

About:

The purpose of the notebook is to demonstrate how Lasso Regression can be utilized for feature selection. The data is related to music, the target variable will be the song's popularity and there are 10 features. The first step is to import all relevant libraries.

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5
        6 from sklearn.model_selection import train_test_split
        7 from sklearn.linear_model import Lasso
        8 from sklearn.linear_model import LinearRegression
        9
       10 from sklearn.model_selection import GridSearchCV
```

```
In [2]: 1 df1 = pd.read_csv('music.csv')
```

```
In [3]: 1 display(df1.head())
```

	popularity	acousticness	danceability	duration_ms	energy	instrumentalness	liveness	loudn
0	60	0.896000	0.726	214547	0.177	0.000002	0.1160	-14.
1	63	0.003840	0.635	190448	0.908	0.083400	0.2390	-4.
2	59	0.000075	0.352	456320	0.956	0.020300	0.1250	-3.
3	54	0.945000	0.488	352280	0.326	0.015700	0.1190	-12.
4	55	0.245000	0.667	273693	0.647	0.000297	0.0633	-7.

```
In [4]: 1 df1.rename(columns = {'popularity' : 'popularity_target'}, inplace=True)
```

In [5]:

```
1 df1.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 11 columns):
#   Column                Non-Null Count  Dtype
---  -
0   popularity_target     1000 non-null   int64
1   acousticness          1000 non-null   float64
2   danceability          1000 non-null   float64
3   duration_ms           1000 non-null   int64
4   energy                1000 non-null   float64
5   instrumentalness       1000 non-null   float64
6   liveness              1000 non-null   float64
7   loudness              1000 non-null   float64
8   speechiness           1000 non-null   float64
9   tempo                 1000 non-null   float64
10  valence                1000 non-null   float64
dtypes: float64(9), int64(2)
memory usage: 86.1 KB
```

In [6]:

```
1 X = df1.drop(columns = 'popularity_target')
2 y = df1['popularity_target']
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30,
```

In part I below I instigate a lasso regression to be integrated with the GRidSearchCV library to ensure the best value for the hyperparameter alpha is selected.

In part II below, I instigate a lasso regression model with the prescribed values for alpha determined by GridSearchCV in part I. X and y are fitted. The coefficients for each feature, sorted by most important -highest absolute value - are lastly printed:

In [7]:

```

1  #part I
2  models = {'LassoReg' : Lasso()}
3  hyperparameters = {'LassoReg' : {'alpha' : np.linspace(0.05, 0.15, 20)}}
4
5  grid_lasso_regression = GridSearchCV(
6  estimator = list(models.values())[0],
7  param_grid = list(hyperparameters.values())[0],
8  scoring='neg_mean_squared_error',
9  n_jobs=4,
10 cv = 3,
11 refit=True,
12 return_train_score=True)
13 grid_lasso_regression.fit(X, y)
14
15 #find the Root mean squared error below:
16 rmse = np.sqrt(-grid_lasso_regression.best_score_)
17 display('RMSE: ', rmse)
18 display('BEST VALUE FOR ALPHA: ', grid_lasso_regression.best_params_)
19
20 #Part II
21 LassoReg = Lasso(alpha = list(grid_lasso_regression.best_params_.values()))
22 lasso_coef = LassoReg.fit(X, y)
23 lasso_coef = LassoReg.fit(X, y).coef_
24 feature_coefficient_dict = dict(zip(X.columns, abs(lasso_coef)))
25 feature_coefficient_dict = {k: v for k, v in sorted(feature_coefficient_dict.items(),
26 key=lambda item: item[1], reverse=True)}
27 keyssortedbyimportance = list(feature_coefficient_dict.keys())
28 display('Features Sorted by Importance, Most Important on Top: '.upper(),

```

'RMSE: '

13.851349024427424

'BEST VALUE FOR ALPHA: '

{'alpha': 0.09210526315789475}

'FEATURES SORTED BY IMPORTANCE, MOST IMPORTANT ON TOP: '

```

['instrumentalness',
 'liveness',
 'acousticness',
 'danceability',
 'loudness',
 'tempo',
 'duration_ms',
 'energy',
 'speechiness',
 'valence']

```

Below: Conduct Linear regression with sklearn, starting with the most important feature ('instrumentalness') for the first iteration, then adding the second most important feature for the second iteration, and so on until every feature is used. record the training RMSE and testing RMSE.

```

In [8]: 1 train_scores = []
        2 test_scores = list()
        3
        4 i = 1
        5 while i < len(keyssortedbyimportance) + 1:
        6     X = df1.drop(columns='popularity_target')
        7     X = df1[keyssortedbyimportance[:i]]
        8
        9     y = df1['popularity_target']
       10     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
       11
       12     LinReg = LinearRegression()
       13
       14     parameters = {'fit_intercept' : [True, False], 'positive' : [True, Fal
       15
       16     grid_rf_regression = GridSearchCV(
       17                                     estimator = LinReg,
       18                                     param_grid = parameters,
       19                                     scoring='neg_mean_squared_error',
       20                                     #scoring='r2',
       21                                     n_jobs=4,
       22                                     cv = 3,
       23                                     refit=True,
       24                                     return_train_score=True)
       25     grid_rf_regression.fit(X_train, y_train)
       26
       27     train_scores.append(np.sqrt(-grid_rf_regression.score(X_train, y_train
       28     test_scores.append(np.sqrt(-grid_rf_regression.score(X_test, y_test)))
       29
       30     i += 1

```

Plot the results:

```

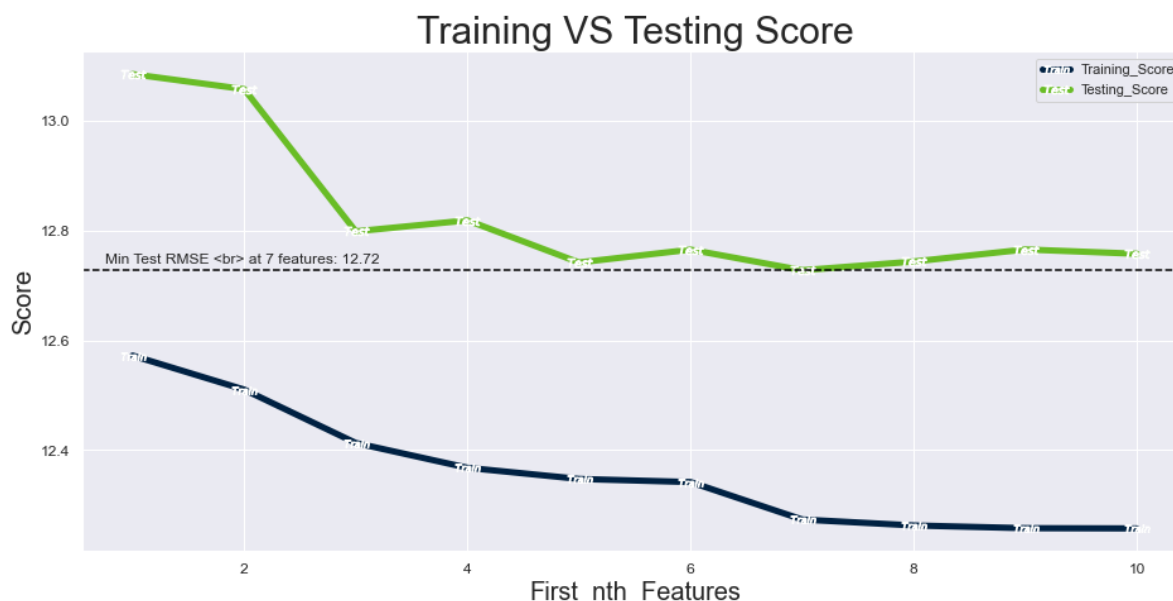
In [9]: 1 sns.set_style('darkgrid')
        2 sns.set(rc={'figure.figsize':(15.7,7.27)})

```

```

In [10]: 1 training = sns.lineplot(x=range(1, len(keyssortedbyimportance)+1), y=train
2           linewidth = 5, color = (0, 0.13, 0.26), label = 'T
3           marker = '$Train$', ms=20, markeredgecolor = 'whit
4
5 testing = sns.lineplot(x=range(1, len(keyssortedbyimportance)+1), y=test_s
6           linewidth = 5, color = (0.41, 0.74, 0.15), label =
7           marker = '$Test$', ms=20, markeredgecolor = 'white'
8
9 plt.title('Training VS Testing Score', fontsize=30)
10 training.set_ylabel('Score', fontsize=20)
11 training.set_xlabel('First_nth_Features', fontsize=20)
12 plt.tick_params(axis='both', labelsize=12)
13 plt.axhline(y = np.min(test_scores), color = 'Black', linestyle = '--')
14 plt.text(0.75, 12.74, 'Min Test RMSE <br> at 7 features: 12.72')
15 plt.show()

```



From the chart above, the test score of RMSE is minimized by using the first seven most essential features:

['instrumentalness', 'liveness', 'acousticness', 'danceability', 'loudness', 'tempo', 'duration_ms']
 once eight or more features are added, the linear regression model begins to show signs of overfitting - the testing score begins to rise, or degrade as the smaller the RMSE the better at predicting the model, while the training score continues to fall.

Note: The rise in testing RMSE is not that drastic in this example however, it is the best I could find within my timeframe. This same notebook should be able to do the same task on any other dataframe with only minor adjustments. Also, in previous versions, I preprocessed the data by scaling but decided not to move forward, as it made a minimal overall impact and was beyond the scope of what I intended for this project.