

```
In [1]: 1 #import all the libraries
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 from sklearn import preprocessing,svm
7 from sklearn.model_selection import train_test_split
8 from sklearn.linear_model import LinearRegression
```

```
In [2]: 1 #Reading the files
2 df=pd.read_csv(r"C:\Users\HP\Downloads\bottle.csv.zip")
3 df
```

C:\Users\HP\AppData\Local\Temp\ipykernel_20064\1808516658.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\HP\Downloads\bottle.csv.zip")

Out[2]:

	Cst_Cnt	Btl_Cnt	Sta_ID	Depth_ID	Depthm	T_degC	Salnty	O2ml_L	STheta	O2Sat	...	R_PHAEO	R_PRE	R_SAMP	DIC1	DIC2
0	1	1	054.0 056.0	19-4903CR-HY-060-0930-05400560-0000A-3	0	10.500	33.4400	NaN	25.64900	NaN	...	NaN	0	NaN	NaN	NaN
1	1	2	054.0 056.0	19-4903CR-HY-060-0930-05400560-0008A-3	8	10.460	33.4400	NaN	25.65600	NaN	...	NaN	8	NaN	NaN	NaN
2	1	3	054.0 056.0	19-4903CR-HY-060-0930-05400560-0010A-7	10	10.460	33.4370	NaN	25.65400	NaN	...	NaN	10	NaN	NaN	NaN
3	1	4	054.0 056.0	19-4903CR-HY-060-0930-05400560-0019A-3	19	10.450	33.4200	NaN	25.64300	NaN	...	NaN	19	NaN	NaN	NaN
4	1	5	054.0 056.0	19-4903CR-HY-060-0930-05400560-0020A-7	20	10.450	33.4210	NaN	25.64300	NaN	...	NaN	20	NaN	NaN	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	108.74	...	0.18	0	NaN	NaN	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	108.74	...	0.18	2	4.0	NaN	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	108.46	...	0.18	5	3.0	NaN	NaN
864861	34404	864862	093.4 026.4	20-1611SR-MX-310-2239-09340264-0010A-3	10	18.161	33.4062	5.816	24.01426	107.74	...	0.31	10	2.0	NaN	NaN
864862	34404	864863	093.4 026.4	20-1611SR-MX-310-2239-09340264-0015A-3	15	17.533	33.3880	5.774	24.15297	105.66	...	0.61	15	1.0	NaN	NaN

864863 rows × 74 columns

```
In [3]: 1 df=df[['Salnty','T_degC']]
        2 df.columns=['Sal','Temp']
```

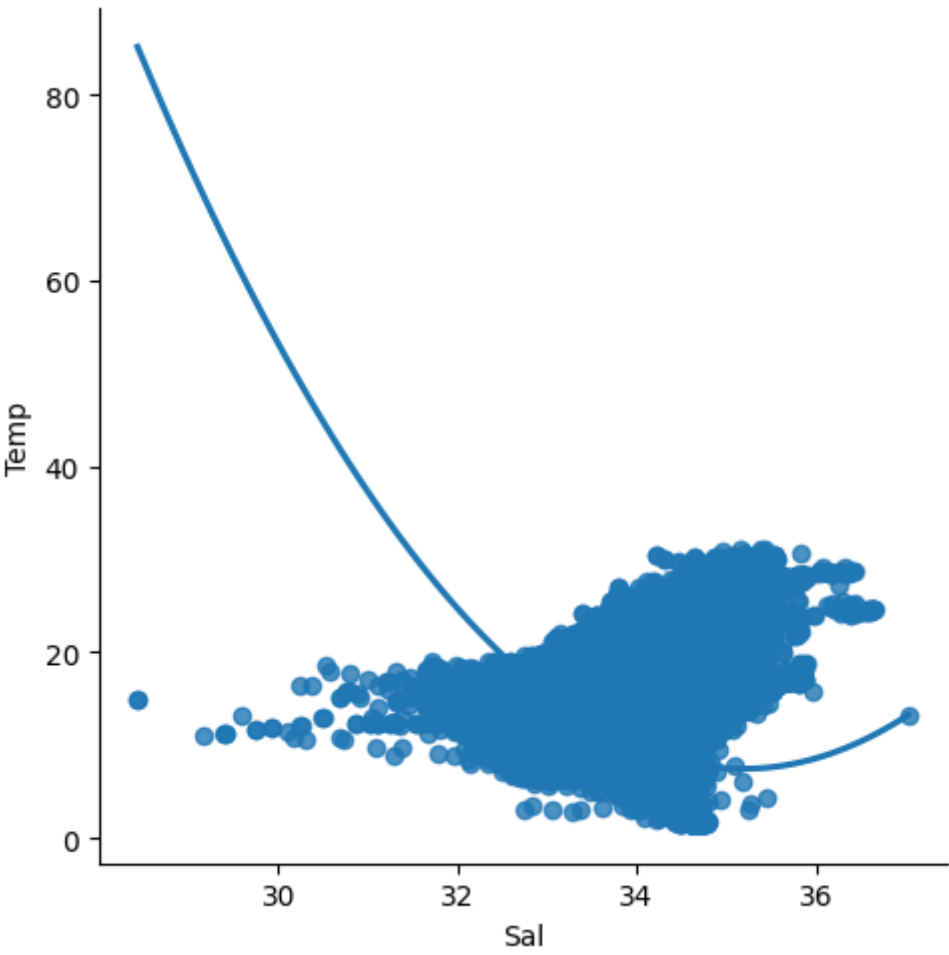
```
In [4]: 1 df.head(10)
```

Out[4]:

	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

```
In [5]: 1 #Exploring the data scatter_plotting the data scatter
        2 sns.lmplot(x = "Sal", y = "Temp", data = df, order = 2, ci = None)
```

Out[5]: <seaborn.axisgrid.FacetGrid at 0x2875f78e260>



```
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 *#data cleaning eliminating Nan or misssing input numbers*
2 df.fillna(method='ffill',inplace=True)

C:\Users\HP\AppData\Local\Temp\ipykernel_20064\3955359400.py:2: 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(method='ffill',inplace=True)
```

In [9]: 1 x=np.array(df['Sal']).reshape(-1,1)
2 y=np.array(df['Temp']).reshape(-1,1)

In [10]: 1 df.dropna(inplace=True)

C:\Users\HP\AppData\Local\Temp\ipykernel_20064\1379821321.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

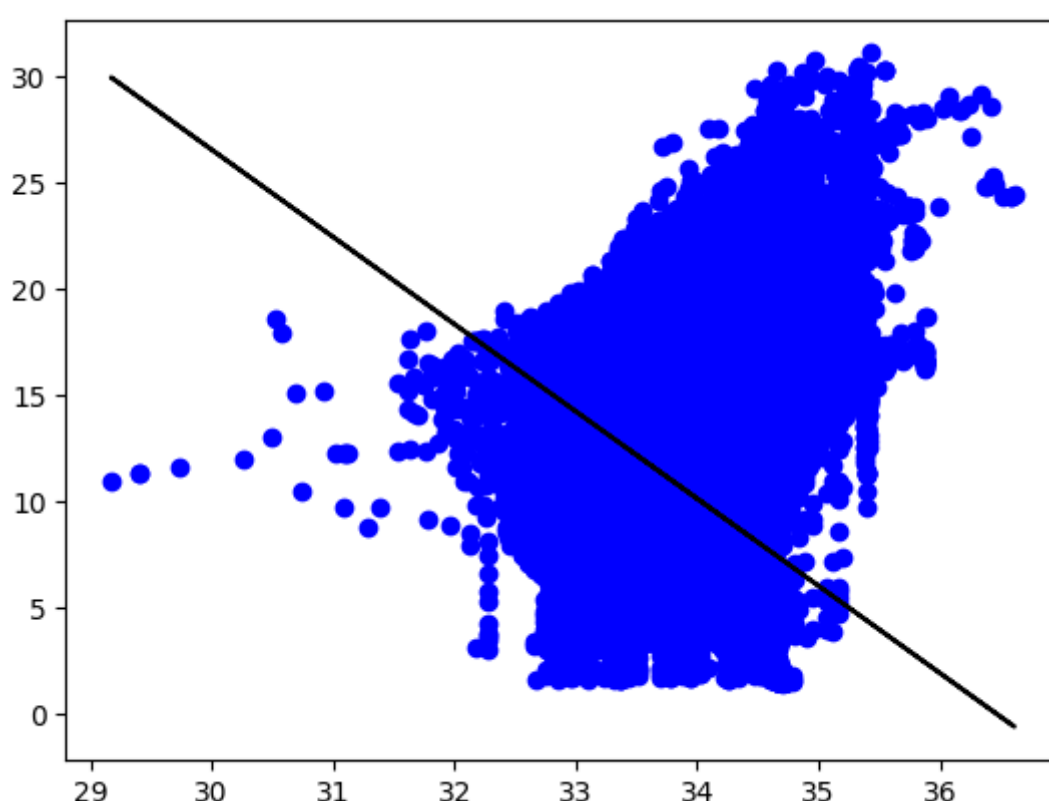
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

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

In [11]: 1 x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.25)
2 regr = LinearRegression()
3 regr.fit(x_train,y_train)
4 print(regr.score(x_test,y_test))

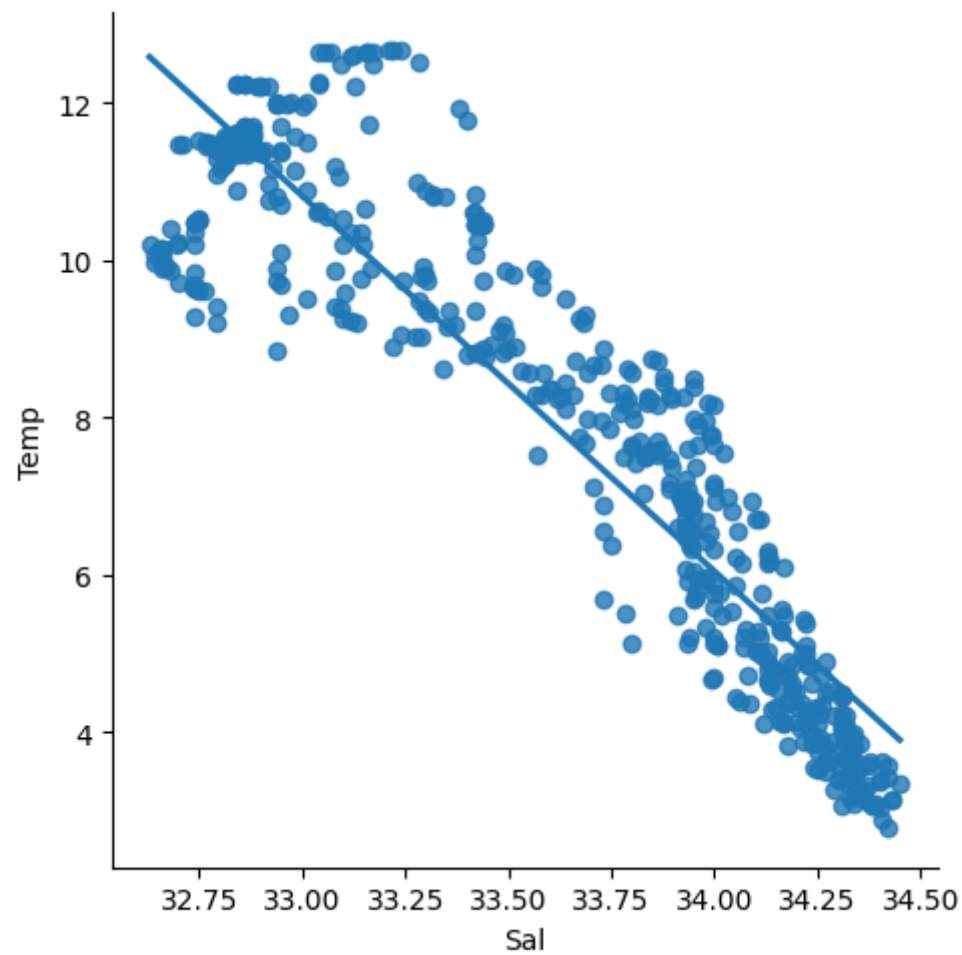
0.20231303641684095

In [12]: 1 *#Exploring our results*
2 y_pred = regr.predict(x_test)
3 plt.scatter(x_test,y_test,color='b')
4 plt.plot(x_test,y_pred,color='k')
5 plt.show()



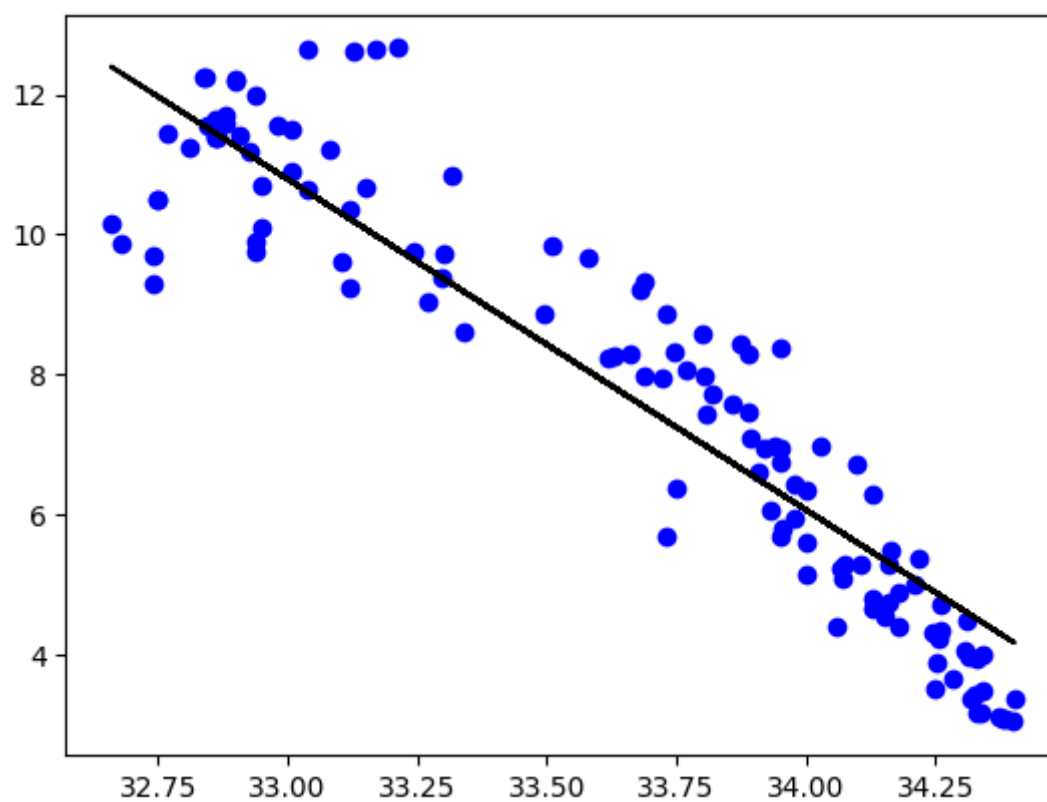
```
In [13]: 1 #Working with a smaller Dataset
2 df500 = df[:][:500]
3 sns.lmplot(x = 'Sal', y = 'Temp', data = df500, order=1, ci=None)
```

Out[13]: <seaborn.axisgrid.FacetGrid at 0x28760a982b0>



```
In [14]: 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 regr=LinearRegression()
7 regr.fit(X_train,y_train)
8 print("Regression: ",regr.score(x_test,y_test))
9 y_pred=regr.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: 0.873118880343903



In [15]:

```
1 #Evaluatevating a model
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import r2_score
4 model=LinearRegression()
5 model.fit(X_train,y_train)
6 y_pred=model.predict(x_test)
7 r2=r2_score(y_test,y_pred)
8 print("R2_Score: ",r2)
```

R2_Score: 0.873118880343903

Conclusion:-

Dataset we have taken is poor for linear Model but with smaller data works well with linear model

In []:

1