

# Final Project of Quantecon

Yalin Liu

June 8, 2023



# Contents

1 Dataset

2 Models

3 Dimension Reduction

4 Prospect

# Dataset

# Introduction

- I spilt the dataset into 10 subsets by the city.
- This is because different cities tend to have different factors which are invisible.
- Create new variables through 'amenities' feature.
- Using ColumnTransformer to handle categorical variables.
- Through Google map, maybe we can define more categorical variables.

# Introduction

In [54]:

```
# Encoding categorical features

transformer = ColumnTransformer([
    ('One Hot', OneHotEncoder(drop='first'), ['neighbourhood_group', 'room_type']),
    ('Binary', ce.BinaryEncoder(), ['neighbourhood'])
], remainder='passthrough')
```

# Example: Hong Kong

```
In [32]: 1 HK.copy=HK
2 X=HK.drop(['amenities','longitude','latitude','city','price','room_type'],axis=1)
3 y=HK['price']
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=111)
```

```
In [33]: 1 X.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 3772 entries, 9244 to 268762
Data columns (total 17 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   host_is_superhost                     3772 non-null   float64
1   host_has_profile_pic                  3772 non-null   float64
2   host_identity_verified                 3772 non-null   float64
3   accommodates                          3772 non-null   int64
4   bedrooms                             3772 non-null   float64
5   review_scores_rating                   3772 non-null   float64
6   review_scores_accuracy                 3772 non-null   float64
7   review_scores_cleanliness              3772 non-null   float64
8   review_scores_checkin                  3772 non-null   float64
9   review_scores_communication            3772 non-null   float64
10  review_scores_location                  3772 non-null   float64
11  review_scores_value                     3772 non-null   float64
12  instant_bookable                       3772 non-null   int64
13  Entire place                           3772 non-null   uint8
14  Hotel room                             3772 non-null   uint8
15  Private room                           3772 non-null   uint8
16  Shared room                            3772 non-null   uint8
dtypes: float64(11), int64(2), uint8(4)
memory usage: 427.3 KB
```

# Models

# OLS

- Serves as the baseline model.
- It may exist the problem of overfit.



## OLS

```

1 #baseline 需要进一步对数据进行处理
2 from sklearn import linear_model
3 reg = linear_model.LinearRegression()
4 reg.fit(X_train,y_train)
5 y_pred=reg.predict(X_test)
6 print("Coefficients: \n", reg.coef_)
7 print("Mean squared error: %.2f\n" % mean_squared_error(y_test, y_pred))
8 y_pred=reg.predict(X_train)
9 print("Mean squared error: %.2f\n" % mean_squared_error(y_train, y_pred))

```

Coefficients:

```

[-8.58703951e+00  1.40810313e+02 -9.50472837e+00  6.21455917e+01
 3.52025946e+02  1.40227260e+02 -7.49295317e+01  4.86927775e+01
 8.87018198e+01 -7.83583805e+01  8.66776198e+00 -6.10255374e+01
 5.69284038e+01 -2.03156397e+15 -2.03156397e+15 -2.03156397e+15
 -2.03156397e+15]

```

Mean squared error: 1581285.09

Mean squared error: 3698527.64

## OLS

```

1 result = permutation_importance(reg, X_train, y_train, n_repeats=10, random_state=111)
2 print(result['importances_mean']/result['importances_mean'].sum())
3 print(result.importances_std)

```

```

[ 1.45472005e-29 -3.90821045e-30  5.45735952e-29  8.87483221e-27
 2.97142225e-26  1.15572685e-26  4.56234827e-27  1.45527443e-27
 3.77779131e-27  4.59452028e-27 -1.15220088e-29  3.01900193e-27
 2.17696319e-28  4.18617900e-01  6.77377087e-02  4.26251923e-01
 8.73924685e-02]
[5.07199345e-05 3.81612755e-05 7.90721919e-05 2.89723164e-03
 5.00467957e-03 2.51501569e-03 1.14501987e-03 8.24247360e-04
 1.27834817e-03 1.10881389e-03 9.79333368e-05 1.02571612e-03
 4.74785145e-04 8.71269686e+21 1.98690864e+21 6.04809251e+21
 1.33276730e+21]

```

# LASSO

Using GridSearchCV to get the best hyperparameter, again we analyse the importance of different features.

# LASSO

## LASSO

```
In [38]: 1 clf = linear_model.Lasso(alpha=1.0)
2         clf.fit(X_train, y_train)
3         result = permutation_importance(clf, X_train, y_train, n_repeats=10, random_state=111)
4         result['importances_mean']/result['importances_mean'].sum()
```

```
Out[38]: array([ 1.20404843e-04,  0.00000000e+00,  3.97131691e-04,  1.19904276e-01,
  4.02580948e-01,  1.47213475e-01,  5.70993151e-02,  1.77505647e-02,
  4.68008266e-02,  5.73138648e-02, -1.67339312e-04,  3.60617009e-02,
  2.89900293e-03,  5.66509240e-02,  4.22766031e-02,  4.74845962e-05,
  1.30508175e-02])
```

```
In [39]: 1 alpha_range = np.logspace(-6, 6, 13)
2         param_grid = {'alpha': alpha_range}
3         grid_search = GridSearchCV(clf, param_grid=param_grid, cv=6, scoring='neg_root_mean_squared_error')
4         grid_search.fit(X_train, y_train)
5         best_alpha_lasso = grid_search.best_params_['alpha']
6         clf = linear_model.Lasso(alpha=best_alpha_lasso)
7         lasso_score = -1.0 * cross_val_score(clf, X_train, y_train, cv=6, scoring='neg_root_mean_squared_error').mean()
8         lasso_score
```

```
Out[39]: 1751.7979783357844
```

# Random Forest

Using Optuna to get the best hyperparameter, again we analyse the importance of different features.

Random Forest:

```
1 import optuna
2 from sklearn.ensemble import RandomForestRegressor
3
4 def objective(trial):
5     n_estimators = trial.suggest_int("n_estimators", 50, 400)
6     clf = RandomForestRegressor(n_estimators=n_estimators, random_state=111)
7     scores = cross_validate(clf, X_train, y_train, cv=5, scoring='neg_mean_squared_error')
8     rmse=scores['test_score'].mean()
9     return rmse
10
11 if __name__ == "__main__":
12     study = optuna.create_study(direction="maximize")
13     study.optimize(objective, n_trials=100)
14     print(study.best_trial)
```

# XGBoost

Using  
BayesSearch  
CV to get  
the best  
hyperparameter.

```

1  from xgboost import XGBRegressor
2  estimator = XGBRegressor()
3
4  fit_params = {
5      'early_stopping_rounds': 10,
6      'eval_set': [(X_train, y_train)],
7      'verbose': False,
8  }
9
10 search_space = {
11     'max_depth': (0, 50),
12     'n_estimators': (0, 1000),
13     'learning_rate': (0.01, 1.0, 'log-uniform'),
14     'gamma': (1e-9, 0.5, 'log-uniform'),
15     'scale_pos_weight': (1e-6, 500, 'log-uniform'),
16 }
17
18
19 opt = BayesSearchCV(
20     estimator=estimator,
21     search_spaces=search_space,
22     fit_params=fit_params,
23     cv=5,
24     scoring="neg_mean_squared_error",
25     random_state=111,
26     n_iter=3,
27     verbose=1,
28 )

```

## ANN

```

1 import keras
2
3 import optuna
4
5 # 1. Define an objective function to be maximized.
6 def objective(trial):
7     model = Sequential()
8
9     # 2. Suggest values of the hyperparameters using a trial object.
10    model.add(
11        Conv2D(filters=trial.suggest_categorical('filters', [32, 64]),
12              kernel_size=trial.suggest_categorical('kernel_size', [3, 5]),
13              strides=trial.suggest_categorical('strides', [1, 2]),
14              activation=trial.suggest_categorical('activation', ['relu', 'linear']),
15              input_shape=input_shape))
16    model.add(Flatten())
17    model.add(Dense(CLASSES, activation='softmax'))
18
19    # We compile our model with a sampled learning rate.
20    lr = trial.suggest_float('lr', 1e-5, 1e-1, log=True)
21    model.compile(loss='sparse_categorical_crossentropy', optimizer=RMSprop(lr=lr), metrics=['accuracy'])
22    return accuracy
23
24 # 3. Create a study object and optimize the objective function.
25 study = optuna.create_study(direction='maximize')
26 study.optimize(objective, n_trials=100)

```

# Dimension Reduction



# PCA

```

1 from sklearn.decomposition import PCA
2 pca = PCA(n_components='mle',svd_solver='full')
3 pca.fit_transform(X)

```

```

array([[ -1.86033700e+00,  1.29128088e+00, -3.77213121e-01, ...,
        -1.79817912e-01,  2.74769947e-02, -5.01513618e-03],
       [ -8.17786060e-01,  1.24206669e+00,  5.52618142e-01, ...,
        -1.16381018e-01, -7.15575971e-02,  4.82734668e-03],
       [ -4.95193573e-01,  1.48852959e-01,  1.25546679e+00, ...,
        -8.18328595e-03, -2.15117652e-02, -3.87388513e-04],
       ...,
       [ -2.32611096e+00, -6.02576701e-01, -3.83352489e-01, ...,
        -1.37331844e-01, -1.27181554e-03, -4.88760937e-03],
       [  2.67988881e-01,  5.41140443e+00,  1.89909331e-01, ...,
        -6.50847068e-02, -3.22408445e-02,  1.74959715e-03],
       [ -4.01702240e-01, -1.07898766e+00,  1.73712161e-01, ...,
        -1.10682938e-01, -8.02502432e-03, -4.34183628e-03]])

```

# autoencoder

```

1 latent_dim = 28
2
3 class Autoencoder(Model):
4     def __init__(self, latent_dim):
5         super(Autoencoder, self).__init__()
6         self.latent_dim = latent_dim
7         self.encoder = tf.keras.Sequential([
8             layers.Flatten(),
9             layers.Dense(latent_dim, activation='relu'),
10        ])
11        self.decoder = tf.keras.Sequential([
12            layers.Dense(784, activation='sigmoid'),
13            layers.Reshape((28, 28))
14        ])
15
16    def call(self, x):
17        encoded = self.encoder(x)
18        decoded = self.decoder(encoded)
19        return decoded
20
21 autoencoder = Autoencoder(latent_dim)

```

```

1 autoencoder.compile(optimizer='adam', loss=losses.MeanSquaredError())

```

```

1 autoencoder.fit(X_train, X_train,
2                 epochs=10,
3                 shuffle=True,
4                 validation_data=(X_test, X_test))

```

# Prospect

# Further Improvement

- Some further Improvements may include:
- Also add some error analysis in the following.

## Further Improvement

You've implemented regularized linear regression on housing prices

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

But it makes unacceptably large errors in predictions. What do you try next?

Get more training examples

Try smaller sets of features

Try getting additional features

Try adding polynomial features ( $x_1^2, x_2^2, x_1x_2, etc$ )

Try decreasing  $\lambda$

Try increasing  $\lambda$

# Thank you!