

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from sklearn import preprocessing, svm
6 from sklearn.model_selection import train_test_split
7 from sklearn.linear_model import LinearRegression
8 from sklearn.linear_model import Lasso
9 from sklearn.linear_model import Ridge
```

In [2]:

```
1 #Step-2:Reading the dataset
2 df=pd.read_csv(r"C:\Users\91955\Desktop\Data Analysis with Python\bottle.csv")
3 df
```

C:\Users\91955\AppData\Local\Temp\ipykernel_4240\292355879.py:2: DtypeWarning: Columns (47,73) have mixed types. Specify dtype option on import or set low_memory=False.

```
df=pd.read_csv(r"C:\Users\91955\Desktop\Data Analysis with Python\bottle.csv")
```

Out[2]:

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

In [3]:

```
1 df=df[['Salnty','T_degC']]
2 #Taking only the selected attributes from the dataset
3 df.columns=['Sal','Temp']
4 #Renaming the columns for easier writing of the code
```

In [4]:

```
1 df.head(10)
2 #Displaying only the 1st 10 rows
```

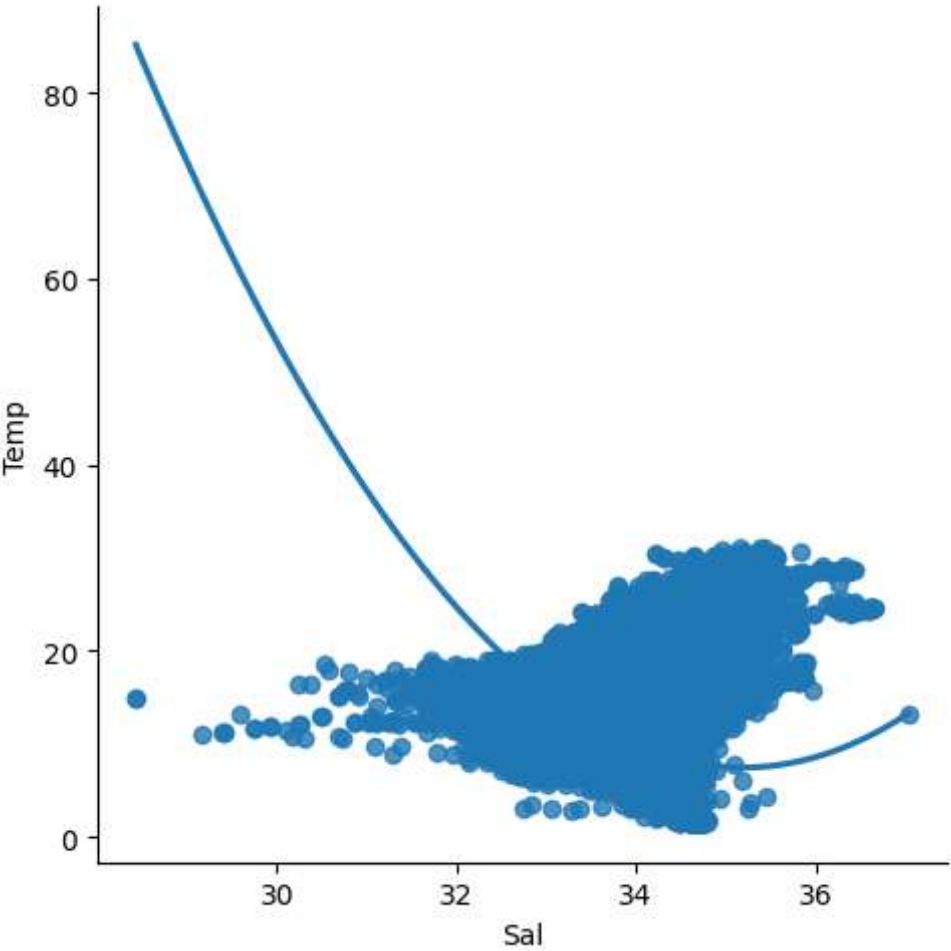
Out[4]:

Out[4]:	1	5	054.0 056.0	HY-060- 0930- 05400560- 0020A-7	20	10.450	33.4210	NaN	25.64300
	Sal	Temp							
0	33.440	10.50
1	33.440	10.46		20- 1611SR- MX-310- 2239-					
2 864858	33.437 34404	10.46 34404	864859	093.4 026.4	0	18.744	33.4083	5.805	23.87055
3	33.420	10.45		09340264- 0000A-7					
4	33.421	10.45							
5	33.431	10.45		20- 1611SR- MX-310- 2239-					
6 864859	33.440 34404	10.45 34404	864860	093.4 026.4	2	18.744	33.4083	5.805	23.87072
7	33.424	10.24		09340264- 0002A-3					
8	33.420	10.06							
9	33.494	9.86		20- 1611SR- MX-310- 2239-					
864860	34404	864861	093.4 026.4	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

```
In [5]: Cst_Cnt  Btl_Cnt  Sta_ID  Depth_ID  Depthm  T_degC  Salnty  O2ml_L  STheta
```

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

```
Out[5]: 34404 864863 093.4 20-1611SR- 15 17.533 33.3880 5.774 24.15297
026.4 2239-
09340264-
0015A-3
<seaborn.axisgrid.FacetGrid at 0x235b87ac880>
```



```
In [6]:
```

```
1 df.describe()
```

```
Out[6]:
```

	Sal	Temp
count	817509.000000	853900.000000
mean	33.840350	10.799677
std	0.461843	4.243825
min	28.431000	1.440000
25%	33.488000	7.680000
50%	33.863000	10.060000
75%	34.196900	13.880000
max	37.034000	31.140000

In [7]:

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 864863 entries, 0 to 864862
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  -
0    Sal      817509 non-null    float64
1    Temp     853900 non-null    float64
dtypes: float64(2)
memory usage: 13.2 MB
```

In [8]:

```
1 df.fillna(method='ffill')
```

Out[8]:

	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 [9]:

```
1 df.isnull().sum()
```

Out[9]:

```
Sal      47354
Temp     10963
dtype: int64
```

In [10]:

```
1 df.fillna(value = 0,  
2           inplace = True)  
3
```

C:\Users\91955\AppData\Local\Temp\ipykernel_4240\3904320923.py:1: SettingWithCopyWarning:

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-docs/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.fillna(value = 0,
```

In [11]:

```
1 df.isna().any()
```

Out[11]:

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

In [12]:

```
1 x=np.array(df['Sal']).reshape(-1,1)
```

In [13]:

```
1 y=np.array(df['Temp']).reshape(-1,1)
```

In [14]:

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

C:\Users\91955\AppData\Local\Temp\ipykernel_4240\1379821321.py:1: SettingWithCopyWarning:

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-docs/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)
```

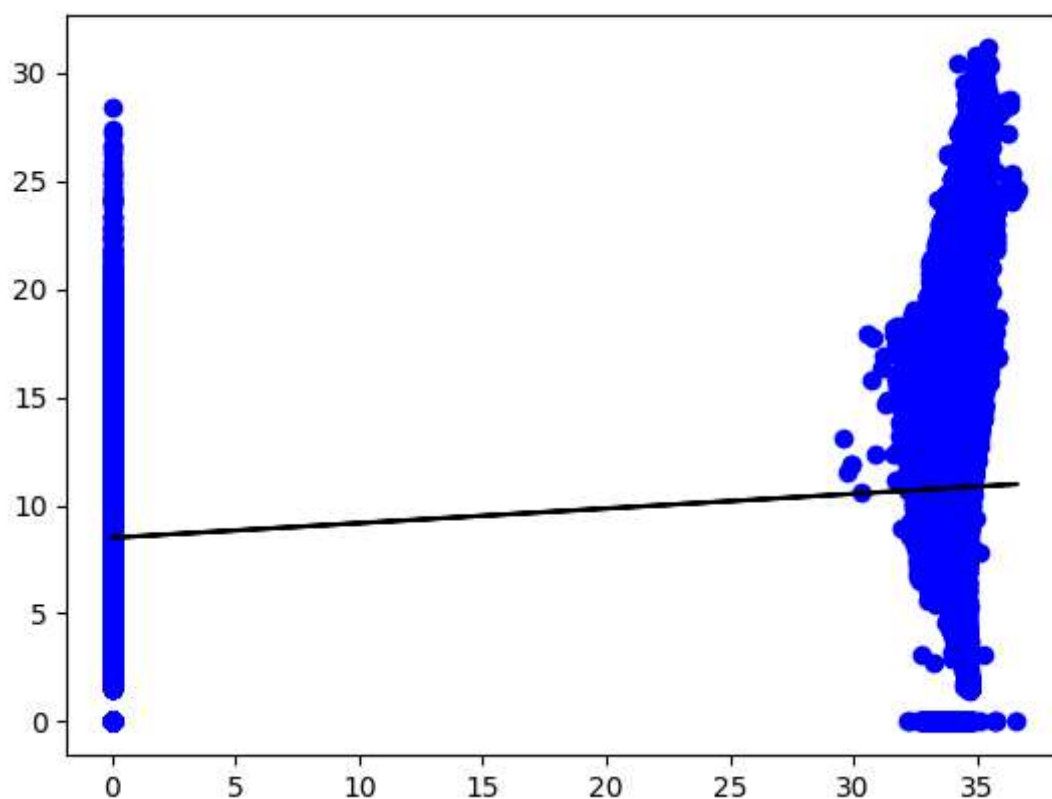
In [15]:

```
1 X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.25)  
2 reg=LinearRegression()  
3 reg.fit(X_train,y_train)  
4 print(reg.score(X_test,y_test))
```

0.014074386054016674

In [16]:

```
1 y_pred=reg.predict(X_test)
2 plt.scatter(X_test,y_test,color='b')
3 plt.plot(X_test,y_pred,color='k')
4 plt.show()
```

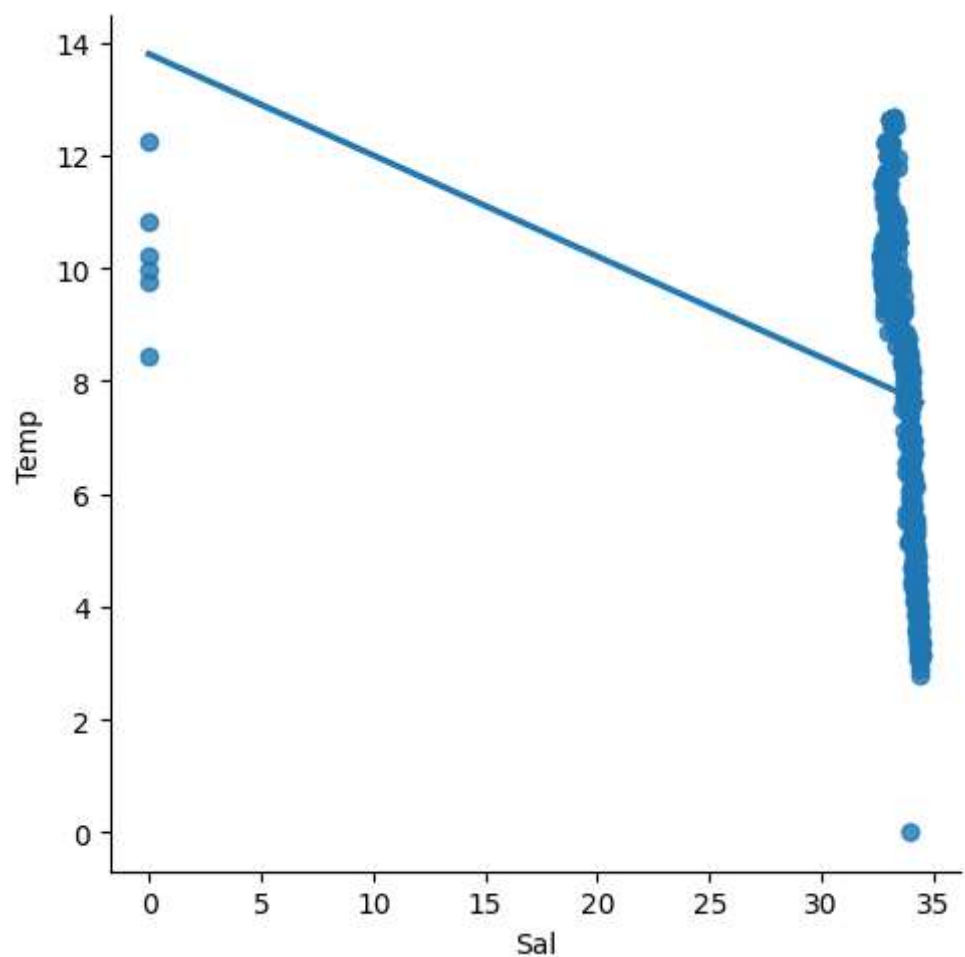


In [17]:

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

Out[17]:

<seaborn.axisgrid.FacetGrid at 0x239b87adc60>



In [18]:

```

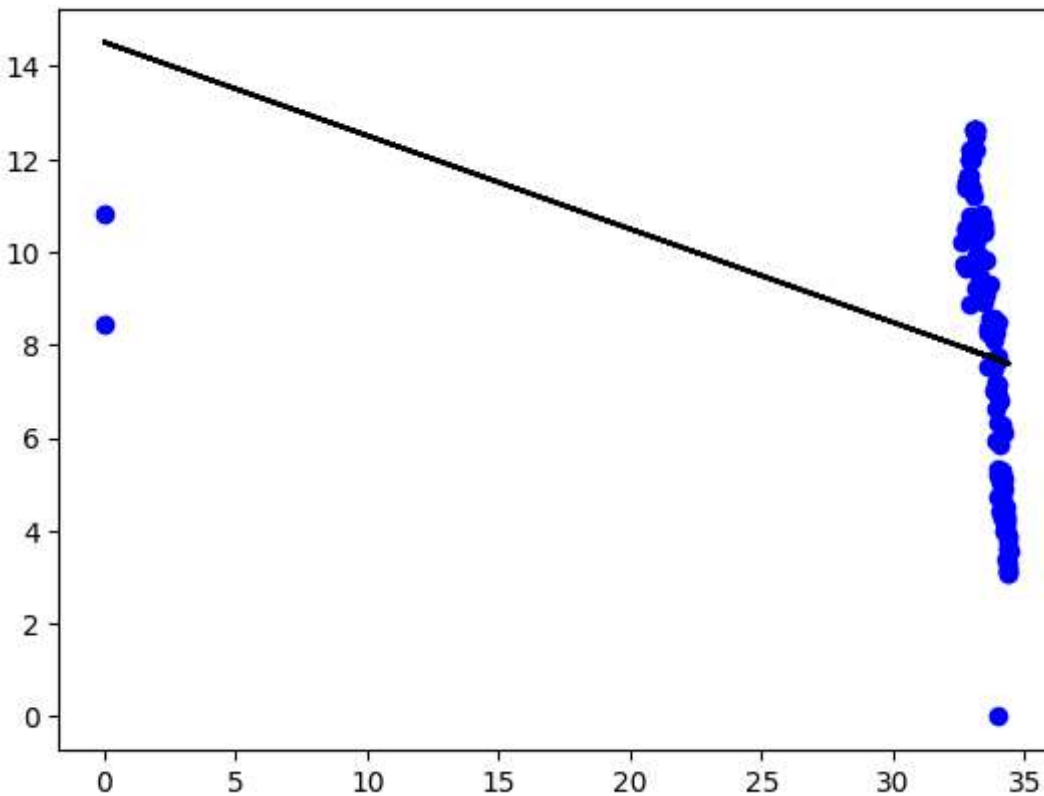
1 df500.fillna(method='ffill',inplace=True)
2 X=np.array(df500['Sal']).reshape(-1,1)
3 y=np.array(df500['Temp']).reshape(-1,1)
4 df500.dropna(inplace=True)
5 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.25)
6 reg=LinearRegression()
7 reg.fit(X_train,y_train)
8 print("Regression:",reg.score(X_test,y_test))
9 y_pred=reg.predict(X_test)
10 plt.scatter(X_test,y_test,color='b')
11 plt.plot(X_test,y_pred,color='k')
12 plt.show

```

Regression: 0.027406409205778193

Out[18]:

<function matplotlib.pyplot.show(close=None, block=None)>



In [19]:

```

1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score
3 model=LinearRegression()
4 model.fit(X_train,y_train)
5 y_pred=model.predict(X_test)
6 r2=r2_score(y_test,y_pred)
7 print("R2 score: ",r2)

```

R2 score: 0.027406409205778193

In [20]:

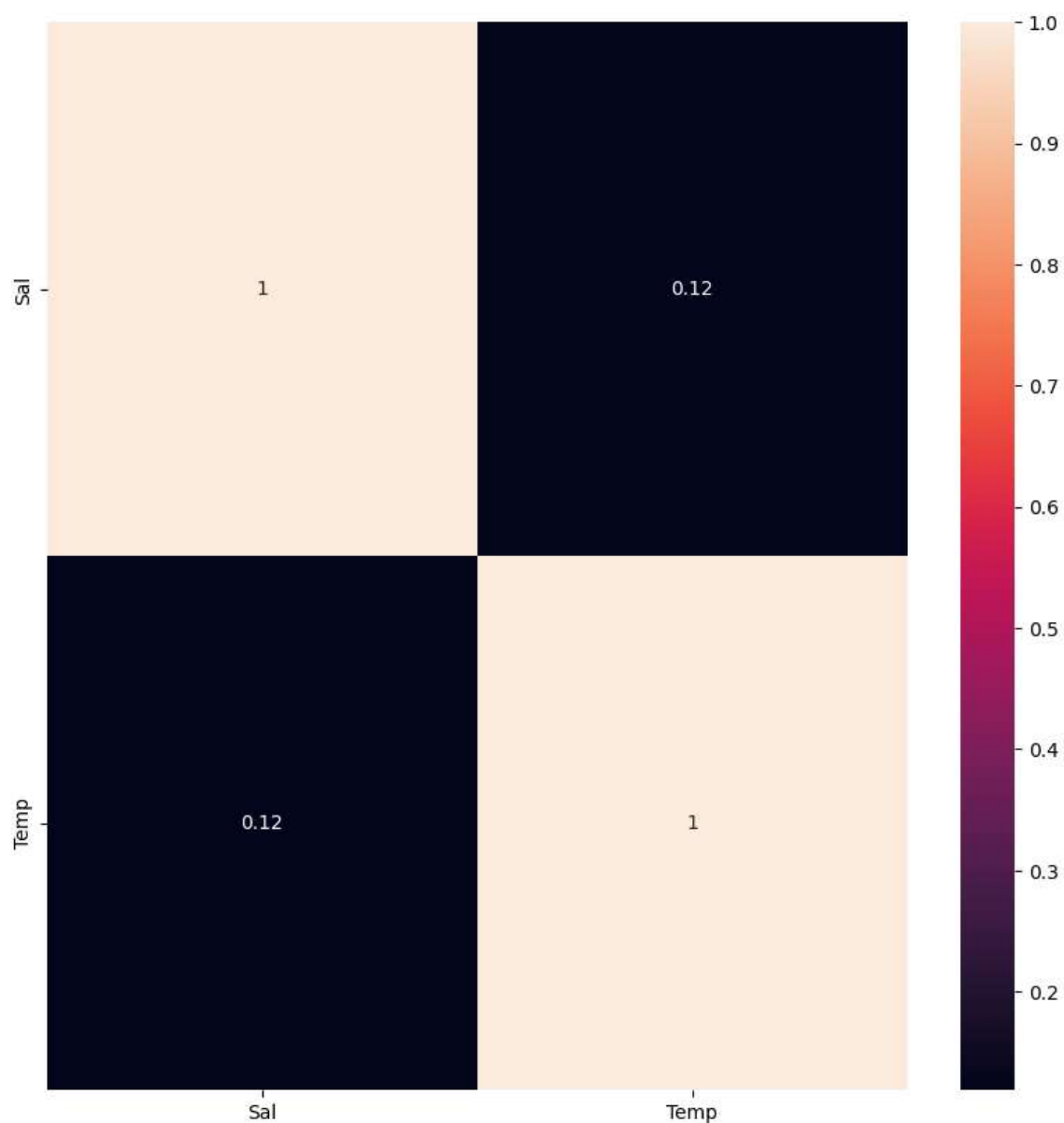
```
1 #conclusion:Linear regression is the best fit for the model
```

In [21]:

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

Out[21]:

<Axes: >



Ridge and Lasso

In [22]:

```
1 features = df.columns[0:2]
2 target = df.columns[-1]
3 #X and y values
4 X = df[features].values
5 y = df[target].values
6 #split
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_stat
8 print("The dimension of X_train is {}".format(X_train.shape))
9 print("The dimension of X_test is {}".format(X_test.shape))
10 #Scale features
11 from sklearn.preprocessing import StandardScaler
12 scaler = StandardScaler()
13 X_train = scaler.fit_transform(X_train)
14 X_test = scaler.transform(X_test)
```

The dimension of X_train is (605404, 2)

The dimension of X_test is (259459, 2)

In [23]:

```
1 #Model
2 lr = LinearRegression()
3 #Fit model
4 lr.fit(X_train, y_train)
5 #predict
6 #prediction = lr.predict(X_test)
7 #actual
8 actual = y_test
9 train_score_lr = lr.score(X_train, y_train)
10 test_score_lr = lr.score(X_test, y_test)
11 print("\nLinear Regression Model:")
12 print("The train score for lr model is {}".format(train_score_lr))
13 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 [24]:

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

The train score for ridge model is 0.9999999986797505

The train score for ridge model is 0.9999999986778121

In [25]:

```
1 ridgeReg = Ridge(alpha=10)
2 ridgeReg.fit(X_train,y_train)
3 #train and test scorefor ridge regression
4 train_score_ridge = ridgeReg.score(X_train, y_train)
5 test_score_ridge = ridgeReg.score(X_test, y_test)
6 print("\nRidge Model:")
7 print("The train score for ridge model is {}".format(train_score_ridge))
8 print("The test score for ridge model is {}".format(test_score_ridge))
```

Ridge Model:

The train score for ridge model is 0.99999999723243

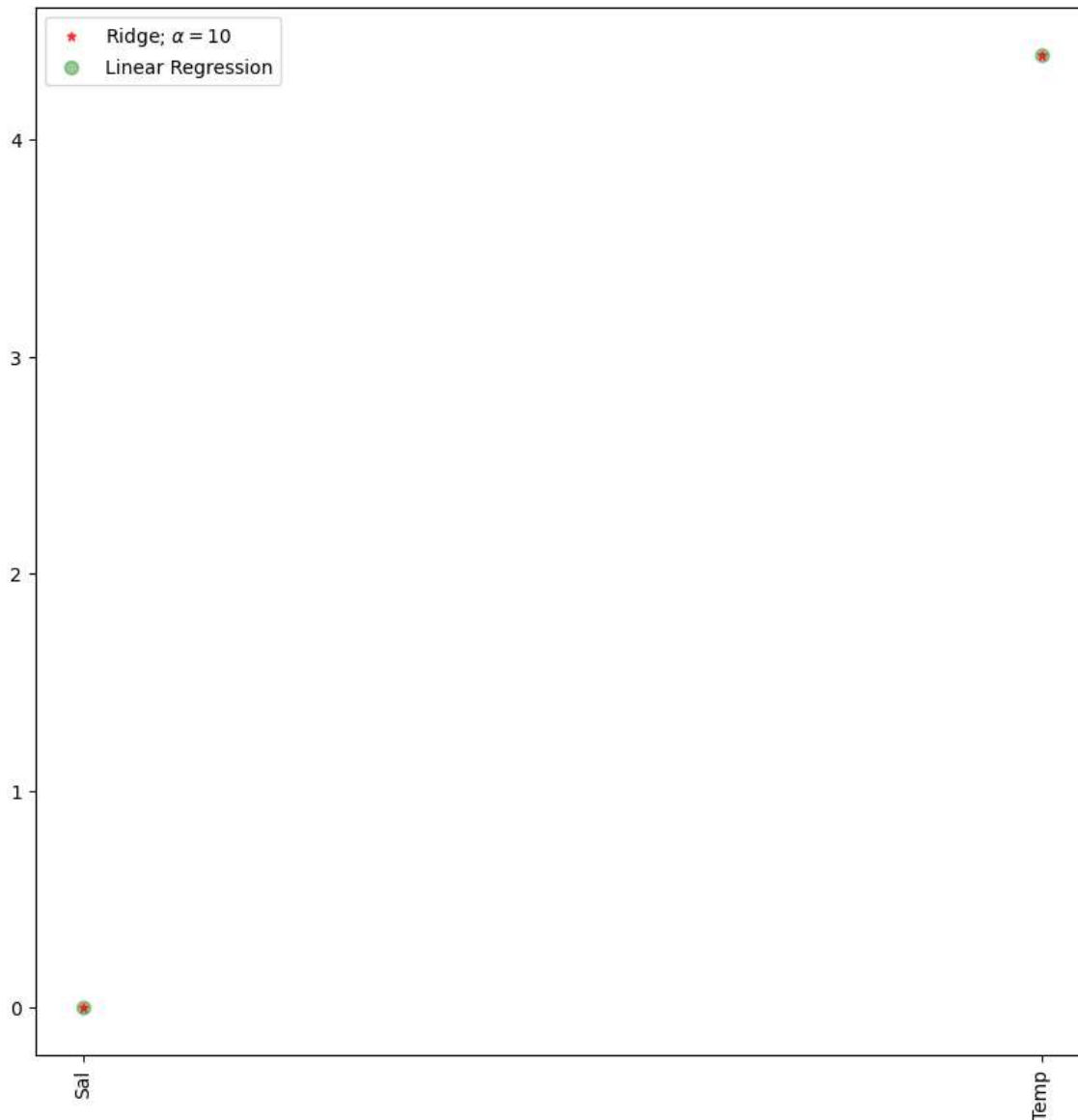
The test score for ridge model is 0.999999997231402

In [26]:

```

1 plt.figure(figsize = (10, 10))
2 plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,
3
4 plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color=
5 plt.xticks(rotation = 90)
6 plt.legend()
7 plt.show()

```



In [27]:

```

1 #Using the linear CV model
2 from sklearn.linear_model import LassoCV
3 #Lasso Cross validation
4 lasso_cv = LassoCV(alphas = [0.0001, 0.001,0.01, 0.1, 1, 10], random_state=0).fit(X_
5 #score
6 print(lasso_cv.score(X_train, y_train))
7 print(lasso_cv.score(X_test, y_test))

```

0.9999999994806811

0.9999999994806712

In [28]:

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

Lasso Model:

The train score for ls model is 0.0

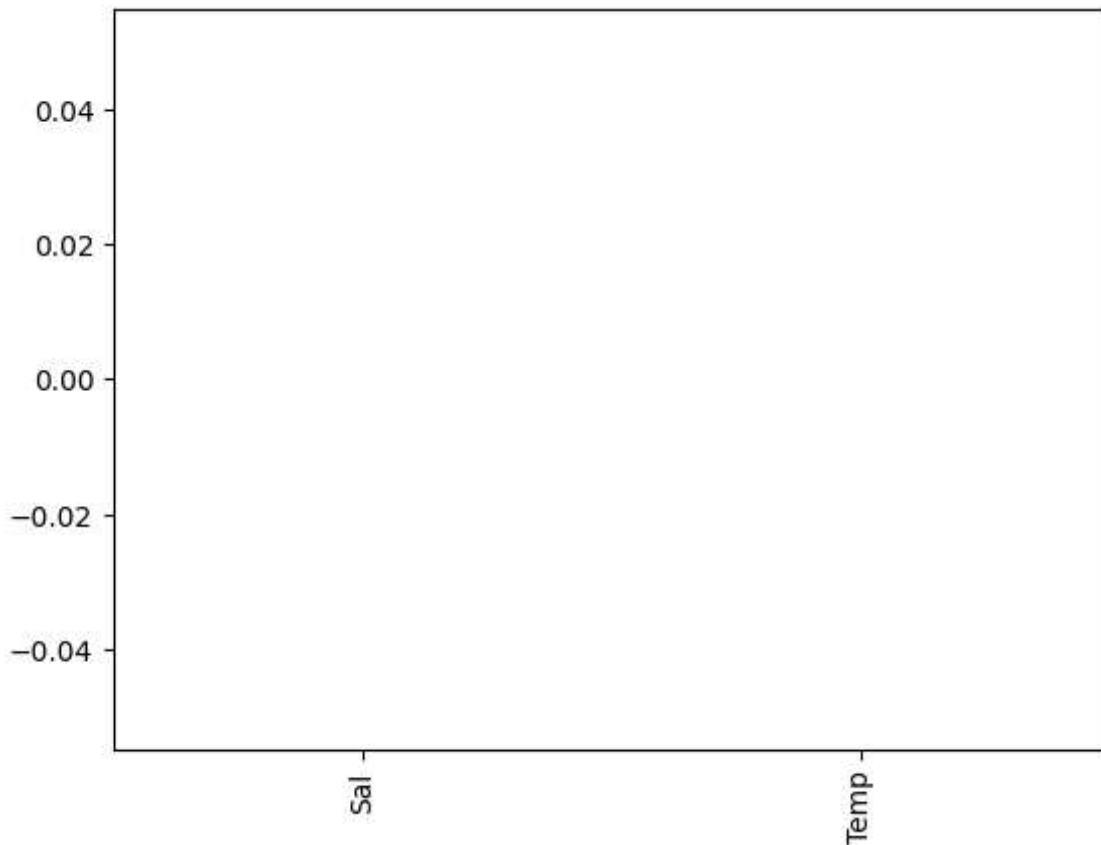
The test score for ls model is -1.9031696447013857e-05

In [29]:

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

Out[29]:

<Axes: >

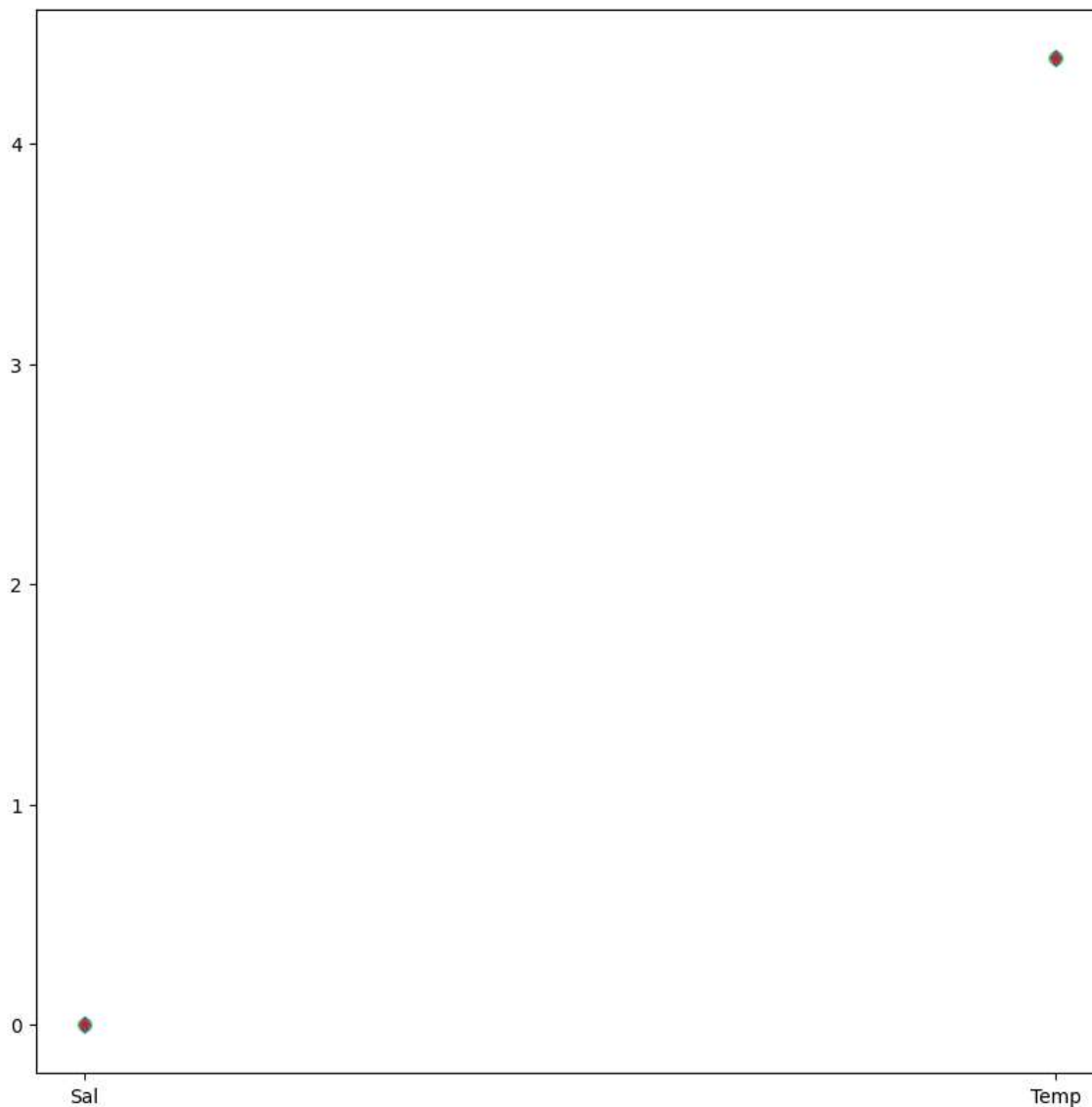


In [30]:

```
1 #plot size
2 plt.figure(figsize = (10, 10))
3 #add plot for ridge regression
4 plt.plot(features,ridgeReg.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,
5 #add plot for lasso regression
6 plt.plot(lasso_cv.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='b')
7 #add plot for linear model
8 plt.plot(features,lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='r')
```

Out[30]:

[<matplotlib.lines.Line2D at 0x23986b02950>]



Elastic Regression

In [31]:

```
1 from sklearn.linear_model import ElasticNet
2 regr=ElasticNet()
3 regr.fit(X,y)
4 print(regr.coef_)
5 print(regr.intercept_)
```

```
[0.          0.94934511]
0.540121963106797
```

In [32]:

```
1 y_pred_elastic=regr.predict(X_train)
2 mean_squared_error=np.mean((y_pred_elastic-y_train)**2)
3 print("Mean squared error on test set",mean_squared_error)
```

Mean squared error on test set 114.40984808659212

In []:

```
1
```