

Data-Driven Marketing 021 1222

The P-values Team

PDF file for the Final Project

Before we start:

- Unfortunately, this file is not converted directly from the .rmd file. we met many issues when copying and pasting each python line into the .rmd file. In addition, r studio was unstable, and many codes which can run in other IDE couldn't run smoothly in the .rmd file. Thus, this file was converted into PDF from multiple VS code files.
- Purpose of the ML, in this case, we are trying to understand the Dynamic aspects of the price today by exploring the Mercari dataset.
- In this project, as mentioned above, the top objective is to understand dynamic pricing. ML here is more about helping us to enhance our understanding rather than predicting the price accurately.
- I intended to do the EDA part in Python. However, I found Tableau is surprisingly helpful in this case. Thus, all the charts in the final presentation are done by Tableau.
- Models Considered: Catboost | LGBM | XGB | RNN | CNN
- Models used: LGBM | XGB | RNN | CNN (Catboost was dropped since it requires more than 2 hours to run)
- For RNN model, I got energy from this link: [a-simple-nn-solution-with-keras](#)
- For LGBM model, energy: [ridge-lgbm](#), [ftrl-fm-lgb](#), and [price-recommendation-tds](#)
- Got some EDA inspiration from: [mercari-price-prediction](#)
- I only have limited time to develop these models. There's still a lot of space for improvements, especially in the optimization part.
- Closing thoughts can be found at the very end of this file.
- I hope you will enjoy it, cheers.

Content (LGBM)

1.
Why LightGBM
2.
How LightGBM Work
3.
Dataset Prep
4.
Model
5.
Model Evaluation

Why LightGBM

[LightGBM](#) is a gradient boosting framework that uses tree based learning algorithms. It is designed to be distributed and efficient with the following advantages:

- Faster training speed and higher efficiency
- Lower memory usage
- Better accuracy
- Parallel and GPU learning supported
- Capable of handling large-scale data

Therefore, we are going to give it a try.

How LightGBM Work

LightGBM is a supported decision tree algorithm, it splits the tree leaf wise with the simplest fit whereas other boosting algorithms split the tree depth wise or level wise instead of leaf-wise.

When growing on an equivalent leaf in Light GBM, the leaf-wise algorithm can reduce more loss than the level-wise algorithm and hence leads to far better accuracy which may rarely be achieved by any of the prevailing boosting algorithms.

It produces far more complex trees by following leaf wise split approach instead of a level-wise approach which is that the main think about achieving higher accuracy. However, it can sometimes cause overfitting which may be avoided by setting the `max_depth` parameter.

Reference: [LightGBM Website](#)

Dataset Prep

PS: Only Train Dataset from [Mercari Train Dataset \(Kaggle\)](#)

In [1]: `# All libraries required for the Base Model: LGBM are listed below:`

```
#!/pip install numpy
#!/pip install pandas
#!/pip install matplotlib
#!/pip install scipy
#!/pip install sklearn
#!/pip install lightgbm

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from scipy.sparse import csr_matrix, hstack
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.preprocessing import LabelBinarizer
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn import metrics
from sklearn.metrics import mean_squared_error
import lightgbm as lgb
```

In [2]: `df = pd.read_csv('train.tsv', sep='\t')
df.head()
#df.shape`

Out[2]:

	train_id	name	item_condition_id	category_name	brand_name	price	shipping	item_description
0	0	MLB Cincinnati Reds T Shirt Size XL	3	Men/Tops/T-shirts	NaN	10.0	1	No description ye
1	1	Razer BlackWidow Chroma Keyboard	3	Electronics/Computers & Tablets/Components & P...	Razer	52.0	0	This keyboard is in great condition and works ..
2	2	AVA-VIV Blouse	1	Women/Tops & Blouses/Blouse	Target	10.0	1	Adorable top with a hint of lace and a key hol..
3	3	Leather Horse Statues	1	Home/Home Décor/Home Décor Accents	NaN	35.0	1	New with tags Leather horses Retail for [rm]..
4	4	24K GOLD plated rose	1	Women/Jewelry/Necklaces	NaN	44.0	0	Complete with certificate of authenticity

From EDA we know there're missing values in many of the cols

In [3]: `print('Missing value counts:')
for i in df:
 if df[i].isnull().sum()>0:
 print(i, df[i].isnull().sum())`

Missing value counts:
category_name 6327
brand_name 632682
item_description 4

General Setting

In [4]: `num_params` 4000

```
NUM_CATEGORIES = 1000
NAME_MIN_DF = 10
MAX_FEATURES_ITEM_DESCRIPTION = 50000
```

Functions: Replace the null with 'missing'/Cut dataset/convert cat

```
In [5]: def handle_missing_inplace(df):
        df['category_name'].fillna(value='missing', inplace=True)
        df['brand_name'].fillna(value='missing', inplace=True)
        df['item_description'].replace('No description yet', 'missing', inplace=True)
        df['item_description'].fillna(value='missing', inplace=True)
        def cutting(df):
            pop_brand = df['brand_name'].value_counts().loc[lambda x: x.index != 'missing'].index
            df.loc[~df['brand_name'].isin(pop_brand), 'brand_name'] = 'missing'
            pop_category = df['category_name'].value_counts().loc[lambda x: x.index != 'missing'].index
            df.loc[~df['category_name'].isin(pop_category), 'category_name'] = 'missing'
        def to_categorical(df):
            df['category_name'] = df['category_name'].astype('category')
            df['brand_name'] = df['brand_name'].astype('category')
            df['item_condition_id'] = df['item_condition_id'].astype('category')
```

Split dataset into train and test and remove price of 0

```
In [6]: msk = np.random.rand(len(df)) < 0.8
        train = df[msk]
        test = df[~msk]
        test_new = test.drop('price', axis=1)
        y_test = np.log1p(test["price"])
        train = train[train.price != 0].reset_index(drop=True)
        train.shape
```

Out[6]: (1185898, 8)

Merge train and new test data

```
In [7]: nrow_train = train.shape[0]
        y = np.log1p(train["price"])
        merge: pd.DataFrame = pd.concat([train, test_new])
```

Implement previous functions to dataset

```
In [8]: handle_missing_inplace(merge)
        cutting(merge)
        to_categorical(merge)
```

Model

```
In [9]: cv = CountVectorizer(min_df=NAME_MIN_DF)
        X_name = cv.fit_transform(merge['name'])
        cv = CountVectorizer()
        X_category = cv.fit_transform(merge['category_name'])
```

```
In [10]: tv = TfidfVectorizer(max_features=MAX_FEATURES_ITEM_DESCRIPTION, ngram_range=(1, 3), stop_words='english')
        X_description = tv.fit_transform(merge['item_description'])
```

```
In [11]: lb = LabelBinarizer(sparse_output=True)
```

```
X_brand = lb.fit_transform(merge['brand_name'])
```

```
In [12]: X_dummies = csr_matrix(pd.get_dummies(merge[['item_condition_id', 'shipping']], sparse=True))
```

```
In [13]: sparse_merge = hstack((X_dummies, X_description, X_brand, X_category, X_name)).tocsr()
```

```
In [14]: mask = np.array(np.clip(sparse_merge.getnnz(axis=0) - 1, 0, 1), dtype=bool)
sparse_merge = sparse_merge[:, mask]
```

```
In [15]: X = sparse_merge[:nrow_train]
X_test = sparse_merge[nrow_train:]
```

```
In [16]: train_X = lgb.Dataset(X, label=y)
```

Set parameters

```
In [17]: params = {
    'learning_rate': 0.75,
    'application': 'regression',
    'max_depth': 3,
    'num_leaves': 100,
    'verbosity': -1,
    'metric': 'RMSE',
}
```

Training

```
In [18]: gbm = lgb.train(params, train_set=train_X, num_boost_round=3200, verbose_eval=100)
```

```
C:\Users\Abram\Anaconda3\lib\site-packages\lightgbm\engine.py:239: UserWarning: 'verbose_eval' argument is deprecated and will be removed in a future release of LightGBM. Pass 'log_evaluation()' callback via 'callbacks' argument instead.
  _log_warning("'verbose_eval' argument is deprecated and will be removed in a future release of LightGBM. ")
```

Prediction

```
In [19]: y_pred = gbm.predict(X_test, num_iteration=gbm.best_iteration)
```

RMSE

```
In [21]: from sklearn.metrics import mean_squared_error
print('RMSE:', mean_squared_error(y_test, y_pred) ** 0.5)
```

RMSE: 0.458683501078225

Closing Thoughts:

- Many pages of model files could cause redundancy in this report.
- I intended to build a content-based recommendation engine for this project. However, this may not be related to our objective, plus not sufficient time available.
- Optimization of different models could be the crucial part of this project because we can further simulate price changes in different scenarios, for example, price changes in B2B/B2C, Season Promotion, Inflation, supply and demand situations, and so on. Maybe I will search for more datasets to let this feature happen in the future.
- We chose XGBoost because it took less time to generate results, even though LGBM can offer us better evaluation scores, most optimization service providers (like Intel Optimization Environment) now do not support LGBM. Thus, if this is a real company case, for better future development considerations, XGB is a better choice.
- The library used in the XGBoost model part is xgboost. Unfortunately, this is not ideal since sklearn's xgboost do compatible with more services. For example, SHAP analysis, some evaluation dashboards and interactive web libraries do not support xgboost but sklearn's xgboost. This is one of my regrets after all this finished.
- There's a lot to improve in both the EDA and model parts. Interactive dashboards were built but due to the limitations of the presentation that we can't directly show in class, but hopefully we can be better in the future.

Cheers