**~\Desktop\class\ml\data processing pipeline\data processing pipeline.py**

```python
1   import numpy as np
2   import matplotlib.pyplot as plt
3   import pandas as pd
4
5   data=pd.read_csv(r'C:\Users\Admin\Desktop\class\ml\simple linear regression
    pipeline\Salary_Data.csv')
6
7   x=data.iloc[:,:-1]
8   y=data.iloc[:,-1]
9
10  from sklearn.model_selection import train_test_split
11  x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.20,random_state=0)
12
13  x_train=x_train.values.reshape(-1,1) #to convert values into array
14  x_test=x_test.values.reshape(-1,1)
15
16  from sklearn.linear_model import LinearRegression
17
18  regressor=LinearRegression()
19  regressor.fit(x_train,y_train)
20
21  y_predict=regressor.predict(x_test)
22
23  plt.scatter(x_test, y_test, color='red')
24  plt.plot(x_train,regressor.predict(x_train))
25  plt.title('salary vs experience (test set)')
26  plt.xlabel('years of experience')
27  plt.ylabel('salary')
28
29  m_slope=regressor.coef_ #for slope (m)
30  print(m_slope)
31
32  c_intercept=regressor.intercept_ #for constant (c)
33  print(c_intercept)
34
35  y_15=m_slope*15+c_intercept #y^
36  print(y_15)
37
38  comparsion=pd.DataFrame({'actual':y_test,'predicted':y_predict})
39  print(comparsion)
40
41  data.mean()
42
43  data.std()
44
45  data['Salary'].mean()
46
47  data.median()
48
49
50  data['Salary'].median()
51
```

```python
52  data.mode()
53
54  data['Salary'].mode()
55
56  data.var()
57
58  data['Salary'].var()
59
60  from scipy.stats import variation #coff variation
61
62  variation(data.values)
63
64  variation(data['Salary'])
65
66  data.corr()
67
68  data['Salary'].corr(data['YearsExperience'])
69
70  data.skew()
71
72  data['Salary'].skew()
73
74  import scipy.stats as stats
75
76  data.apply(stats.zscore)
77
78  y_mean=np.mean(y)
79  stats.zscore(data['Salary'])
80  ssr=np.sum((y_predict-y_mean)**2)
81  print(ssr)
82
83  y=y[0:6]
84  sse=np.sum((y-y_predict)**2)
85  print(sse)
86
87
88  mean_total=np.mean(data.values)
89  sst=np.sum((data.values-mean_total)**2)
90  print(sst)
91
92
93
94  rsquare=1-(ssr/sst)
95  print(rsquare)
96
97  import pickle
98  filename = 'linear_regression_model.pkl'
99  with open(filename, 'wb') as file:
100     pickle.dump(regressor, file)
101 print("Model has been pickled and saved as linear_regression_model.pkl")
102
103 import os
104 print(os.getcwd())
105
```

```
106
107
```