In [1]:

```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn import preprocessing,svm
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Lasso
from sklearn.linear_model import Ridge
```

In [2]:

```
a=pd.read_csv(r"C:\Users\Gowthami\Downloads\bottle.csv.zip")
a
```

```
C:\Users\Gowthami\AppData\Local\Temp\ipykernel_22952\1376288322.py:1: Dtyp
eWarning: Columns (47,73) have mixed types. Specify dtype option on import
or set low_memory=False.
  a=pd.read_csv(r"C:\Users\Gowthami\Downloads\bottle.csv.zip")
```

Out[2]:

| | Cst_Cnt | Btl_Cnt | Sta_ID | Depth_ID | Depthm | T_degC | Salnty | O2ml_L | STheta |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0000A-3 | 0 | 10.500 | 33.4400 | NaN | 25.64900 |
| 1 | 1 | 2 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0008A-3 | 8 | 10.460 | 33.4400 | NaN | 25.65600 |

In [3]:

```
a.head()
```

| 2 | 1 | 3 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0010A-7 | 10 | 10.460 | 33.4370 | NaN | 25.65400 |

Out[3]:

| 3 | 1 | 4 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0019A-3 | 19 | 10.450 | 33.4200 | NaN | 25.64300 |

| | Cst_Cnt | Btl_Cnt | Sta_ID | Depth_ID | Depthm | T_degC | Salnty | O2ml_L | STheta | O2Sat | . |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0000A-3 | 0 | 10.50 | 33.440 | NaN | 25.649 | NaN | . |
| 4 | 1 | 5 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0020A-7 | 20 | 10.450 | 33.4210 | NaN | 25.64300 | | |
| 1 | 1 | 2 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0008A-3 | 8 | 10.46 | 33.440 | NaN | 25.656 | NaN | . |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | . |
| 864858 | 34404 | 864859 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0000A-7 | 0 | 18.744 | 33.4083 | 5.805 | 23.87055 | | |
| 2 | 1 | 3 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0010A-7 | 10 | 10.46 | 33.437 | NaN | 25.654 | NaN | . |
| 864859 | 34404 | 864860 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0002A-3 | 2 | 18.744 | 33.4083 | 5.805 | 23.87072 | | |
| 3 | 1 | 4 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0019A-3 | 19 | 10.45 | 33.420 | NaN | 25.643 | NaN | . |
| 864860 | 34404 | 864861 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0005A-3 | 5 | 18.692 | 33.4150 | 5.796 | 23.88911 | | |
| 4 | 1 | 5 | 054.0 056.0 | 19-4903CR-HY-060-0930-05400560-0020A-7 | 20 | 10.45 | 33.421 | NaN | 25.643 | NaN | . |
| 864861 | 34404 | 864862 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264- | 10 | 18.161 | 33.4062 | 5.816 | 24.01426 | | |

5 rows × 74 columns

In [4]:

```
a.info()
```

| | Cst_Cnt | Btl_Cnt | Sta_ID | Depth_ID | Depthm | T_degC | Salnty | O2ml_L | STheta |
|---|---|---|---|---|---|---|---|---|---|
| **864862** | 34404 | 864863 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0015A-3 | 15 | 17.533 | 33.3880 | 5.774 | 24.15297 |

864863 rows × 74 columns

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 864863 entries, 0 to 864862
Data columns (total 74 columns):
 #   Column       Non-Null Count    Dtype
---  ------       --------------    -----
 0   Cst_Cnt      864863 non-null   int64
 1   Btl_Cnt      864863 non-null   int64
 2   Sta_ID       864863 non-null   object
 3   Depth_ID     864863 non-null   object
 4   Depthm       864863 non-null   int64
 5   T_degC       853900 non-null   float64
 6   Salnty       817509 non-null   float64
 7   O2ml_L       696201 non-null   float64
 8   STheta       812174 non-null   float64
 9   O2Sat        661274 non-null   float64
 10  Oxy_µmol/Kg  661268 non-null   float64
 11  BtlNum       118667 non-null   float64
 12  RecInd       864863 non-null   int64
 13  T_prec       853900 non-null   float64
 14  T_qual       23127 non-null    float64
 15  S_prec       817509 non-null   float64
 16  S_qual       74914 non-null    float64
 17  P_qual       673755 non-null   float64
 18  O_qual       184676 non-null   float64
 19  SThtaq       65823 non-null    float64
 20  O2Satq       217797 non-null   float64
 21  ChlorA       225272 non-null   float64
 22  Chlqua       639166 non-null   float64
 23  Phaeop       225271 non-null   float64
 24  Phaqua       639170 non-null   float64
 25  PO4uM        413317 non-null   float64
 26  PO4q         451786 non-null   float64
 27  SiO3uM       354091 non-null   float64
 28  SiO3qu       510866 non-null   float64
 29  NO2uM        337576 non-null   float64
 30  NO2q         529474 non-null   float64
 31  NO3uM        337403 non-null   float64
 32  NO3q         529933 non-null   float64
 33  NH3uM        64962 non-null    float64
 34  NH3q         808299 non-null   float64
 35  C14As1       14432 non-null    float64
 36  C14A1p       12760 non-null    float64
 37  C14A1q       848605 non-null   float64
 38  C14As2       14414 non-null    float64
 39  C14A2p       12742 non-null    float64
 40  C14A2q       848623 non-null   float64
 41  DarkAs       22649 non-null    float64
 42  DarkAp       20457 non-null    float64
 43  DarkAq       840440 non-null   float64
 44  MeanAs       22650 non-null    float64
 45  MeanAp       20457 non-null    float64
 46  MeanAq       840439 non-null   float64
 47  IncTim       14437 non-null    object
 48  LightP       18651 non-null    float64
 49  R_Depth      864863 non-null   float64
 50  R_TEMP       853900 non-null   float64
 51  R_POTEMP     818816 non-null   float64
 52  R_SALINITY   817509 non-null   float64
 53  R_SIGMA      812007 non-null   float64
 54  R_SVA        812092 non-null   float64
 55  R_DYNHT      818206 non-null   float64
```

```
 56   R_O2                696201 non-null   float64
 57   R_O2Sat             666448 non-null   float64
 58   R_SIO3              354099 non-null   float64
 59   R_PO4               413325 non-null   float64
 60   R_NO3               337411 non-null   float64
 61   R_NO2               337584 non-null   float64
 62   R_NH4               64982 non-null    float64
 63   R_CHLA              225276 non-null   float64
 64   R_PHAEO             225275 non-null   float64
 65   R_PRES              864863 non-null   int64
 66   R_SAMP              122006 non-null   float64
 67   DIC1                1999 non-null     float64
 68   DIC2                224 non-null      float64
 69   TA1                 2084 non-null     float64
 70   TA2                 234 non-null      float64
 71   pH2                 10 non-null       float64
 72   pH1                 84 non-null       float64
 73   DIC Quality Comment  55 non-null      object
dtypes: float64(65), int64(5), object(4)
memory usage: 488.3+ MB
```

In [5]:

```
a.tail()
```

Out[5]:

| | Cst_Cnt | Btl_Cnt | Sta_ID | Depth_ID | Depthm | T_degC | Salnty | O2ml_L | STheta |
|---|---|---|---|---|---|---|---|---|---|
| **864858** | 34404 | 864859 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0000A-7 | 0 | 18.744 | 33.4083 | 5.805 | 23.87055 |
| **864859** | 34404 | 864860 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0002A-3 | 2 | 18.744 | 33.4083 | 5.805 | 23.87072 |
| **864860** | 34404 | 864861 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0005A-3 | 5 | 18.692 | 33.4150 | 5.796 | 23.88911 |
| **864861** | 34404 | 864862 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0010A-3 | 10 | 18.161 | 33.4062 | 5.816 | 24.01426 |
| **864862** | 34404 | 864863 | 093.4 026.4 | 20-1611SR-MX-310-2239-09340264-0015A-3 | 15 | 17.533 | 33.3880 | 5.774 | 24.15297 |

5 rows × 74 columns

In [6]:

```python
a.describe()
```

Out[6]:

|       | Cst_Cnt       | Btl_Cnt       | Depthm        | T_degC        | Salnty        | O        |
|-------|---------------|---------------|---------------|---------------|---------------|----------|
| count | 864863.000000 | 864863.000000 | 864863.000000 | 853900.000000 | 817509.000000 | 696201.0 |
| mean  | 17138.790958  | 432432.000000 | 226.831951    | 10.799677     | 33.840350     | 3.3      |
| std   | 10240.949817  | 249664.587269 | 316.050259    | 4.243825      | 0.461843      | 2.0      |
| min   | 1.000000      | 1.000000      | 0.000000      | 1.440000      | 28.431000     | -0.0     |
| 25%   | 8269.000000   | 216216.500000 | 46.000000     | 7.680000      | 33.488000     | 1.3      |
| 50%   | 16848.000000  | 432432.000000 | 125.000000    | 10.060000     | 33.863000     | 3.4      |
| 75%   | 26557.000000  | 648647.500000 | 300.000000    | 13.880000     | 34.196900     | 5.5      |
| max   | 34404.000000  | 864863.000000 | 5351.000000   | 31.140000     | 37.034000     | 11.1     |

8 rows × 70 columns

In [7]:

```python
a.isna().any()
```

Out[7]:

```
Cst_Cnt                False
Btl_Cnt                False
Sta_ID                 False
Depth_ID               False
Depthm                 False
                       ...
TA1                     True
TA2                     True
pH2                     True
pH1                     True
DIC Quality Comment     True
Length: 74, dtype: bool
```

In [8]:

```python
a.isnull().sum()
```

Out[8]:

```
Cst_Cnt                  0
Btl_Cnt                  0
Sta_ID                   0
Depth_ID                 0
Depthm                   0
                       ...
TA1                 862779
TA2                 864629
pH2                 864853
pH1                 864779
DIC Quality Comment 864808
Length: 74, dtype: int64
```

In [9]:

```python
a.loc[:,['Salnty','T_degC']]
```

Out[9]:

|        | Salnty  | T_degC |
|--------|---------|--------|
| 0      | 33.4400 | 10.500 |
| 1      | 33.4400 | 10.460 |
| 2      | 33.4370 | 10.460 |
| 3      | 33.4200 | 10.450 |
| 4      | 33.4210 | 10.450 |
| ...    | ...     | ...    |
| 864858 | 33.4083 | 18.744 |
| 864859 | 33.4083 | 18.744 |
| 864860 | 33.4150 | 18.692 |
| 864861 | 33.4062 | 18.161 |
| 864862 | 33.3880 | 17.533 |

864863 rows × 2 columns

In [10]:

```python
a=a[['Salnty','T_degC']]
a.columns=['Sal','Temp']
```

In [11]:

```python
a.head(20)
```

Out[11]:

|    | Sal    | Temp  |
|----|--------|-------|
| 0  | 33.440 | 10.50 |
| 1  | 33.440 | 10.46 |
| 2  | 33.437 | 10.46 |
| 3  | 33.420 | 10.45 |
| 4  | 33.421 | 10.45 |
| 5  | 33.431 | 10.45 |
| 6  | 33.440 | 10.45 |
| 7  | 33.424 | 10.24 |
| 8  | 33.420 | 10.06 |
| 9  | 33.494 | 9.86  |
| 10 | 33.510 | 9.83  |
| 11 | 33.580 | 9.67  |
| 12 | 33.640 | 9.50  |
| 13 | 33.689 | 9.32  |
| 14 | 33.847 | 8.76  |
| 15 | 33.860 | 8.71  |
| 16 | 33.876 | 8.53  |
| 17 | NaN    | 8.45  |
| 18 | 33.926 | 8.26  |
| 19 | 33.980 | 7.96  |

In [12]:

```python
sns.lmplot(x='Sal',y='Temp',data=a,order=2,ci=None)
```

Out[12]:

```
<seaborn.axisgrid.FacetGrid at 0x17ff59ebc10>
```

In [13]:

```python
a.fillna(method='ffill')
```

Out[13]:

|        | Sal     | Temp   |
|--------|---------|--------|
| 0      | 33.4400 | 10.500 |
| 1      | 33.4400 | 10.460 |
| 2      | 33.4370 | 10.460 |
| 3      | 33.4200 | 10.450 |
| 4      | 33.4210 | 10.450 |
| ...    | ...     | ...    |
| 864858 | 33.4083 | 18.744 |
| 864859 | 33.4083 | 18.744 |
| 864860 | 33.4150 | 18.692 |
| 864861 | 33.4062 | 18.161 |
| 864862 | 33.3880 | 17.533 |

864863 rows × 2 columns

In [14]:

```python
a.fillna(value=0,inplace=True)
```

```
C:\Users\Gowthami\AppData\Local\Temp\ipykernel_22952\3015722605.py:1: Sett
ingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-doc
s/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://
pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-
view-versus-a-copy)
  a.fillna(value=0,inplace=True)
```

In [15]:

```python
x=np.array(a['Sal']).reshape(-1,1)
y=np.array(a['Temp']).reshape(-1,1)
```

In [16]:

```python
a.isna().any()
```

Out[16]:

```
Sal      False
Temp     False
dtype: bool
```

In [17]:

```python
a.isnull().sum()
```

Out[17]:

```
Sal     0
Temp    0
dtype: int64
```

In [18]:

```python
a.dropna(inplace=True)
```

```
C:\Users\Gowthami\AppData\Local\Temp\ipykernel_22952\2317726482.py:1: Sett
ingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-doc
s/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://
pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-
view-versus-a-copy)
  a.dropna(inplace=True)
```

In [19]:

```python
 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
```

In [20]:

```python
regr=LinearRegression()
regr.fit(x_train,y_train)
print(regr.score(x_test,y_test))
```

```
0.014809436284115463
```

In [21]:

```python
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
```

In [22]:

```
a500=a[:][:500]
sns.lmplot(x='Sal',y='Temp',data=a500,order=1,ci=None)
```

Out[22]:

```
<seaborn.axisgrid.FacetGrid at 0x17ff475fd90>
```

In [23]:

```python
a500.fillna(method='ffill',inplace=True)
x=np.array(a500['Sal']).reshape(-1,1)
y=np.array(a500['Temp']).reshape(-1,1)
a500.dropna(inplace=True)
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
regr=LinearRegression()
regr.fit(x_train,y_train)
print("Regression:",regr.score(x_test,y_test))
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
```

Regression: 0.022412040359021446



In [24]:

```python
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
model=LinearRegression()
model.fit(x_train,y_train)
y_pred=model.predict(x_test)
r2=r2_score(y_test,y_pred)
print("R2 score:",r2)
```

R2 score: 0.022412040359021446

#conclusion: Linear regression is best fit for the model

# Ridge and Lasso regression

In [25]:

```python
plt.figure(figsize = (10, 10))
sns.heatmap(a.corr(), annot = True)
```

Out[25]:

```
<Axes: >
```

In [26]:

```python
features = a.columns[0:2]
target = a.columns[-1]
#X and y values
X = a[features].values
y = a[target].values
#splot
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=17
print("The dimension of X_train is {}".format(X_train.shape))
print("The dimension of X_test is {}".format(X_test.shape))
#Scale features
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```
The dimension of X_train is (605404, 2)
The dimension of X_test is (259459, 2)
```

In [27]:

```python
#Model
lr = LinearRegression()
#Fit model
lr.fit(X_train, y_train)
#predict
#prediction = lr.predict(X_test)
#actual
actual = y_test
train_score_lr = lr.score(X_train, y_train)
test_score_lr = lr.score(X_test, y_test)
print("\nLinear Regression Model:\n")
print("The train score for lr model is {}".format(train_score_lr))
print("The test score for lr model is {}".format(test_score_lr))
```

```
Linear Regression Model:

The train score for lr model is 1.0
The test score for lr model is 1.0
```

In [28]:

```python
#Using the linear CV model
from sklearn.linear_model import RidgeCV
#Ridge Cross validation
ridge_cv = RidgeCV(alphas = [0.0001, 0.001,0.01, 0.1, 1, 10]).fit(X_train, y_train)
#score
print("The train score for ridge model is {}".format(ridge_cv.score(X_train, y_train)))
print("The train score for ridge model is {}".format(ridge_cv.score(X_test, y_test)))
```

```
The train score for ridge model is 0.9999999981135502
The train score for ridge model is 0.99999999811206
```

In [29]:

```python
ridgeReg = Ridge(alpha=10)
ridgeReg.fit(X_train,y_train)
#train and test scorefor ridge regression
train_score_ridge = ridgeReg.score(X_train, y_train)
test_score_ridge = ridgeReg.score(X_test, y_test)
print("\nRidge Model:\n")
print("The train score for ridge model is {}".format(train_score_ridge))
print("The test score for ridge model is {}".format(test_score_ridge))
```
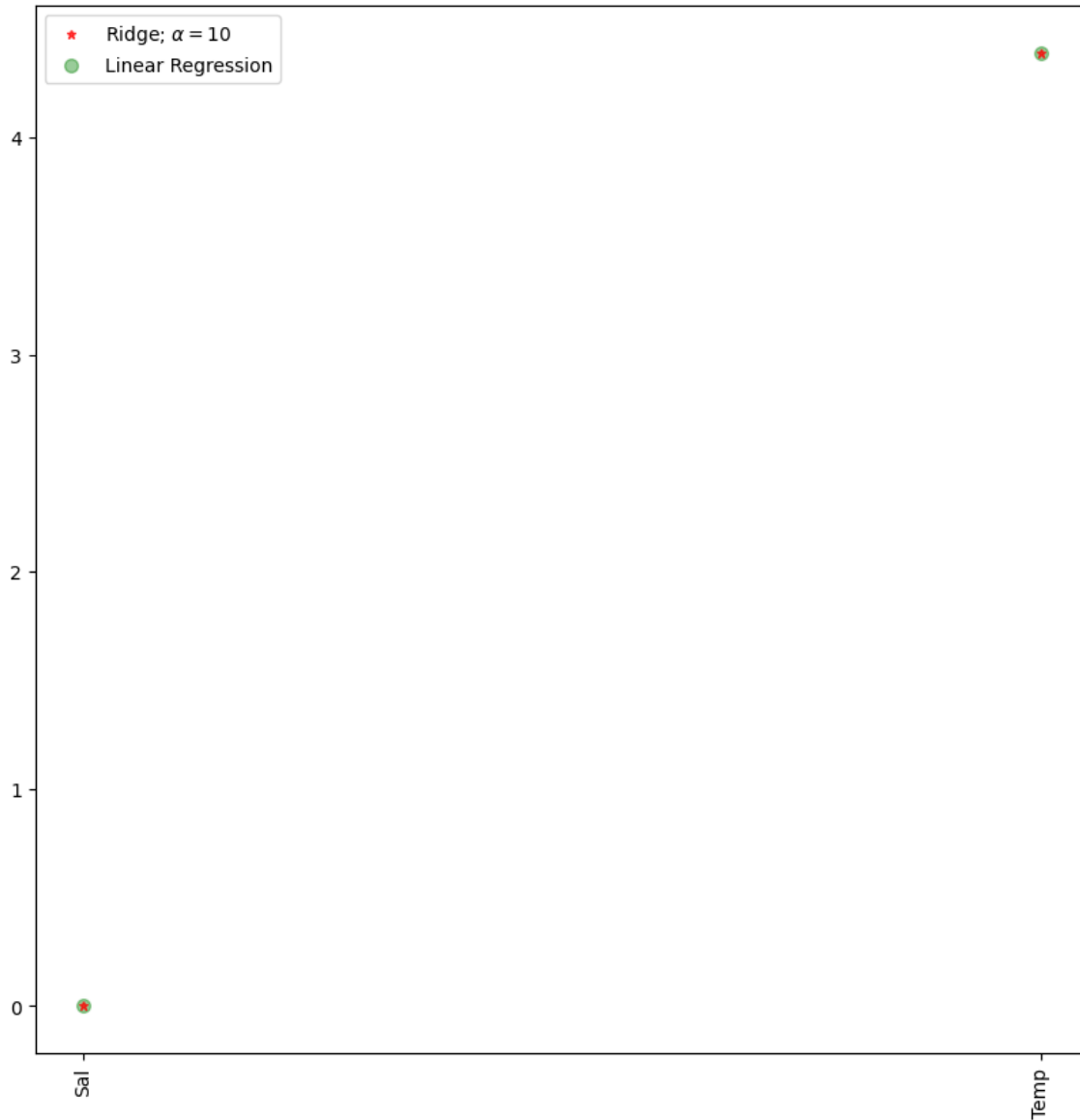
```
Ridge Model:

The train score for ridge model is 0.999999999723243
The test score for ridge model is 0.9999999997231402
```

In [30]:

```python
plt.figure(figsize = (10, 10))
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,colo

plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='gre
plt.xticks(rotation = 90)
plt.legend()
plt.show()
```



In [31]:

```python
#Using the linear CV model
from sklearn.linear_model import LassoCV
#Lasso Cross validation
lasso_cv = LassoCV(alphas = [0.0001, 0.001,0.01, 0.1, 1, 10], random_state=0).fit(X_trai
#score
print(lasso_cv.score(X_train, y_train))
print(lasso_cv.score(X_test, y_test))
```

```
0.9999999994806811
0.9999999994806712
```

In [32]:

```python
#Lasso regression model
print("\nLasso Model: \n")
lasso = Lasso(alpha = 10)
lasso.fit(X_train,y_train)
train_score_ls =lasso.score(X_train,y_train)
test_score_ls =lasso.score(X_test,y_test)
print("The train score for ls model is {}".format(train_score_ls))
print("The test score for ls model is {}".format(test_score_ls))
```

```
Lasso Model:

The train score for ls model is 0.0
The test score for ls model is -1.9031696447013857e-05
```

In [33]:

```python
pd.Series(lasso.coef_, features).sort_values(ascending = True).plot(kind = "bar")
```

Out[33]:

```
<Axes: >
```

In [34]:

```python
#plot size
plt.figure(figsize = (10, 10))
#add plot for ridge regression
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,colo
#add plot for lasso regression
plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue',
#add plot for linear model
plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='gre
```

Out[34]:

[<matplotlib.lines.Line2D at 0x17fd2e42b10>]



# Elastic net

In [35]:

```python
from sklearn.linear_model import ElasticNet
regr=ElasticNet()
regr.fit(X,y)
print(regr.coef_)
print(regr.intercept_)
```

```
[0.         0.94934511]
0.5401219631063316
```

In [36]:

```python
y_pred_elastic=regr.predict(X_train)
```

In [37]:

```python
mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
print("Mean Squared Error on test set",mean_squared_error)
```

```
Mean Squared Error on test set 114.40984808660129
```

# Vehicle selection

In [38]:

```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn import preprocessing,svm
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Lasso
from sklearn.linear_model import Ridge


df=pd.read_csv(r"C:\Users\Gowthami\Downloads\fiat500_VehicleSelection_Dataset.csv")
df
```

Out[38]:

| | ID | model | engine_power | age_in_days | km | previous_owners | lat | |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | lounge | 51 | 882 | 25000 | 1 | 44.907242 | 8.611 |
| 1 | 2 | pop | 51 | 1186 | 32500 | 1 | 45.666359 | 12.241 |
| 2 | 3 | sport | 74 | 4658 | 142228 | 1 | 45.503300 | 11.417 |
| 3 | 4 | lounge | 51 | 2739 | 160000 | 1 | 40.633171 | 17.634 |
| 4 | 5 | pop | 73 | 3074 | 106880 | 1 | 41.903221 | 12.495 |
| ... | ... | ... | ... | ... | ... | ... | ... | |
| 1533 | 1534 | sport | 51 | 3712 | 115280 | 1 | 45.069679 | 7.704 |
| 1534 | 1535 | lounge | 74 | 3835 | 112000 | 1 | 45.845692 | 8.666 |
| 1535 | 1536 | pop | 51 | 2223 | 60457 | 1 | 45.481541 | 9.413 |
| 1536 | 1537 | lounge | 51 | 2557 | 80750 | 1 | 45.000702 | 7.682 |
| 1537 | 1538 | pop | 51 | 1766 | 54276 | 1 | 40.323410 | 17.568 |

1538 rows × 9 columns

In [39]:

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1538 entries, 0 to 1537
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   ID               1538 non-null   int64
 1   model            1538 non-null   object
 2   engine_power     1538 non-null   int64
 3   age_in_days      1538 non-null   int64
 4   km               1538 non-null   int64
 5   previous_owners  1538 non-null   int64
 6   lat              1538 non-null   float64
 7   lon              1538 non-null   float64
 8   price            1538 non-null   int64
dtypes: float64(2), int64(6), object(1)
memory usage: 108.3+ KB
```

In [40]:

```python
df.head()
```

Out[40]:

| | ID | model | engine_power | age_in_days | km | previous_owners | lat | lon | |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | lounge | 51 | 882 | 25000 | 1 | 44.907242 | 8.611560 | |
| **1** | 2 | pop | 51 | 1186 | 32500 | 1 | 45.666359 | 12.241890 | |
| **2** | 3 | sport | 74 | 4658 | 142228 | 1 | 45.503300 | 11.417840 | |
| **3** | 4 | lounge | 51 | 2739 | 160000 | 1 | 40.633171 | 17.634609 | |
| **4** | 5 | pop | 73 | 3074 | 106880 | 1 | 41.903221 | 12.495650 | |

In [41]:

```python
df.tail()
```

Out[41]:

| | ID | model | engine_power | age_in_days | km | previous_owners | lat | l |
|---|---|---|---|---|---|---|---|---|
| **1533** | 1534 | sport | 51 | 3712 | 115280 | 1 | 45.069679 | 7.704! |
| **1534** | 1535 | lounge | 74 | 3835 | 112000 | 1 | 45.845692 | 8.666! |
| **1535** | 1536 | pop | 51 | 2223 | 60457 | 1 | 45.481541 | 9.413! |
| **1536** | 1537 | lounge | 51 | 2557 | 80750 | 1 | 45.000702 | 7.682; |
| **1537** | 1538 | pop | 51 | 1766 | 54276 | 1 | 40.323410 | 17.568; |

In [42]:

```
df.info
```

Out[42]:

```
<bound method DataFrame.info of          ID    model   engine_power   age_in_d
ays    km  previous_owners
0        1   lounge            51      882    25000              1  \
1        2      pop            51     1186    32500              1
2        3    sport            74     4658   142228              1
3        4   lounge            51     2739   160000              1
4        5      pop            73     3074   106880              1
...    ...      ...           ...      ...      ...            ...
1533  1534    sport            51     3712   115280              1
1534  1535   lounge            74     3835   112000              1
1535  1536      pop            51     2223    60457              1
1536  1537   lounge            51     2557    80750              1
1537  1538      pop            51     1766    54276              1

            lat        lon  price
0     44.907242   8.611560   8900
1     45.666359  12.241890   8800
2     45.503300  11.417840   4200
3     40.633171  17.634609   6000
4     41.903221  12.495650   5700
...         ...        ...    ...
1533  45.069679   7.704920   5200
1534  45.845692   8.666870   4600
1535  45.481541   9.413480   7500
1536  45.000702   7.682270   5990
1537  40.323410  17.568270   7900

[1538 rows x 9 columns]>
```

In [43]:

```
df.describe()
```

Out[43]:

|       | ID          | engine_power | age_in_days | km            | previous_owners | l...        |
|-------|-------------|--------------|-------------|---------------|-----------------|-------------|
| count | 1538.000000 | 1538.000000  | 1538.000000 | 1538.000000   | 1538.000000     | 1538.00000  |
| mean  | 769.500000  | 51.904421    | 1650.980494 | 53396.011704  | 1.123537        | 43.54136    |
| std   | 444.126671  | 3.988023     | 1289.522278 | 40046.830723  | 0.416423        | 2.13351     |
| min   | 1.000000    | 51.000000    | 366.000000  | 1232.000000   | 1.000000        | 36.85583    |
| 25%   | 385.250000  | 51.000000    | 670.000000  | 20006.250000  | 1.000000        | 41.80299    |
| 50%   | 769.500000  | 51.000000    | 1035.000000 | 39031.000000  | 1.000000        | 44.39409    |
| 75%   | 1153.750000 | 51.000000    | 2616.000000 | 79667.750000  | 1.000000        | 45.46796    |
| max   | 1538.000000 | 77.000000    | 4658.000000 | 235000.000000 | 4.000000        | 46.79561    |

In [44]:

```python
df.isna().any()
```

Out[44]:

```
ID               False
model            False
engine_power     False
age_in_days      False
km               False
previous_owners  False
lat              False
lon              False
price            False
dtype: bool
```

In [45]:

```python
df.isnull().sum()
```

Out[45]:

```
ID               0
model            0
engine_power     0
age_in_days      0
km               0
previous_owners  0
lat              0
lon              0
price            0
dtype: int64
```

In [46]:

```
df.isnull()
```

Out[46]:

| | ID | model | engine_power | age_in_days | km | previous_owners | lat | lon | price |
|---|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False | False | False |
| 1 | False | False | False | False | False | False | False | False | False |
| 2 | False | False | False | False | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False | False |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 1533 | False | False | False | False | False | False | False | False | False |
| 1534 | False | False | False | False | False | False | False | False | False |
| 1535 | False | False | False | False | False | False | False | False | False |
| 1536 | False | False | False | False | False | False | False | False | False |
| 1537 | False | False | False | False | False | False | False | False | False |

1538 rows × 9 columns

In [47]:

```
df.loc[:11,["ID","price"]]
```

Out[47]:

| | ID | price |
|---|---|---|
| 0 | 1 | 8900 |
| 1 | 2 | 8800 |
| 2 | 3 | 4200 |
| 3 | 4 | 6000 |
| 4 | 5 | 5700 |
| 5 | 6 | 7900 |
| 6 | 7 | 10750 |
| 7 | 8 | 9190 |
| 8 | 9 | 5600 |
| 9 | 10 | 6000 |
| 10 | 11 | 8950 |
| 11 | 12 | 10990 |

In [48]:

```python
df=df[["engine_power","price"]]
df.columns=["Engine","pri"]
```

In [49]:

```python
sns.lmplot(x='Engine',y='pri',data=df,order=2,ci=None)
```

Out[49]:

```
<seaborn.axisgrid.FacetGrid at 0x17fd2e74350>
```

In [50]:

```python
df.describe()
```

Out[50]:

|  | Engine | pri |
|---|---|---|
| count | 1538.000000 | 1538.000000 |
| mean | 51.904421 | 8576.003901 |
| std | 3.988023 | 1939.958641 |
| min | 51.000000 | 2500.000000 |
| 25% | 51.000000 | 7122.500000 |
| 50% | 51.000000 | 9000.000000 |
| 75% | 51.000000 | 10000.000000 |
| max | 77.000000 | 11100.000000 |

In [51]:

```python
df.fillna(method="ffill")
```

Out[51]:

|  | Engine | pri |
|---|---|---|
| 0 | 51 | 8900 |
| 1 | 51 | 8800 |
| 2 | 74 | 4200 |
| 3 | 51 | 6000 |
| 4 | 73 | 5700 |
| ... | ... | ... |
| 1533 | 51 | 5200 |
| 1534 | 74 | 4600 |
| 1535 | 51 | 7500 |
| 1536 | 51 | 5990 |
| 1537 | 51 | 7900 |

1538 rows × 2 columns

In [52]:

```python
x=np.array(df['Engine']).reshape(-1,1)
y=np.array(df['pri']).reshape(-1,1)
```

In [53]:

```python
df.dropna(inplace=True)
```

```
C:\Users\Gowthami\AppData\Local\Temp\ipykernel_22952\1379821321.py:1: Sett
ingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-doc
s/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://
pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-
view-versus-a-copy)
  df.dropna(inplace=True)
```

In [54]:

```python
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
```

In [55]:

```python
regr=LinearRegression()
regr.fit(x_train,y_train)
print(regr.score(x_test,y_test))
```

```
0.046186912889935816
```

In [56]:

```python
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
```

In [57]:

```
df500=df[:][:500]
sns.lmplot(x='Engine',y='pri',data=df500,order=1,ci=None)
```

Out[57]:

```
<seaborn.axisgrid.FacetGrid at 0x17fd2e13710>
```



#df500.fillna(method='ffill',inplace=True) x=np.array(df500['Engine']).reshape(-1,1)
y=np.array(df500['pri']).reshape(-1,1) df500.dropna(inplace=True)
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25) regr=LinearRegression()
regr.fit(x_train,y_train) print("Regression:",regr.score(x_test,y_test)) y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='b') plt.plot(x_test,y_pred,color='k') plt.show()

In [58]:

```
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
model=LinearRegression()
model.fit(x_train,y_train)
y_pred=model.predict(x_test)
r2=r2_score(y_test,y_pred)
print("R2 score:",r2)
```

```
R2 score: 0.04618691289935816
```

#conclusion: Linear regression is not fit for the model

# Ridge and Lasso regression

In [59]:

```python
features = df.columns[0:2]
target = df.columns[-1]
#X and y values
X = df[features].values
y = df[target].values
#splot
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=17
print("The dimension of X_train is {}".format(X_train.shape))
print("The dimension of X_test is {}".format(X_test.shape))
#Scale features
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```
The dimension of X_train is (1076, 2)
The dimension of X_test is (462, 2)
```

In [60]:

```python
#Model
lr = LinearRegression()
#Fit model
lr.fit(X_train, y_train)
#predict
#prediction = lr.predict(X_test)
#actual
actual = y_test
train_score_lr = lr.score(X_train, y_train)
test_score_lr = lr.score(X_test, y_test)
print("\nLinear Regression Model:\n")
print("The train score for lr model is {}".format(train_score_lr))
print("The test score for lr model is {}".format(test_score_lr))
```

```
Linear Regression Model:

The train score for lr model is 1.0
The test score for lr model is 1.0
```
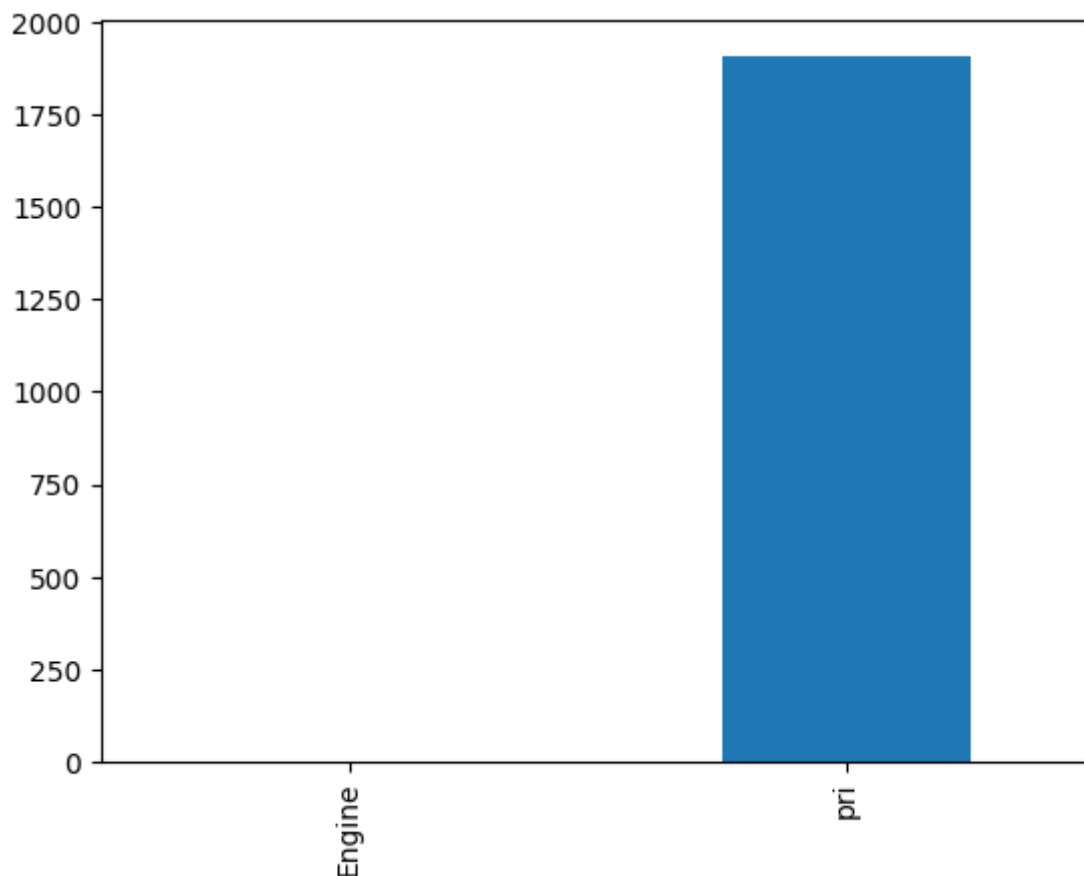
In [61]:

```python
#Using the linear CV model
from sklearn.linear_model import RidgeCV
#Ridge Cross validation
ridge_cv = RidgeCV(alphas = [0.0001, 0.001,0.01, 0.1, 1, 10]).fit(X_train, y_train)
#score
print("The train score for ridge model is {}".format(ridge_cv.score(X_train, y_train)))
print("The train score for ridge model is {}".format(ridge_cv.score(X_test, y_test)))
```

```
The train score for ridge model is 0.9999999999999897
The train score for ridge model is 0.9999999999999899
```

In [62]:

```python
ridgeReg = Ridge(alpha=10)
ridgeReg.fit(X_train,y_train)
#train and test scorefor ridge regression
train_score_ridge = ridgeReg.score(X_train, y_train)
test_score_ridge = ridgeReg.score(X_test, y_test)
print("\nRidge Model:\n")
print("The train score for ridge model is {}".format(train_score_ridge))
print("The test score for ridge model is {}".format(test_score_ridge))
```

```
Ridge Model:

The train score for ridge model is 0.9999088581979684
The test score for ridge model is 0.9999100853681022
```

In [63]:

```python
plt.figure(figsize = (10, 10))
sns.heatmap(df.corr(), annot = True)
```

Out[63]:

```
<Axes: >
```

In [64]:

```python
plt.figure(figsize = (10, 10))
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,colo

plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='gre
plt.xticks(rotation = 90)
plt.legend()
plt.show()
```



In [65]:

```python
#Using the linear CV model
from sklearn.linear_model import LassoCV
#Lasso Cross validation
lasso_cv = LassoCV(alphas = [0.0001, 0.001,0.01, 0.1, 1, 10], random_state=0).fit(X_trai
#score
print(lasso_cv.score(X_train, y_train))
print(lasso_cv.score(X_test, y_test))
```

```
0.999999999501757
0.999999999638806
```

In [66]:

```python
#Lasso regression model
print("\nLasso Model: \n")
lasso = Lasso(alpha = 10)
lasso.fit(X_train,y_train)
train_score_ls =lasso.score(X_train,y_train)
test_score_ls =lasso.score(X_test,y_test)
print("The train score for ls model is {}".format(train_score_ls))
print("The test score for ls model is {}".format(test_score_ls))
```

```
Lasso Model:

The train score for ls model is 0.9999728562194999
The test score for ls model is 0.9999728508562553
```

In [67]:

```python
pd.Series(lasso.coef_, features).sort_values(ascending = True).plot(kind = "bar")
```

Out[67]:

```
<Axes: >
```

In [68]:

```python
#plot size
plt.figure(figsize = (10, 10))
#add plot for ridge regression
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,colo
#add plot for lasso regression
plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue',
#add plot for linear model
plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='gre
```

Out[68]:

`[<matplotlib.lines.Line2D at 0x17fd2c63290>]`



# Elastic net

In [69]:

```python
from sklearn.linear_model import ElasticNet
regr=ElasticNet()
regr.fit(X,y)
print(regr.coef_)
print(regr.intercept_)
```

```
[-0.         0.99999973]
0.002280249860632466
```

In [70]:

```python
y_pred_elastic=regr.predict(X_train)
```

In [71]:

```python
mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
print("Mean Squared Error on test set",mean_squared_error)
```

```
Mean Squared Error on test set 77371869.93693778
```

# Data

In [72]:

```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn import preprocessing,svm
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Lasso
from sklearn.linear_model import Ridge
```

In [73]:

```python
de=pd.read_csv(r"C:\Users\Gowthami\Downloads\data.csv")
de
```

Out[73]:

| | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront |
|---|---|---|---|---|---|---|---|---|
| 0 | 2014-05-02 00:00:00 | 3.130000e+05 | 3.0 | 1.50 | 1340 | 7912 | 1.5 | 0 |
| 1 | 2014-05-02 00:00:00 | 2.384000e+06 | 5.0 | 2.50 | 3650 | 9050 | 2.0 | 0 |
| 2 | 2014-05-02 00:00:00 | 3.420000e+05 | 3.0 | 2.00 | 1930 | 11947 | 1.0 | 0 |
| 3 | 2014-05-02 00:00:00 | 4.200000e+05 | 3.0 | 2.25 | 2000 | 8030 | 1.0 | 0 |
| 4 | 2014-05-02 00:00:00 | 5.500000e+05 | 4.0 | 2.50 | 1940 | 10500 | 1.0 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 4595 | 2014-07-09 00:00:00 | 3.081667e+05 | 3.0 | 1.75 | 1510 | 6360 | 1.0 | 0 |
| 4596 | 2014-07-09 00:00:00 | 5.343333e+05 | 3.0 | 2.50 | 1460 | 7573 | 2.0 | 0 |
| 4597 | 2014-07-09 00:00:00 | 4.169042e+05 | 3.0 | 2.50 | 3010 | 7014 | 2.0 | 0 |
| 4598 | 2014-07-10 00:00:00 | 2.034000e+05 | 4.0 | 2.00 | 2090 | 6630 | 1.0 | 0 |
| 4599 | 2014-07-10 00:00:00 | 2.206000e+05 | 3.0 | 2.50 | 1490 | 8102 | 2.0 | 0 |

4600 rows × 18 columns

In [74]:

```python
de=de[["sqft_living","sqft_lot"]]
de.columns=["Liv","Lot"]
```

In [75]:

```python
sns.lmplot(x='Liv',y='Lot',data=de,order=2,ci=None)
```

Out[75]:

```
<seaborn.axisgrid.FacetGrid at 0x17fd2c2f7d0>
```



In [76]:

```python
de.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4600 entries, 0 to 4599
Data columns (total 2 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   Liv     4600 non-null   int64
 1   Lot     4600 non-null   int64
dtypes: int64(2)
memory usage: 72.0 KB
```

In [77]:

```
de.describe()
```

Out[77]:

|  | Liv | Lot |
|---|---|---|
| **count** | 4600.000000 | 4.600000e+03 |
| **mean** | 2139.346957 | 1.485252e+04 |
| **std** | 963.206916 | 3.588444e+04 |
| **min** | 370.000000 | 6.380000e+02 |
| **25%** | 1460.000000 | 5.000750e+03 |
| **50%** | 1980.000000 | 7.683000e+03 |
| **75%** | 2620.000000 | 1.100125e+04 |
| **max** | 13540.000000 | 1.074218e+06 |

In [78]:

```
de.fillna(method='ffill')
```

Out[78]:

|  | Liv | Lot |
|---|---|---|
| **0** | 1340 | 7912 |
| **1** | 3650 | 9050 |
| **2** | 1930 | 11947 |
| **3** | 2000 | 8030 |
| **4** | 1940 | 10500 |
| **...** | ... | ... |
| **4595** | 1510 | 6360 |
| **4596** | 1460 | 7573 |
| **4597** | 3010 | 7014 |
| **4598** | 2090 | 6630 |
| **4599** | 1490 | 8102 |

4600 rows × 2 columns

In [79]:

```
x=np.array(de['Liv']).reshape(-1,1)
y=np.array(de['Lot']).reshape(-1,1)
```

In [80]:

```python
de.dropna(inplace=True)
```

```
C:\Users\Gowthami\AppData\Local\Temp\ipykernel_22952\836337131.py:1: Setti
ngWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-doc
s/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://
pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-
view-versus-a-copy)
  de.dropna(inplace=True)
```

In [81]:

```python
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
```

In [82]:

```python
regr=LinearRegression()
regr.fit(x_train,y_train)
print(regr.score(x_test,y_test))
```

```
0.10640487967838053
```

In [83]:

```python
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
```

In [84]:

```
de500=de[:][:500]
sns.lmplot(x='Liv',y='Lot',data=de500,order=1,ci=None)
```
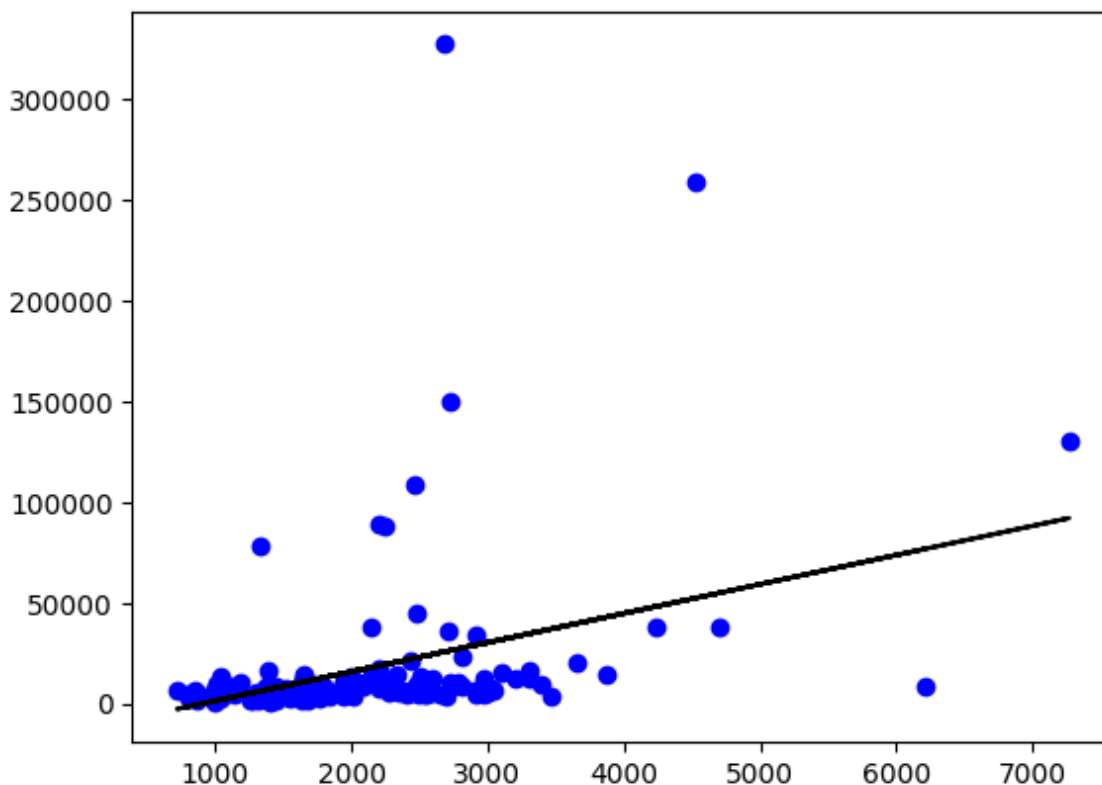
Out[84]:

```
<seaborn.axisgrid.FacetGrid at 0x17fd2e50650>
```

In [85]:

```python
de500.fillna(method='ffill',inplace=True)
x=np.array(de500['Liv']).reshape(-1,1)
y=np.array(de500['Lot']).reshape(-1,1)
de500.dropna(inplace=True)
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
regr=LinearRegression()
regr.fit(x_train,y_train)
print("Regression:",regr.score(x_test,y_test))
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='b')
plt.plot(x_test,y_pred,color='k')
plt.show()
```

Regression: 0.11354254617812931



In [86]:

```python
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
model=LinearRegression()
model.fit(x_train,y_train)
y_pred=model.predict(x_test)
r2=r2_score(y_test,y_pred)
print("R2 score:",r2)
```

R2 score: 0.11354254617812931

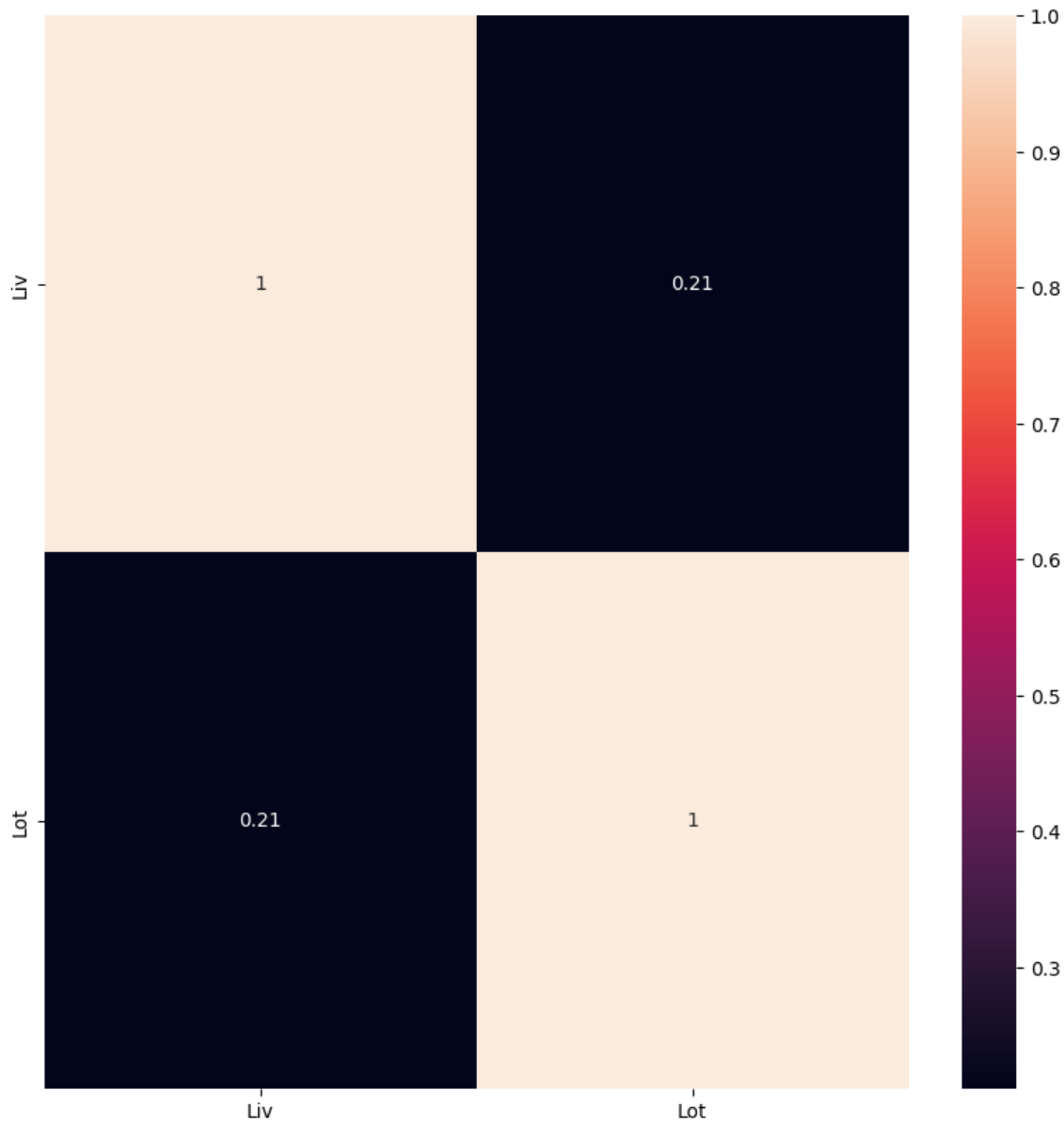#conclusion: Linear regression is fit for the model

## Ridge and lasso

In [87]:

```python
plt.figure(figsize = (10, 10))
sns.heatmap(de.corr(), annot = True)
```

Out[87]:

`<Axes: >`

In [88]:

```python
features = de.columns[0:2]
target = de.columns[-1]
#X and y values
X = de[features].values
y = de[target].values
#splot
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=17
print("The dimension of X_train is {}".format(X_train.shape))
print("The dimension of X_test is {}".format(X_test.shape))
#Scale features
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```
The dimension of X_train is (3220, 2)
The dimension of X_test is (1380, 2)
```

In [89]:

```python
#Model
lr = LinearRegression()
#Fit model
lr.fit(X_train, y_train)
#predict
#prediction = lr.predict(X_test)
#actual
actual = y_test
train_score_lr = lr.score(X_train, y_train)
test_score_lr = lr.score(X_test, y_test)
print("\nLinear Regression Model:\n")
print("The train score for lr model is {}".format(train_score_lr))
print("The test score for lr model is {}".format(test_score_lr))
```

```
Linear Regression Model:

The train score for lr model is 1.0
The test score for lr model is 1.0
```

In [90]:

```python
#Using the Linear CV model
from sklearn.linear_model import RidgeCV
#Ridge Cross validation
ridge_cv = RidgeCV(alphas = [0.0001, 0.001,0.01, 0.1, 1, 10]).fit(X_train, y_train)
#score
print("The train score for ridge model is {}".format(ridge_cv.score(X_train, y_train)))
print("The train score for ridge model is {}".format(ridge_cv.score(X_test, y_test)))
```

```
The train score for ridge model is 0.9999999999999997
The train score for ridge model is 0.9999999999999996
```

In [91]:

```python
ridgeReg = Ridge(alpha=10)
ridgeReg.fit(X_train,y_train)
#train and test scorefor ridge regression
train_score_ridge = ridgeReg.score(X_train, y_train)
test_score_ridge = ridgeReg.score(X_test, y_test)
print("\nRidge Model:\n")
print("The train score for ridge model is {}".format(train_score_ridge))
print("The test score for ridge model is {}".format(test_score_ridge))
```

```
Ridge Model:

The train score for ridge model is 0.9999900245012017
The test score for ridge model is 0.9999902306741419
```
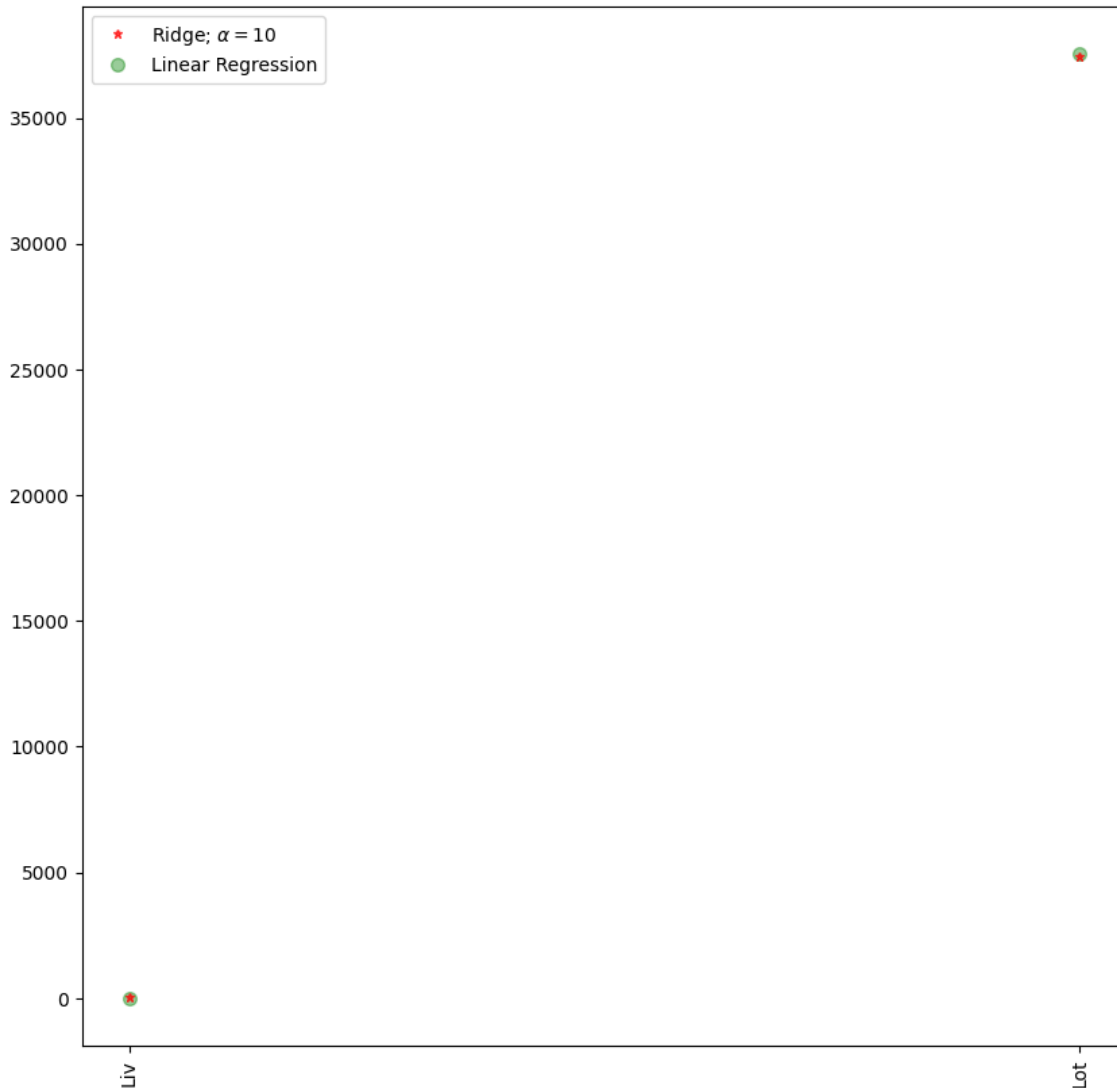
In [92]:

```
plt.figure(figsize = (10, 10))
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,colo

plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='gre
plt.xticks(rotation = 90)
plt.legend()
plt.show()
```



In [93]:

```
#Using the Linear CV model
from sklearn.linear_model import LassoCV
#Lasso Cross validation
lasso_cv = LassoCV(alphas = [0.0001, 0.001,0.01, 0.1, 1, 10], random_state=0).fit(X_trai
#score
print(lasso_cv.score(X_train, y_train))
print(lasso_cv.score(X_test, y_test))
```

```
0.999999999997705
0.999999999996928
```

In [94]:

```python
#Lasso regression model
print("\nLasso Model: \n")
lasso = Lasso(alpha = 10)
lasso.fit(X_train,y_train)
train_score_ls =lasso.score(X_train,y_train)
test_score_ls =lasso.score(X_test,y_test)
print("The train score for ls model is {}".format(train_score_ls))
print("The test score for ls model is {}".format(test_score_ls))
```

```
Lasso Model:

The train score for ls model is 0.9999999291360324
The test score for ls model is 0.9999999291231869
```
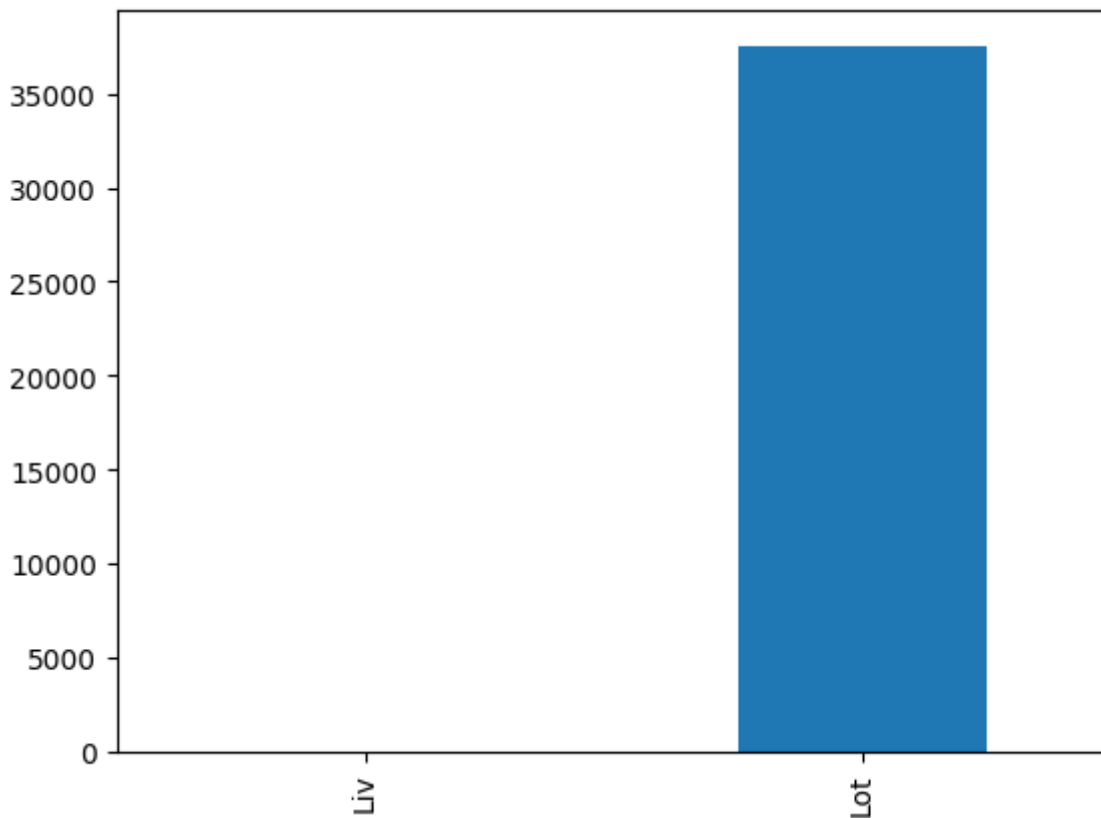
In [95]:

```python
pd.Series(lasso.coef_, features).sort_values(ascending = True).plot(kind = "bar")
```
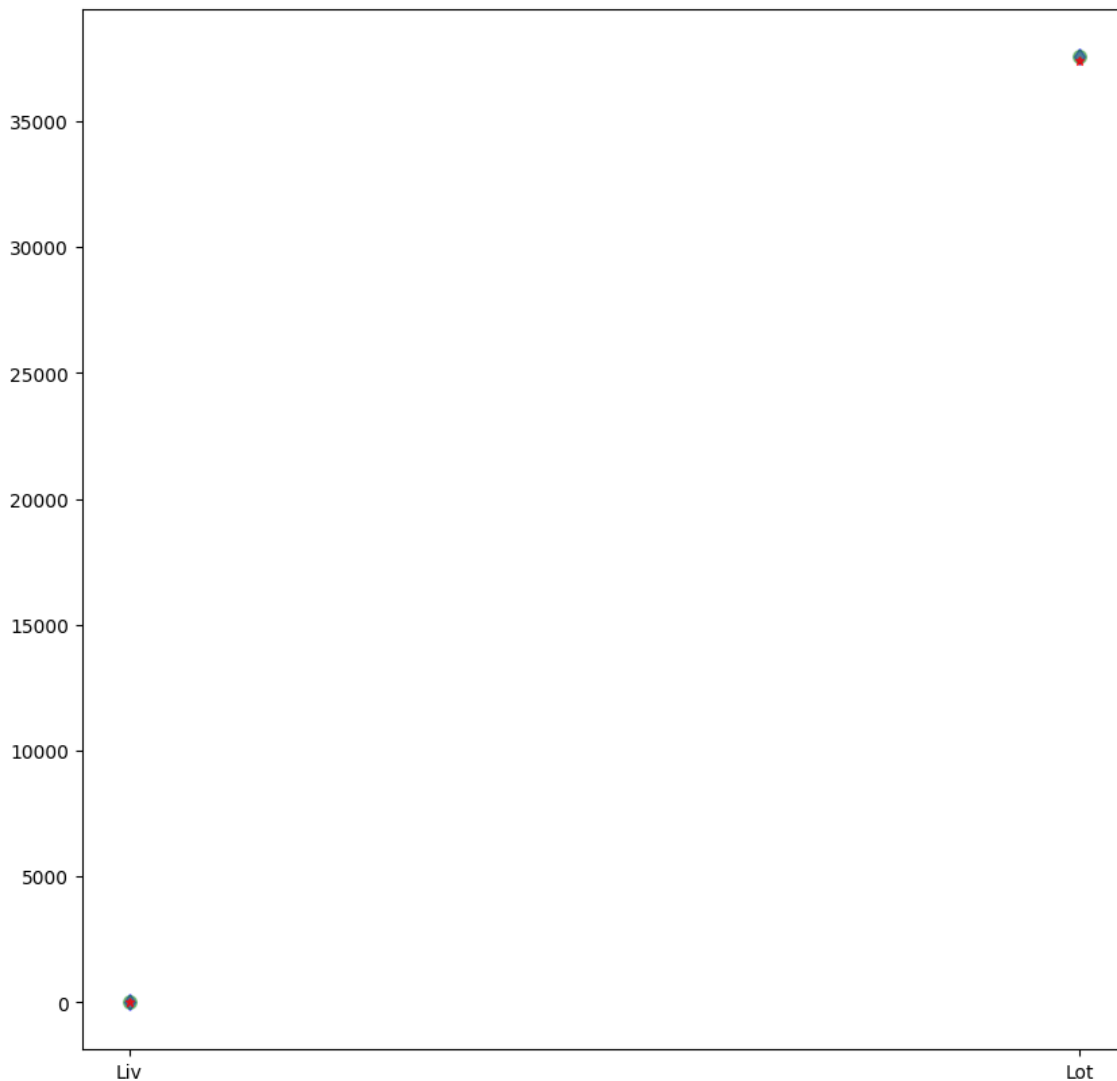
Out[95]:

```
<Axes: >
```

In [96]:

```python
#plot size
plt.figure(figsize = (10, 10))
#add plot for ridge regression
plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,colo
#add plot for lasso regression
plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue',
#add plot for linear model
plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='gre
```

Out[96]:

[<matplotlib.lines.Line2D at 0x17ffa835350>]



# Elastic net

In [97]:

```python
from sklearn.linear_model import ElasticNet
regr=ElasticNet()
regr.fit(X,y)
print(regr.coef_)
print(regr.intercept_)
```

```
[7.84570096e-07 9.99999995e-01]
-0.001601077769009862
```

In [98]:

```python
y_pred_elastic=regr.predict(X_train)
```

In [99]:

```python
mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
print("Mean Squared Error on test set",mean_squared_error)
```

```
Mean Squared Error on test set 1635484254.240867
```

In [ ]: