

# Data Analysis and Interpretation Capstone

## Title: Report on Net Migration causing factors ( population, industry, urbanization ) in year 2012

### Introduction to the Research Question

This report is examining factors that cause human migration in year 2012. We will examine and check any relationships between predictor variables (FDIs,Health,Facilities,Populations) in regards to target variable which is Net Migration.

The dataset is provided by World Bank.

### Data Dictionary

	Field	Description
	x131_2012	FOREIGN DIRECT INVESTMENT, NET INFLOWS (% OF GDP)
	x142_2012	GDP PER CAPITA (CURRENT US\$)
	x150_2012	HEALTH EXPENDITURE, TOTAL (% OF GDP)
	x155_2012	IMPROVED SANITATION FACILITIES (% OF POPULATION WITH ACCESS)
	x156_2012	IMPROVED WATER SOURCE (% OF POPULATION WITH ACCESS)
	x1_2012	ACCESS TO ELECTRICITY (% OF POPULATION)
	x258_2012	RURAL POPULATION (% OF TOTAL POPULATION)
	x283_2012	URBAN POPULATION (% OF TOTAL)
	x81_2012	AGRICULTURAL LAND (% OF LAND AREA)
	x195_2012	NET MIGRATION

### Data Preparation

```
In [1]: import pandas as pd

In [2]: df = pd.read_csv('worldbank.csv')
df.head()
```

	country	x1_2012	x2_2012	x9_2012	x11_2012	x12_2012	x14_2012	x15_2012	x16_2012	x18_2012	...	x244_2013	
0	AFGHANISTAN	43.000000	19.480962	1.852684e+10	623.26804	0.677345	8.965464	1.833022e+09	1.600000	0.125552	...	0.0	
1	ALBANIA	104.000000	62.086412	1.037147e+10	3575.76667	0.386039	12.381013	1.521033e+09	2.842804	1.2730370	...	13.1	
2	ALGERIA	100.000000	99.990000	1.480000e+11	3942.202841	0.620760	7.019903	1.417333e+10	4.467196	19.550386	...	2.4	
3	AMERICAN_SAMOA	59.32891	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	
4	ANDORRA	100.000000	100.000000	NaN	NaN	NaN	0.152983	9.237132	2.919443e+08	3.100000	0.000000	...	NaN

5 rows x 163 columns

```
In [3]: df2 = df[['x131_2012','x142_2012','x150_2012','x155_2012','x156_2012','x1_2012','x258_2012','x283_2012','x31_2012',
df2.head()
```

	x131_2012	x142_2012	x150_2012	x155_2012	x156_2012	x1_2012	x258_2012	x283_2012	x31_2012	x195_2012
0	0.299592	690.842629	8.479199	30.5	51.6	43.000000	74.532	25.468	58.067580	473007.0
1	7.468318	4247.485437	5.627439	92.1	95.4	100.000000	45.670	54.330	43.843066	-91750.0
2	0.717733	5583.616160	6.007278	87.0	84.9	100.000000	31.130	68.870	17.381700	-143268.0
3	NaN	NaN	NaN	62.5	100.0	59.32891	12.587	87.413	24.500000	NaN
4	NaN	39666.369210	8.320219	100.0	100.0	100.000000	13.292	86.708	42.978723	NaN

```
In [5]: df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 248 entries, 0 to 247
Data columns (total 10 columns):
# Column Non-Null Count Dtype
---
0 x131_2012 219 non-null float64
1 x142_2012 225 non-null float64
2 x150_2012 221 non-null float64
3 x155_2012 226 non-null float64
4 x156_2012 227 non-null float64
5 x1_2012 245 non-null float64
6 x258_2012 245 non-null float64
7 x283_2012 245 non-null float64
8 x31_2012 240 non-null float64
9 x195_2012 127 non-null float64
dtypes: float64(10)
memory usage: 19.5 KB
```

```
In [6]: df2.columns = ['fdi','gdp','health','sanitation','water','electricity','ruralpop','urbanpop','agri','migration']
```

```
In [7]: df2
```

	fdi	gdp	health	sanitation	water	electricity	ruralpop	urbanpop	agri	migration
0	0.299592	690.842629	8.479199	30.5	51.6	43.000000	74.532	25.468	58.067580	473007.0
1	7.468318	4247.485437	5.627439	92.1	95.4	100.000000	45.670	54.330	43.843066	-91750.0
2	0.717733	5583.616160	6.007278	87.0	84.9	100.000000	31.130	68.870	17.381700	-143268.0
3	NaN	NaN	NaN	62.5	100.0	59.32891	12.587	87.413	24.500000	NaN
4	NaN	39666.369210	8.320219	100.000000	100.000000	100.000000	13.292000	86.708000	42.978723	NaN
...	...	...	...	...	...	...	...	...	...	...
243	0.559266	2782.95026	NaN	92.2000	65.100000	97.697830	25.422000	74.577000	43.354582	-43750.0
244	2.537712	10460.113770	9.940811	65.93004	89.487651	84.583874	47.551252	52.448806	37.785216	0.0
245	-0.044394	1289.034078	5.553370	53.30000	54.900000	48.406710	67.126000	32.874000	44.604807	-50000.0
246	6.942853	1686.618024	4.754491	43.20000	63.000000	22.062560	60.413000	39.587000	32.063923	-34490.0
247	3.223668	850.827694	NaN	37.3000	77.500000	40.462560	67.166000	32.834000	41.876696	-19922.0

248 rows x 10 columns

```
In [8]: df2.to_csv('wbresearch.csv', index=False)
```

```
In [9]: import seaborn as sns
import matplotlib.pyplot as plt
df = pd.read_csv('wbresearch.csv')
df.head()
```

	fdi	gdp	health	sanitation	water	electricity	ruralpop	urbanpop	agri	migration
0	0.299592	690.842629	8.479199	30.5	51.6	43.000000	74.532	25.468	58.067580	473007.0
1	7.468318	4247.485437	5.627439	92.1	95.4	100.000000	45.670	54.330	43.843066	-91750.0
2	0.717733	5583.616160	6.007278	87.0	84.9	100.000000	31.130	68.870	17.381700	-143268.0
3	NaN	NaN	NaN	62.5	100.0	59.32891	12.587	87.413	24.500000	NaN
4	NaN	39666.369210	8.320219	100.000000	100.000000	100.000000	13.292000	86.708000	42.978723	NaN
...	...	...	...	...	...	...	...	...	...	...
243	0.559266	2782.95026	NaN	92.2000	65.100000	97.697830	25.422000	74.577000	43.354582	-43750.0
244	2.537712	10460.113770	9.940811	65.93004	89.487651	84.583874	47.551252	52.448806	37.785216	0.0
245	-0.044394	1289.034078	5.553370	53.30000	54.900000	48.406710	67.126000	32.874000	44.604807	-50000.0
246	6.942853	1686.618024	4.754491	43.20000	63.000000	22.062560	60.413000	39.587000	32.063923	-34490.0
247	3.223668	850.827694	NaN	37.3000	77.500000	40.462560	67.166000	32.834000	41.876696	-19922.0

248 rows x 10 columns

```
In [8]: df2.to_csv('wbresearch.csv', index=False)
```

```
In [9]: import seaborn as sns
import matplotlib.pyplot as plt
df = pd.read_csv('wbresearch.csv')
df.head()
```

	fdi	gdp	health	sanitation	water	electricity	ruralpop	urbanpop	agri	migration
0	0.299592	690.842629	8.479199	30.5	51.6	43.000000	74.532	25.468	58.067580	473007.0
1	7.468318	4247.485437	5.627439	92.1	95.4	100.000000	45.670	54.330	43.843066	-91750.0
2	0.717733	5583.616160	6.007278	87.0	84.9	100.000000	31.130	68.870	17.381700	-143268.0
3	NaN	NaN	NaN	62.5	100.0	59.32891	12.587	87.413	24.500000	NaN
4	NaN	39666.369210	8.320219	100.000000	100.000000	100.000000	13.292	86.708	42.978723	NaN
...	...	...	...	...	...	...	...	...	...	...
243	0.559266	2782.95026	NaN	92.2000	65.100000	97.697830	25.422000	74.577000	43.354582	-43750.0
244	2.537712	10460.113770	9.940811	65.93004	89.487651	84.583874	47.551252	52.448806	37.785216	0.0
245	-0.044394	1289.034078	5.553370	53.30000	54.900000	48.406710	67.126000	32.874000	44.604807	-50000.0
246	6.942853	1686.618024	4.754491	43.20000	63.000000	22.062560	60.413000	39.587000	32.063923	-34490.0
247	3.223668	850.827694	NaN	37.3000	77.500000	40.462560	67.166000	32.834000	41.876696	-19922.0

248 rows x 10 columns

```
In [8]: df2.to_csv('wbresearch.csv', index=False)
```

```
In [9]: import seaborn as sns
import matplotlib.pyplot as plt
df = pd.read_csv('wbresearch.csv')
df.head()
```

	fdi	gdp	health	sanitation	water	electricity	ruralpop	urbanpop	agri	migration
0	0.299592	690.842629	8.479199	30.5	51.6	43.000000	74.532	25.468	58.067580	473007.0
1	7.468318	4247.485437	5.627439	92.1	95.4	100.000000	45.670	54.330	43.843066	-91750.0
2	0.717733	5583.616160	6.007278	87.0	84.9	100.000000	31.130	68.870	17.381700	-143268.0
3	NaN	NaN	NaN	62.5	100.0	59.32891	12.587	87.413	24.500000	NaN
4	NaN	39666.369210	8.320219	100.000000	100.000000	100.000000	13.292	86.708	42.978723	NaN
...	...	...	...	...	...	...	...	...	...	...
243	0.559266	2782.95026	NaN	92.2000	65.100000	97.697830	25.422000	74.577000	43.354582	-43750.0
244	2.537712	10460.113770	9.940811	65.93004	89.487651	84.583874	47.551252	52.448806	37.785216	0.0
245	-0.044394	1289.034078	5.553370	53.30000	54.900000	48.406710	67.126000	32.874000	44.604807	-50000.0
246	6.942853	1686.618024	4.754491	43.20000	63.000000	22.062560	60.413000	39.587000	32.063923	-34490.0
247	3.223668	850.827694	NaN	37.3000	77.500000	40.462560	67.166000	32.834000	41.876696	-19922.0

248 rows x 10 columns

```
In [8]: df2.to_csv('wbresearch.csv', index=False)
```

```
In [9]: import seaborn as sns
import matplotlib.pyplot as plt
df = pd.read_csv('wbresearch.csv')
df.head()
```

	fdi	gdp	health	sanitation	water	electricity	ruralpop	urbanpop	agri	migration
0	0.299592	690.842629	8.479199	30.5	51.6	43.000000	74.532	25.468	58.067580	473007.0
1	7.468318	4247.485437	5.627439	92.1	95.4	100.000000	45.670	54.330	43.843066	-91750.0
2	0.717733	5583.616160	6.007278	87.0	84.9	100.000000	31.130	68.870	17.381700	-143268.0
3	NaN	NaN	NaN	62.5	100.0	59.32891	12.587	87.413	24.500000	NaN
4	NaN	39666.369210	8.320219	100.000000	100.000000	100.000000	13.292	86.708	42.978723	NaN
...	...	...	...	...	...	...	...	...	...	...
243	0.559266	2782.95026	NaN	92.2000	65.100000	97.697830	25.422000	74.577000	43.354582	-43750.0
244	2.537712	10460.113770	9.940811	65.93004	89.487651	84.583874	47.551252	52.448806	37.785216	0.0
245	-0.044394	1289.034078	5.553370	53.30000	54.900000	48.406710	67.126000	32.874000	44.604807	-50000.0
246	6.942853	1686.618024	4.754491	43.20000	63.000000	22.062560	60.413000	39.587000	32.063923	-34490.0
247	3.223668	850.827694	NaN	37.3000	77.500000	40.462560	67.166000	32.834000	41.876696	-19922.0

248 rows x 10 columns

```
In [8]: df2.to_csv('wbresearch.csv', index=False)
```

```
In [9]: import seaborn as sns
import matplotlib.pyplot as plt
df = pd.read_csv('wbresearch.csv')
df.head()
```

```

import xgboost as xgb

#from xgboost import XGBClassifier, XGBRegressor
#from xgboost import plot_importance

#from sklearn.experimental import enable_hist_gradient_boosting
#from sklearn.linear_model import ElasticNet, Lasso, LinearRegression, LogisticRegression, Ridge
#from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor, ExtraTreeClassifier, Extra
#from sklearn.ensemble import GradientBoostingClassifier, GradientBoostingRegressor, HistGradientBoostingClassifier

%matplotlib inline
#sets the default autoreave frequency in seconds
%autoreave 60
sns.set_style('dark')
sns.set(font_scale=1.2)

```



```
In [38]: (array([[ 0.29959, 690.8462629, 8.47919927, ..., 74.532,
[ 25.461, 58.06175957],
[ 7.46831756, 42.84748843, 5.62743862, ..., 45.67,
[ 54.33, 43.8430569],
[ 6.00727769, ..., 31.13,
[ 66.87, 17.38169993],
...,
[ -0.04439374, 1289.034079, 5.55337022, ..., 67.126,
[ 32.874, 44.60480708],
[ 6.94285335, 1686.618024, 4.75449068, ..., 60.413,
[ 30.587, 22.063923, 6.2333354, ..., 67.166,
[ 3.22567999, 850.827694, 6.2333354, ...,
array([ 473007, -917500, 143268, -10000, -10000,
[ 102322, -56, -1168750, 30000, -9876,
[ 3253, 147089, 2182, 120535, 269988,
[ 29915, -2226481, 2180, 120535, 269988,
[ 7594, -10000, -10000, 10000, -161794,
[ -2306, 20000, 15924, 2102, 50000,
[ -125000, 40000, -11052, -149999, -60000,
[ 1175963, -121221, -1000, -175013, -5589,
[ -10000, 8654, 201289, -180000, -144998,
[ -1000, -95920, -60000, 19658, 50000,
[ -2000, -78903, -10000, 15310, 29898,
[ 96839, -15998, -10000, -153010, -1457489,
[ -728059, -36003, -21581, 240415, 20000,
[ -160001, -118501, -60001, 157503, 547385,
[ 1003278, 2324066, -10000, -28720, 107409,
[ -405505, 193461, 200003, -25001, -10000,
[ -296323, 1249998, -49999, -136299, -10000,
[ -4274, 8, -120001, -10000, 12476,
[ -20278, -150000, -6072487, 1648636, 4722878,
[ 11735448, -80000, 150000, 23999, -378,
[ -259821, -700000, -300001, 58666, -140001,
[ -1000, 1987, 52828, -2700, 350000,
[ 229617, 159807, -50000, -2130, 0,
[ 30000, -10000, 15350, -113963, -117700,
[ -2081948, -2120734, -73442, -4741314, 1250000,
[ -19998, -20000, -501692, -10000, -169529,
[ -1644928, -1166621, -6002193, 46104, 35000,
[ -4999, -5000, -30000, 450000, -53,
[ -302449, 6252, -10000, 20000, 0,
[ -52385, -8192, -21302, -361566, 15283807,
[ -9529, -10000, -15001, -2412, -310624,
[ -23500, -475276, -137, -373399, 110066,
[ 5660, 7265, -135000, -28497, -300000,
[ 6183750, -10000, 235665, -10000, 1321866,
[ 1211010, -95513, -71075, -1081818, -10000,
[ 28105, 0, -86700, -240000, -70000,
[ 13803, -140000, -10337, 363500, -437201,
[ 1117884, -75001, -12630, -1000, -5589,
[ 85000, -99998, -99999, -1551, -21000,
[ 397936, -10000, 0, 119, -10000,
[ -11868, -400000, 600000, -6288064, 865000,
[ -593069, -484772, -10000, 40, -10000,
[ -5000, 1684894, -1707443, -800069, -28779,
[ -6000, 272626, 382267, -4023996, -117382,
[ -199999, 100000, -50004, -9994, -8078,
[ -5000, -13941, 200003, -25001, -10000,
[ -1000, -150000, 195000, 495000, 900000,
[ 5007807, 71584, -30000, -195003, 803,
[ -69121, -200002, -3604, -43750, 0,
[ -50000, -34490, -21992,]])
```

```
In [39]: X_train, X_test, y_train, y_test = train_test_split(X.values, y.values, test_size=0.2, random_state=0)
```

```
In [40]: X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

```
Out[40]: ((198, 9), (50, 9), (198,), (50,))
```

## Feature Selection

### Recursive Feature Elimination with Cross Validation (Test each ML model)

```
In [41]: model1 = LinearRegression()
```

```
In [42]: model2 = Lasso(random_state=0)
```

```
In [43]: model3 = DecisionTreeRegressor(random_state=0)
```

```
In [44]: model4 = RandomForestRegressor(random_state=0)
```

```
In [45]: rfecv = RFECV(estimator=model2, cv=5, scoring='neg_log_mean_squared_error', verbose=1, step=1, n_jobs=-1)
```

```
In [46]: rfecv.fit(X_train, y_train)
```

```
Fitting estimator with 9 features.
Fitting estimator with 8 features.
Fitting estimator with 7 features.
```

```
Out[46]: RFECV(cv=5, estimator=Lasso(random_state=0), n_jobs=-1,
scoring='neg_log_mean_squared_error', verbose=1)
```

```
In [47]: print("Optimal no of features:", rfecv.n_features_)
```

```
Optimal no of features: 6
```

```
In [48]: print("Best features:", rfecv.support_)
```

```
Best features: [False False True True True True False True]
```

```
In [49]: print("Selector Ranking:", rfecv.ranking_)
```

```
Selector Ranking: (2 4 1 1 1 1 1 3 1)
```

```
In [50]: X.columns
```

```
Out[50]: Index(['fai', 'gdp', 'health', 'sanitation', 'water', 'electricity', 'ruralpop', 'urbanpop', 'agri'], dtype='<ob
ject')
```

Model 2 (Lasso) selects 6 out of 9 variables for our new model.

## Feature Scaling

```
In [51]: X_train
```

```
Out[51]: array([[ 2.80294079, 5445.894718, 5.63399128, ...,
[ 45.875, 54.102, 40.99722992,
[ 2.30169913, 5721.834908, 5.00823663, ...,
[ 24.533, 75.467, 43.34928813],
[ 1.85941408, 28647.83524, 3.29736719, ...,
[ 21.098, 79.802, 53.86137822],
...,
[ -0.65024907, 50903.9046, 2.56544882, ...,
[ 1.706, 98.294, 8.52974186],
[ 1.753259992, 3191.164299, 31.04538799],
[ 35.9, 64.1, 3.9750746, ...,
[ 5.1142608, 3537.281329, 5.41421679, ...,
[ 63.23772629, 36.76227371, 13.41227797]])
```

```
In [52]: scaler = StandardScaler()
```

```
In [53]: X_train_scaled = scaler.fit_transform(X_train)
```

```
In [54]: X_test_scaled = scaler.transform(X_test)
```

```
In [55]: X_train_scaled
```

```
Out[55]: array([[-0.25294346, -0.40498985, -0.38721389, ..., 0.1606332,
[-0.1608367, 0.2115446,
[-0.29434692, 0.34441394, -0.64022407, ..., -0.74778921,
[0.74778909, 0.33683209],
[0.33072528, 0.69618709, 1.04025178, ..., -0.89401371,
[0.89401177, 0.84200913],
...,
[-0.34796829, 1.75245308, -1.59729721, ..., -1.71951171,
[1.71951099, -1.37957819],
[0.81236308, -0.5119784, -1.0448789, ..., -0.26390741,
[0.26390454, -0.27615399],
[0.06240442, -0.49551116, -0.48115696, ..., 0.89983219,
[0.89983679, -1.14031518]])
```

```
In [56]: X_test_scaled
```

```
Out[56]: array([[-0.21825984, -0.62030009, -0.16022902, -1.28579295, -0.82190573,
[-1.25588248, 1.06705486, -1.06705971, 0.25465381,
[-0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.38194575,
[0.60396636, -0.0827876, -0.08282204, 0.9048102, 0.80982068,
[-0.35264338, 1.27531468, 1.9464014, 0.9048102, 0.80982068,
[0.73574973, -0.89052021, -0.0592211, 0.9048102, 0.80982068,
[-0.16422809, -0.563021, 2.01438977, 0.08795037, -0.06768553,
[0.73574973, -0.55524781, -0.3552519, 1.87125621,
[0.34635479, -0.20291116, -0.02107303, 0.38914491, 0.45446692,
[0.70564314, -0.87294206, 0.87294009, 0.89245961,
[0.06963358, -0.44747586, -1.0420818, 0.63371262, 0.49072751,
[0.62479885, -0.23863636, -0.22880994, -0.65732531,
[-0.32862649, -0.62782024, -0.04154747, -1.39595223, 0.09186106,
[0.29121933, -1.215989, -0.21642405, 1.76420991,
[1.903031087, 1.68462375, 2.36509099, 0.87270654, 0.80982068,
[0.73574973, -1.30578362, 1.30578229, 0.879061991,
[0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.38194575,
[0.61072578, 0.0613623, -0.06136566, 1.97666951,
[-0.18004802, -0.58453845, 0.01199471, -1.7772895, -2.41737156,
[1.52013139, -1.04657916, -0.04658399, 0.52693281,
[-0.00903671, -0.54976066, 1.10055631, 0.23776745, 0.03384412,
[0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.38194575,
[-0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.38194575,
[-1.51528784, 0.84063037, 0.84063489, 1.64615591,
[0.42618575, 0.38765552, 1.0290463, 0.9551145, 0.80982068,
[0.73574973, -0.81304748, 0.81304542, 1.29949261,
[-0.51871214, -0.56133493, -0.89533907, -1.9411448, -3.54870187,
[-0.00372003, -1.9121346, -1.9121407, -1.68842591,
[-0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.80982068,
[0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.80982068,
[-0.388044, -0.30303866, 0.89257829, -0.3079948, 0.22965129,
[0.24734402, -0.22866028, 0.22865735, 2.114677031,
[0.41336675, -0.07196214, -0.70399952, 0.5601387, 0.80982068,
[0.73574973, -1.24478225, 1.24478083, 1.233717371,
[-0.11271842, -0.49253313, -0.16022902, 0.38414084, 0.38194575,
[0.73574973, -0.0827876, -0.08282204, 0.9048102,
[-0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.38194575,
[0.40303713, -1.215989, -0.21642405, 1.76420991,
[-0.05081498, -0.62401392, -0.19182689, -1.692671, -1.86799726,
[-1.63417797, 1.04999959, -1.05000452, 0.144663541,
[0.40303713, -0.58203758, -1.14020109, -1.00317957, 0.12948642,
[-0.00157853, 1.11192834, -1.111944, 0.976246181,
[-0.13981349, -0.43991816, 0.09907943, 0.63034594, 0.76630797,
[-0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.38194575,
[-0.00903671, -0.54976066, 1.10055631, 0.23776745, 0.03384412,
[-0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.38194575,
[-0.42705302, -0.63689946, -0.18871417, -1.89334252, -1.10473831,
[-0.58202486, 0.46729769, -0.46726169, 1.89334252, -1.10473831,
[0.14661933, -0.28207558, 0.01358699, 0.0264016, 0.17169773,
[0.63539445, -0.75357859, 0.75357644, 0.084845281,
[0.59592229, -0.01596136, -0.0430434, 0.38414084, 0.68603468,
[0.45051713, -1.10191807, -1.10203305, -0.6667715, 1.
[-0.24513354, -0.38023812, -0.16022902, 0.38414084, 0.62126662,
[0.40303713, -1.79233442, -1.7923381, 1.26460776, 0.30942458,
[0.06982529, -0.08820446, -0.91170651, 0.86200552,
[0.73574973, 0.18911186, -0.18911541, 2.143815751,
[-0.32020343, -0.47130468, 0.02635486, -0.02635959, 0.20064282,
[0.52500363, -0.31443681, 0.31443401, 1.920542391,
[-0.48325316, 1.55195309, 1.42540926, 0.95318216, 0.80982068,
[0.73574973, -1.44741068, 1.44740956, 1.18809981,
[0.73574973, -0.0827876, -0.08282204, 0.9048102,
[0.73574973, -0.95769578, 0.95769979, 0.9048102,
[0.27499168, 0.01798811, 0.04318401, 0.79066385, 0.70829103,
[-0.72905938, -1.42319713, 1.42319704, 0.8793761,
[0.21548824, -0.59280778, 0.59280778, -1.38843806, 0.59225715,
[-0.58684726, -0.22866028, 0.22865937, 0.08946062, 0.17169773,
[0.43926392, 0.44134733, -1.22701103, 0.0558467, -2.38303671,
[-0.40161014, -0.78484495, -0.78484941, -1.3012385, 1.
[0.37327087, 0.01726353, 0.0081711, 0.87977969, 0.90523175,
[0.73574973, -0.37050018, 0.37049744, 0.24630361,
[-0.39714659, -0.40089374, -1.24130902, -0.86465827, -2.9757846,
[-1.37171121, 0.6898607, -0.68986499, 0.522179,
[-0.41321029, 0.3924013, -0.3924013, 0.38414084, 0.38194575,
[0.73574973, -1.79233442, 1.7923381, 1.233717371,
[-0.28812662, -0.60843361, 2.15623542, -1.97073885, -0.55307739,
[1.30137253, -1.36831795, 1.36813325, 1.893136061,
[-0.37838348, 1.63002369, 1.74523224, 0.95318216, 0.74455162,
[0.73574973, -0.3899912, 0.38999376, 0.9048102,
[-0.45248809, -0.60339678, -1.03786663, -0.5362875, -0.61884645,
[0.52168646, 0.87149289, -0.87149745, 0.495044281,
[-0.04129731, -0.58015705, -0.58015705, -1.2242017, 0.3601894,
[0.70229797, 1.11669031, -1.11669524, -0.082931561,
[-0.20010359, -0.4646006, 0.15682686, 0.60517603, 0.008500504,
[0.73574973, -1.35645241, 0.35641967, 1.32224591,
[-0.48387208, 0.99027874, 0.99718102, 0.93334679, 0.80982068,
[0.73574973, -0.45375434, 0.45371625, 0.9048102,
[-0.33679102, -0.59096826, -1.9027389, 0.21636501,
[-0.85780653, 1.08289255, -1.08289742, -0.85831961,
[-0.29251898, -0.58731087, -0.58739436, -1.0303351,
[-1.17775303, 0.76397826, -0.76398266, -0.104967911,
[0.02598216, -0.36021753, -0.36082779, 0.01304182, -0.24173634,
[-0.41314623, -0.85920398, 0.85920796, 0.86603691,
[-0.37838348, 1.77868661, 0.40767973, 0.95318216, 0.74455162,
[0.73574973, -0.87920029, 0.8791912, 0.9048102,
[-0.27022302, -0.47551109, -0.3457602, -0.51925251, -0.11447803,
[0.09864518, 0.51396272, -0.51396599, 0.32999471,
[-0.32656111, 1.58689734, 0.58045752, -0.78484941, 0.80982068,
[0.73574973, -1.10325247, 1.10325284, -0.429307831,
[-0.45167216, -0.46830136, -0.86275124, 0.60484676, 0.76630797,
[0.59746307, 1.46447302, -1.46448446, 0.31243871])]
```

## Model Training

### Using Regression or Classification Models

#### K-Fold Cross-Validation (Generalization Performance)

```
In [57]: lasso = Lasso(random_state=0)
```

```
In [58]: kf = KFold(n_splits=5, shuffle=True, random_state=0)
```

```
In [59]: lasso_cv = cross_validate(estimator=lasso, X=X_train_scaled, y=y_train, scoring='neg_log_mean_squared_error',
cv=kf, n_jobs=-1, return_train_score=True)
```

```
In [60]: lasso_cv
```

```
{ 'fit_time': (0.00099778, 0.00199533, 0.00498724, 0.00099826, 0.00588478)),
[ 'score_time': array([1.318096, 6.956839, -24.904011, 117.8899, -2053996.62388626,
[-2.687320, 11.584607, -0.958145, 50.324661],
[ 'train_score': array([-2.078470, 47.000646, -2.299170, 975.0352, -2.392965, 577330008,
[-2.241988, 45.163901, -25.64741, 64.225141])]
```

```
In [61]: np.mean(lasso_cv['train_score']), np.std(lasso_cv['train_score'])
```

```
Out[61]: (-2315467.423246431, 161230.20156291054)
```

```
In [62]: np.mean(lasso_cv['test_score']), np.std(lasso_cv['test_score'])
```

```
Out[62]: (-2247045.07026801, 773720.771476094)
```

Cross-validation scores for both train and test acceptable.

```
In [63]: lasso_model = Lasso(random_state=0)
```

```
In [64]: lasso_model.fit(X_train_scaled, y_train)
```

```
Out[64]: Lasso(random_state=0)
```

```
In [65]: lasso_pred = lasso_model.predict(X_test_scaled)
```

```
In [66]: lasso_pred
```

```
Out[66]: array([-603435.19120308, -334056.55528799, 1987713.28014251,
[0.89601, 15283336, -2345019.012506, -846566.86178877,
[-1587641.71599513, 2235017.67116915, -1277578.7963206,
[-69211.84567274, 720727.6174006, -192446.88547985,
[899073.96022103, 20973.19275671, 97360.44207589,
[41557.53744135, 2119691.39370042, -429106.99727561,
[598065.20102145, -302576.59632366, -2128832.74470593,
[-368510.18863117, 160753.01769488, -967119.48110553,
[117916.85658373, -511335.09822084, 32388.26686068,
[1104842.04324412, -619345.35974784, 1944631.42016344,
[-412697.85298733, 310369.06215966, -630600.1272115,
[518008.87434103, 34716.63649429, -41901.05256173,
[1851939.54882276, 61945.35974784, 1944631.42016344,
[-0.72905938, -1.42319713, 1.42319704, 0.8793761,
[0.21548824, -0.59280778, 0.59280778, -1.38843806, 0.59225715,
[-0.58684726, -0.22866028, 0.22865937, 0.08946062, 0.17169773,
[0.43926392, 0.44134733, -1.22701103, 0.0558467, -2.38303671,
[-0.40161014, -0.78484495,
```