

Outlet Sales Forecasting With Ensemble Modeling





CONTENTS

Part 01 Exploratory Data Analysis

Part 02 Data Preprocessing & Feature Engineering

Part 03 Model Building & Validation

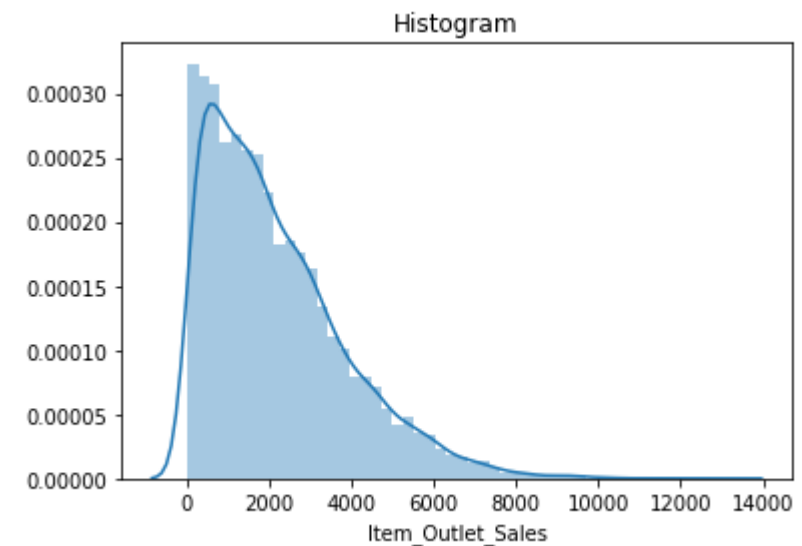
Part 04 Prediction & Future Improvements

EDA

What is the target variable?

And what are the features we have?

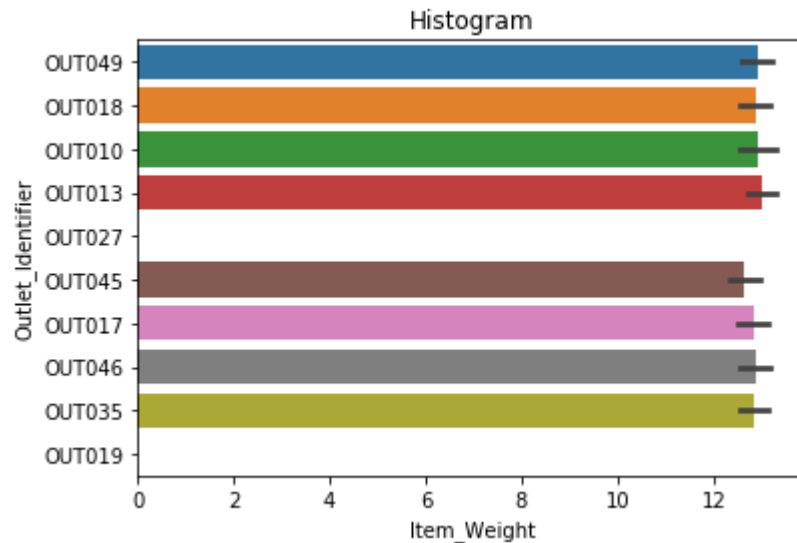
- From the descriptive analysis we could find out the distribution of target variables is not normalized.





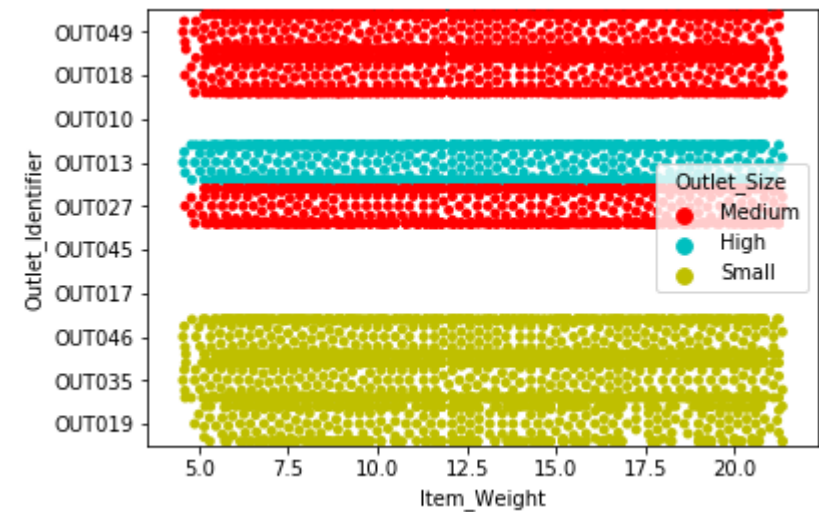
Item Weight

The missing values in this feature are limited to Outlet ID : Out019 and Out027.



Outlet Size

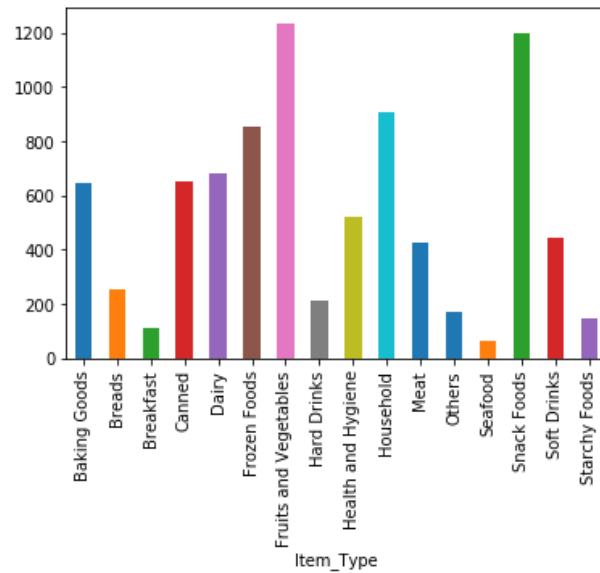
Outlet with ID Out010, Out017 and Out045 have missing values in this feature.





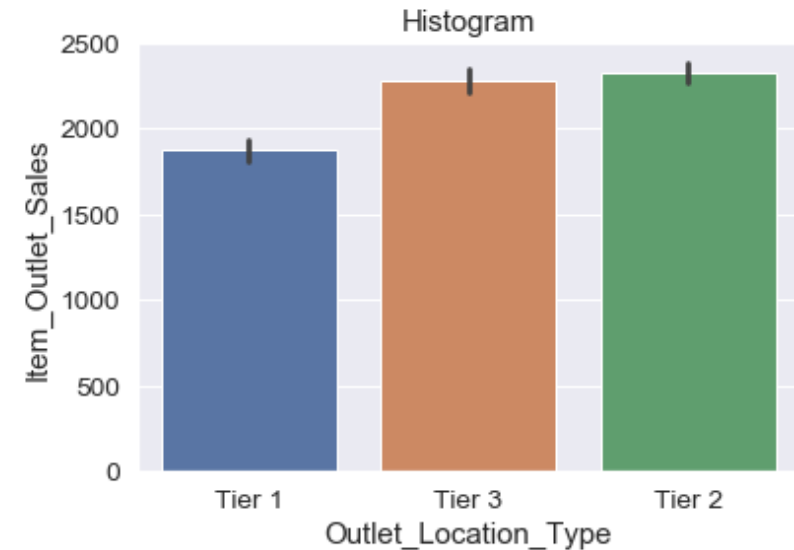
Item Type

Different item types fall into different groups with regard to sales figure.



Outlet Location

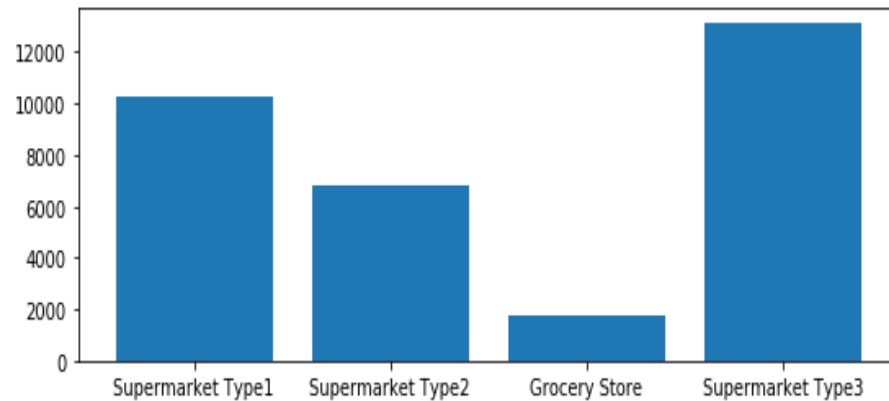
Location Tier 2 & 3 generally have higher sales than Tier 1.





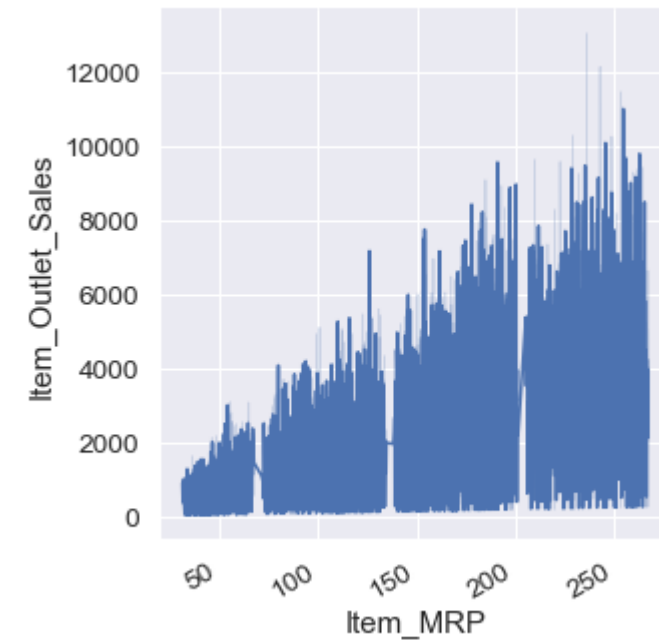
Outlet Type

Grocery Store have a lot lower sales number than Supermarket.



Item MRP

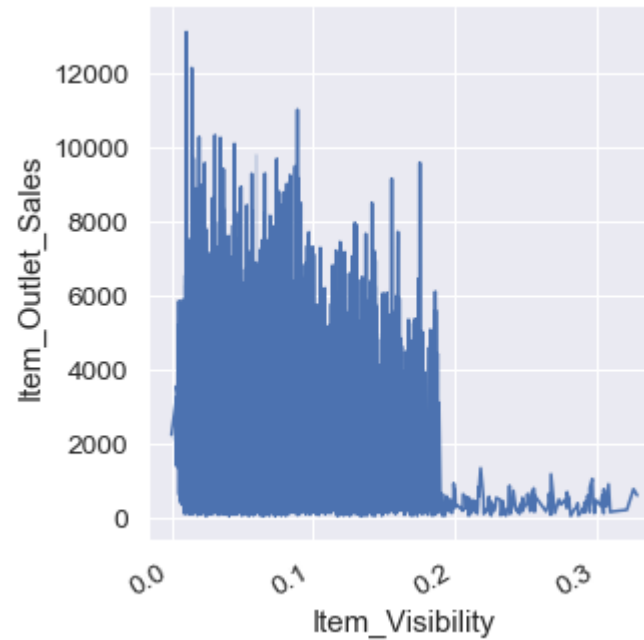
Increasing of sales number demonstrated a gradual pattern with the growth of Item MRP in different levels.





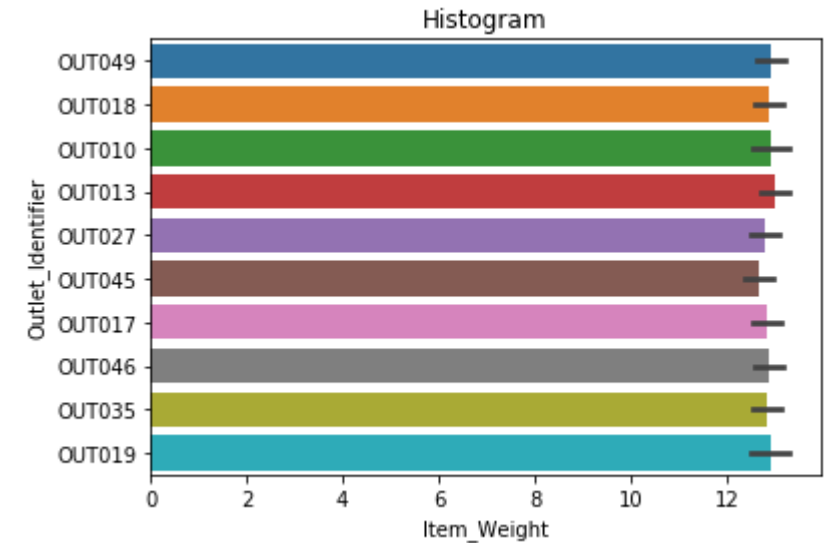
Item Visibility

There are some spurious outliers in the Item visibility feature.



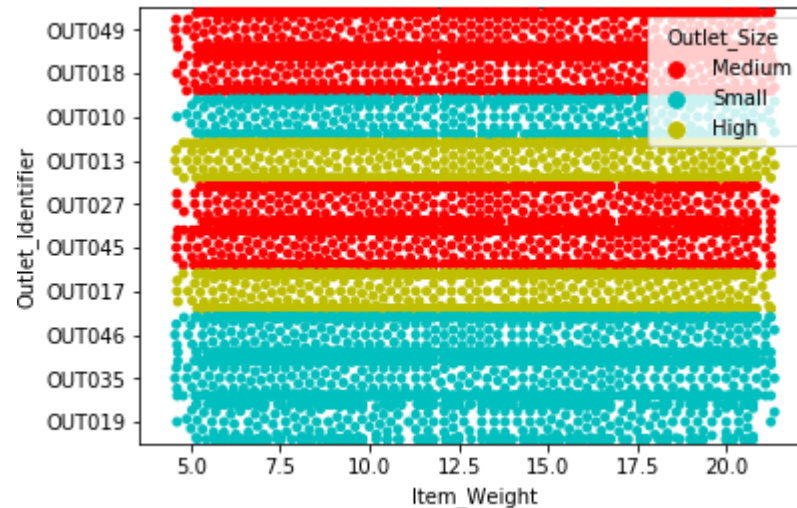
Data Preprocessing

Data imputation for the Item Weight feature would be conducted based on The findings that only outlet with ID 019 and 027 have missing values, we randomly chose item weights from those two kinds of outlet to fill in as the actual item weight to minimize the effects of preexisting missing values.



Item Weight After Imputation

Data Preprocessing



Outlet Size Scatterplot After Imputation.

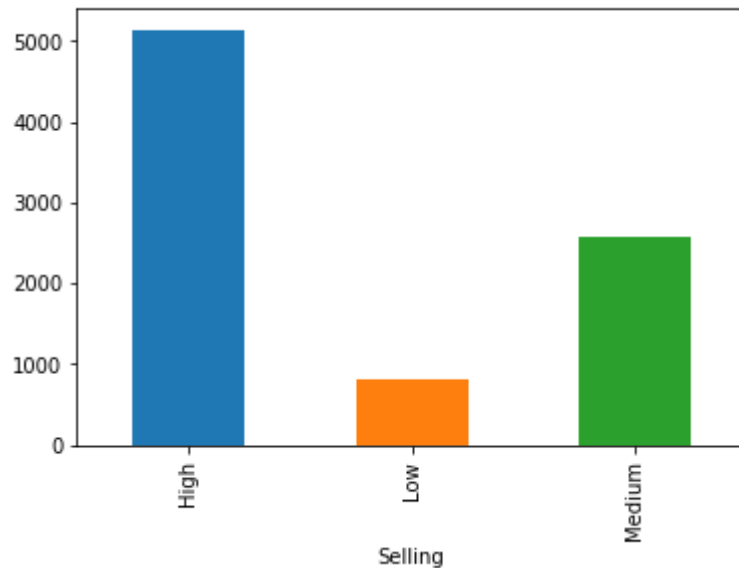
Missing values for the categorical variables would be imputed based on the specific outlet ids and its corresponding Outlet Size and Item sales.

For example, We find out that Out010 corresponds to grocery store therefore the missing values for Outlet Size should be small.



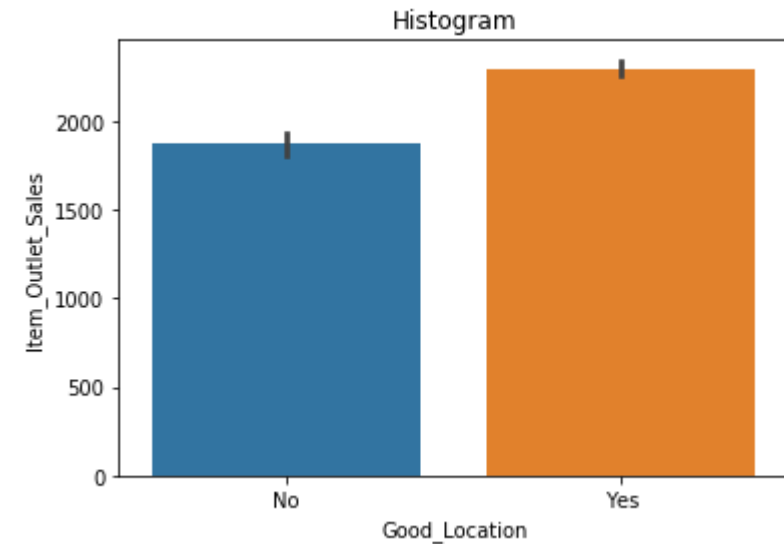
Selling Level

Different Selling level are created to represent the difference in sales.



Location Level

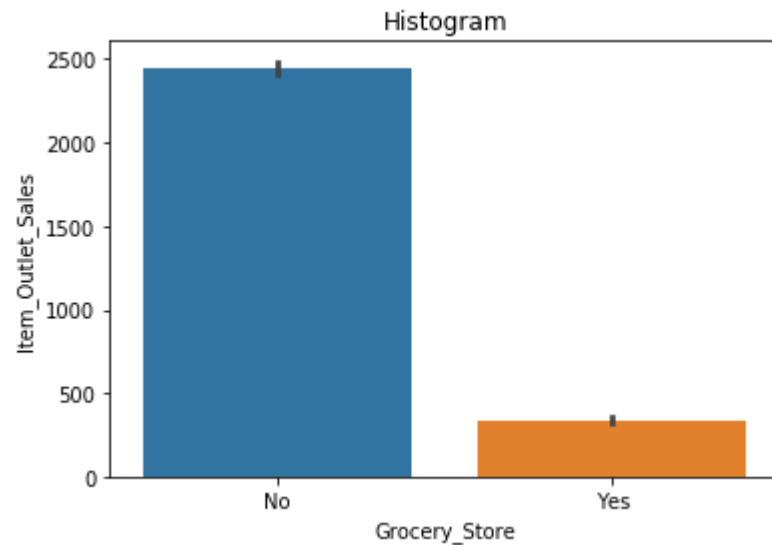
Location levels are applied to show the sales gap between outlets in various locations.





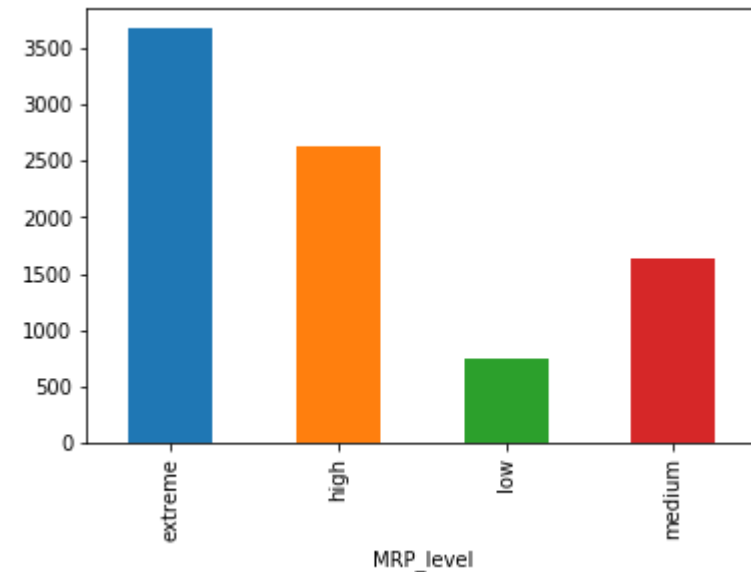
Outlet level

Sales gap between Supermarket and Grocery Store was shown in this feature.



MRP Level

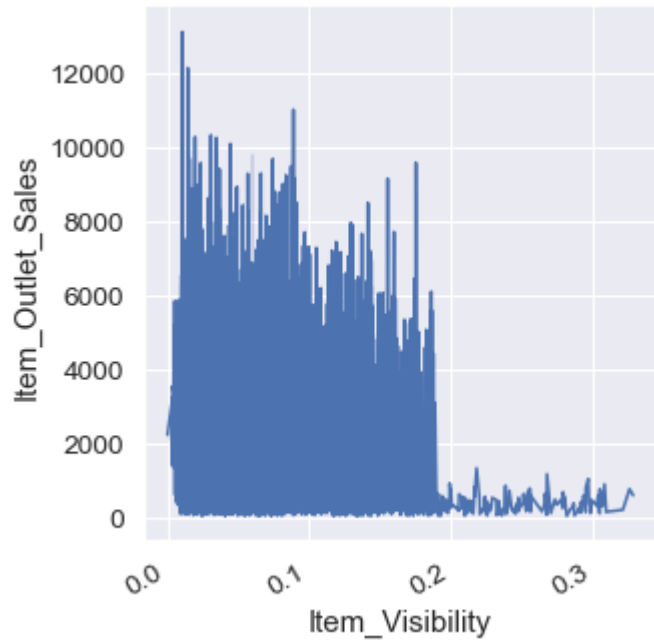
Gradual pattern in sales growth regarding MRP level were reflected in this feature.





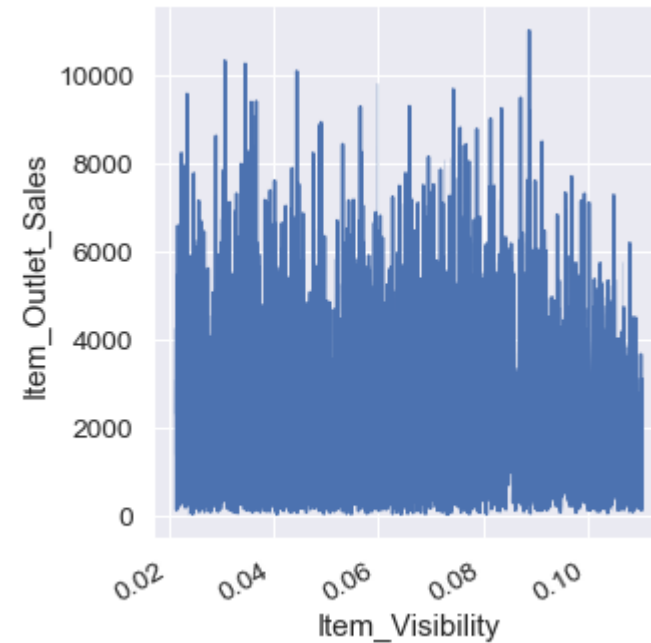
Item Visibility

We decided to apply winsorization to the feature. Set a limit to the extreme values to reduce the effects of spurious outliers.



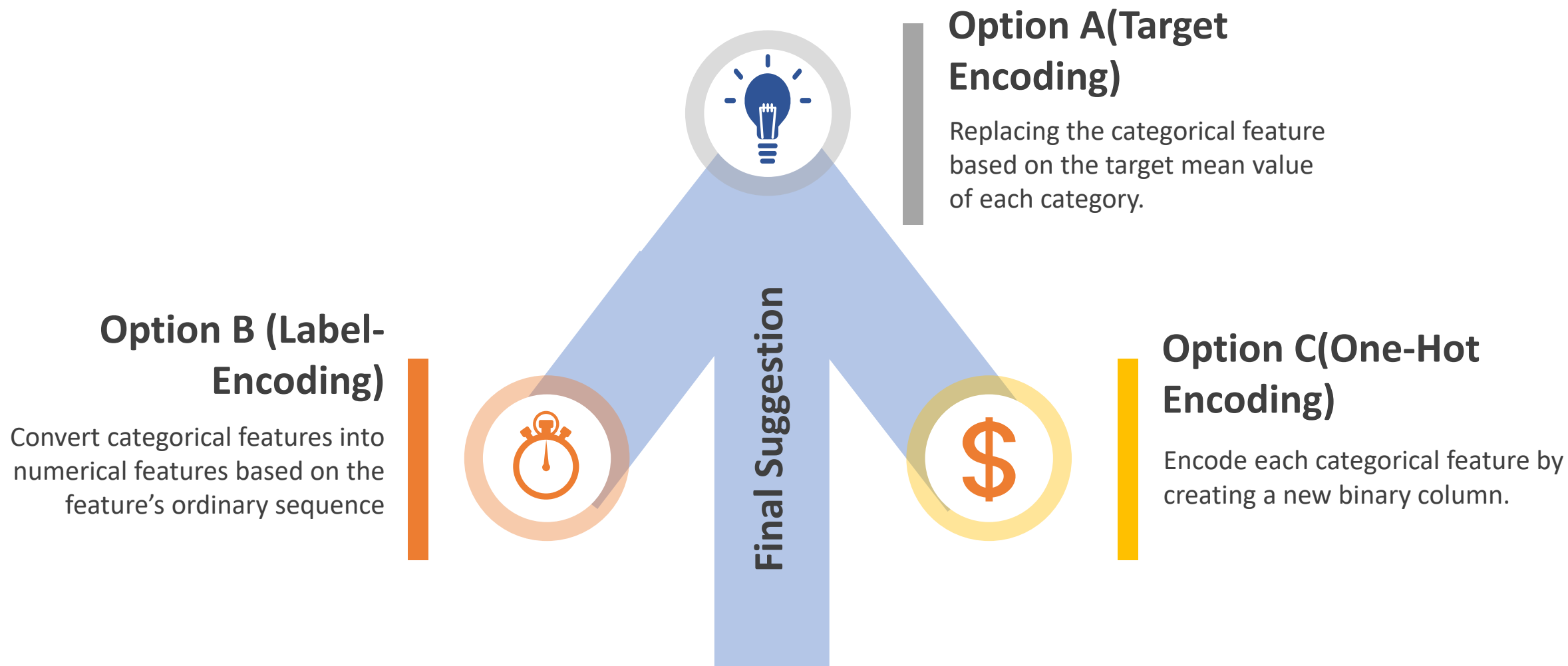
After Winsorization

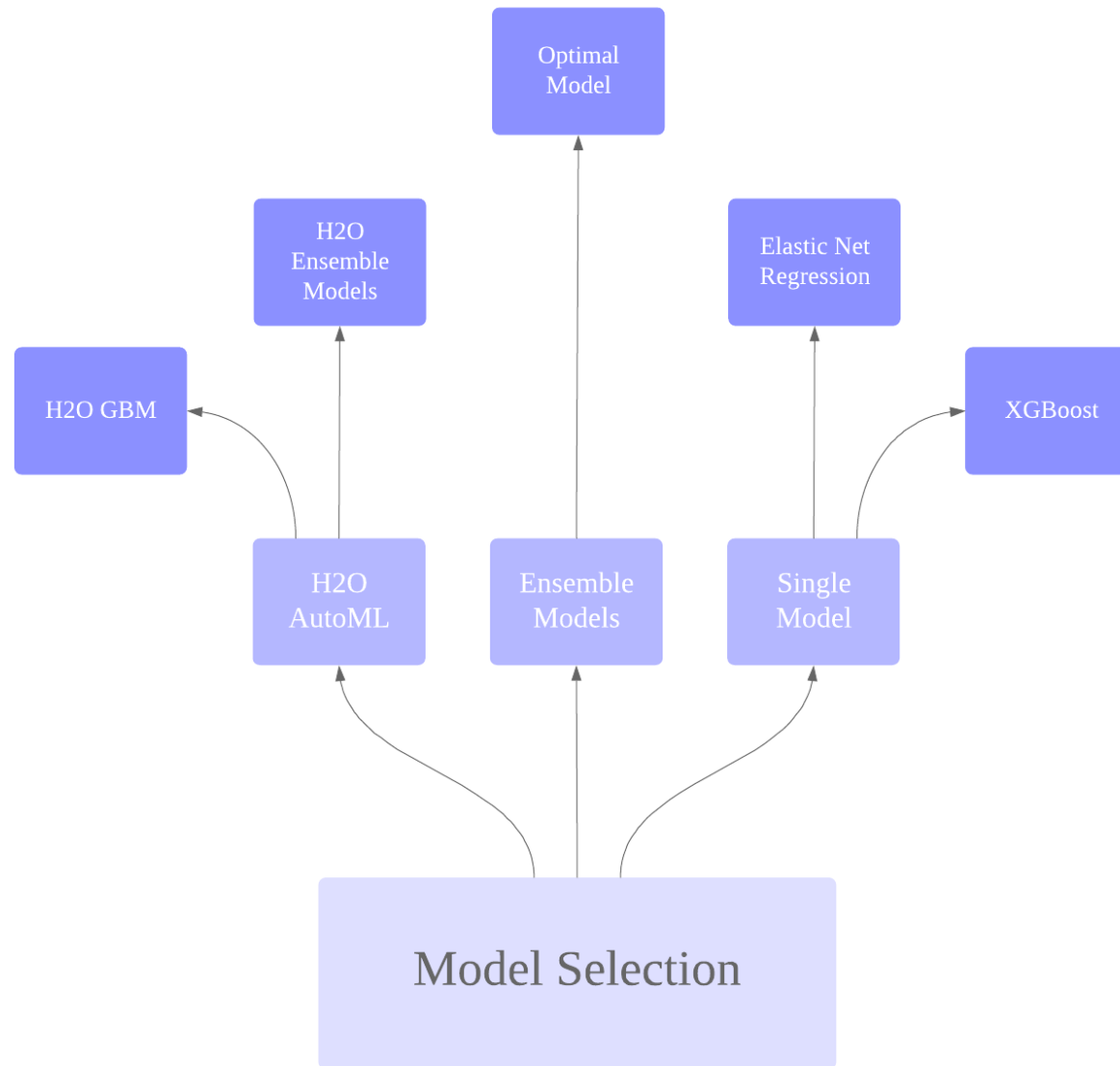
The extreme values has successfully been transformed within the range and the feature looks fine.



Feature Engineering

Encoding





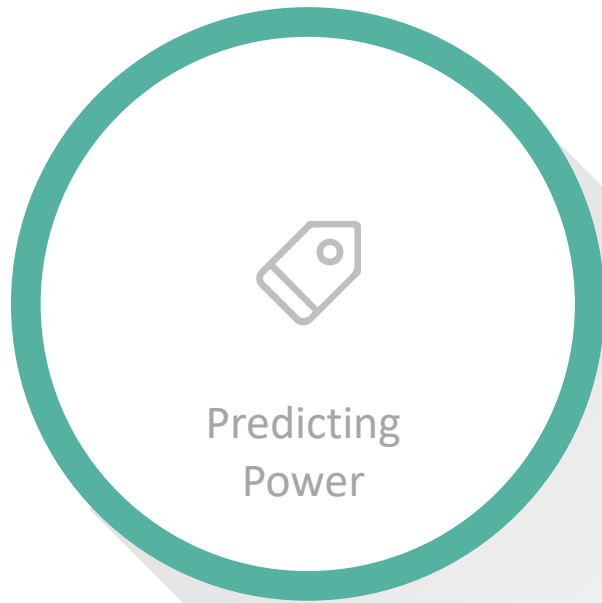
Optimal Model

```
lr = LinearRegression()
svr_lin = SVR(kernel='linear')
ridge = Ridge(random_state=1)
lasso = Lasso(random_state=1)
reg = xgb.XGBRegressor()
svr_rbf = SVR(kernel='rbf')
regressors = [svr_lin, lr, ridge, lasso, reg]
streg = StackingRegressor(regressors=regressors,
                           meta_regressor=svr_rbf)

params = {'lasso__alpha': [0.1, 0.5, 1.0],
          'ridge__alpha': [0.1, 0.5, 1.0],
          'svr__C': [0.1, 0.5, 1.0],
          #'reg__max_depth': [3, 4, 5],
          #'reg__n_estimators': [10, 20, 30],
          #'reg__learning_rate': [0.1, 0.3, 0.5],
          'meta-svr__C': [0.1, 1.0, 5.0],
          'meta-svr__gamma': [0.1, 1.0, 5.0]}

bestmodel4 = RandomizedSearchCV(estimator=streg,
                                param_distributions=params,
                                cv=5,
                                n_iter=10, n_jobs=-1, random_state=1, scoring='r2')
bestmodel4.fit(train_x, train_y1)
```

Notebook Display



■ Goodness of Fit

75.65 %

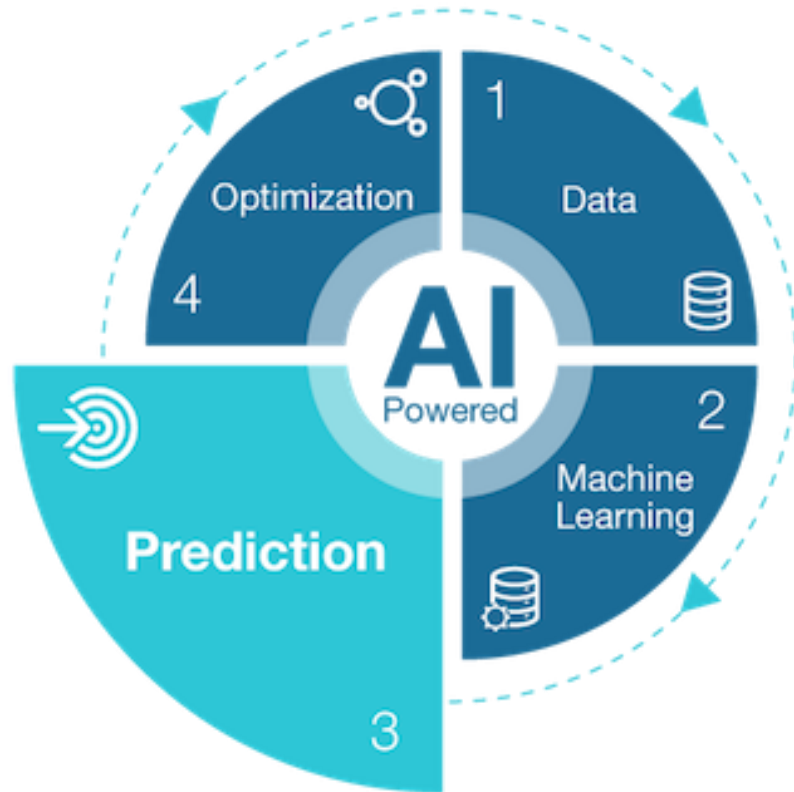
We used cross-validation to evaluate our model since test.csv did not provide the label. Our model finally achieved a **75.65%** goodness of fit score(R^2) with an RMSE of 1.76. Although it is not our ultimate performance since parameter tuning is not finished due to time limit and complexity, it still indicated that our model performed well.

```
streg.predict(train_x)
print("Mean Squared Error: %.4f"
      % np.mean((streg.predict(train_x) - train_y1) ** 2))
print('Variance Score: %.4f' % streg.score(train_x, train_y1))
```

```
Mean Squared Error: 3.1240
Variance Score: 0.7565
```

```
np.sqrt(mean_squared_error(train_y1, streg.predict(train_x)))
```

```
1.767470448092271
```



Further Improvement:

- Impute the missing values for Item weight based on the relationship with other column values.
- Customized feature encoding for every categorical variable.
- Parameter tuning for the final ensemble model.

 mengqi0128@126.com

A person with long hair, wearing a light-colored sweater, is sitting at a wooden desk and typing on a laptop. To the left of the laptop is a dark-colored mug. In the foreground, there is a notebook with handwritten notes and a black pen. The background is slightly blurred, showing a desk with various items. The overall lighting is warm and soft.

THANK YOU

2019.04.03