

```
In [1]: 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
```

```
In [2]: dt=pd.read_csv(r"C:\Users\91903\Downloads\bottle.csv")
dt
```

C:\Users\91903\AppData\Local\Temp\ipykernel_25376\3720528792.py:1: DtypeWarning: Columns (47,73) have mixed types. Specify dtype option on import or set low_memory=False.

```
dt=pd.read_csv(r"C:\Users\91903\Downloads\bottle.csv")
```

Out[2]:

	Cst_Cnt	Btl_Cnt	Sta_ID	Depth_ID	Depthm	T_degC	Salnty	O2ml_L	STheta	O2S
0	1	1	054.0 056.0	19- 4903CR- HY-060- 0930- 05400560- 0000A-3	0	10.500	33.4400	NaN	25.64900	Ni
1	1	2	054.0 056.0	19- 4903CR- HY-060- 0930- 05400560- 0008A-3	8	10.460	33.4400	NaN	25.65600	Ni
2	1	3	054.0 056.0	19- 4903CR- HY-060- 0930- 05400560- 0010A-7	10	10.460	33.4370	NaN	25.65400	Ni
3	1	4	054.0 056.0	19- 4903CR- HY-060- 0930- 05400560- 0019A-3	19	10.450	33.4200	NaN	25.64300	Ni
4	1	5	054.0 056.0	19- 4903CR- HY-060- 0930- 05400560- 0020A-7	20	10.450	33.4210	NaN	25.64300	Ni
...
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.
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.
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.
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.

	Cst_Cnt	Btl_Cnt	Sta_ID	Depth_ID	Depthm	T_degC	Salnty	O2ml_L	STheta	O2S
864862	34404	864863	093.4026.4	20-1611SR-MX-310-2239-09340264-0015A-3	15	17.533	33.3880	5.774	24.15297	105.

864863 rows × 74 columns

```
In [3]: dt=dt[['Salnty','T_degC']]
dt.columns=['Sal','Temp']
```

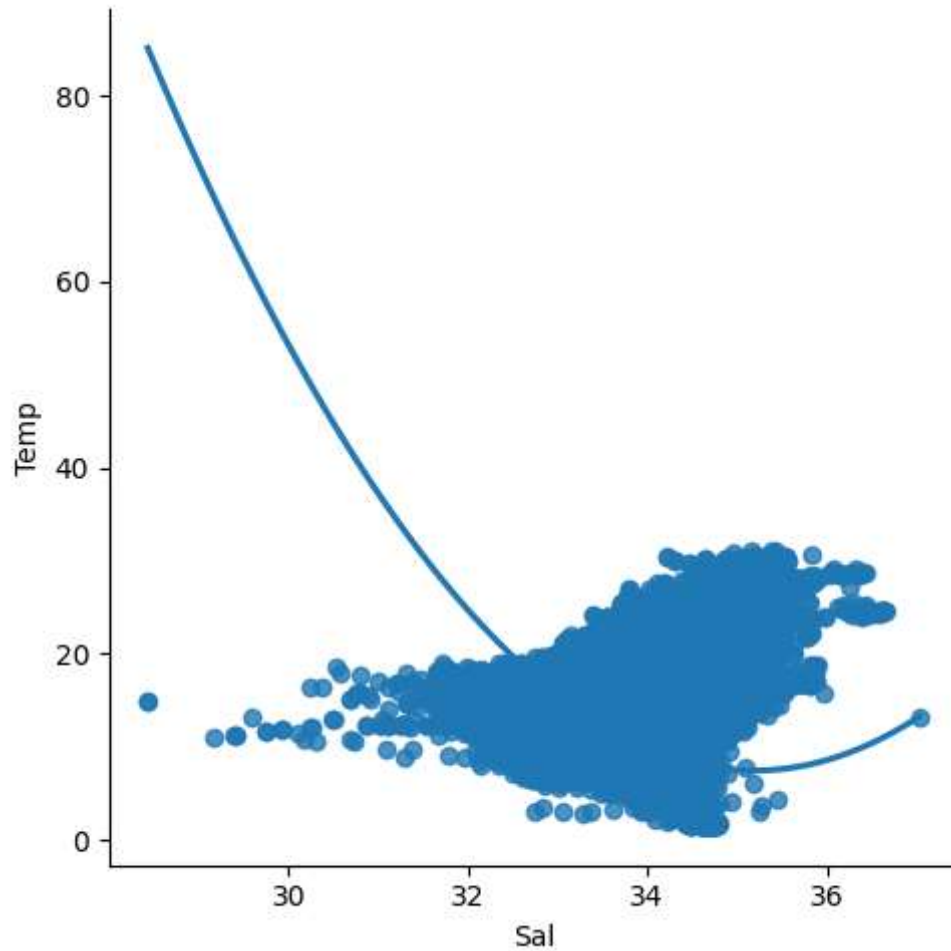
```
In [4]: dt.head(10)
```

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 [16]: sns.lmplot(x='Sal',y='Temp',data=dt,order=2,ci=None)
```

```
Out[16]: <seaborn.axisgrid.FacetGrid at 0x2915c8a4970>
```



```
In [17]: dt.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 [18]: dt.describe()
```

```
Out[18]:
```

	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 [19]: dt.fillna(method='ffill')
```

```
Out[19]:
```

	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 [20]: x=np.array(dt['Sal']).reshape(-1,1)
          y=np.array(dt['Temp']).reshape(-1,1)
```

```
In [21]: dt.dropna(inplace=True)
```

C:\Users\91903\AppData\Local\Temp\ipykernel_25376\735218168.py:1: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame

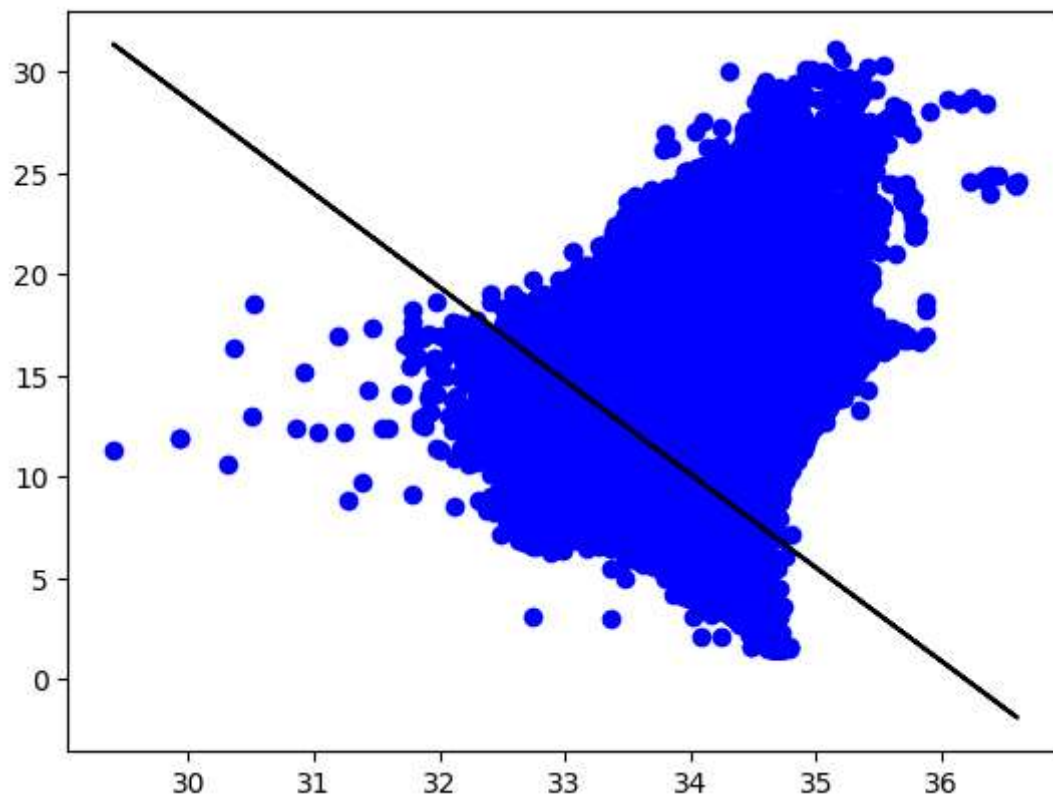
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

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

```
In [22]: X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.25)
reg=LinearRegression()
reg.fit(X_train,y_train)
print(reg.score(X_test,y_test))
```

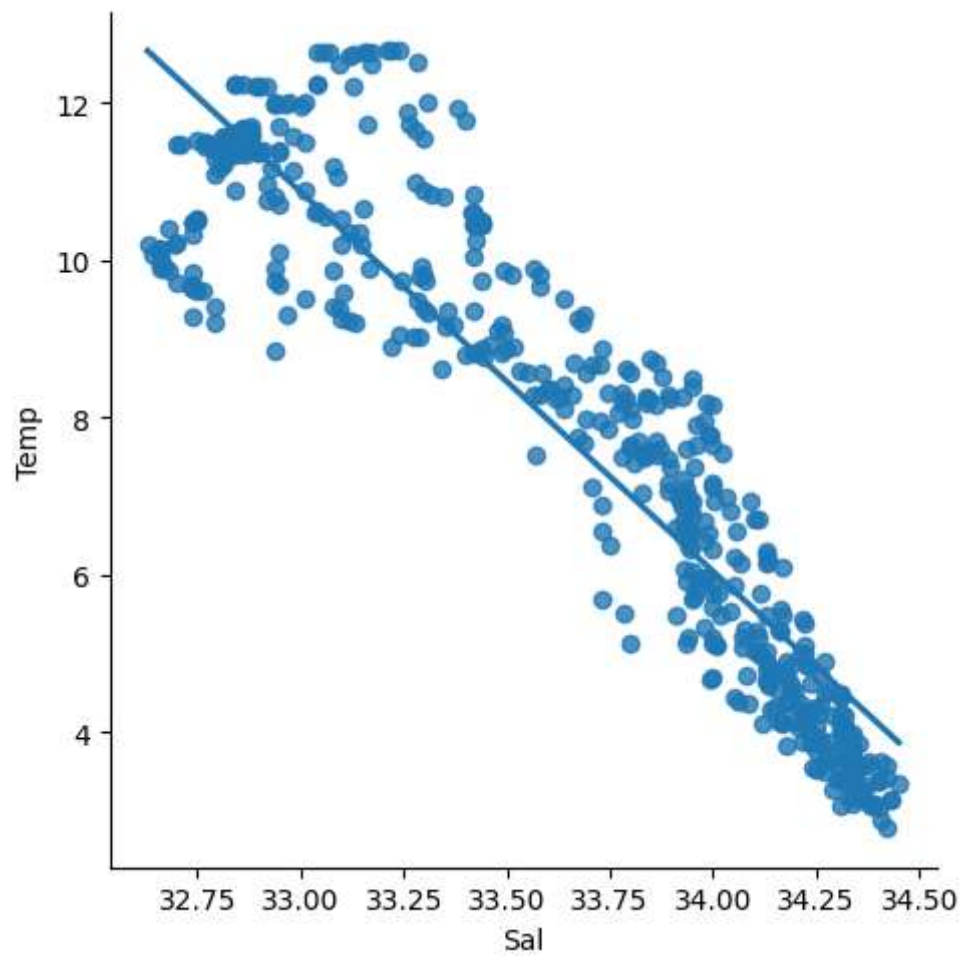
0.25894156633387044

```
In [23]: y_pred=reg.predict(X_test)
plt.scatter(X_test,y_test,color='b')
plt.plot(X_test,y_pred,color='k')
plt.show()
```



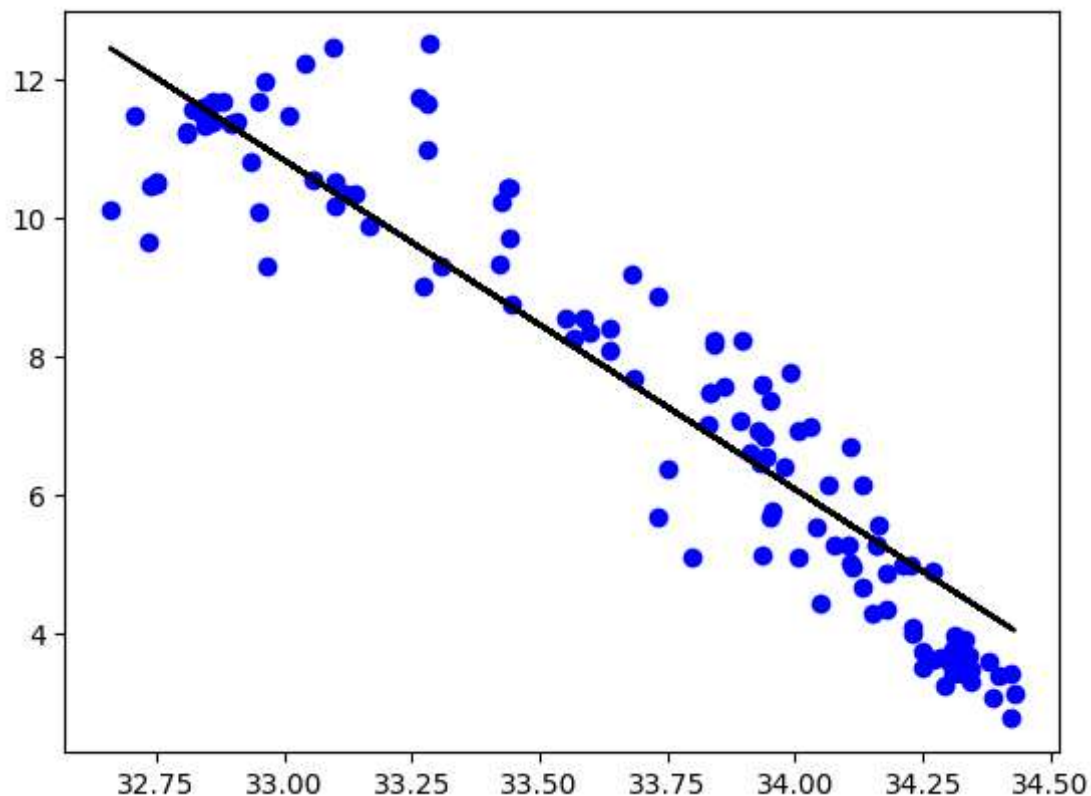
```
In [24]: dt500=dt[:][:500]  
sns.lmplot(x="Sal",y="Temp",data=dt500,order=1,ci=None)
```

```
Out[24]: <seaborn.axisgrid.FacetGrid at 0x2913035d810>
```




```
In [27]: dt500.fillna(method='ffill',inplace=True)
X=np.array(dt500['Sal']).reshape(-1,1)
y=np.array(dt500['Temp']).reshape(-1,1)
dt500.dropna(inplace=True)
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.25)
reg=LinearRegression()
reg.fit(X_train,y_train)
print("Regression:",reg.score(X_test,y_test))
y_pred=reg.predict(X_test)
plt.scatter(X_test,y_test,color="b")
plt.plot(X_test,y_pred,color='k')
plt.show()
```

Regression: 0.8846089899631949



```
In [30]: 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: 0.8846089899631949

#conclusion : Linear regression is best fit for the model

In []: