

08 Lasso, Ridge, and Elastic-Net Regression in Supervised Learning

Linear regression is a fundamental tool in statistics and machine learning. However, when dealing with high-dimensional data, or when multicollinearity (correlation between predictors) is present, plain linear regression can have limitations. To address these challenges we use penalized regression techniques such as Lasso, Ridge, and Elastic-Net regression.

Penalized regression models:

- handle feature selection
- prevent overfitting from multiclinear predictors

Penalized regression's feature selection is not a substitute for subject matter knowledge driven feature selection.

2. Experiment with different values of λ (from very small to very large). What do you observe in terms of coefficient shrinkage?

whenever the lambda seerds to increase, my coefficient numbers would decrease/shrink.

3. Using the same dataset, split it into a training and test set. Evaluate the performance of ridge, lasso, and elastic-net on the test set. Which method performs the best and why? Looking at the last chart to compare all three, it seems that the LassoCV performed the best out of all three. This is because it has the lowest score of all methods, meaning it had better performance

4. Reflect on the importance of feature selection in real-world applications. Discuss scenarios where you'd prefer lasso over ridge regression and vice-versa. Feature selection is important for when you do not need added variables that make the data more 'messy'. This allows for a 'cleaner' seeming dataset so it doesn't get too messy with deleting or changing variables often. I would say you use Lasso for when you would like an easier sort of model that has most coefficients (ying at 0, or close to it. You would use Ridge Regression when there is a high amount of predictor variables and that number is bigger than the amount of observations you have.

Load our data

```
In [32]: import pandas as pd
import numpy as np

df = pd.read_csv(
    filepath_or_buffer = "D:\\output\\data_prepared.csv",
    engine = 'python',
    dtype = {}

    'cat_id': str,
    'symboling': float,
    'displacement': float,
    'wheelbase': float,
    'carlength': float,
    'carwidth': float,
    'carheight': float,
    'curbweight': float,
    'enginehorsepower': float,
    'enginesize': float,
    'compressionratio': float,
    'horsepower': float,
    'peakrpm': float,
    'citympg': float,
    'highwaympg': float,
    'trains': float,
    'carname_1_low': str,
    'carname_2_moderate': int,
    'carname_3_high': int,
    'drivewheel_1_low': int,
    'drivewheel_2_moderate': int,
    'drivewheel_3_high': int,
    'enginelocation_2_moderate': int,
    'enginelocation_3_high': int,
    'aspiration_2_moderate': int,
    'aspiration_3_high': int,
    'fuelsystem_1_low': int,
    'fuelsystem_2_moderate': int,
    'fuelsystem_3_high': int,
    'fueltype_2_moderate': int,
    'fueltype_3_high': int,
    'fueltype_1_low': int,
    'carbody_1_low': int,
    'carbody_2_moderate': int,
    'carbody_3_high': int
)

<SyntaxWarning: invalid escape sequence '\o'
C:\Users\andri\AppData\Local\Temp\ipykernel_20160123\6787289.py:41: SyntaxWarning: invalid escape sequence '\o'
filepaths_or_buffers = "D:\\output\\data_prepared.csv"
```

```
In [33]: list_full = [
    'symboling',
    'displacement',
    'wheelbase',
    'carlength',
    'carwidth',
    'carheight',
    'curbweight',
    'enginehorsepower',
    'enginesize',
    'compressionratio',
    'horsepower',
    'peakrpm',
    'citympg',
    'highwaympg',
    'carname_1_low',
    'carname_2_moderate',
    'carname_3_high',
    'drivewheel_1_low',
    'drivewheel_2_moderate',
    'drivewheel_3_high',
    'enginelocation_2_moderate',
    'enginelocation_3_high',
    'aspiration_2_moderate',
    'aspiration_3_high',
    'fuelsystem_1_low',
    'fuelsystem_2_moderate',
    'fuelsystem_3_high',
    'fueltype_2_moderate',
    'fueltype_3_high',
    'fueltype_1_low',
    'carbody_1_low',
    'carbody_2_moderate',
    'carbody_3_high'
]

list_reduced = [
    'symboling',
    'displacement',
    'wheelbase',
    'carlength',
    'carwidth',
    'carheight',
    'curbweight',
    'enginehorsepower',
    'enginesize',
    'compressionratio',
    'horsepower',
    'peakrpm',
    'citympg',
    'highwaympg',
    'carname_1_low',
    'carname_2_moderate',
    'carname_3_high',
    'drivewheel_1_low',
    'drivewheel_2_moderate',
    'drivewheel_3_high',
    'enginelocation_2_moderate',
    'enginelocation_3_high',
    'aspiration_2_moderate',
    'aspiration_3_high',
    'fuelsystem_1_low',
    'fuelsystem_2_moderate',
    'fuelsystem_3_high',
    'fueltype_2_moderate',
    'fueltype_3_high',
    'fueltype_1_low',
    'carbody_1_low',
    'carbody_2_moderate',
    'carbody_3_high'
]

list_predictions = list_full + list_reduced

list_float = [
    'symboling',
    'displacement',
    'wheelbase',
    'carlength',
    'carwidth',
    'carheight',
    'curbweight',
    'enginehorsepower',
    'enginesize',
    'compressionratio',
    'horsepower',
    'peakrpm',
    'citympg',
    'highwaympg',
    'carname_1_low',
    'carname_2_moderate',
    'carname_3_high',
    'drivewheel_1_low',
    'drivewheel_2_moderate',
    'drivewheel_3_high',
    'enginelocation_2_moderate',
    'enginelocation_3_high',
    'aspiration_2_moderate',
    'aspiration_3_high',
    'fuelsystem_1_low',
    'fuelsystem_2_moderate',
    'fuelsystem_3_high',
    'fueltype_2_moderate',
    'fueltype_3_high',
    'fueltype_1_low',
    'carbody_1_low',
    'carbody_2_moderate',
    'carbody_3_high'
]

In [34]: df[list_float].describe().loc[['mean','std']]
```

```
Out[34]:
```

	symboling	displacement	wheelbase	carlength	carwidth	carheight	curbweight	cylindernumber	engineize	borebore	stroke	compressionratio	horsepower	peakrpm	citympg	highwaympg
mean	-2.95947e-17	4.332576e-18	0.001478	4.332576e-17	0.004605	4.332576e-17	0.000309	0.018917	0.002056	-5.190909e-17	-0.000507	-3.890200e-17	-1.733031e-17	-0.000021	6.498866e-18	2.166289e-18
std	1.001222e+00	1.001222e+00	0.006660	1.001222e+00	0.004605	1.001222e+00	1.000024	0.089748	0.002056	1.001222e+00	0.000507	1.001222e+00	1.001222e+00	1.001169	1.001222e+00	1.001222e+00

```
In [35]: from sklearn.preprocessing import StandardScaler
StandardScaler().fit(
    X = df[list_float]
)
y = df['price']

X = StandardScaler().fit_transform(
    X = df[list_float]
)

df[list_float].describe().loc[['mean','std']]
```

```
Out[35]:
```

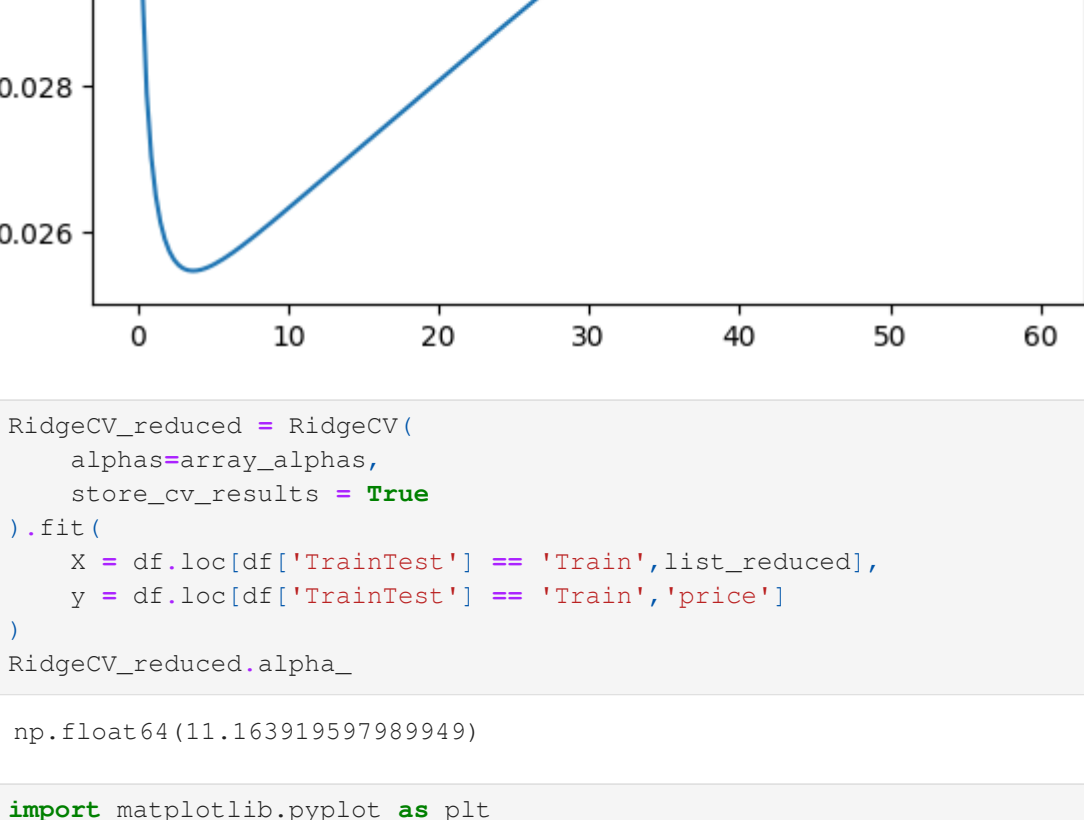
	symboling	displacement	wheelbase	carlength	carwidth	carheight	curbweight	cylindernumber	engineize	borebore	stroke	compressionratio	horsepower	peakrpm	citympg	highwaympg
mean	-2.95947e-17	4.332576e-18	2.166289e-18	-1.733031e-17	1.299773e-17	1.733031e-17	5.190909e-17	2.707861e-17	4.332576e-18	-6.196090e-17	-1.733031e-17	-1.299773e-17	6.498866e-18	-6.661506e-18	6.498866e-18	2.166289e-18
std	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00	1.001222e+00

```
In [36]: array_alpha = np.linspace(0.01,0.200)
from sklearn.linear_model import RidgeCV
RidgeCV_full = RidgeCV(
    alpha=array_alpha,
    store_cv_results = True
)
X = df.loc[df['train/test'] == 'Train',list_full],
y = df.loc[df['train/test'] == 'Train','price']
RidgeCV_full.alpha_
```

Out[36]: np.float64(13.6274874371855296)

```
In [37]: import matplotlib.pyplot as plt
plt.plot(RidgeCV_full.alpha, RidgeCV_full.cv_results_.mean(axis = 0))
plt.show()
```

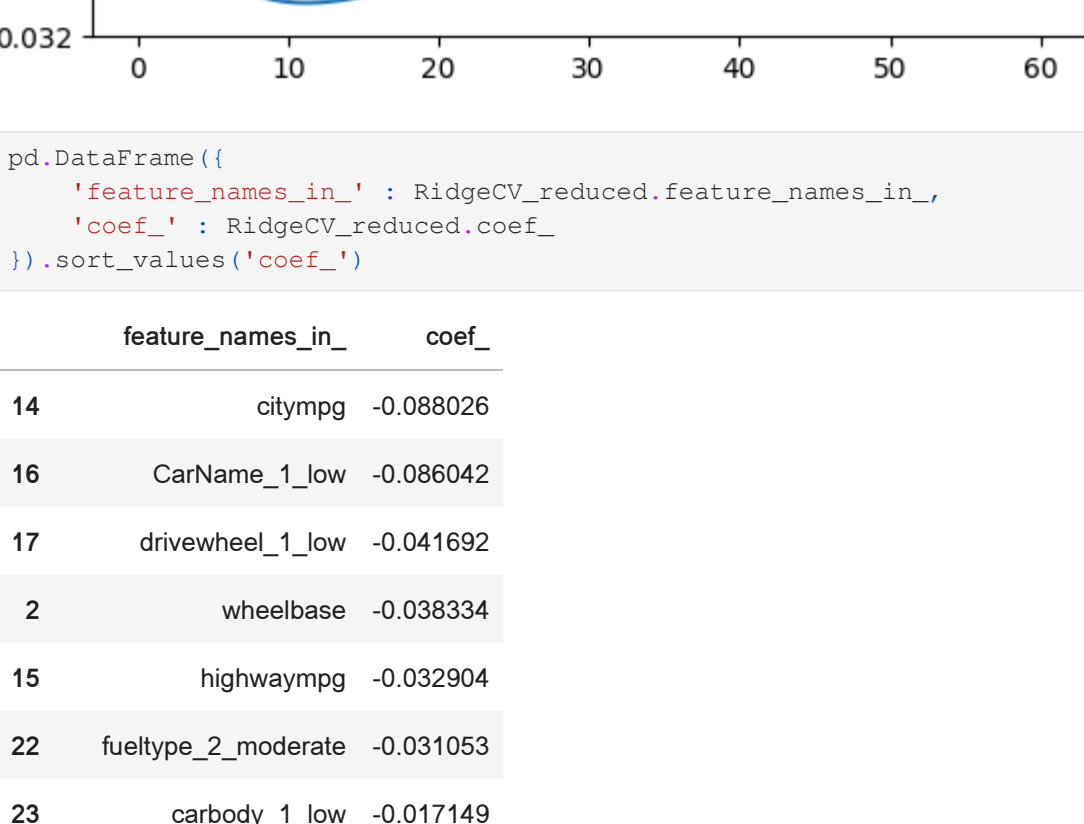
DeprecationWarning: Attribute 'cv_results_' is deprecated in version 1.5 and will be removed in 1.7. Use 'cv_results_' instead.



```
In [38]: RidgeCV_reduced = RidgeCV(
    alpha=array_alpha,
    store_cv_results = True
)
X = df.loc[df['train/test'] == 'Train',list_reduced],
y = df.loc[df['train/test'] == 'Train','price']
RidgeCV_reduced.alpha_
```

Out[38]: np.float64(11.163919597989943)

```
In [39]: import matplotlib.pyplot as plt
plt.plot(RidgeCV_reduced.alpha, RidgeCV_reduced.cv_results_.mean(axis = 0))
plt.show()
```



```
In [40]: pd.DataFrame(
    {'feature_names_in_': RidgeCV_reduced.feature_names_in_,
    'coef_': RidgeCV_reduced.coef_,
    })

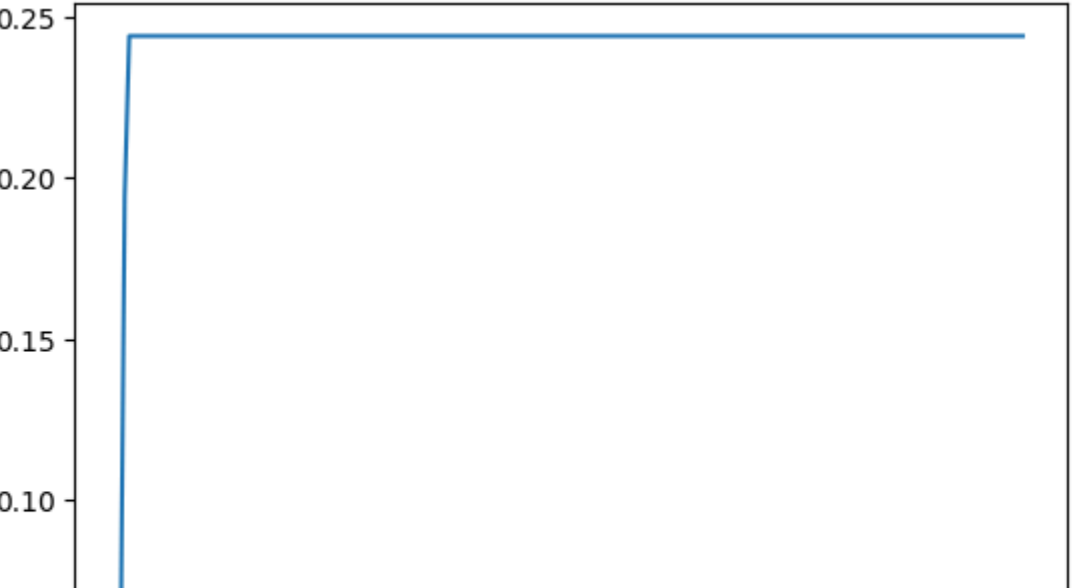
Out[40]:
```

	feature_names_in_	coef_
14	citympg	-0.080206
16	CarName_1_low	-0.08042
17	drivewheel_1_low	-0.041682
2	wheelbase	-0.038334
15	highwaympg	-0.032904
22	fueltype_2_moderate	-0.031053
23	carbody_1_low	-0.017449
9	borebore	-0.014830
21	enginehorsepower_2_moderate	-0.011476
0	symboling	-0.011089
10	stroke	-0.003085
20	fuelsystem_3_high	0.003338
1	displacement	0.003391
3	carlength	0.010208
19	aspiration_3_high	0.015682
5	carheight	0.020559
13	peakrpm	0.023336
18	enginelocation_3_high	0.044834
8	enginesize	0.053351
4	carwidth	0.020468
11	compressionratio	0.069765
12	horsepower	0.089130
6	curbweight	0.101649
7	cylindernumber	0.107086

```
In [41]: from sklearn.linear_model import LassoCV
LassoCV_full = LassoCV(
    alpha=array_alpha
)
X = df.loc[df['train/test'] == 'Train',list_full],
y = df.loc[df['train/test'] == 'Train','price']
LassoCV_full.alpha_
```

Out[41]: np.float64(0.01)

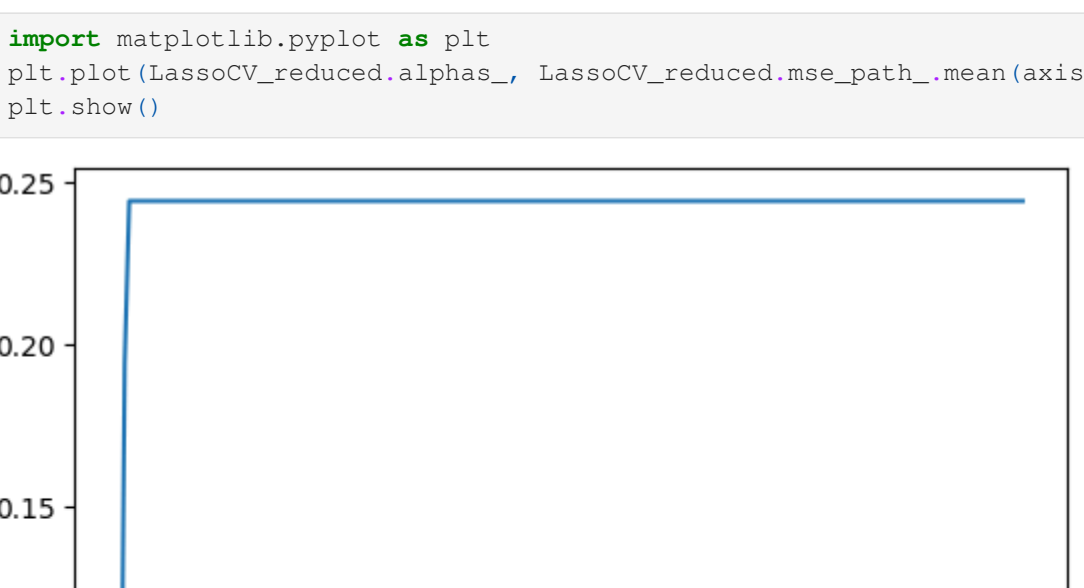
```
In [42]: import matplotlib.pyplot as plt
plt.plot(LassoCV_full.alpha, LassoCV_full.mse_path_.mean(axis = 1))
plt.show()
```



```
In [43]: LassoCV_reduced = LassoCV(
    alpha=array_alpha
)
X = df.loc[df['train/test'] == 'Train',list_reduced],
y = df.loc[df['train/test'] == 'Train','price']
LassoCV_reduced.alpha_
```

Out[43]: np.float64(0.01)

```
In [44]: import matplotlib.pyplot as plt
plt.plot(LassoCV_reduced.alpha, LassoCV_reduced.mse_path_.mean(axis = 1))
plt.show()
```



```
In [45]: pd.DataFrame(
    {'feature_names_in_': LassoCV_reduced.feature_names_in_,
    'coef_': LassoCV_reduced.coef_,
    })

Out[45]:
```

	feature_names_in_	coef_
14	citympg	-0.112572
16	CarName_1_low	-0.081983
2	wheelbase	-0.019779
1	displacement	0.000000
5	carheight	0.000000
10	stroke	-0.000000
8	enginesize	0.000000
0	symboling	-0.000000
15	highwaympg	-0.000000
9	borebore	-0.000000
21	enginehorsepower_2_moderate	-0.000000
20	fuelsystem_3_high	0.000000
19	aspiration_3_high	0.000000
18	enginelocation_3_high	0.000000
17	drivewheel_1_low	-0.000000
3	carlength	0.000000
22	fueltype_2_moderate	-0.000000
23	carbody_1_low	-0.000000
13	peakrpm	0.006255
4	carwidth	0.019680
11	compressionratio	0.084713
12	horsepower	0.103955
6	cylindernumber	0.124934
7	curbweight	0.178139

In [46]: array_l1_ratio = np.linspace(0.01, 1, 100)

```
In [47]: from sklearn.linear_model import ElasticNetCV
import numpy as np
ElasticNetCV_full = ElasticNetCV(
    alpha=array_alpha,
    l1_ratio=array_l1_ratio,
    max_iter = 1000
)
X = df.loc[df['train/test'] == 'Train',list_full],
y = df.loc[df['train/test'] == 'Train','price']
```

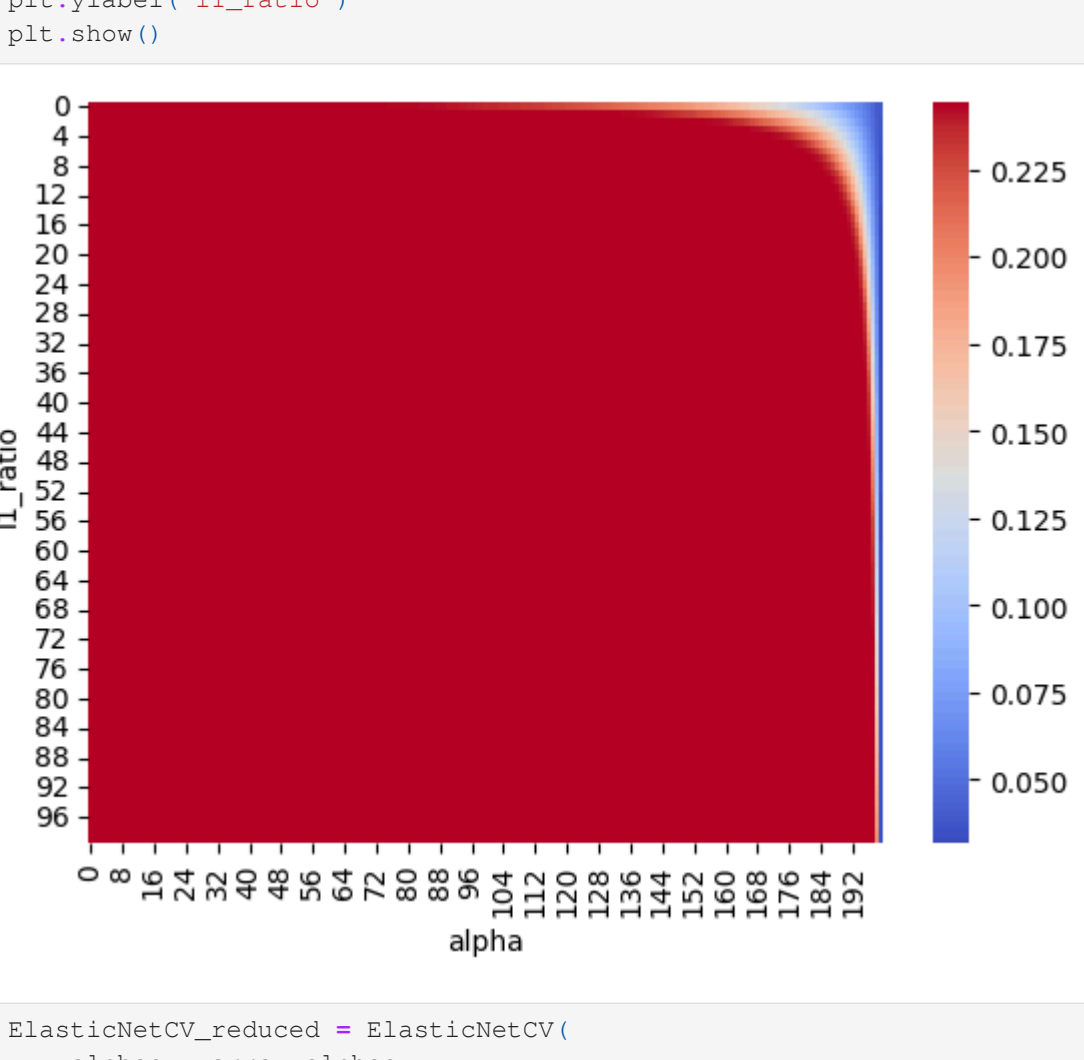
Out[47]: ElasticNetCV_full.alpha_

Out[48]: np.float64(0.01)

```
In [49]: ElasticNetCV_full.l1_ratio_
```

Out[49]: np.float64(0.25)

```
In [50]: import seaborn as sns
sns.heatmap(ElasticNetCV_full.mse_path_.mean(axis = 2),cmap="coolwarm")
plt.xlabel('alpha')
plt.ylabel('l1_ratio')
plt.show()
```



```
In [51]: ElasticNetCV_reduced = ElasticNetCV(
    alpha=array_alpha,
    l1_ratio=array_l1_ratio,
    max_iter = 1000
)
X = df.loc[df['train/test'] == 'Train',list_reduced],
y = df.loc[df['train/test'] == 'Train','price']
```

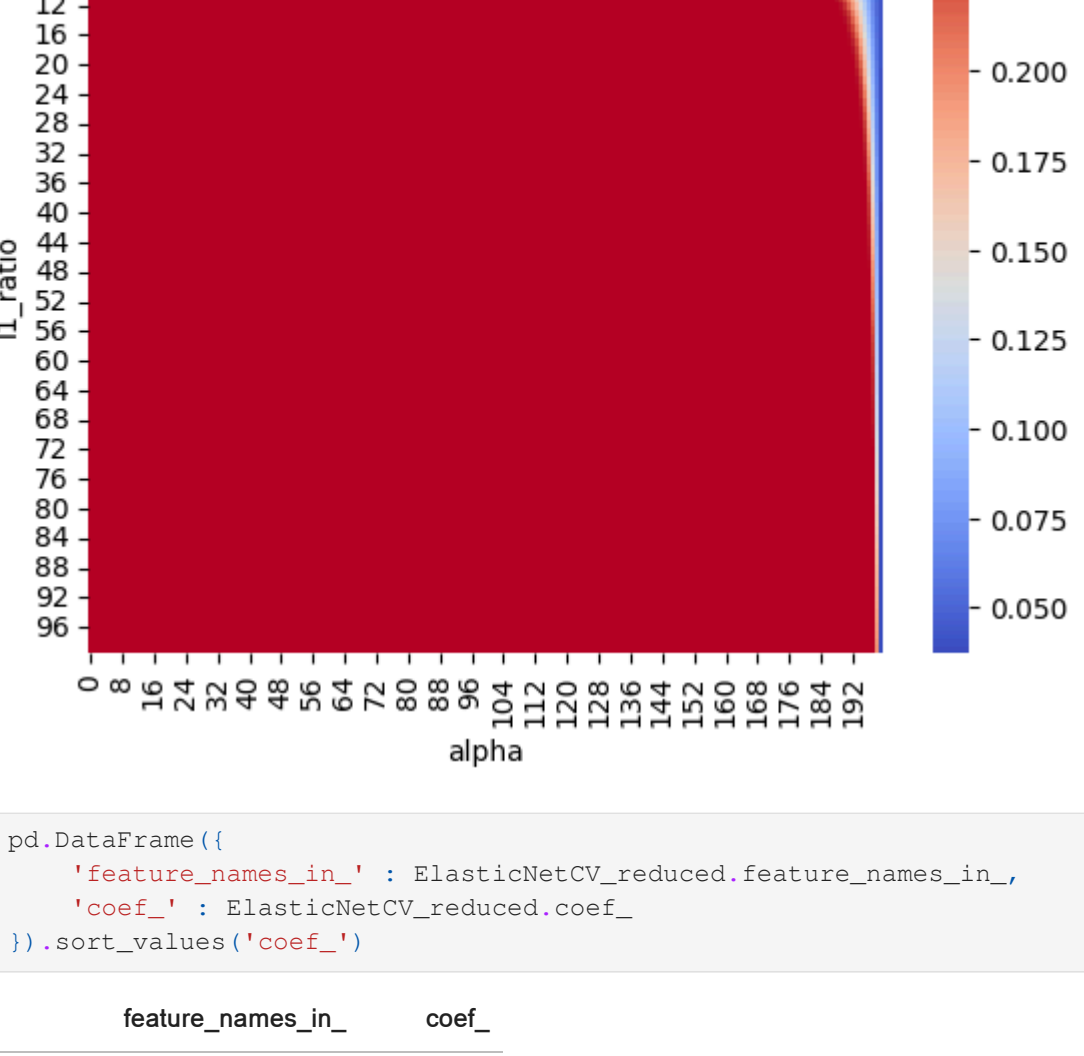
Out[51]: ElasticNetCV_reduced.alpha_

Out[52]: np.float64(0.01)

Out[53]: ElasticNetCV_reduced.l1_ratio_

Out[53]: np.float64(0.89)

```
In [54]: import seaborn as sns
sns.heatmap(ElasticNetCV_reduced.mse_path_.mean(axis = 2),cmap="coolwarm")
plt.xlabel('alpha')
plt.ylabel('l1_ratio')
plt.show()
```



```
In [55]: pd.DataFrame(
    {'feature_names_in_': ElasticNetCV_reduced.feature_names_in_,
    'coef_': ElasticNetCV_reduced.coef_,
    })

Out[55]:
```

	feature_names_in_	coef_
14	citympg	-0.112567
16	CarName_1_low	-0.089057
2	wheelbase	-0.017487
1	displacement	0.000000
9	borebore	-0.000000
10	stroke	-0.000000
8	enginesize	0.000000
0	symboling	-0.000000
15	highwaympg	-0.000000
22	fueltype_2_moderate	-0.000000
21	enginehorsepower_2_moderate	0.000000
20	fuelsystem_3_high	-0.000000
19	aspiration_3_high	0.000000
18	enginelocation_3_high	0.000000
17	drivewheel_1_low	-0.000000
3	carlength	0.000000
23	carbody_1_low	-0.000000
5	carheight	0.002777
13	peakrpm	0.008146
4	carwidth	0.023912
11	compressionratio	0.086222
12	horsepower	0.106068
7	cylindernumber	0.130758
6	curbweight	0.175916

```
In [56]: list_cv = [RidgeCV_full,RidgeCV_reduced,LassoCV_full,LassoCV_reduced,ElasticNetCV_full,ElasticNetCV_reduced]
df_predict = pd.DataFrame(
    CV_predict.X = df[CV_feature_names_in_] = df['price']**2 for CV in list_cv
)
df_predict.columns = ['LassoCV_full','LassoCV_reduced','ElasticNetCV_full','ElasticNetCV_reduced']
df_predict.groupby(df['train/test']).mean()
```

```
Out[56]:
```

	RidgeCV_full	RidgeCV_reduced	LassoCV_full	LassoCV_reduced	ElasticNetCV_full	ElasticNetCV_reduced
Score	NaN	NaN	NaN	NaN	NaN	NaN
Test	0.053390	0.056342	0.051899	0.053177	0.053512	0.053135
Train	0.012117	0.019963	0.017819	0.020717	0.012229	0.020262
Validation	0.018458	0.026295	0.018843	0.027990	0.017515	0.026076

