

Analyze Salary with Linear Regression

Python Automation; Scott Schmidt; Illinois State University

The main purpose of this assignment is to do linear regression without the use of imports so that one gains a full understanding of statistics and Python.

Section Summaries

1. The dataset will have two columns of data which are years and salary.
2. Split the data using 75% training and 25% testing.
3. Calculate Covariance and Variance to find the coefficients.
4. Evaluate the model performance using Root Mean Square Error (RMSE).
5. R_Squared using no imports
6. Full Linear Regression process with sklearn imports
7. Plot Graph data with line of best fit

Math Functions

Three functions called `sumList()`, `meanList()`, and `stdevList()` will be used to calculate the mean and standard deviation of any given list. Each function has a list as an input, and returns the corresponding result (i.e., sum, mean, min, max, stdev, etc.). Please note: a) `meanList()` should use `sumList()` to calculate the summation of a list; b) `stdevList()` should use `meanList()` to calculate the mean value; c) `minMaxList()` should return both min and max values.

```
In [1]: 1  #Return the total:
2  def sumList(numbers):
3      count=0
4      numbers=list(numbers)
5      for n in numbers:
6          count=count+float(n)
7      total=round(count,2)
8      return total
9
10
11 #Return the mean:
12 def meanList(numbers):
13     numbers=list(numbers)
14     total=sumList(numbers)
15     length=len(numbers)
16     mean=total/length
17     return round(mean,2)
18
19 #Return the standard deviation:
20 def stdevList(mean):
21     var=0 #variance
22     for d in data:
23         n=d-mean
24         n2=n*n
25         var=var+n2
26     std=var/len(data)
27     std=round(std**0.5,2)
28     return std
```

1)View Data

Read the csv file and display the data below. There is no missing data.

```
In [2]: 1 import numpy as np # Linear algebra
2 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
3 from sklearn.model_selection import train_test_split, cross_validate
4 from sklearn.metrics import mean_absolute_error, mean_squared_error
5
6 # sFile=r'/kaggle/input/salary/salary.csv'
7 sFile=r'C:\Users\sschm\Desktop\IT170\salary.csv'
8 sDF = pd.read_csv(sFile, header=None, names=['years', 'salary'])
9 sDF
```

Out[2]:

	years	salary
0	1.1	39343
1	1.3	46205
2	1.5	37731
3	2.0	43525
4	2.2	39891
5	2.9	56642
6	3.0	60150
7	3.2	54445
8	3.2	64445
9	3.7	57189
10	3.9	63218
11	4.0	55794
12	4.0	56957
13	4.1	57081
14	4.5	61111
15	4.9	67938
16	5.1	66029
17	5.3	83088
18	5.9	81363
19	6.0	93940
20	6.8	91738
21	7.1	98273
22	7.9	101302
23	8.2	113812
24	8.7	109431
25	9.0	105582
26	9.5	116969
27	9.6	112635
28	10.3	122391

	years	salary
29	10.5	121872

2) Split Data

Function will be called `train_test_split` that will split the dataset List into a training set and a test based on 75% for training and 25% for testing. Normally, one can use `from sklearn.model_selection import train_test_split` to do this but in this case the `train_test_split` function will do this manually. A list sent to this function will get split by a ratio as a way to use training data to make future predictions on the test data.

```
In [3]: 1 def train_test_split(aList, ratio):
2         elements = len(aList)
3         middle = int(elements * ratio)
4         trainSet=[aList[:middle]]
5         testSet=[aList[middle:]]
6         return list(trainSet), list(testSet)
```

3)Calculate Variance and Coefficients

Some basic math functions used in this section will be:

```
mean=meanList(data)
stdev=stdevList(mean)
```

Variance Function

Variance is a numeric measurement of how far a dataset is more a mean. In this section we will calculate covariance and variance from the training dataset.

```
In [4]: 1 # Calculate the variance of a list of numbers:
2 def varianceList(aList):
3     # calculate mean
4     m = meanList(aList)
5
6     # calculate variance using a list comprehension
7     var = sum((xi - m) ** 2 for xi in aList) / len(aList)
8     return var
```

Covariance Function

Covariance measures the relationships between two assets. A positive correlation would be found between the Nadsaq and Apple. When the tech sector does well during a duration of time, normally, so does Apple.

```

In [5]: 1 #Calculate the covariance of two lists of numbers:
2 def covarianceList(x, y):
3     # Finding the mean of the series x and y
4     mean_x = meanList(x)
5     mean_y = meanList(y)
6
7     # Subtracting mean from the individual elements
8     sub_x = [i - mean_x for i in x]
9     sub_y = [i - mean_y for i in y]
10
11     numerator = sum([sub_x[i]*sub_y[i] for i in range(len(sub_x))])
12     denominator = len(x)-1
13     cov = numerator/denominator
14     return cov

```

Calculate Coefficients Function

Calculate the w_0 and w_1 values

Hint: Calculate the covariance of two lists of numbers:

The input is a list with each item a list [x, y]

The diagram shows the equation $4x - 7 = 5$. Arrows point from labels to the corresponding parts of the equation: 'Coefficient' points to '4', 'Variable' points to 'x', 'Operator' points to '-', 'Constants' points to '7', and another 'Constants' points to '5'.

!

```

In [6]: 1 def coefficeints(trainingDatasetList):
2     x=list(trainingDatasetList[0])
3     y=list(trainingDatasetList[1])
4
5     w1 = covarianceList(x, y) / float(varianceList(x))
6
7     # Coefficient w0 = mean of y_readings - ( w1 * the mean of the x_reading
8     w0 = meanList(y) - (w1 * meanList(x))
9     return w0, w1

```

```

In [7]: 1 #GET DATA IN FORM OF X AND Y:
2 x=sDF['years']
3 y=sDF['salary']
4 x=list(x)
5 y=list(y)
6
7 #SPLIT DATA:
8 x_train, x_test=train_test_split(x, .25)
9 y_train, y_test=train_test_split(y, .25)
10
11 #Turn double bracket list into single bracket list:
12 x_train, y_train,x_test, y_test=x_train[0], y_train[0],x_test[0], y_test[0]
13
14 #COVARIANCE:
15 cov=round(covarianceList(x_train,y_train),4)
16 print("The covariance is: ", cov)
17
18 #VARIANCE:
19 var=round(varianceList(x_train),4)
20 print("The variance is: ", var)
21
22 #COEFFICIENTS:
23 trainingDatasetList= [x_train, y_train]
24 coef=coefficients(trainingDatasetList)
25 print("The coefficients are: ", coef)
26
27 w0=coef[0]
28 w1=coef[1]
29 print(w0)
30 print(w1)

```

The covariance is: 5414.7167

The variance is: 0.4857

The coefficients are: (23916.537843137256, 11147.946078431372)

23916.537843137256

11147.946078431372

Make Prediction

Based on the coefficients between year and salary, we can make a prediction of what an individual would make based on their years of service. Y is the final prediction and the formula is below:

$$y = w_0 + w_1 \cdot x_1$$

Normally, the code with numpy and linear regression import would look like this:

```

x_new = np.array([0, 1, 2, 3])
y_prediction = lr.intercept_ + x_new*lr.coef_[0]
lr.predict(x_new.reshape(-1,1))

```

w_0 and w_1 are the two coefficients, where w_0 is the intercept (of the y-axis), and w_1 is the slope of the line. w_1 shows the impact of the independent variable x_1 on y . For example, when $w_1 = 0$, there's no impact of x_1 on y since ($0 \cdot x_1 = 0$).

In simple terms, linear regression is an algorithm that finds the best values of w_0 and w_1 to fit the training dataset.

A. Predictions with numpy

Both numpy and the manual loop prediction get the same prediction results.

```
In [8]: 1 x_new=np.array([x_test])
        2 predictions = np.round_(w0 + x_new*w1, decimals=2)
        3 print(predictions)
```

```
[[ 59589.97  59589.97  65163.94  67393.53  68508.32  68508.32  69623.12
   74082.3   78541.47  80771.06  83000.65  89689.42  90804.21  99722.57
  103066.95 111985.31 115329.7   120903.67 124248.05 129822.03 130936.82
  138740.38 140969.97]]
```

B. Predictions using a loop

```
In [9]: 1 predictions=[]
        2 for x in x_test:
        3     predict = round(w0 + x*w1,2)
        4     predictions.append(predict)
        5 print(predictions)
```

```
[59589.97, 59589.97, 65163.94, 67393.53, 68508.32, 68508.32, 69623.12, 74082.3,
78541.47, 80771.06, 83000.65, 89689.42, 90804.21, 99722.57, 103066.95, 111985.3
1, 115329.7, 120903.67, 124248.05, 129822.03, 130936.82, 138740.38, 140969.97]
```

4)Calculate RMSE

Root Mean Squared Error (RMSE) is the standard deviation of the prediction errors. In simple terms, it tells an individual how far the predictions are away from the line of best fit. Here is the official formula below:

$$RMSE = \sqrt{\frac{\sum (f - o)^2}{N}}$$

Σ = summation ("add up")

$(z_{fi} - z_{oi})^2$ = differences, squared

N = sample size.

Specifically, one must do the following to get the RMSE:

1. Squaring the residuals.
2. Finding the average of the residuals.
3. Taking the square root of the result.

Calculate RMSE function

The Root Mean Square Error is 10023.7. Both methods get the same answer to verify the calculation was done correctly. This means that on average, a salary can be predicted on average within 10,023.70 dollars based on the current data.

```
In [10]: 1 def calcRMSE(predictions, targets):
2         difference=[]
3         zip_object = zip(predictions, targets)
4         for predictions, targets in zip_object:
5             difference.append(predictions-targets)
6         final=[]
7
8         for d in difference:
9             sq=d*d
10            result=sq**.5
11            final.append(result)
12        return meanList(final)
13
14 RMSE=calcRMSE(predictions, y_test)
15 print("Root Mean Square Error:", RMSE)
```

Root Mean Square Error: 10023.7

Calulcate RMSE with numpy

```
In [11]: 1 def evalRMSE(predictions, targets):
2         diff=np.subtract(predictions,targets)
3         square=np.square(diff)
4         MSE=square.mean()
5         RMSE=np.sqrt(MSE)
6         RMSE=np.round_(RMSE, decimals=2)
7         return RMSE
8
9 rmse=evalRMSE(predictions, y_test)
10 print("Root Mean Square Error:", rmse)
```

Root Mean Square Error: 10023.7

Calculate r_squared no imports

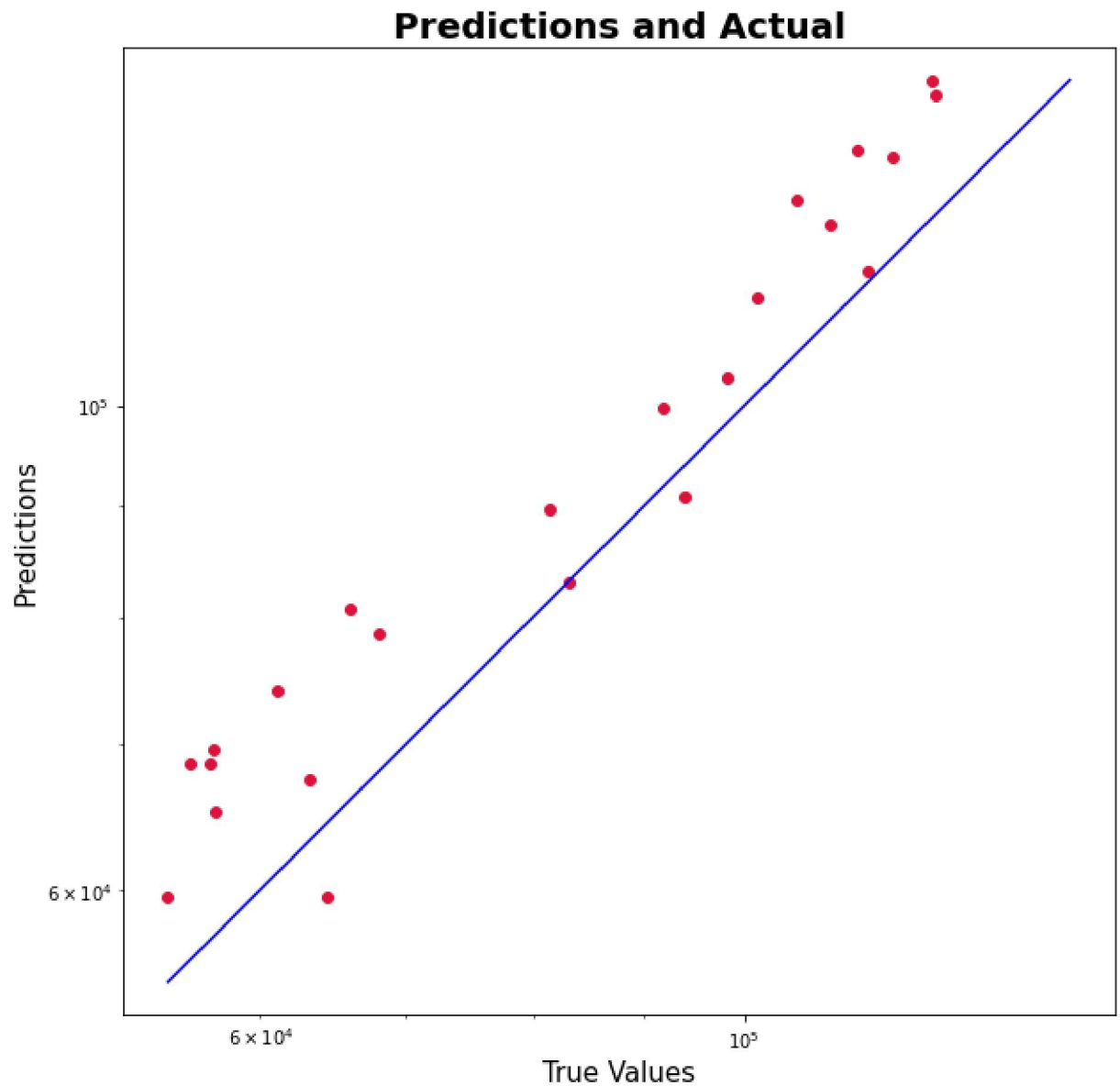
To make sure that the assignment is done correctly, I made a function to calculate r_squared.


```
In [12]: 1 def calcR2(predicted, actual):
2         diff=np.subtract(predicted,actual)
3         ssres=sum(np.square(diff))
4         sstot=sum(np.square(actual-np.mean(actual)))
5
6         #ssres = sum((actual - predicted)**2)
7         #sstot = sum((actual-np.mean(actual))**2)
8         r2_m = 1-(ssres/sstot)
9         return r2_m
10
11 r2=calcR2(predictions, y_test)
12 print("R_Squared: ", round(r2,4))
```

R_Squared: 0.7749

Plot Predictions vs Actual

```
In [13]: 1 import matplotlib.pyplot as plt
2 plt.figure(figsize=(10,10))
3 plt.scatter(y_test, predictions, c='crimson')
4 plt.yscale('log')
5 plt.xscale('log')
6
7 p1 = max(max(predictions), max(y_test))
8 p2 = min(min(predictions), min(y_test))
9
10 plt.title("Predictions and Actual", size=20,fontweight="bold")
11 plt.plot([p1, p2], [p1, p2], 'b-')
12 plt.xlabel('True Values', fontsize=15)
13 plt.ylabel('Predictions', fontsize=15)
14 plt.axis('equal')
15 plt.show()
```



Linear Regression with imports

Extra section to verify that the linear regression prediction and error are correct by using the official imports from Sklearn. One of the most basic and used types of predictive modeling is linear regression. Linear regression can help find the strength and effect independent variables have on dependent variables. I will verify the results using sklearn imports, instead of doing it by hand.

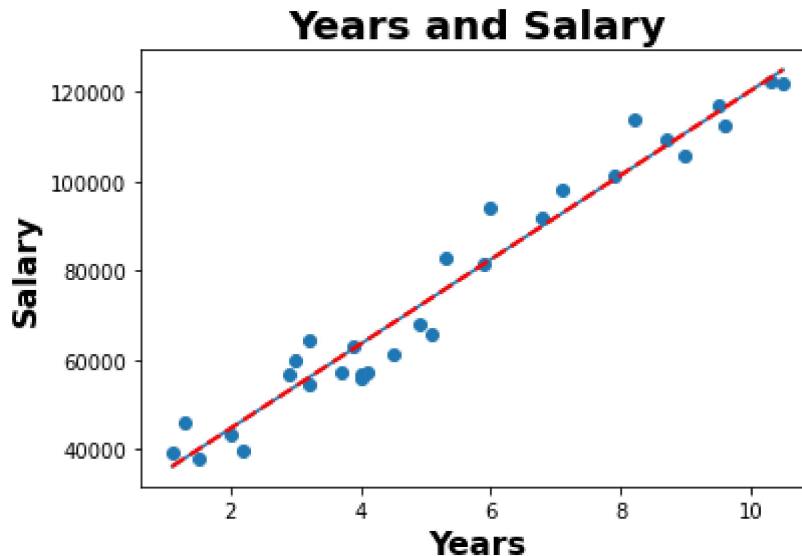
```
In [14]: 1 import pandas as pd
2 from sklearn.linear_model import LinearRegression
3 from sklearn.model_selection import train_test_split
4 import sklearn.metrics as metrics
5 from sklearn.metrics import r2_score, mean_squared_error
6
7 #Get data as x and y:
8 x=sDF[['years']]
9 x=sDF[['salary']]
10 X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.25, ra
11
12 #Fit and predict:
13 lrModel = LinearRegression()
14 lrModel.fit(X_train, y_train)
15 lrPredict = lrModel.predict(X_test)
16 print(lrPredict)
17
18 #Linear Metrics:
19 r2 = r2_score(y_test, lrPredict)
20 print("Linear regression r2 score: ", r2.round(4))
21
22 #PRINT COEFFICIENT:
23 lrCoef=lrModel.coef_
24 print("Linear coef: ", lrCoef)
25
26 #PRINT MEAN SQUARED ERROR:
27 lrRMSE=np.sqrt(metrics.mean_squared_error(y_test, lrPredict))
28 print("Linear regression MSE: ", round(lrRMSE,4))
```

```
[101302.  67938. 116969.  46205. 105582.  64445.  54445. 121872.]
Linear regression r2 score:  1.0
Linear coef:  [1.]
Linear regression MSE:  0.0
```

Plot Data

The graph clearly shows a trend that salary and years are correlated. Meaning, the more years experience, the salary tends to be higher.

```
In [15]: 1 import matplotlib.pyplot as plt
2 x=sDF['years']
3 y=sDF['salary']
4
5 #find line of best fit
6 a, b = np.polyfit(x, y, 1)
7
8 #add points to plot
9 plt.scatter(x, y)
10
11 #add line of best fit to plot
12 plt.plot(x, a*x+b)
13
14 #add line of best fit to plot
15 plt.plot(x, a*x+b, color='red', linestyle='--', linewidth=2)
16 plt.title("Years and Salary", size=20,fontweight="bold")
17 plt.xlabel('Years',size=16, fontweight="bold")
18 plt.ylabel('Salary',size=16,fontweight="bold")
19 plt.show()
```



References

1. <https://www.urbanpro.com/machine-learning/linear-regression-without-any-libraries>
(<https://www.urbanpro.com/machine-learning/linear-regression-without-any-libraries>)
2. <https://dataaspirant.com/simple-linear-regression-python-without-any-machine-learning-libraries/> (<https://dataaspirant.com/simple-linear-regression-python-without-any-machine-learning-libraries/>)

[learning-libraries/](#)).