

In [2]:

```
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
```

In [3]:

```
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_2828\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[3]:

	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 [4]:

```
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 [5]:

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

Out[5]:

	Sal	Temp								
0	33.440	10.50
1	33.440	10.46								
2	33.437	10.46								
864858	33.440	10.46	864859	093.4	20-1611SR-MX-310-2239-09340264-0000A-7	0	18.744	33.4083	5.805	23.87055
3	33.420	10.45								
4	33.421	10.45								
5	33.431	10.45								
864859	33.440	10.45	864860	093.4	20-1611SR-MX-310-2239-09340264-0002A-3	2	18.744	33.4083	5.805	23.87072
7	33.424	10.24								
8	33.420	10.06								
9	33.494	9.86								
864860	34404	10.05	864861	093.4	20-1611SR-MX-310-2239-09340264-0005A-3	5	18.692	33.4150	5.796	23.88911
864861	34404	10.05	864862	093.4	20-1611SR-MX-310-2239-09340264-0010A-3	10	18.161	33.4062	5.816	24.01426

In [24]:

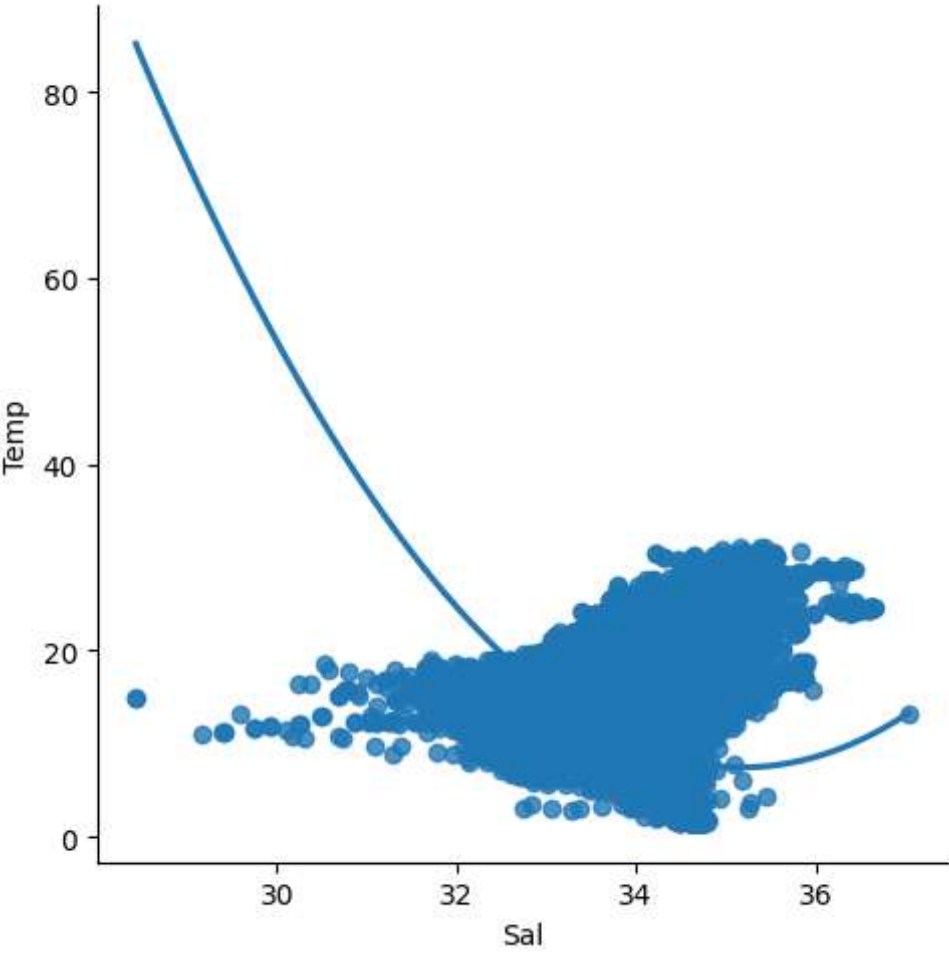
	Cst_Cnt	Btl_Cnt	Sta_ID	Depth_ID	Depthm	T_degC	Salnty	O2ml_L	STheta
--	---------	---------	--------	----------	--------	--------	--------	--------	--------

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

Out[24]:

34404	864863	093.4	MX-310-	15	17.533	33.3880	5.774	24.15297
026.4		2239-						
09340264-								
0015A-3								

<seaborn.axisgrid.FacetGrid at 0x276a76a4580>



In [25]:

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

Out[25]:

	Sal	Temp
count	814247.000000	814247.000000
mean	33.841337	10.860287
std	0.461636	4.224930
min	28.431000	1.440000
25%	33.489000	7.750000
50%	33.866000	10.110000
75%	34.197000	13.930000
max	37.034000	31.140000

In [26]:

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

```
<class 'pandas.core.frame.DataFrame'>
Index: 814247 entries, 0 to 864862
Data columns (total 2 columns):
#   Column   Non-Null Count  Dtype
---  -
0    Sal      814247 non-null  float64
1    Temp     814247 non-null  float64
dtypes: float64(2)
memory usage: 18.6 MB
```

In [27]:

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

Out[27]:

	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

814247 rows × 2 columns

In [29]:

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

In [30]:

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

In [31]:

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

C:\Users\91955\AppData\Local\Temp\ipykernel_2828\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)
```

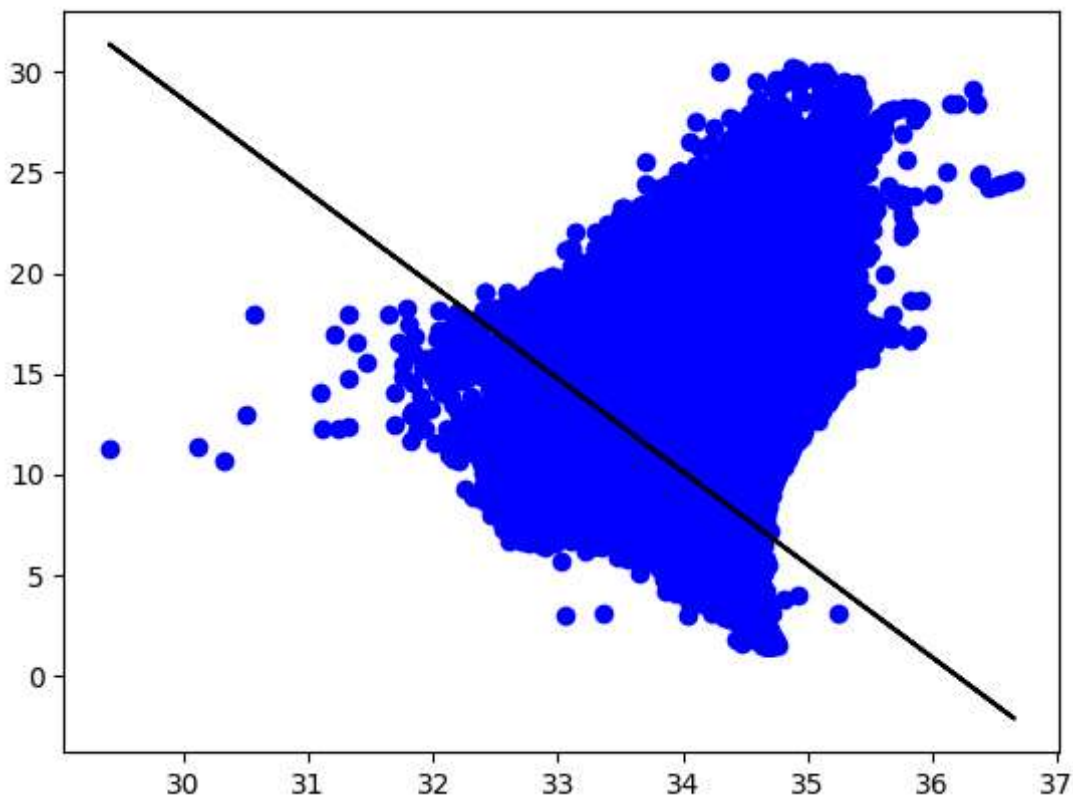
In [41]:

```
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))
```

0.2544304702117517

In [43]:

```
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()
```

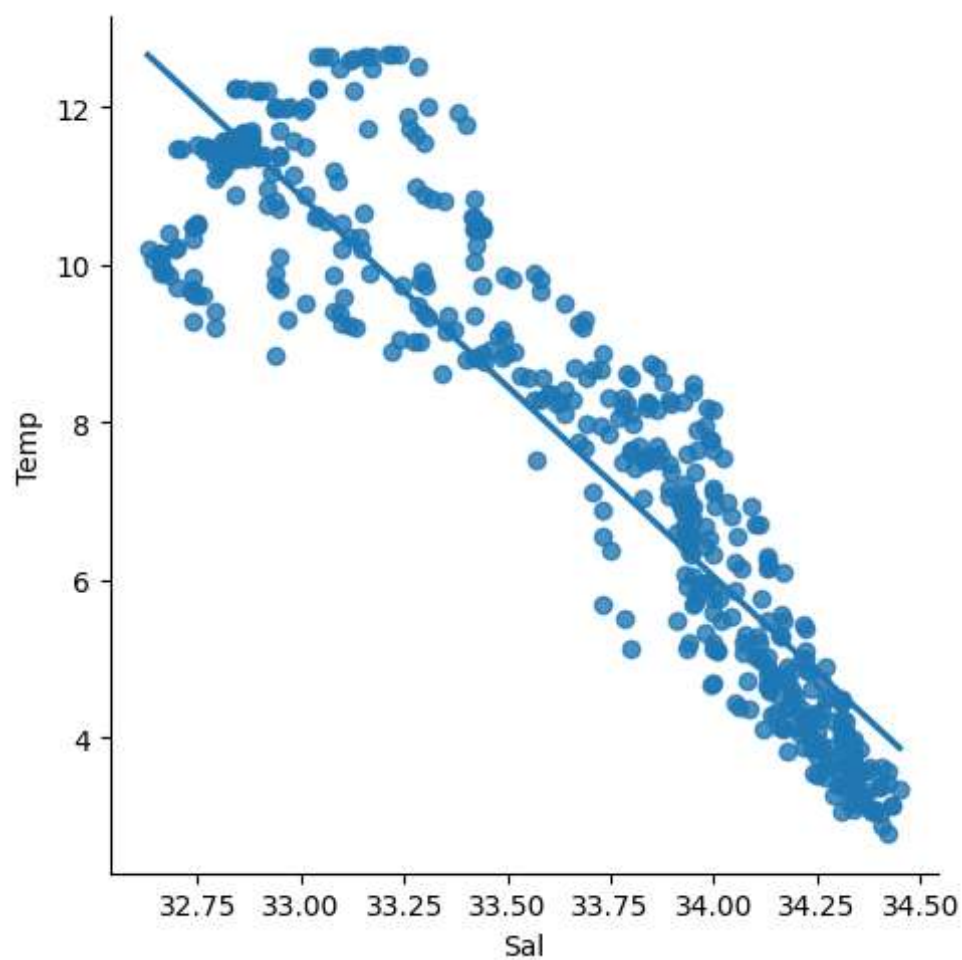


In [44]:

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

Out[44]:

<seaborn.axisgrid.FacetGrid at 0x276a7840cd0>



In [46]:

```

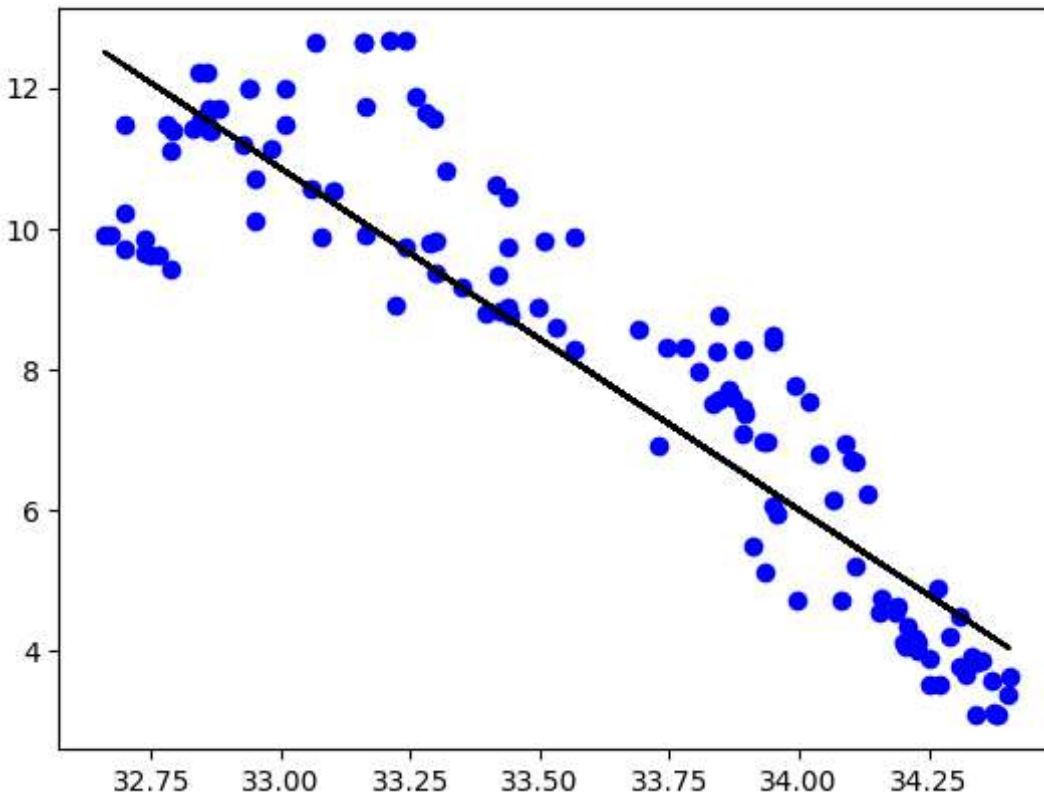
1 df500.fillna(method='ffill',inplace=True)
2 X=np.array(df500['Sal']).reshape(-1,1)
3 y=np.array(df500['Temp']).reshape(-1,1)
4 df500.dropna(inplace=True)
5 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.25)
6 reg=LinearRegression()
7 reg.fit(X_train,y_train)
8 print("Regression:",reg.score(X_test,y_test))
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: 0.8180243753722324

Out[46]:

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



In [48]:

```

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: 0.8180243753722324

In []:

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