Out of 84 item categories there are handful of categories that have significant number of items. There are 5000 items in category_id 40.



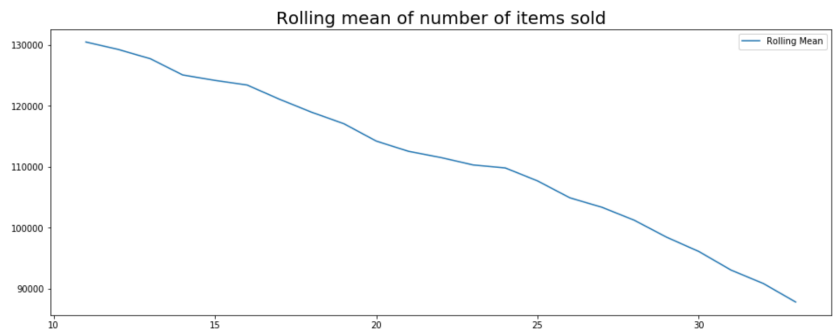
Number of sold items trend shows

1. Declining sales.
2. There are peaks in Nov-Dec and some seasonal behaviors during June-July-August. Seasonality maybe likely due to some national events such as – national holidays, or some other seasonal influence during these times of the year.



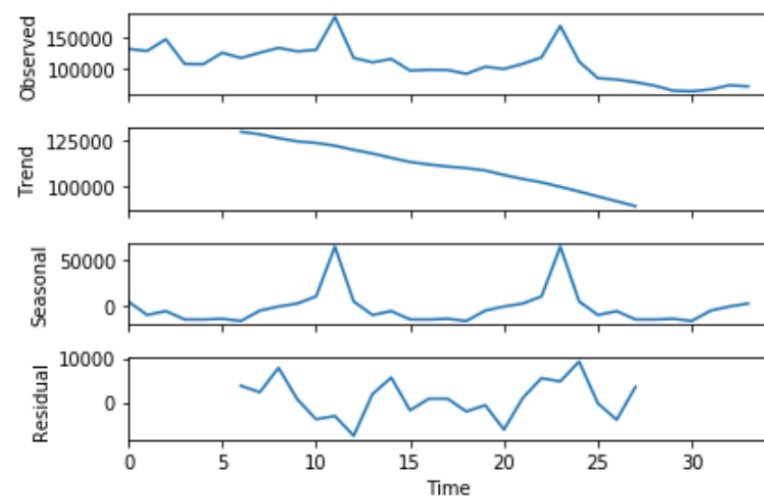
Rolling mean of number of sold items indicates.

1. Declining sales for this business over the last 3 years.



Here I have decomposed sales data into Trend, seasonality and residuals - using Statsmodels, a python library available for statisticians.

1. We see obvious "seasonality" trend
2. Sales peak around a certain time of the year - right around October and November.



Statsmodel provide additive as well as multiplicative model.

- The additive model is useful when the seasonal variation is relatively constant over time.
- The multiplicative model is useful when the seasonal variation increases over time.



Item categories & shops heatmap

Benchmark

We will use following 2 benchmark models.

1. Our base benchmark model will predict next month sales based on the rolling average sales of the previous similar months from the dataset.
2. Our target benchmark model to beat will be based on Prophet, a Facebook forecasting tool available in Python.

Prophet forecasting tool is an additive regression model with four main components:

- Prophet automatically detects changes in trends by selecting changepoints from the data.
- Prophet has yearly seasonal component modeled using Fourier series.
- Prophet has weekly seasonal component using dummy variables.
- Prophet has weekly can accept user-provided list of important holidays.

Comparing our results with above 2 benchmark models will allow us to figure out how well is our prediction model.

Limitations of Prophet forecasting tool –

- a. It can only perform univariate analysis.
- b. It takes 2 variable inputs – a datetime series (ds) and time sequence output y (corresponding numeric output).
- c. For our project, we will aggregate all the monthly sales count as y and all date_block_num into sequential months starting with 0 to 33.

Implementation of benchmark model are provided in the associated jupyter notebook – [PredictSales – Baseline.ipynb](#)

Methodology

Data consists of daily sales for 60 stores, that sell 22,140 unique items in 84 different categories. We are given 34 months of daily sales data ~2.9 million rows.

Contrary to most established forecasting models used in the industry which are based on Moving Averages such as ARIMA (Auto Regressive Integrated Moving Average), I

have attempted to use an almost fully automated approach based on deep neural networks.

The architecture I have used is a simple deep neural network that takes in a sequence of sales data aggregated monthly for every single shop and item pair for all months and predicts the sales for every shop-item pair during the next month in sequence.

My approach could easily be adapted to other applications in which the goal is to predict a fixed-length output from a variable-length sequence.

The key to this methodology is to systematically engineer lagging sales indicator (item_sales_cnt for previous month) for each shop-item pair for every month as input data for deep neural network. Key idea here is to present neural network sales trend for every entry of shop_id and item_id pair and letting network figure out all the hidden patterns.

Data Preprocessing

- I. Daily sales data that is made available to train our model looks as below –

	date	date_block_num	shop_id	item_id	item_price	item_cnt_day
0	02.01.2013	0	59	22154	999.00	1.0
1	03.01.2013	0	25	2552	899.00	1.0
2	05.01.2013	0	25	2552	899.00	-1.0
3	06.01.2013	0	25	2554	1709.05	1.0
4	15.01.2013	0	25	2555	1099.00	1.0

- II. For data exploration and analysis, we will break down the date feature into individual feature components to explore daily, weekly, monthly, quarterly and yearly seasonality.

	date_block_num	shop_id	item_id	item_price	item_cnt_day	year	month	day	dayofyear	weekofyear	dayofweek	quarter
7554	0	19	18976	399.0	1.0	2013	1	1	1	1	1	1
7644	0	19	18284	199.0	1.0	2013	1	1	1	1	1	1
7646	0	19	18320	199.0	1.0	2013	1	1	1	1	1	1
7647	0	19	18329	299.0	1.0	2013	1	1	1	1	1	1
7694	0	19	19367	399.0	1.0	2013	1	1	1	1	1	1

- III. Since we have to predict monthly sales numbers, we will aggregate daily sales data into monthly sales numbers before providing it as input to deep neural net model. We will also aggregate shop-item pairs for every month (date_block_num) and calculate the mean price of the item at that shop during that month(date_block_num).

	shop_id	item_id	date_block_num	month	year	item_cnt_day	item_price
103792	5	13334	33	10	2015	4.0	359.0
103791	5	13331	33	10	2015	1.0	399.0
103790	5	13331	33	7	2015	1.0	399.0
723715	25	13332	33	10	2015	6.0	1399.0
1760607	50	20742	33	10	2015	1.0	499.0

IV. We have used two separate deep neural network model implementations.

- a. DNN with categorical entity embeddings.
- b. DNN with monthly sales count lag features.

V. Categorical entity embeddings

In this approach we will take daily sales data which has date feature broken down into individual features – month, day, dayofyear, weekofyear, dayofweek and quarter and submit it to a sequential model with first layer being the embedding layer.

Our input will look as below and item_sales_cnt is our target.

	date_block_num	shop_id	item_id	item_price	item_cnt_day	year	month	day	dayofyear	weekofyear	dayofweek	quarter
7554	0	19	18976	399.0	1.0	2013	1	1	1	1	1	1
7644	0	19	18284	199.0	1.0	2013	1	1	1	1	1	1
7646	0	19	18320	199.0	1.0	2013	1	1	1	1	1	1
7647	0	19	18329	299.0	1.0	2013	1	1	1	1	1	1
7694	0	19	19367	399.0	1.0	2013	1	1	1	1	1	1

VI. Monthly sales count lag features

Lag features are values at prior time steps. We will generate lag features based on 'item_cnt_day' and grouped by 'shop_id' and 'item_id'. We will generate item sales count for every shop-id, item_id pair for all the months.

These lag features are the key engineered features we want our model to learn from.

Lag feature explanation –

For the very first month (date_block_num=0) – all the lag features (0,1,2,3,...33) will be 0 because at that point in time we will have no prior sales count information. So all rows of date_block_num=0 will have lag features set to 0. For the second month rows (date_block=1), we will populate the lag feature – '0' which indicates the previous month sales count, we will set the current month lag feature "1" and all future months (2 to 33) to 0 since we have no information about it.

Finally by 33rd month – we should have item_sales_cnt in 0,1,2,3,...32 for all rows where date_block_num=32 and column '33' lag feature will be 0.

Sample input data (X) to deep neural net should look as below –

	shop_id	item_id	date_block_num	month	year	item_price	0	1	2	3	...	24	25	26	27	28	29	30	31	32	33
0	10	4885	0	1	2013	716.0	0	0	0	0	...	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1	10	4885	1	2	2013	716.0	1	0	0	0	...	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2	10	4885	3	4	2013	716.0	1	2	0	0	...	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	10	4885	3	7	2013	716.0	1	2	0	0	...	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	10	4885	4	5	2013	716.0	1	2	0	2	...	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Our target variable will be item_sales_cnt - sales count for the next month in sequence.

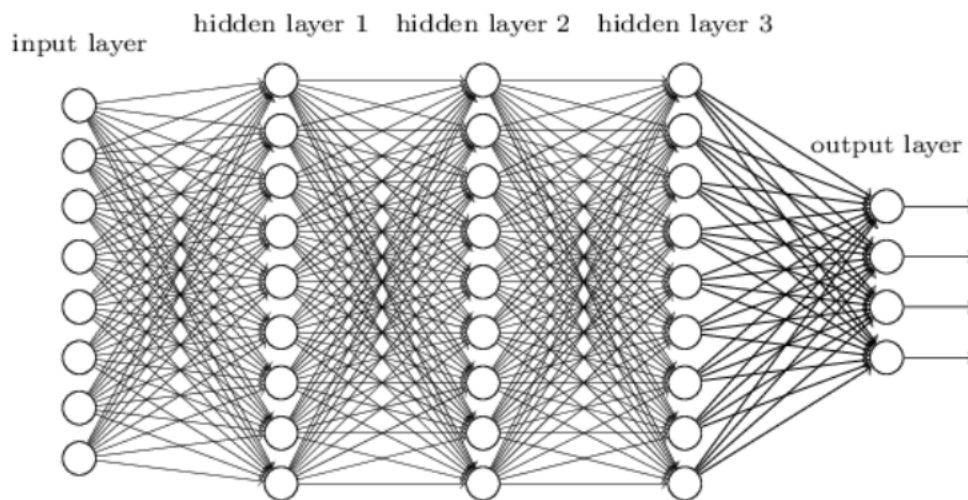
Algorithms and Techniques

Common techniques used in the industry for demand forecasting involve techniques including both informal methods, such as educated guesses, and quantitative methods, such as the use of historical sales data and statistical techniques or current data. For the purpose of this project we will focus on quantitative method specifically using deep neural network technique to build our forecasting model.

Generally, such a problem is modeled using linear regression techniques and since deep neural network can be thought of as a linear regression at its very basic perceptron level, I would like to build an automated forecasting tool that can learn from historical data patterns to predict demand for next time step in sequence.

Also there are many Kaggle submissions with different linear-regression techniques but my goal was to try DNN for sparse tabular data.

Deep neural network



Source - <http://neuralnetworksanddeeplearning.com/chap5.html>

Implementation

Implementation of benchmark model are provided in the associated jupyter notebook – PredictSales – DeepNeuralNetwork.ipynb

Implementation of Neural Network with Embeddings model are provided in the associated jupyter notebook – PredictSales – DNN-CategoricalEmbeddings.ipynb

Implementation of Neural Network with Embeddings model are provided in the associated jupyter notebook – PredictSales – DeepNeuralNetwork.ipynb

Refinement

Of the two neural network models, we tried, neural network with lagging sales indicator appeared to be comparable to benchmark and none of the methods or various tuning of neural network improved the RMSE score.

To further refine the model,

1. We can add Russian holidays as additional one hot encoded features.
2. Since new shop_id's start showing up randomly, we can try to identify that as a separate feature to see if it is important for the model. Correspondingly we can try to identify shops that have not generated any sales for more than 6 months and try to mark them as closed.

3. We can translate shop names and item names to English and perform “most favored shop” analysis based on most to least revenue contribution by shop to company’s total revenue.

Results

Model Evaluation and Validation

Below we have evaluated models based on the predicted sales count and actual sales count for last 6 months of data in the set.

Model	Root Mean Square Error	Mean Square Error	Absolute Error
Moving Average	n/a	n/a	4546
Prophet	n/a	n/a	15753
DNN /w Entity Embeddings	9.75	95.16	16240
DNN /w Lag Features	6.067	38.81	4390

Justification

DNN with lag features fared better compared to DNN with Categorical embeddings, Moving Average and Prophet. Moving average was the quickest and least cumbersome.

The reason DNN better fared maybe due to monthly lag features captured the sales pattern for shop-item over 33 months better than just date features embedded in daily sales data.

Conclusion

Reflection

In this project, my goal was to apply techniques that are generally popular in image recognition and classification problems with high dimensional feature set. In this project we tried extracting features from historical sales figures. Our goal was to come with feature to capture trends and seasonality etc. into every row that go as input to NN. Our final model fell below my expectations, but I believe there is lot of room for improvement. If we can engineer features that capture knowledge and intuition of people responsible for forecasting, we can come up with a NN based forecasting tool that will come very close to something that has both qualitative as well quantitative expertise to predict accurately.

My number one challenge was the size of the data and using pandas to perform complex join, grouping, aggregation. Framing the data, the way I was intending to be was way more difficult than I anticipated. Coming from database and data background I found working with pandas challenging and cumbersome but at the same time it was very rewarding. I now have an additional toolset under my belt. I discovered the whole process of preparing data and feature engineering as more than half the battle for such problems.

My ultimate goal was to try emerging techniques on age old problem in the industry using structured data. Enterprises today have tons of structured data that are hiding patterns that we can uncover using powerful technique such as NN.

Improvement

There were two particular improvements we can try with this project.

1. Define a model that will concatenate categorical embedding layer with dense layer.

With 60 shops and 22k items and various categorical date features made feature set huge and model very slow but it did not improve the model's prediction accuracy. I believe categorical embedding layer on entire input set wasn't as helpful as I initially anticipated.

Instead we should try to build a categorical embedding layer of just shop_id, item_id, item_price and dayofyear and then concatenate rest of the features before submitting as input to dense layers.

2. Neural network model with just lagging sales count indicator doesn't capture all the nuances, additional features such holidays, shops popularity and items popularity etc., elapsed time between when an item is sold at the same shop again may also improve the NN model.
3. LSTM/RNN based network model may better capture the time-series pattern.

References

- An introduction to Deep Learning for Tabular Data - <http://www.fast.ai/2018/04/29/categorical-embeddings/>
- An Introductory Study of Time-Series Modeling and Forecasting - <https://arxiv.org/pdf/1302.6613>
- Understanding LSTM Networks - <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

- Avoiding Look Ahead Bias in Time Series Modelling
- <https://www.datasciencecentral.com/profiles/blogs/avoiding-look-ahead-bias-in-time-series-modelling-1>
- Prophet: forecasting at scale - <https://research.fb.com/prophet-forecasting-at-scale/>
- An Introduction to Time Series Forecasting with Prophet Package in Exploratory
- <https://blog.exploratory.io/an-introduction-to-time-series-forecasting-with-prophet-package-in-exploratory-129ed0c12112>