

## Elo Model Building and Training

This notebook focuses on building different regression models for

```
import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings('ignore')
```

```
from prettypable import PrettyTable
from sklearn.model_selection import StratifiedFoldSplit
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import Ridge
from sklearn.neighbors import KNeighborsRegressor
import lightgbm as lgb
import xgboost as xgb
import optuna

Function to reduce the memory usage by any pandas dataframe variable

In [2]: https://www.kaggle.com/c/champs-scalar-coupling/discussion/96655
def reduce_mem_usage(df, verbose=True):

    numerics = ['int16', 'int32', 'int64', 'float16', 'float32', 'float64']
    start_mem = df.memory_usage().sum() / 1024**2
    for col in df.columns:
        col_type = df[col].dtypes
        if col_type in numerics:
            c_min = df[col].min()
            c_max = df[col].max()
            if str(col_type)[:3] == 'int':
                if c_min > np.iinfo(np.int8).min and c_max < np.iinfo(np.int8).max:
                    df[col] = df[col].astype(np.int8)
                elif c_min > np.iinfo(np.int16).min and c_max < np.iinfo(np.int16).max:
                    df[col] = df[col].astype(np.int16)
                elif c_min > np.iinfo(np.int32).min and c_max < np.iinfo(np.int32).max:
```

```
df[col] = df[col].astype(np.int32)
elif c_min > np.iinfo(np.int64).min and c_max < np.iinfo(np.int64).m
    df[col] = df[col].astype(np.int64)
else:
    if c_min > np.finfo(np.float16).min and c_max < np.finfo(np.float16).
        df[col] = df[col].astype(np.float16)
    elif c_min > np.finfo(np.float32).min and c_max < np.finfo(np.float32).
        df[col] = df[col].astype(np.float32)
    else:
        df[col] = df[col].astype(np.float64)
```

```

elif c_min > np.inf or float32_c_min and c_max < np.float32_c_min:
    df[col] = df[col].astype(np.float32)
else:
    df[col] = df[col].astype(np.float64)
end_mem = df.memory_usage().sum() / 1024**2
if Verbose: print('Mem. usage decreased to {:5.2f} Mb ({:.1f}% reduction)'.format(end_mem, 100 *
                                                                              (start_mem - end_mem) / start_mem))
return df

```

## Loading featureized train and test dataset

```

In [3]: train = pd.read_csv('data/featureized_train.csv')
        test = pd.read_csv('data/featureized_test.csv')

In [4]: train = reduce_mem_usage(train)
        test = reduce_mem_usage(test)

Mem. usage decreased to 176.00 Mb (67.3% reduction)
Mem. usage decreased to 113.53 Mb (65.3% reduction)

```

We will use the target column in train set as Y train and outlier column as Outlier array to stratify while splitting data. Then we will drop the card id, target and outlier from train and card id from test.

```

In [5]: Y_train = train['target']
        Outlier = train['outlier']
        X_train = train.drop(columns = ['card_id', 'target', 'outlier'])
        X_test = test.drop(columns = ['card_id'])

```

## Model Training

We will be training different models for our problem and analyze which model performs better. There is large variation in values of different features since the values are not normalized, so we will be using non-linear regression models for our problem.

### Baseline Model

The Baseline model will predict the mean of the target values of train as target for all test points. The RSME score of the Baseline model will provide us with benchmark from where we have to increase our model performance. Each of our model should have RMSE score less than the RSME score of the Baseline model.

```

In [6]: Y_train_mean = pd.DataFrame(Y_train, columns = ['target'])
        Y_train_mean['target'] = Y_train.astype('float').mean()

```

```
print("RSME score for Baseline Model: ", rsme_baseline)
```

```

First non-linear model we will use is KNN Regressor which is the simplest of the regression models. We will use Grid Search to find the optimum value of number of neighbors and then train the final model with that number of neighbors.

In [7]: knn_model = KNeighborsRegressor(algorithm='kd_tree')
         parameters = [{"n_neighbors": [1, 2, 5, 10, 50, 100, 200]}]
         knn_data = StratifiedFoldIn splits = 4, random_state = 9, shuffle = True).split(X_train, Outlier.values)
         regressor_KNN = GridSearchCV(knn_model, parameters, cv = knn_data, scoring = 'neg_mean_squared_error', n_jobs=

In [8]: regressor_KNN.fit(X_train, Y_train)

Out[8]: GridSearchCV(cv=generator object BaseFold.split at 0x000001C8B388E3B0,
                    estimator=KNeighborsRegressor(algorithm='kd_tree', n_neighbors=1, n_jobs=-1,
                    param_grid=[{'n_neighbors': [1, 2, 5, 10, 50, 100, 200]}],
                    scoring='neg_mean_squared_error')

In [9]: best_params = regressor_KNN.best_params_
         print("RMSE : ", np.sqrt(abs(regressor_KNN.best_score_)))
         print("Best Hyperparameter")
         print("-" * 20)
         for hyperparameter, value in best_params.items():
             print(hyperparameter, ' : ', value)

RMSE : 3.8250306371322047
Best Hyperparameter
-----
n_neighbors : 200

```

```
Y_train_pred = np
knn_folds = Strat
```

```
for fold, (train_idx, val_idx) in enumerate(knn_folds.split(X_train, Outlier.values)):
    print("Training for fold {}".format(fold + 1))
    regressor_knn = KNeighborsRegressor(n_neighbors = 20)
    regressor_knn.fit(X[train_idx].loc[train_idx != 1], Y_train.loc[train_idx != 1])
    Y_train_pred[val_idx] = regressor_knn.predict(X[train_idx].loc[val_idx])

    Training for fold 1.....
    Training for fold 2.....
    Training for fold 3.....
    Training for fold 4.....

mse_knn = np.sqrt(mean_squared_error(Y_train, Y_train_pred))
print("Mean square error for KNN Regressor model")
```

RMSE score for KNN Regressor model: 3.824844001316022

The KNNRegressor model with RMSE score of 3.8248 shows a small improvement over the Baseline model. We will now build some complex models for our problem.

## XGBoost Model

Now we will build a XGBoost model for our problem. But before we build the final XGBoost model we will use Optuna for finding the optimum hyperparameters for the model.

```
In [7]: def objective(trial):
parameters = {
    'objective':      : 'reg:squarederror',
    'learning_rate'   : 0.01,
    'eval_metric'     : 'rmse',
    'tree_method'     : 'gpu_hist',
```

```
'random_state'
'verbosity'
'max_depth'
'subsample'
```

```

    'min_split_loss': trial.suggest_uniform('min_split_loss', 0, 10),
    'min_split_weight': trial.suggest_uniform('min_split_weight', 0.0, 0.2),
    'min_out_lambda': trial.suggest_uniform('min_out_lambda', 0.1, 1.0),
    'reg_lambda': trial.suggest_uniform('reg_lambda', 0.1, 1.0)
)

X_train_pred = np.zeros(len(X_train))
xgb folds = StratifiedKFold(n_splits = 4, shuffle = True, random_state = 9)

for fold, (train_idx, val_idx) in enumerate(xgb.folds.split(X_train, Outlier.values())):
    val_data = xgb.Matrix(X_train.iloc[train_idx], label = Y_train.iloc[train_idx])
    val_data = xgb.Matrix(X_train.iloc[val_idx], label = Y_train.iloc[val_idx])

```

```

        evals = {'train_data': train_data, 'val_data': val_data, 'evals': 10, 'num_boost_round': 10}
        early_stopping_rounds = 500, verbose_eval=3,
        iteration_report_callback=callback,
        iteration_report_callback = (0, regressor_XGB.best_iteration))

    return np.sqrt(mean_squared_error(Y_train, Y_train_pred))

[In 8]: study = optuna.create_study(
study_optimize(objective, n_trials = 20)

[study_optimize object at 0x15912323374] A new study created in memory with name: name-303e1ad4-a906-4c5b-9c4a-2c7f6da24241

[2022-04-04 15:25:27.693] Trial 0 finished with value: 3.66139386128575 and parameters: {'max_depth': 4, 'n_estimators': 0.9286454084973134, 'colsample_bytree': 0.614241327865513, 'min_split_loss': 6.71781814212325, 'n_child_weight': 29.62485731, 'colsample_bytree': 0.5902700700000001, 'reg_lambda': 5.66396265787472}. Best trial 0 with value: 3.66139386128575.

[2022-04-04 15:35:32.200] Trial 1 finished with value: 3.6567087315942403 and parameters: {'max_depth': 6, 'n_estimators': 0.9478008000000001, 'colsample_bytree': 0.703681842358705, 'min_split_loss': 0.478213188239679, 'n_child_weight': 24.5126125159515, 'reg_lambda': 0.9372726007798, 'reg_lambda': 2.611671219322766}. Best trial 1 with value: 3.6567087315942403.

[2022-04-04 15:45:32.154] Trial 2 finished with value: 3.66179940717565 and parameters: {'max_depth': 4, 'n_estimators': 0.923039261562233, 'colsample_bytree': 0.599451658687579, 'min_split_loss': 0.039124612570901, 'n_child_weight': 30.790700410321156, 'reg_lambda': 0.8127952825779, 'reg_lambda': 9.58173819060805}. Best trial 2 with value: 3.6567087315942403.

[2022-04-04 15:55:16.126] Trial 3 finished with value: 3.69238031173204 and parameters: {'max_depth': 2, 'n_estimators': 0.3389583982587135, 'colsample_bytree': 0.553013536136286, 'min_split_loss': 8.336412223436559, 'n_child_weight': 3.422300146098364, 'reg_lambda': 1.884149193673174, 'reg_lambda': 1.5014639334415}. Best trial 3 with value: 3.6567087315942403.

[2022-04-04 16:07:152.306] Trial 4 finished with value: 3.688077278472054 and parameters: {'max_depth': 4, 'n_estimators': 0.9286454084973134, 'colsample_bytree': 0.614241327865513, 'min_split_loss': 6.71781814212325, 'n_child_weight': 29.62485731, 'colsample_bytree': 0.5902700700000001, 'reg_lambda': 5.66396265787472}. Best trial 4 with value: 3.6567087315942403.

```

```

n_child_weight: 0.149677494373413813, 'reg_alpha': 7.95981811659342, 'reg_lambda': 7.8076242165406285, 'Best
| train | 1 | 0.5150000000000000 | 0.5150000000000000 | 0.5150000000000000 | 0.5150000000000000 | 0.5150000000000000 |
| [ 2022-04-06 16:19:38.272] Train 1 finished with value: 6.568919664289907 and parameters: {'max_depth': 8,
| n_child_weight': 0.6264273905441251, 'csc_feature': 'none', 'csc_lambda': 0.7813574403989164, 'min_split_loss': 3.883280248762573,
| 'min_child_weight': 0.149677494373413813, 'reg_alpha': 8.8307646403961604, 'reg_lambda': 4.78882921075339, 'Best
| train | 1 with value: 6.3567087315942403.
| [ 2022-04-06 16:25:39.503] Train 4 finished with value: 6.598962342343767 and parameters: {'max_depth': 5,
| n_child_weight': 0.6150000000000000, 'csc_feature': 'none', 'csc_lambda': 0.7813574403989164, 'min_split_loss': 3.883280248762573,
| 'min_child_weight': 0.149677494373413813, 'reg_alpha': 7.003260369141438, 'reg_lambda': 0.3982322115697671, 'Best
| train | 1 with value: 6.3567087315942403.
| [ 2022-04-06 16:31:40.000] Train 5 finished with value: 6.613616416145252 and parameters: {'max_depth': 5,
| n_child_weight': 0.5662326242796118, 'csc_feature': 'none', 'csc_lambda': 0.94631373787876, 'min_split_loss': 3.33880581090044,
| 'n_child_weight': 0.7664864879982074, 'reg_alpha': 1.6701454863960626, 'reg_lambda': 2.790671934700588, 'Best
| train | 1 with value: 6.3567087315942403.
| [ 2022-04-06 16:43:21.758] Train 8 finished with value: 6.660426020620123 and parameters: {'max_depth': 8,
| n_child_weight': 0.38598955575724, 'csc_feature': 'none', 'csc_lambda': 0.9358780695619144, 'min_split_loss': 6.93954709338007,
| 'min_child_weight': 0.149677494373413813, 'reg_alpha': 2.84232024356516, 'reg_lambda': 0.3670849382799191, 'Best
| train | 1 with value: 6.3567087315942403.
| [ 2022-04-06 16:52:16.042] Train 9 finished with value: 6.573103329789856 and parameters: {'max_depth': 7,
| n_child_weight': 0.048431944, 'csc_feature': 'none', 'csc_lambda': 0.9358780695619144, 'min_split_loss': 0.02140956417,
| 'min_child_weight': 0.7481582482, 'reg_alpha': 2.561472742829594, 'reg_lambda': 0.03963685219174, 'Best
| train | 1 with value: 6.3567087315942403.
| [ 2022-04-06 16:58:35.723] Train 10 finished with value: 6.376413262004045 and parameters: {'max_depth': 6,
| n_child_weight': 0.029999999999999998, 'csc_feature': 'none', 'csc_lambda': 0.12555862829676016, 'min_split_loss': 0.02140956417,
| 0914, 'min_child_weight': 0.1671758155670189, 'reg_alpha': 4.630714944040244, 'reg_lambda': 0.37066795064217,
| 'Best | train | 1 with value: 6.3567087315942403.
| [ 2022-04-06 17:04:17.000] Train 11 finished with value: 6.58025614142854 and parameters: {'max_depth': 8,
| n_child_weight': 0.1170734675866, 'csc_feature': 'none', 'csc_lambda': 0.7549047838853784, 'min_split_loss': 3.81890563828661,
| 'min_child_weight': 0.149677494373413813, 'reg_alpha': 9.956747430902987, 'reg_lambda': 4.9222689280582661,
| 'Best | train | 1 with value: 6.3567087315942403.

```

```

    'subsample': 0.714561031369386, 'coalesce_bytree': 0.87486110015906, 'min_split_loss': 2.26857348380745,
    'min_child_weight': 1.7235937389902428, 'reg_alpha': 0.3649511648210071, 'reg_lambda': 3.474818069961041,
    'max_depth': 12, 'min_child_weight': 0.63549842961493846,
    'min_split_loss': 1.7235937389902428, 'coalesce_bytree': 0.87486110015906, 'min_split_loss': 2.26857348380745,
    'min_child_weight': 0.7546894736891103, 'coalesce_bytree': 0.35935055040456897, 'min_split_loss': 1.743007273688249,
    'min_child_weight': 0.20715810259551817, 'reg_alpha': 0.754543704302853, 'reg_lambda': 2.391766901471961,
    'min_split_loss': 1.743007273688249, 'min_child_weight': 0.63549842961493846,
    'min_split_loss': 1.743007273688249, 'min_child_weight': 0.63549842961493846, 'coalesce_bytree': 0.87486110015906,
    'min_split_loss': 0.378778070370372, 'coalesce_bytree': 0.338046217432016, 'min_split_loss': 1.9645302789072023,
    'min_child_weight': 19.93997754531223, 'reg_alpha': 4.65622704266928, 'reg_lambda': 3.749576140612131,
    'min_split_loss': 1.743007273688249, 'min_child_weight': 0.63549842961493846,
    'min_split_loss': 1.743007273688249, 'min_child_weight': 0.63549842961493846, 'coalesce_bytree': 0.87486110015906,
    'min_split_loss': 0.342013728655751, 'coalesce_bytree': 0.2751347587450241, 'min_split_loss': 6.1318482086646,
    'min_child_weight': 7.669472157459031, 'reg_alpha': 6.97139349475482, 'reg_lambda': 7.236383840382761,
    'min_split_loss': 1.743007273688249, 'min_child_weight': 0.63549842961493846,
    'min_split_loss': 1.743007273688249, 'min_child_weight': 0.63549842961493846, 'coalesce_bytree': 0.87486110015906,
    'min_split_loss': 0.4955085084340094, 'coalesce_bytree': 0.40492851902133645, 'min_split_loss': 1.76557193763,
    'min_child_weight': 16.69824788006556, 'reg_alpha': 0.4164545975927003, 'reg_lambda': 3.686565714144346

```

```
[1] 2022-04-08 18:06:48, 'trial 11 finished with value: 6.6554928083918505 and parameters: {'max_depth': 7, 'subsample': 0.31369999999999999, 'min_split_loss': 0.1383819498849513, 'min_split_loss': 0.3813072182273142}
```

```
5. 'min_child_weight': 20.72825773141414, 'reg_alpha': 0.663841413637069, 'reg_lambda': 1.894686811750832, 'trial 12 with value: 6.35484296143846.
```

```
[1] 2022-04-08 18:11:41, 'trial 18 finished with value: 7.34257797121652 and parameters: {'max_depth': 1, 'subsample': 0.4306889904681825, 'colsample_bytree': 0.4634213525307946, 'min_split_loss': 0.305879428400515}
```

```
6. 'min_child_weight': 25.125964232381426, 'reg_alpha': 7.383607233605659, 'reg_lambda': 3.7416948958354593, 'trial 12 with value: 6.35484296143846.
```

```
[1] 2022-04-06 18:37:46, 'trial 19 finished with value: 6.6552223922320303 and parameters: {'max_depth': 6, 'subsample': 0.8036160540143876, 'colsample_bytree': 0.3078380360846494, 'min_split_loss': 1.286284285331473}
```

```
7. 'min_child_weight': 16.57978389902428, 'reg_alpha': 5.93240734545893, 'reg_lambda': 3.677070253937143, 'trial 12 with value: 6.35484296143846.
```

```
In [9]: b_trial = study.best_trial
        print('RMSE :', b_trial.value)
        best_params = b_trial.params
        print('best hyperparameters:')
        print('-----')
        for hyperparameter, value in best_params.items():
            print(hyperparameter, ':', value)

        RMSE : 6.35484296143846
        Best Hyperparameters:
        -----
        max_depth : 7
        subsample : 0.7145610311969036
        colsample_bytree : 0.36484961010159906
        min_split_loss : 2.2695747848075492
        min_child_weight : 16.57978389902428
        reg_alpha : 9.874518468120071
```

```

n [10]: with open('data/XGB_parameters', 'wb') as df_file:
        pickle.dump(best_params, df_file)

We will use the best hyperparameters obtained from optuna trial for building the final XGBoost model.

In [9]: with open('data/XGB_parameters', 'rb') as df_file:
        best_params_xgb = pickle.load(df_file)

n [10]: Parameters = {
    'objective':      'reg:squarederror',
    'learning_rate':  0.01,
    'eval_metric':    'rmse',
    'tree_method':     'gpu_hist',
    'predicator':      'gpu_predictor',
    'random_state':   0,
    'verbosity':       0,
    'max_depth':      best_params_xgb.get('max_depth'),
    'subsample':       best_params_xgb.get('subsample'),
    'colsample_bytree': best_params_xgb.get('colsample_bytree'),
    'min_split_loss':  best_params_xgb.get('min_split_loss'),
    'min_child_weight': best_params_xgb.get('min_child_weight'),
    'reg_alpha':       best_params_xgb.get('reg_alpha'),
    'reg_lambda':      best_params_xgb.get('reg_lambda')
}

n [11]: Y_train_pred_xgb = np.zeros(X_train)
xgb folds = StratifiedFoldIn(splits = 4, shuffle = True, random state = 9)

```

```

for fold, (train_idx, val_idx) in enumerate(xgb folds.split(X_train, Outlier.values)):
    print('Training fold 1.....')
    print('train data = %g' % train_idx)
    train_data = xgb.DMatrix(X_train.iloc[train_idx], label = Y_train.iloc[train_idx])
    val_data = xgb.DMatrix(X_train.iloc[val_idx], label = Y_train.iloc[val_idx])
    regressor_XGB = xgb.train(params=parameters, dtrain=train_data, evals=[(train_data, 'train')],
                              num_boost_round=10000, early_stopping_rounds=500, verbose_eval=1000)
    Y_train_pred_xgb[val_idx] = regressor_XGB.predict(xgb.DMatrix(X_train.iloc[val_idx]),
                                                         iteration_range=(0, regressor_XGB.best_iteration))

Training for fold 1.....
[0] train-rmse:3.98423          eval-rmse:3.83439
[1000] train-rmse:3.28844      eval-rmse:3.53679
[1369] train-rmse:3.16492      eval-rmse:3.53823
Training for fold 2.....
[0] train-rmse:3.93100          eval-rmse:3.96622
[1000] train-rmse:3.21182      eval-rmse:3.72947
[1268] train-rmse:3.13924      eval-rmse:3.73045
Training for fold 3.....
[0] train-rmse:3.93599          eval-rmse:3.98138
[1000] train-rmse:3.21848      eval-rmse:3.69194
[1355] train-rmse:3.12377      eval-rmse:3.69378
Training for fold 4.....
[0] train-rmse:3.93726          eval-rmse:3.97762
[1000] train-rmse:3.22834      eval-rmse:3.67231
[1589] train-rmse:3.07012      eval-rmse:3.67334

n [12]: rms_e_xgb = np.sqrt(mean_squared_error(Y_train, Y_train_pred_xgb))
print('RMSE score for XGBoost model: ', rms_e_xgb)

RMSE score for XGBoost model: 3.6579771246304467

```

The XGBoost model with RMSE score of 3.679 has much better performance than the KNNRegressor and shows high improvement in RMSE score from the Baseline model.

### LightGBM Model

We will also build a LightGBM model for the problem to analyze if it performs better than the XGBoost model. Similar to XGBoost, we will use Optuna to find the best hyperparameters before building the final LightGBM model.

```
n [33]: def objective(trial):
        parameters = {
            'objective':      : 'regression',
            'metric':         : 'rmse',
            'boosting_type'   : 'gbdt',
            'learning_rate'   : 0.01,
            'device'          : 'cpu',
            'n_jobs'          : -1,
            'verbosity'       : -1,
            'random_state'    : 0,
            'bagging_freq'    : 1,
            'bagging_seed'    : 9,
            'max_depth'       : trial.suggest_int('max_depth', 1, 16),
            'num_leaves'      : trial.suggest_int('num_leaves', 16, 128),
            'min_data_in_leaf': trial.suggest_int('min_data_in_leaf', 8, 64),
            'min_child_weight': trial.suggest_uniform('min_child_weight', 0.32),
            'feature_fraction': trial.suggest_uniform('feature_fraction', 0.1, 1.0),
            'bagging_fraction': trial.suggest_uniform('bagging_fraction', 0.1, 1.0),
            'min_split_gain'   : trial.suggest_uniform('min_split_gain', 0.10)
```

[illegible][illegible][illegible]

```

38946, 'req_lambda': 1.800502708613534, 'Best is trial 9 with value: 3.657823429966506,
[1] 2022-04-07 03:46:34,076] Trial 18 finished with value: 3.6506641942500763 and parameters: {'max_depth': 6,
'min_data_in_leaf': 64, 'min_child_weight': 1.54274471514698, 'feature_fraction': 0.656}
4889426578, 'bagging_fraction': 0.6959253014980463, 'min_split_gain': 3.709719854929109, 'req_alpha': 4.6100
38385073, 'req_lambda': 9.820770171501771, 'Best is trial 18 with value: 3.6506641942500763,
[1] 2022-04-07 03:46:34,133] Trial 18 finished with value: 3.657049402472997 and parameters: {'max_depth':
'min_data_in_leaf': 128, 'min_child_weight': 29, 'min_split_gain': 3.709719854929109, 'feature_fraction': 0.656}
0735667374, 'bagging_fraction': 0.8581991444249893, 'min_split_gain': 7.128286954687689, 'req_alpha': 4.17141
30038237, 'req_lambda': 9.206931489860473, 'Best is trial 18 with value: 3.6506641942500763,
[1] 2022-04-07 03:46:34,190] Trial 18 finished with value: 3.6506641942500763 and parameters: {'max_depth': 6,
b_trial = study.best_trial
print('Name: ', b_trial.value)
best_params = b_trial.params
print("Best Hyperparameters")
print("--")
for hyperparameter, value in best_params.items():
    print(hyperparameter, ': ', value)

RMSE : 3.6506641942500763
Best Hyperparameters:
max_depth : 6
num_leaves : 108
min_data_in_leaf : 64
min_child_weight : 1.54274471514698
feature_fraction : 0.646394894206278
bagging_fraction : 0.6959253014980463
min_split_gain : 3.709719854929109
req_alpha : 4.6100383504763505

```

```
[n 13]: with open('data/LGB_parameters', 'wb') as df_file:
pickle.dump(best_params, df_file)

Now the best hyperparameters obtained from Optuna will be used to build the final LightGBM model.

[n 14]: with open('data/LGB_parameters', 'rb') as df_file:
best_params_lgb = pickle.load(df_file)

parameters = {
    'objective': 'regression',
    'metric': 'rmse',
    'boosting_type': 'gbdt',
    'learning_rate': 0.01,
    'device': 'cpu',
    'n_jobs': 1,
    'verbosity': -1,
    'random_state': 9,
    'bagging_freq': 5,
    'bagging_seed': 9,
    'max_depth': best_params_lgb.get('max_depth'),
    'num_leaves': best_params_lgb.get('num_leaves'),
    'min_data_in_leaf': best_params_lgb.get('min_data_in_leaf'),
    'min_child_weight': best_params_lgb.get('min_child_weight'),
    'feature_fraction': best_params_lgb.get('feature_fraction'),
    'bagging_fraction': best_params_lgb.get('bagging_fraction'),
    'min_split_gain': best_params_lgb.get('min_split_gain'),
    'reg_alpha': best_params_lgb.get('reg_alpha'),
```

```

In [15]: Y_train_pred_lgb = np.zeros(len(X_train))
lgb_folds = StratifiedFold(n_splits = 4, shuffle = True, random_state = 9)

for fold, (train_idx, val_idx) in enumerate(lgb_folds.split(X_train, Outlier.values)):
    print("Training for fold 1.....")
    train_data = lgb.Dataset(X_train.iloc[train_idx], label = Y_train.iloc[train_idx])
    val_data = lgb.Dataset(X_train.iloc[val_idx], label = Y_train.iloc[val_idx])
    regressor_LGB = lgb.LGBMRegressor(params = train_data, valid_sets = [train_data, val_data],
                                     num_boost_round = 10000, verbose_exe = 1000, early_stopping_rounds = 50)
    Y_train_pred_lgb[val_idx] = regressor_LGB.predict(X_train.iloc[val_idx], num_iteration = regressor_LGB.best_iteration)

Training for fold 1.....
Training until validation scores don't improve for 500 rounds
(1000) training's rmse: 3.51008      valid_1's rmse: 3.53725
(2000) training's rmse: 3.39051     valid_1's rmse: 3.53737
Early stopping, best iteration is:
(1608) training's rmse: 3.435_val_1's rmse: 3.53617
Training for fold 2.....
Training until validation scores don't improve for 500 rounds
(1000) training's rmse: 3.44442     valid_1's rmse: 3.72588
Early stopping, best iteration is:
(1038) training's rmse: 3.43184     valid_1's rmse: 3.72561
Training for fold 3.....
Training until validation scores don't improve for 500 rounds
(1000) training's rmse: 3.45567     valid_1's rmse: 3.68775
Early stopping, best iteration is:
(1229) training's rmse: 3.42736     valid_1's rmse: 3.68719
Training for fold 4.....

```

```
[1000] training's rmse: 3.46697          valid_1's rmse: 3.67038
[2000] training's rmse: 3.39138          valid_1's rmse: 3.66958
Early stopping, best iteration is:
[1514] training's rmse: 3.40552          valid_1's rmse: 3.66951

rmse_lgb = np.sqrt(mean_squared_error(Y_train, Y_train_pred_lgb))
print("RMSE score for LGBM model is ", rmse_lgb)

RMSE score for LGBM model: 3.6556194192500763
```

The LightGBM model with RMSE score of 3.6551 is the best of all the models with XGBoost not far behind. To improve the RMSE score even more, we will use stacked models of LightGBM and XGBoost with different weights.

## Stacked Model

We will use different weightage for the XGBoost and LightGBM models to decide the weightage for our final Stacked model.

```
[17]: for i in range(100,1010):
      Y_train_pred_stack = ((i / 100) * Y_train_pred_xgb) + ((1000 - i) / 100) * Y_train_pred_lgb
      print("RMSE score for Stacked Model with (%i weightage to XGBoost and (%i weightage to LightGBM : "(i, 1000-i)

RMSE score for Stacked Model with 10% weightage to XGBoost and 90% weightage to LightGBM : 3.6548067300148395
RMSE score for Stacked Model with 20% weightage to XGBoost and 80% weightage to LightGBM : 3.654671232631392
RMSE score for Stacked Model with 30% weightage to XGBoost and 70% weightage to LightGBM : 3.65467175665005
RMSE score for Stacked Model with 40% weightage to XGBoost and 60% weightage to LightGBM : 3.654746630470889
RMSE score for Stacked Model with 50% weightage to XGBoost and 50% weightage to LightGBM : 3.655396918787539
RMSE score for Stacked Model with 60% weightage to XGBoost and 40% weightage to LightGBM : 3.655348997481543
RMSE score for Stacked Model with 70% weightage to XGBoost and 30% weightage to LightGBM : 3.655823258989213
```

```

RMS_E score for Stacked Model with 30% weightage to XGBoost and 10% weightage to LightGBM : 3.657137423989734
[ 26 ]: meta_stack = np.sqrt((mean_squared_error(Y_train, ((0.3 * Y_train_pred_xgb) + (0.7 * Y_train_pred_lgb))))
print("RMS_E score for Stacked model: ", rme_stack)
RMS_E score for Stacked model: 3.654657715656053

The Stacked Model with 30% weightage to XGBoost and 70% weightage to LightGBM provides a small improvement in RMS_E score from our previous best of LightGBM.

Stacked Model using Meta Learner

Now we will build another Stacked model but with a meta learner and see if it is better than the rest of our models. We will be using Ridge Regressor with predictions of XGBoost and LightGBM as input and the Y train as target values.

[ 19 ]: meta_train = np.vstack([Y_train_pred_xgb, Y_train_pred_lgb]).transpose()

[ 20 ]: meta_model = Ridge()
        parameters = ["alpha", 0.0001, 0.001, 0.01, 0.1, 1.0]
        meta_folds = StratifiedKFold(n_splits = 3, random_state = 9, shuffle = True).split(meta_train, Outlier.values)
        regressor_meta = GridSearchCV(meta_model, parameters, cv = meta_folds, scoring = 'neg_mean_squared_error',)

[ 21 ]: regressor_meta.fit(meta_train, Y_train)

[ 22 ]: GridSearchCV(generator object _BaseKFold.split at 0x000001f5627bb300,
estimator=Ridge(),

```

```

scoring='neg_mean_squared_error')

n [22]: best_params = regressor_meta.best_params_
        print("Best Hyperparameters")
        print("-" * 20)
        for hyperparameter, value in best_params.items():
            print(hyperparameter, ': ', value)

Best Hyperparameters
-----
alpha : 1.0

n [24]: Y_train_pred_meta = np.zeros(len(meta_train))
        meta_folds = StratifiedKFold(n_splits = 4, random_state = 9, shuffle = True)

        for fold, (train_idx, val_idx) in enumerate(meta_folds.split(meta_train, Outlier.values)):
            print("Training for fold {}".format(fold + 1))
            regressor_meta = Ridge(alpha = 1.0)
            regressor_meta.fit(meta_train[train_idx], Y_train.iloc[train_idx].values)
            Y_train_pred_meta[val_idx] = regressor_meta.predict(meta_train[val_idx])

Training for fold 1 .....
Training for fold 2 .....
Training for fold 3 .....
Training for fold 4 .....

n [25]: rmsc_meta_stack = np.sqrt(mean_squared_error(Y_train, Y_train_pred_meta))
        print("RMSE score for Stacked model with Meta Learner: ", rmsc_meta_stack)

RMSE score for Stacked model with Meta Learner: 3.654871417535515

```

The Meta Learner Stacked model with RSME score of 3.654871 performed little better than LightGBM but poorly compared to Weighted Stacked models.

### Final RSME Scores of all Models

```

In [27]: myTable = PrettyTable(["Model", "RSME Score"])
myTable.add_row(["Baseline", rsme_baseline])
myTable.add_row(["KNNRegressor", rsme_knn])
myTable.add_row(["XGBoost", rsme_xgb])
myTable.add_row(["LightGBM", rsme_lgb])
myTable.add_row(["Simple Stacked", rsme_stack])
myTable.add_row(["Meta-learner Stacked", rsme_meta_stack])
print(myTable)

```

Model	RSME Score
Baseline	3.850440680607971
KNNRegressor	3.824844001316022
XGBoost	3.6579771246304467
LightGBM	3.6550641942500763
Simple Stacked	3.6546577156656053
Meta-learner Stacked	3.654871417535515

## Submission

```

best_params_xgb = pickle.load(df_file)

In [29]: with open('data/LGB_parameters', 'wb') as df_file:
best_params_lgb = pickle.dump(df_file)

In [30]: def predict(X_train, Y_train, X_test, Outlier, param, num_splits = 4, num_round = 10000, model = 'lgb'):
This function predicts and returns the target value of test_data by training model on
train_data. It takes X_train and Y_train as dataframe, parameters for model as dictionary,
number of boosting rounds for model and whether to train LightGBM or XGBoost model as string.
Y_test_pred = np.zeros(len(X_test))
Folds = StratifiedFold(num_splits = num_splits, shuffle = True, random_state = 9)

if model == 'lgb':

    parameters = {
        'objective': 'regression',
        'metric': 'rmse',
        'boosting_type': 'gbdt',
        'learning_rate': 0.01,
        'device': 'cpu',
        'n_jobs': -1,
        'verbose': -1,
        'random_state': 9,
        'bagging_freq': 1,
        'bagging_seed': 9,
        'max_depth': param.get('max_depth'),
        'num_leaves': param.get('num_leaves'),
        'min_data_in_leaf': param.get('min_data_in_leaf'),
    }

```

```

        'feature_fraction': param.get('feature_fraction'),
        'bagging_fraction': param.get('bagging_fraction'),
        'min_split_gain': param.get('min_split_gain'),
        'reg_alpha': param.get('reg_alpha'),
        'reg_lambda': param.get('reg_lambda')
    }

    for fold, (train_idx, val_idx) in enumerate(Folds.split(X_train, Outlier.values)):
        train_data = lgb.Dataset(X_train.iloc[train_idx], label = Y_train.iloc[train_idx])
        val_data = lgb.Dataset(X_train.iloc[val_idx], label = Y_train.iloc[val_idx])
        regressor_LGB = lgb.TrainParams(parameters)
        regressor_LGB = lgb.Train(model=model, train_data=train_data, valid_sets=[train_data, val_data], num_boost_round=num_round, early_stopping_rounds=500, verbose_eval=
            Y_test_pred += (regressor_LGB.predict(X_test, num_iteration = regressor_LGB.best_iteration) /

```

```

        'req_lambda': param.get('req_lambda')
    }

    for fold, (train_idx, val_idx) in enumerate(folds.split(X_train, Outlier.values)):
        train_data = xgb.MatMatrix(X_train.iloc[train_idx], label = Y_train.iloc[train_idx])
        val_data = xgb.MatMatrix(X_train.iloc[val_idx], label = Y_train.iloc[val_idx])
        regressor_XGB = xgb.train(params = parameters, dtrain = train_data,
                                   evals = [(train_data, 'train'), (val_data, 'eval')], num_boost_round =
                                   early_stopping_rounds = 500, verbose_eval = False)
        Y_test_pred += (regressor_XGB.predict(xgb.MatMatrix(X_test,
                                                            iteration_range = (0, regressor_XGB.best_iteration)) / n

    return Y_test_pred

n [33]: Y_test_pred_xgb = predict(X_train, Y_train, X_test, Outlier, best_params_xgb, model = 'xgb')

n [32]: sub_xgb = pd.DataFrame(["card_id":test["card_id"].values])
sub_xgb["target"] = Y_test_pred_xgb
sub_xgb.to_csv("data/submit_xgb.csv", index = False)



n [33]: Y_test_pred_lgb = predict(X_train, Y_train, X_test, Outlier, best_params_lgb, model = 'lgb')

n [34]: sub_lgb = pd.DataFrame(["card_id":test["card_id"].values])
sub_lgb["target"] = Y_test_pred_lgb
sub_lgb.to_csv("data/submit_lgb.csv", index = False)

```

```

LightGBM Model Kaggle Score
n [35]: Y_test_pred_stack = (0.1 * Y_test_pred_xgb) + (0.9 * Y_test_pred_lgb)

n [36]: sub_stack = pd.DataFrame({"card_id":test["card_id"].values})
sub_stack["target"] = Y_test_pred_stack
sub_stack.to_csv("data/submit_stack.csv", index = False)

Stacked Model Kaggle Score
n [37]: meta_train = np.vstack([Y_train_pred_xgb, Y_train_pred_lgb]).transpose()
meta_test = np.vstack([Y_test_pred_xgb, Y_test_pred_lgb]).transpose()

n [40]: Y_train_pred_meta = np.zeros(len(meta_train))
Y_test_pred_meta = np.zeros(len(meta_test))
meta_folds = StratifiedKFold(n_splits = 4, random_state = 9, shuffle = True)

for fold, (train_idx, val_idx) in enumerate(meta_folds.split(meta_train, Outlier.values)):
    regressor_meta = Ridge(alpha = 1.0)
    regressor_meta.fit(meta_train[train_idx], Y_train.iloc[train_idx].values)
    Y_train_pred_meta[val_idx] = regressor_meta.predict(meta_train[val_idx])
    Y_test_pred_meta += (regressor_meta.predict(meta_test) / meta_folds.n_splits)

n [41]: sub_meta = pd.DataFrame({"card_id":test["card_id"].values})
sub_meta["target"] = Y_test_pred_meta
sub_meta.to_csv("data/submit_meta.csv", index = False)

```

Meta Learner Stacked Model Kaggle Score