

Food Demand Forecasting Challenge

Abstract

Demand forecasting is one of the most common time series problems in today's world holding great significance in several business applications and much more. In this project, we take up a food demand forecasting problem where the goal is to predict the demand for the following 10 weeks from the previously given data. For the purpose of this problem multiple models were implemented of which boosting models proved to produce the best results. The LightGBM model gave us the best performance with an RMLSE of 0.51 placing us in the top 50 in the leaderboard of this forecasting challenge.

Introduction

Demand forecasting is a key component to every growing online business. Without proper demand forecasting processes in place, it can be nearly impossible to have the right amount of stock on hand at any given time. A food delivery service has to deal with a lot of perishable raw materials which makes it all the more important for such a company to accurately forecast daily and weekly demand.

Too much inventory in the warehouse means more risk of wastage, and not enough could lead to out-of-stocks – and push customers to seek solutions from your competitors. Therefore, reliable forecasting is highly essential and models must be created to achieve such levels of performance.

Objective

In this challenge we have a client as a meal delivery company which operates in multiple cities, they have various fulfillment centers in these cities for dispatching meal orders to their customers. The client wants us to help these centers with demand forecasting for upcoming weeks so that these centers will plan the stock of raw materials accordingly.

The replenishment of majority of raw materials is done on weekly basis and since the raw material is perishable, the procurement planning is of utmost importance. Secondly, staffing of the centers is also one area wherein accurate demand forecasts are really helpful.

Given the following information, the task is to predict the demand for the next 10 weeks (Weeks: 146-155) for the center-meal combinations in the test set:

- Historical data of demand for a product-center combination (Weeks: 1 to 145)
- Product(Meal) features such as category, sub-category, current price and discount
- Information for fulfillment center like center area, city information etc.

Data Dictionary

1. **Weekly Demand data (train.csv):** Contains the historical demand data for all centers, test.csv contains all the following features except the target variable

Variable	Definition
id	Unique ID
week	Week No
center_id	Unique ID for fulfillment center
meal_id	Unique ID for Meal
checkout_price	Final price including discount, taxes & delivery charges
base_price	Base price of the meal
emailer_for_promotion	Emailer sent for promotion of meal
homepage_featured	Meal featured at homepage
num_orders	(Target) Orders Count

2. **fulfilment_center_info.csv**: Contains information for each fulfilment center

Variable	Definition
center_id	Unique ID for fulfillment center
city_code	Unique code for city
region_code	Unique code for region
center_type	Anonymized center type
op_area	Area of operation (in km ²)

3. **meal_info.csv**: Contains information for each meal being served

Variable	Definition
meal_id	Unique ID for the meal
category	Type of meal (beverages/snacks/soups....)
cuisine	Meal cuisine (Indian/Italian/...)

Evaluation Metric

The evaluation metric for this competition is 100*RMSLE where RMSLE is Root of Mean Squared Logarithmic Error across all entries in the test set.

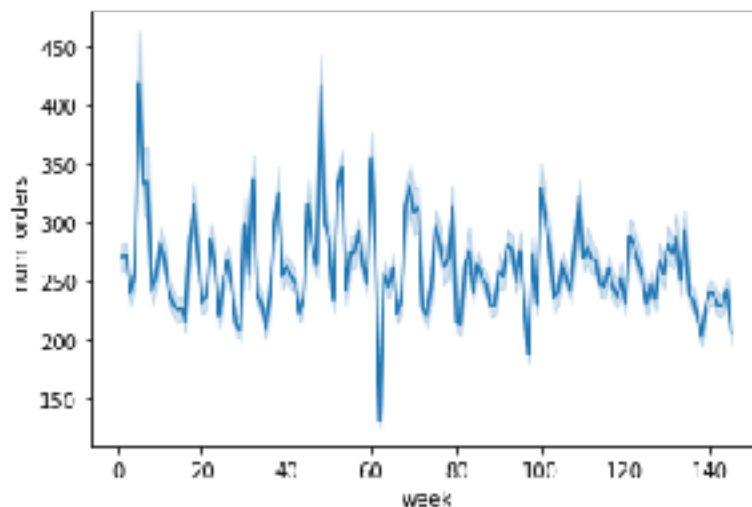
$$\sqrt{\frac{1}{n} \sum_{i=1}^n (\log(x_i+1) - \log(y_i+1))^2}$$

Methodology

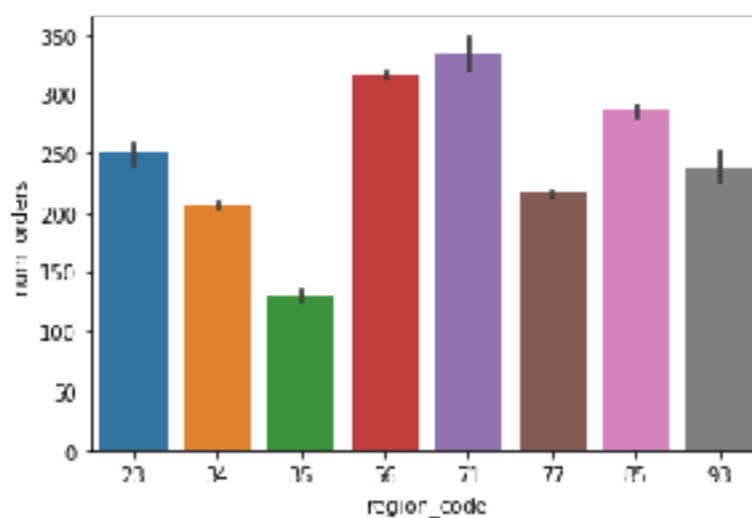
Several different approaches and preprocessing techniques were implemented for the purpose of this challenge some of which helped achieve good results while others didn't turn out to be that useful or rather even degrade performance.

A good amount of preprocessing and exploratory analysis was required to be done on the data before fitting any sort of model. Our exploratory analysis involved various visualizations to identify correlations amongst the various variables and to identify any trends present in the data.

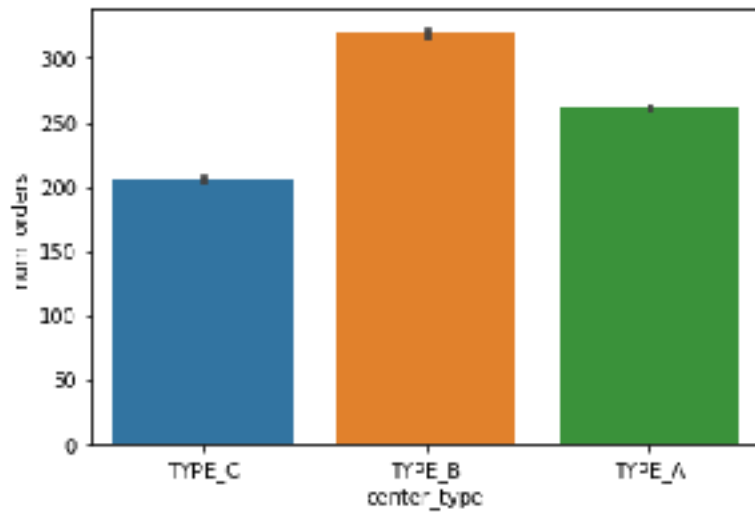
Visualization:



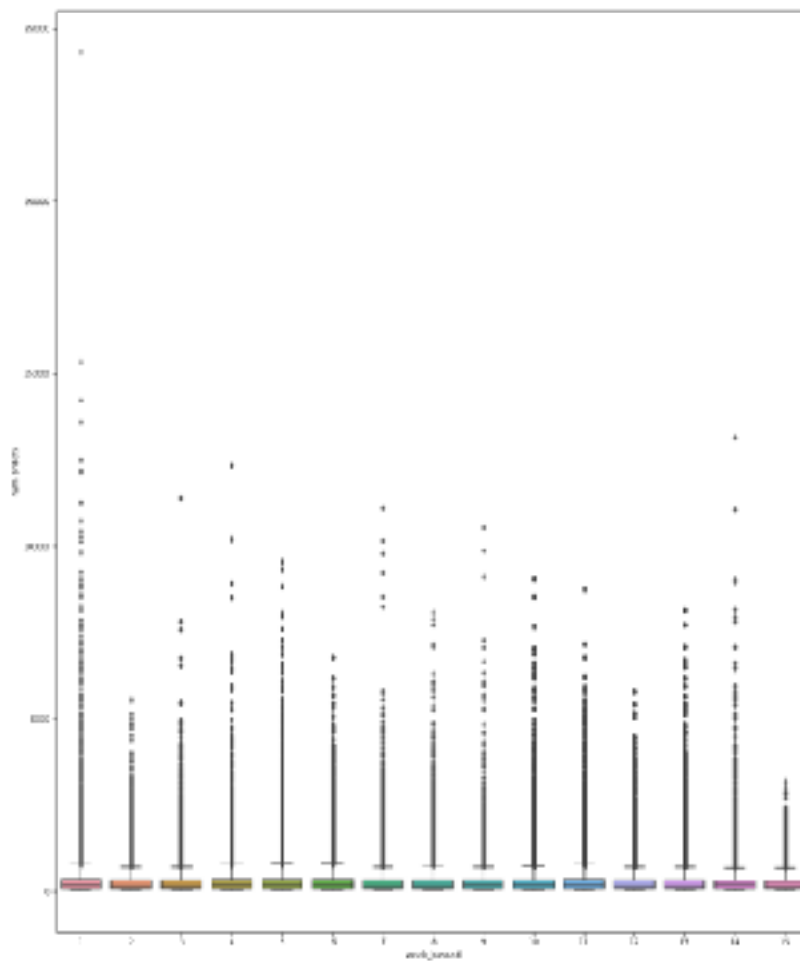
This is a lineplot showing the trend of the number of orders over the weeks of data that was provided to us. We can see a trend similar to a Sin/Cos curve which inspired used to perform the sine transforms on the data.



The above graph shows the num of orders in the different regions.



The graph above shows the number of orders placed in different center types. We can see that type B has the most orders and the centers belonging to center B must have higher stocks of raw materials in comparison to the other centers.



The plot above is the boxplot of the data to understand the distribution. As we can see the data is highly skewed due to the presence of a lot of outliers in it.

Pre-processing:

Outlier detection by means of a plot and removal was performed while ensuring minimal data loss.

The number of orders placed (target variable) was found to be highly right skewed and so we apply a log transformation.

The feature 'price_last_curr_diff' had a few missing values which were filled with the mean.

Categorical variables namely category, cuisine, center_type, center_id, meal_id, city_region, price_increase and profit/loss were all label encoded. A one-hot encoder was used on label encoded features including emailer_for_promotion and homepage_featured. Feature scaling was carried out on the other numerical attributes using Standard Scaler.

Dimensionality reduction techniques including PCA and LDA were applied and tested however they resulted in poorer performance.

Many boosting models such as Catboost and LGBM internally take care of categorical features on their own and therefore our preprocessing is only applied to the appropriate model as and how required.

Feature Engineering:

Apart from the given attributes we have carried out some feature engineering to create a few extra features which helped model the data better. These include profit, pricediff, profit/loss, price_last_curr_diff and price_increase.

As this is a time series problem we found it useful and rather important to create a few lag features. Lag features are variables which contain data from prior time steps. We can convert time-series data into rows where every row contains data about one observation and includes all previous occurrences of that observation. We have created lag features for pricediff and checkout_price for preceding 10, 11 and 12 weeks.

Sine and cos transformations of the weeks attribute was done to create two features called week_sin and week_cos respectively. The purpose of this was to help capture the trend of the time-series more accurately therefore achieve better results when fit in our model.

Train-Test Split:

We use a train test split of 80-20 to evaluate our model. The training was done on data from weeks 1-135 and the testing was done on weeks 136-145. This method of training in the weeks of data and testing on the last few weeks in chronological order as the task at hand is time series forecasting. This method of validation helped us in generalizing the model better.

Models:

A few of the several models we have used for the given problem statement include:

- Random Forest
- CatBoost
- XGBoost
- LightGBM

Random Forest

Random Forest is basically bagging of decision trees. Random Forests provide an improvement over bagged trees by way of small tweak that decorrelates the trees. This reduces the variance when we average the trees. In bagging we build a number of decision trees bootstrapped training samples, but when building these decision trees, each time a split in a tree is considered, a random selection of 'm' predictors is chosen as split candidates from the full set of 'p' predictors. The split allows us to use only one of those m predictors.

A fresh selection of m predictors is taken at each split, and typically the number of predictors considered at each split is approximately equal to the square root of the total number of predictors.

Parameters tuned:

- N estimators = 10
- Random state = 0
- Verbose = 2
- N jobs = 1
- Min samples leaf = 5
- Max depth = 20
- Min samples split = 16
- OOB(Out-of-Bag) Score = true

CatBoost

CatBoost is a recently open-sourced machine learning algorithm from Yandex. It can work with diverse data types to help solve a wide range of problems. To top it up, it provides best-in-class accuracy.

It is especially powerful in two ways:

- It yields state-of-the-art results without extensive data training typically required by other machine learning methods, and
- Provides powerful out-of-the-box support for the more descriptive data formats that accompany many business problems.

“CatBoost” name comes from two words “Category” and “Boosting”.

We are using CatBoost mainly for three reasons:

Performance: CatBoost provides state of the art results and it is competitive with any leading machine learning algorithm on the performance front.

Handling Categorical features automatically: We can use CatBoost without any explicit pre-processing to convert categories into numbers. CatBoost converts categorical values into numbers using various statistics on combinations of categorical features and combinations of categorical and numerical features.

Robust: It reduces the need for extensive hyper-parameter tuning and lower the chances of overfitting also which leads to more generalized models.

In addition to this, CatBoost does not require conversion of data set to any specific format like XGBoost and LightGBM.

Parameters tuned:

- Iterations = 625
- Learning rate = 0.06
- Depth = 8
- l2 leaf reg = 10
- Loss function = 'RMSE'
- Random seed = 2018

XGBoost

XGBoost stands for eXtreme Gradient Boosting. It is a popular ensemble learning technique which implements gradient boosting. XGBoost are also highly efficient in execution speed and over all model performance.

We use XGBoost mainly for three reasons:

Regularization: XGBoost has an option to penalize complex models through both L1 and L2 regularization. Regularization helps in preventing overfitting

Handling sparse data: Missing values or data processing steps like one-hot encoding make data sparse. XGBoost incorporates a sparsity-aware split finding algorithm to handle different types of sparsity patterns in the data

Weighted quantile sketch: Most existing tree-based algorithms can find the split points when the data points are of equal weights (using quantile sketch algorithm). However, they are not equipped to handle weighted data. XGBoost has a distributed weighted quantile sketch algorithm to effectively handle weighted data

Parameters tuned:

- Max depth = 50
- Min child weight = 1
- N estimators = 200
- N jobs = -1
- Verbosity = 2
- Learning rate = 0.16

LightGBM

Light GBM is a fast, distributed, high-performance gradient boosting framework based on decision tree algorithm, used for ranking, classification and many other machine learning tasks.

Since it is based on decision tree algorithms, it splits the tree leaf wise with the best fit whereas other boosting algorithms split the tree depth wise or level wise rather than leaf-wise. So, when growing on the same leaf in Light GBM, the leaf-wise algorithm can reduce more loss than the level-wise algorithm and hence results in much better accuracy which can rarely be achieved by any of the existing boosting algorithms.

We use LightGBM mainly for three reasons:

Faster training speed and higher efficiency: Light GBM use histogram-based algorithm i.e. it buckets continuous feature values into discrete bins which fasten the training procedure.

Better accuracy than any other boosting algorithm: It produces much more complex trees by following leaf wise split approach rather than a level-wise approach which is the main factor in achieving higher accuracy. However, it can sometimes lead to overfitting which can be avoided by setting the max_depth parameter.

Compatibility with Large Datasets: It is capable of performing equally good with large datasets with a significant reduction in training time as compared to XGBOOST.

Parameters tuned:

- Learning rate = 0.003
- N estimators = 40000
- Silent = False
- **g

How each model treats Categorical Variables

- Catboost vs. XGBoost vs. LightGBM

CatBoost

CatBoost has the flexibility of giving indices of categorical columns so that it can be encoded as one-hot encoding using `one_hot_max_size` (Use one-hot encoding for all features with number of different values less than or equal to the given parameter value). If you don't pass anything in `cat_features` argument, CatBoost will treat all the columns as numerical variables.

LightGBM

Similar to CatBoost, LightGBM can also handle categorical features by taking the input of feature names. It does not convert to one-hot coding, and is much faster than one-hot coding. LGBM uses a special algorithm to find the split value of categorical features.

XGBoost

Unlike CatBoost or LGBM, XGBoost cannot handle categorical features by itself, it only accepts numerical values similar to Random Forest. Therefore, one has to perform various encodings like label encoding, mean encoding or one-hot encoding before supplying categorical data to XGBoost.

CODE:

Reading data and merging tables

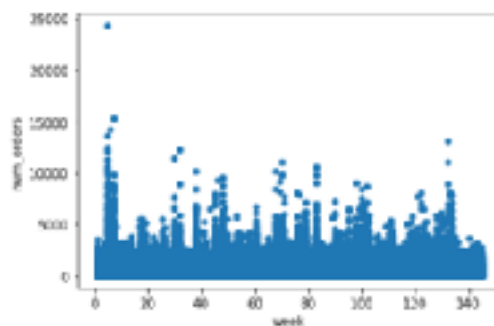
The tables have been read into pandas data frame and have been merged into a single one.

```
In [2]: df.head()
Out[2]:
```

	id	week	order_id	week_id	channel_name	base_price	actual_price	margin	is_promotion	homepage_featured	num_orders	city_code	region_code	order_type
0	1436942	1	90	1002	88.58	131.00		0	0	0	365.0	630	00	TYPE_J
1	1204090	2	90	1002	88.06	131.00		0	0	0	702.0	630	00	TYPE_J
2	1417751	1	90	1002	88.58	131.00		0	0	0	361.0	630	00	TYPE_J
3	1314660	4	90	1002	88.06	131.00		0	0	0	1002.0	630	00	TYPE_J
4	1307020	0	90	1002	88.06	131.00		0	0	0	300.0	630	00	TYPE_J

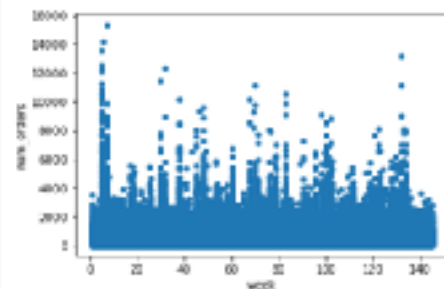
Outlier removal

```
In [234]: df.plot.scatter(x='week', y='num_orders')
Out[235]: <matplotlib.axes._subplots.AxesSubplot at 0x4a2f9502e0>
```



```
In [6]: df.drop(df[df['num_orders'] > 20000].index, inplace = True)
```

```
In [237]: df.plot.scatter(x='week', y='num_orders')
Out[237]: <matplotlib.axes._subplots.AxesSubplot at 0x4a4a6571b>
```



Feature Engineering

Creation of features like profit,pricediff,profit/loss,price_diff,etc using the already existent features.

```

In [8]: df['profit'] = df['checkout_price'] - df['base_price']
df['price_diff'] = df['base_price'] - df['checkout_price']
df['profit_loss'] = (df['profit'] > 0).astype(int)
df['price_last_curr_diff'] = (df['checkout_price'].shift(1) - df['checkout_price']).fillna(1) / df['checkout_price'].shift(1).fillna(1)
df['price_last_curr_diff'] = df['price_last_curr_diff'].fillna(0)
df['price_increase'] = (df['price_last_curr_diff'] < 0).astype(int)

```

Creation of sales lag features for time series forecasting.

```

In [9]: def create_sales_lag_feats(df, gby_cols, target_col, lags):
    gby = df.groupby(gby_cols)
    for i in lags:
        df['_'.join([target_col, 'lag', str(i)])] = \
            gby[target_col].shift(i).values * np.random.normal(scale=1.6, size=[len(df),])
    return df

In [10]: df = create_sales_lag_feats(df, gby_cols=['center_id', 'meal_id'], target_col='price_diff', lags=[10, 11, 12])
df = create_sales_lag_feats(df, gby_cols=['center_id', 'meal_id'], target_col='checkout_price', lags=[10, 11, 12])

```

Merging categorical variables

```

In [16]: df_train['city_region'] = \
    df_train['city_code'].astype('str') + '_' + \
    df_train['region_code'].astype('str')

df_test['city_region'] = \
    df_test['city_code'].astype('str') + '_' + \
    df_test['region_code'].astype('str')

C:\Users\Amrith\Anaconda3\lib\site-packages\ipykernel_launcher.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy
This is separate from the ipykernel package so we can avoid doing imports until
C:\Users\Amrith\Anaconda3\lib\site-packages\ipykernel_launcher.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy
import sys

```

Applying sin and cos transform on weeks

```

In [5]: df_train['week_sin'] = \
    np.sin(2 * np.pi * df_train['week'] / 52.143)
df_train['week_cos'] = \
    np.cos(2 * np.pi * df_train['week'] / 52.143)

df_test['week_sin'] = \
    np.sin(2 * np.pi * df_test['week'] / 52.143)
df_test['week_cos'] = \
    np.cos(2 * np.pi * df_test['week'] / 52.143)

```

Log transformation of the target variable and sorting the final dataframe by weeks

```

In [7]: df_train['num_orders_log1p'] = np.log1p(df_train['num_orders'])

C:\Users\Amrith\Anaconda3\lib\site-packages\ipykernel_launcher.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy
"" entry point for launching an jupyter kernel.

In [8]: df = df_train.append(df_test).reset_index(drop=True)[df_train.columns]
df = df.sort_values(['center_id', 'meal_id', 'week']).reset_index(drop=True)

```

Preprocessing

Handling of null values

```
b1: df['price_last_curr_diff'] = df['price_last_curr_diff'].fillna(df['price_last_curr_diff'].mean())
df['pricediff_lag_10'] = df['pricediff_lag_10'].fillna(0)
df['pricediff_lag_11'] = df['pricediff_lag_11'].fillna(0)
df['pricediff_lag_12'] = df['pricediff_lag_12'].fillna(0)
df['checkout_price_lag_10'] = df['checkout_price_lag_10'].fillna(0)
df['checkout_price_lag_11'] = df['checkout_price_lag_11'].fillna(0)
df['checkout_price_lag_12'] = df['checkout_price_lag_12'].fillna(0)
```

Standard Scaling and Label Encoding

```
14 [20]: from sklearn.preprocessing import LabelEncoder, StandardScaler
        le1=LabelEncoder()
        le2=LabelEncoder()
        le3=LabelEncoder()
        le4=LabelEncoder()
        le5=LabelEncoder()
        le6=LabelEncoder()
        le7=LabelEncoder()
        le8=LabelEncoder()
        se1=StandardScaler()
        se2=StandardScaler()
        se3=StandardScaler()
        se4=StandardScaler()
        se5=StandardScaler()
        se6=StandardScaler()
        se7=StandardScaler()
        fitter=StandardScaler()
        df['category'] = le1.fit_transform(df['category'])
        df['cuisine'] = le2.fit_transform(df['cuisine'])
        df['center_type'] = le3.fit_transform(df['center_type'])
        df['center_id'] = le4.fit_transform(df['center_id'])
        df['meal_id'] = le5.fit_transform(df['meal_id'])
        df['city_region'] = le6.fit_transform(df['city_region'])
        df['price_increase'] = le7.fit_transform(df['price_increase'])
        df['profit/loss'] = le8.fit_transform(df['profit/loss'])
        df['profit'] = se1.fit_transform(df[['profit']])
        df['checkout_price'] = se2.fit_transform(df[['checkout_price']])
        df['base_price'] = se3.fit_transform(df[['base_price']])
        df['op_area'] = se4.fit_transform(df[['op_area']])
        df['price_last_curr_diff'] = se5.fit_transform(df[['price_last_curr_diff']])
        df['pricediff_lag_10'] = se6.fit_transform(df[['pricediff_lag_10']])
        df['pricediff_lag_11'] = se7.fit_transform(df[['pricediff_lag_11']])
        df['pricediff_lag_12'] = se8.fit_transform(df[['pricediff_lag_12']])
        df['checkout_price_lag_10'] = se9.fit_transform(df[['checkout_price_lag_10']])
        df['checkout_price_lag_11'] = se10.fit_transform(df[['checkout_price_lag_11']])
        df['checkout_price_lag_12'] = se11.fit_transform(df[['checkout_price_lag_12']])
```

One-hot encoding

```
15: one_hot = pd.get_dummies(df['emailer_for_promotion'], prefix='emailer_for_promotion_')
df = df.drop('emailer_for_promotion', axis=1)
df = df.join(one_hot)
columns = one_hot.columns
df = df[columns]

one_hot = pd.get_dummies(df['homepage_featured'], prefix='homepage_featured_')
df = df.drop('homepage_featured', axis=1)
df = df.join(one_hot)
columns = one_hot.columns
df = df[columns]

one_hot = pd.get_dummies(df['category'], prefix='category_')
df = df.drop('category', axis=1)
df = df.join(one_hot)
columns = one_hot.columns
df = df[columns]

one_hot = pd.get_dummies(df['cuisine'], prefix='cuisine_')
df = df.drop('cuisine', axis=1)
df = df.join(one_hot)
columns = one_hot.columns
df = df[columns]

one_hot = pd.get_dummies(df['center_type'], prefix='center_type_')
df = df.drop('center_type', axis=1)
df = df.join(one_hot)
columns = one_hot.columns
df = df[columns]
```

```

one_hot = pd.get_dummies(df['price_increase'], prefix='price_increase')
df = df.drop('price_increase', axis=1)
df = df.join(one_hot)
print(one_hot.columns)
print(1)

one_hot = pd.get_dummies(df['profit/loss'], prefix='profit/loss')
df = df.drop('profit/loss', axis=1)
df = df.join(one_hot)
print(one_hot.columns)
print(1)

one_hot = pd.get_dummies(df['city_region'], prefix='city_region')
df = df.drop('city_region', axis=1)
df = df.join(one_hot)
print(one_hot.columns)
print(1)

```

Feature selection

```

In [22]: dfval=['meal_id', 'center_id', 'week', 'week_id', 'week_end', 'profit', 'log_price', 'checkbox_price', 'price_low_corr_diff', 'pricediff',
          'ready_for_protection', 'homepage_featured', 'category', 'cuisine', 'center_type', 'price_increase', 'profit/loss', 'city_region']
          #/usr/local/lib/python3.6/dist-packages/sklearn
          #sklearn.feature
          self.train(dfval)
          self.train(['run_orders_logp'])

```

Train-Test split:

```

In [25]: from sklearn.model_selection import train_test_split
          x_train, x_test, y_train, y_test = train_test_split(X, y,
                                                            test_size=0.02,
                                                            shuffle=False)

```

Modelling

Random forest

```

In [26]: from sklearn.ensemble import RandomForestRegressor
          from sklearn.metrics import mean_squared_log_error
          regressor = RandomForestRegressor(n_estimators=10, random_state=0, verbose=2, n_jobs=4, min_samples_leaf=5, max_depth=20, min_sam
          regressor.fit(X_train, y_train)
          y_pred = regressor.predict(X_test)
          print("Accuracy", 100*mean_squared_log_error(np.exp(y_test), np.exp(y_pred)))

```

building tree 1 of 10building tree 2 of 10building tree 3 of 10
building tree 4 of 10
building tree 5 of 10
building tree 6 of 10
building tree 7 of 10
building tree 8 of 10
building tree 9 of 10
building tree 10 of 10

[Parallel(n_jobs=1)]: Done 10 out of 10 | elapsed: 11.5s finished

Accuracy: 41.538814712235

//anaconda3/lib/python3.7/site-packages/sklearn/ensemble/forest.py:277: UserWarning: Some inputs do not have OOB scores. This probably means too few trees were used to compute any reliable oob estimates.
Warning: Some inputs do not have OOB scores.
[Parallel(n_jobs=4)]: Using backend ThreadedBackend with 4 concurrent workers.
[Parallel(n_jobs=4)]: Done 10 out of 10 | elapsed: 0.0s finished

Catboost Regressor

```
In [48]: import catboost as cb
cat_reg = cb.CatBoostRegressor(iterations=625, learning_rate=0.05, depth=4, l2_leaf_reg=10, loss_function='RMSE', random_seed=70)

In [49]: cat_reg.fit(X_train, y_train, cat_features=cat, verbose=True)
000:   learn: 0.4270463   total: 2m 33s   remaining: 0.74s
001:   learn: 0.4875642   total: 2m 33s   remaining: 0.58s
002:   learn: 0.4265431   total: 2m 34s   remaining: 0.46s
003:   learn: 0.4372575   total: 2m 34s   remaining: 3.00s
004:   learn: 0.4272279   total: 2m 34s   remaining: 2.00s
005:   learn: 0.4873578   total: 2m 34s   remaining: 0.65s
006:   learn: 0.4263238   total: 2m 34s   remaining: 0.48s
007:   learn: 0.4371287   total: 2m 35s   remaining: 2.25s
008:   learn: 0.4273532   total: 2m 35s   remaining: 1.90s
009:   learn: 0.4871685   total: 2m 35s   remaining: 1.76s
010:   learn: 0.4261836   total: 2m 35s   remaining: 1.54s
011:   learn: 0.4369336   total: 2m 36s   remaining: 1.37s
012:   learn: 0.4868522   total: 2m 36s   remaining: 1.1s
013:   learn: 0.4866448   total: 2m 36s   remaining: 880ms
014:   learn: 0.4266188   total: 2m 36s   remaining: 660ms
015:   learn: 0.4368736   total: 2m 37s   remaining: 420ms
016:   learn: 0.4868367   total: 2m 37s   remaining: 280ms
017:   learn: 0.4867644   total: 2m 37s   remaining: 0ms

Out[49]: <catboost.core.CatBoostRegressor at 0xab1e3e668>

In [50]: pred = cat_reg.predict(X_test)
print("Accuracy:", 100*mean_squared_log_error(np.exp(y_test), np.exp(pred)))

Accuracy: 92.68229615488888
```

LightGBM

```
In [32]: from lightgbm import LGBMRegressor

g = {'early_stopping_rounds': 100,
     'min_child_samples': 5,
     'min_child_weight': 0.01}

estimator = LGBMRegressor(learning_rate=0.001, n_jobs=-1, max_depth=20,
                           n_estimators=10000,
                           silent=False,
                           **g)

fit_params = {'early_stopping_rounds': 1000,
              'feature_name': 'feat',
              'categorical_feature': cat,
              'eval_set': [(X_train, y_train), (X_test, y_test)]}

estimator.fit(X_train, y_train, **fit_params)

In [35]: from sklearn.metrics import mean_squared_log_error
pred = estimator.predict(X_test)
print("Accuracy:", 100*mean_squared_log_error(np.exp(y_test), np.exp(pred)))

Accuracy: 90.993407955456927
```

XGBoost

[illegible]

RESULTS:

Serial No	Model	Test Accuracy (RMSLE*100)	Submission Score(RMSLE*100)
1.	Catboost	32.48	0.53
2.	Random Forest	41.53	0.68
3.	XGBoost	41.39	0.65
4.	LightGBM	30.99	0.514

CONCLUSION:

Multiple models were fit with varying features and parameter tuning resulting in different levels of performance each time. The best performances achieved were as follows: Random forest gave us an RMSLE of 0.68. Catboost resulted in an RMSLE of 0.53 while LightGBM produced the best results with an RMSLE of 0.5144.

Finally, the model that placed us in the 50th position of the competition is Light-GBM producing a RMSLE score of 0.5144.

Further hyper parameter tuning and model enhancements could lead to potential performance improvements.