

House price Prediction

ML/DL Basic

한유진

```
lookup.KeyValue  
f.constant(['em  
=tf.constant([0  
ce = tf.lookup.StaticV  
init,  
num_oov_buckets=5)  
  
lookup.StaticVocabular  
initializer,  
num_oov_buckets,  
lookup_key_dtype=None  
name=None,  
experimental_is_open
```

Code Review - Missing Value

`df = df.dropna(axis=1, how='any')`

: 결측치를 삭제하는 작업

```
1 df = df.dropna(axis=1, how='any')
```

`isnull().sum()`

: column별로 결측치의 개수를 계산하는 작업

```
1 df.isnull().sum()
```

Id	0
MSSubClass	0
MSZoning	0
LotFrontage	259
LotArea	0
...	
MoSold	0
YrSold	0
SaleType	0
SaleCondition	0
SalePrice	0
Length: 81, dtype: int64	

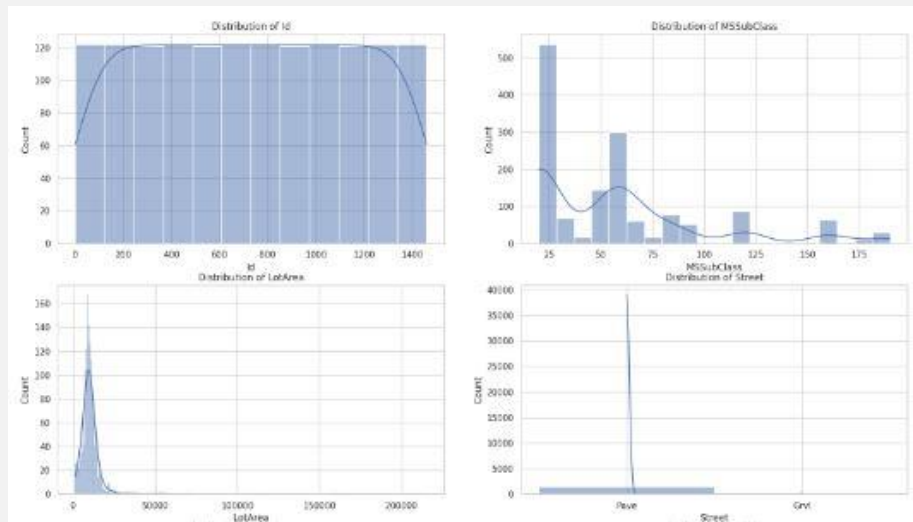
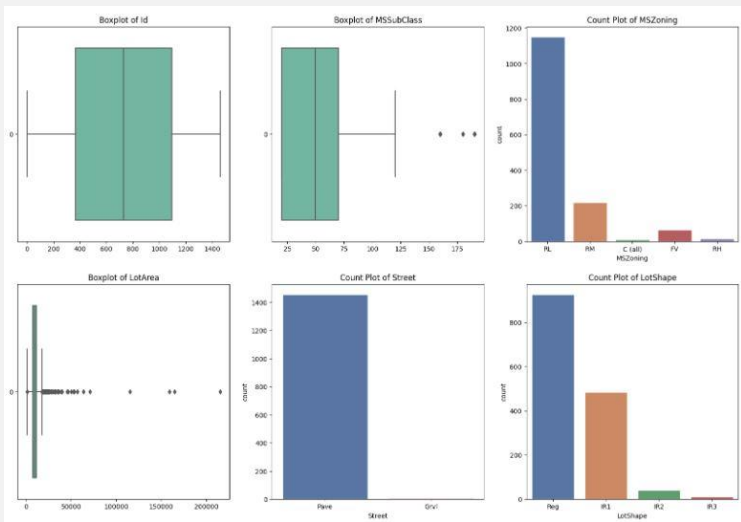


```
1 df.isnull().sum()
```

Id	0
MSSubClass	0
MSZoning	0
LotArea	0
Street	0
...	
MoSold	0
YrSold	0
SaleType	0
SaleCondition	0
SalePrice	0
Length: 62, dtype: int64	

Code Review - Outlier

Column 데이터를 시각화하는 작업



수치형 데이터는 boxplot,
범주형 데이터는 countplot 생성

데이터의 분포 확인

Code Review - Outlier

```

1 def remove_outliers(df, columns, threshold=1.5):
2     for column in columns:
3         Q1 = df[column].quantile(0.25)
4         Q3 = df[column].quantile(0.75)
5         IQR = Q3 - Q1
6         lower_bound = Q1 - threshold * IQR
7         upper_bound = Q3 + threshold * IQR
8         df = df[(df[column] >= lower_bound) & (df[column] <= upper_bound)]
9     return df
10
11 columns_to_remove_outliers = ['MSSubClass', 'OverallQual', 'OverallCond', 'BsmtFinSF1',
12 removed_outlier = remove_outliers(ptr_df, columns_to_remove_outliers)
13 removed_outlier

```

```

1 def factorize_categorical_columns(column):
2     if column.dtype == 'object':
3         column_encoded, _ = pd.factorize(column)
4         return column_encoded
5     return column
6
7 # Apply factorize only to categorical columns
8 df_encoded = removed_outlier.apply(factorize_categorical_columns)

```

Outlier 제거

: IQR을 설정하여 이상치를 정의
($Q1 - \text{threshold} * IQR$ 보다 작거나
 $Q3 + \text{threshold} * IQR$ 보다 큰 값)

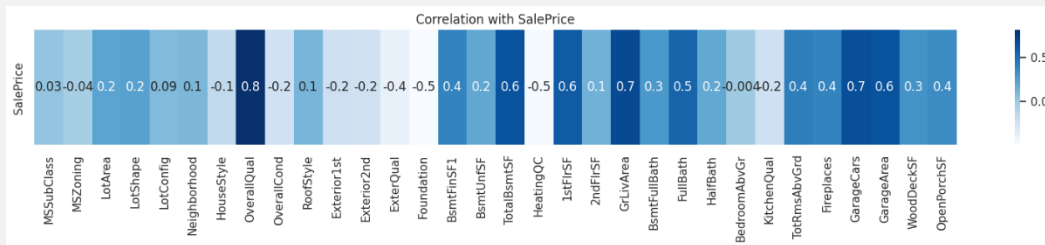
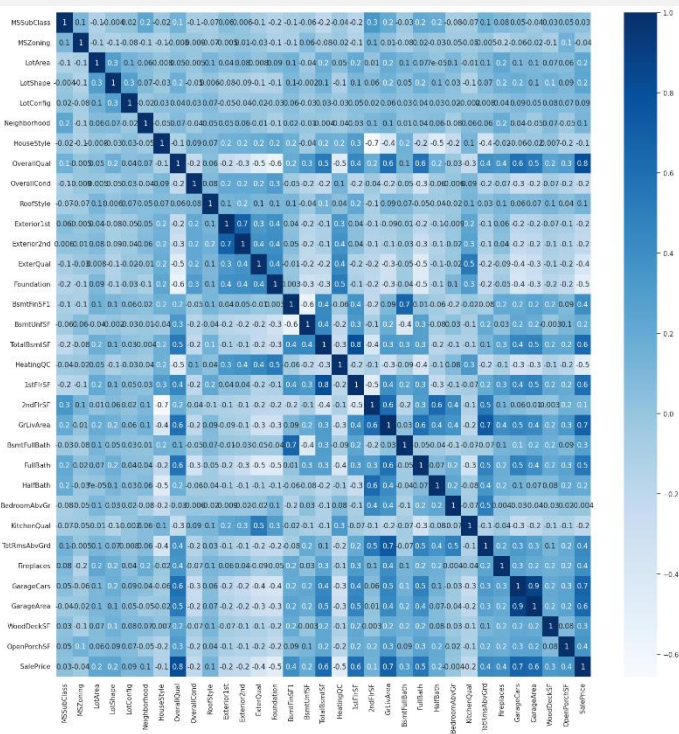
범주형 데이터를 정수로 인코딩

: pd.factorize 함수를 이용하여

데이터 유형이 object인 경우 범주형

데이터로 간주해 수치형 데이터로 변환

Code Review - Correlation Heatmap



Correlation analysis

: 두 변수간의 관계를 상관계수로 표현하는 작업

`final_df.corr():` 기본적으로 pearson 상관계수 계산

Pearson 상관계수: -1과 1 사이의 값을 가짐

Code Review - VIF

VIF(Variance Inflation Factor)

: 분산 팽창 인수로 데이터들이 다중공선성이 있는지 확인하고 100이 넘으면 문제가 있음

다중공선성 - 독립변수들 간에 상관관계를 가짐

```
1 from statsmodels.stats.outliers_influence import variance_inflation_factor
2 vif=pd.DataFrame()
3 vif['VIF']=[variance_inflation_factor(sel_df,i) for i in range(sel_df.shape[1])]
4 vif['features']=sel_df.columns
```

Code Review - Applying ML Algorithm

Linear Regression 진행

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.ensemble import GradientBoostingRegressor
3 from sklearn.metrics import mean_squared_error, r2_score
4 from sklearn.linear_model import LinearRegression
5
6 model = LinearRegression()
7 model.fit(X_train, y_train)
```

```
1 test_pred=model.predict(X_test)
2 end_pred=pd.DataFrame(test_pred,index=df_test.index)
3 end_pred.columns=['SalePrice']
4 end_pred.to_csv('submission.csv',sep=',')
5 end_pred.head()
```

Performance Improvement

GBoost Regression

```
1 GBoost = GradientBoostingRegressor(n_estimators=3000, learning_rate=0.05,  
2                                   max_depth=4, max_features='sqrt',  
3                                   min_samples_leaf=15, min_samples_split=10,  
4                                   loss='huber', random_state =5)
```

```
1 GBoost.fit(train,y_train)  
2 GB_train_pred = GBoost.predict(train)  
3 GB_pred = np.expml(GBoost.predict(test.values))  
4 print(mse(y_train, GB_train_pred))
```


Performance Improvement

XGBoost Regression

```
1 model_xgb = xgb.XGBRegressor(colsample_bytree=0.4603, gamma=0.0468,  
2                               learning_rate=0.05, max_depth=3,  
3                               min_child_weight=1.7817, n_estimators=2200,  
4                               reg_alpha=0.4640, reg_lambda=0.8571,  
5                               subsample=0.5213, silent=1,  
6                               random_state =7, nthread = -1)
```

```
1 model_xgb.fit(train, y_train)  
2 xgb_train_pred = model_xgb.predict(train)  
3 xgb_pred = np.expml(model_xgb.predict(test))  
4 print(mse(y_train, xgb_train_pred))  
_
```

Performance Improvement

LightGBM Regression

```
1 model_lgb = lgb.LGBMRegressor(objective='regression', num_leaves=5,  
2                               learning_rate=0.05, n_estimators=720,  
3                               max_bin = 55, bagging_fraction = 0.8,  
4                               bagging_freq = 5, feature_fraction = 0.2319,  
5                               feature_fraction_seed=9, bagging_seed=9,  
6                               min_data_in_leaf =6, min_sum_hessian_in_leaf = 11)
```

```
1 model_lgb.fit(train, y_train)  
2 lgb_train_pred = model_lgb.predict(train)  
3 lgb_pred = np.expml(model_lgb.predict(test.values))  
4 print(mse(y_train, lgb_train_pred))
```

Performance Improvement

Ensemble

```
1 ensemble1 = xgb_pred*0.25 + lgb_pred*0.25 + GB_pred*0.5
```

```
1 ensemble2 = xgb_pred*0.5 + lgb_pred*0.25 + GB_pred*0.25
```

```
1 ensemble3 = xgb_pred*0.25 + lgb_pred*0.5 + GB_pred*0.25
```

Performance Improvement

	submission_ensemble3.csv Complete · now	0.12533
	submission_ensemble2.csv Complete · 24s ago	0.12662
	submission_ensemble1.csv Complete · 1m ago	0.12565
	submission_xgb.csv Complete · 1m ago	0.13247
	submission_lgb.csv Complete · 1m ago	0.12624
	submission_GB.csv Complete · 2m ago	0.12897
	submission (1).csv Complete · 7h ago	0.53875