

Multiple linear Regression

Multiple linear Regression means you have more than one feature column in regression problem.

Guideline for Multiple linear Regression:

- 1- Check for Missing values using
- 2- Check for Outliers - Only on Feature Columns if performing Inferential Stats else perform for all columns
(Never check outliers for Label Column if performing Supervised Learning)
- 3- Inferential Stats
 - Ensure your data is Complete
 - Ensure your data is Strictly Numeric
- 4- Separate data as features and label and store the same in the form of Numpy Array(Inferential Stats)
- 5- Split our data as Training Set and Testing Set where Training Set is used for training/convergence purpose(Inferential Stats)
 - (whereas Testing set is used for Quality Check Purpose)
 - Decide the Split Ratio (65%:35%)(80:20)
 - Implement the same using Sci-kit Learn package
- 6- Build the model:
 - Regression Algo -----> Linear Regression
 - $y = mx + b$ (Slope Intercept Formula)
 - $\text{profit} = b_0(\text{California}) + b_1(\text{Florida}) + b_2(\text{New York}) + b_3(\text{R\&D Spend}) + b_4(\text{AdminSpend}) + b_5(\text{MarketingSpend}) + \text{intercept}$
 - Goal of this algo is to derive values of $b_0, b_1, b_2, b_3, b_4, b_5, \text{intercept}$ based on the historical data !
- 7- Train the model | Converging Training Set to the Algorithm
- 8- $\text{fit}(\text{features}, \text{label})$
- 9- **Quality Check (Guideline)**
 - 1- Ensure your model that is converged is a Generalized Model
 - **Generalized model is a model that performs well with Known and Unknown data**
 - Technique: **$\text{Accuracy}(\text{Test Data}) > \text{Accuracy}(\text{Train Data})$** -----> Model is **Generalized Model**
 - 2- Ensure the Model's accuracy score must be greater than or equal to CL score (*CL = 1 - SL)
 - Technique: **$\text{Accuracy}(\text{Test Data}) \geq \text{CL}$**
- 10- Deployment Test

Use-case: To create a model that can **predict the profit of the company** based on **company's location** and **company's spending pattern**.

		<pre>data = pd.read_csv('50_Startups.csv') data.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 50 entries, 0 to 49 Data columns (total 5 columns): # Column Non-Null Count Dtype --- - - - - - - - - - - - - - - - - 0 R&D Spend 50 non-null float64 1 Administration 50 non-null float64 2 Marketing Spend 50 non-null float64 3 State 50 non-null object 4 Profit 50 non-null float64 dtypes: float64(4), object(1) memory usage: 2.1+ KB</pre>																																													
		<pre>data.State.unique() array(['New York', 'California', 'Florida'], dtype=object)</pre>																																													
Check for Missing values using		<pre>data.isna().sum() R&D Spend 0 Administration 0 Marketing Spend 0 State 0 Profit 0 dtype: int64</pre>																																													
Check for Outliers - Only on Feature Columns if performing Inferential Stats else perform for all columns Never check outliers for Label Column if performing Supervised Learning		<pre>data.describe()</pre> <table><thead><tr><th></th><th>R&D Spend</th><th>Administration</th><th>Marketing Spend</th><th>Profit</th></tr></thead><tbody><tr><td>count</td><td>50.000000</td><td>50.000000</td><td>50.000000</td><td>50.000000</td></tr><tr><td>mean</td><td>73721.615600</td><td>121344.639600</td><td>211025.097800</td><td>112012.639200</td></tr><tr><td>std</td><td>45902.256482</td><td>28017.802755</td><td>122290.310726</td><td>40306.180338</td></tr><tr><td>min</td><td>0.000000</td><td>51283.140000</td><td>0.000000</td><td>14681.400000</td></tr><tr><td>25%</td><td>39936.370000</td><td>103730.875000</td><td>129300.132500</td><td>90138.902500</td></tr><tr><td>50%</td><td>73051.080000</td><td>122699.795000</td><td>212716.240000</td><td>107978.190000</td></tr><tr><td>75%</td><td>101602.800000</td><td>144842.180000</td><td>299469.085000</td><td>139765.977500</td></tr><tr><td>max</td><td>165349.200000</td><td>182645.560000</td><td>471784.100000</td><td>192261.830000</td></tr></tbody></table>		R&D Spend	Administration	Marketing Spend	Profit	count	50.000000	50.000000	50.000000	50.000000	mean	73721.615600	121344.639600	211025.097800	112012.639200	std	45902.256482	28017.802755	122290.310726	40306.180338	min	0.000000	51283.140000	0.000000	14681.400000	25%	39936.370000	103730.875000	129300.132500	90138.902500	50%	73051.080000	122699.795000	212716.240000	107978.190000	75%	101602.800000	144842.180000	299469.085000	139765.977500	max	165349.200000	182645.560000	471784.100000	192261.830000
	R&D Spend	Administration	Marketing Spend	Profit																																											
count	50.000000	50.000000	50.000000	50.000000																																											
mean	73721.615600	121344.639600	211025.097800	112012.639200																																											
std	45902.256482	28017.802755	122290.310726	40306.180338																																											
min	0.000000	51283.140000	0.000000	14681.400000																																											
25%	39936.370000	103730.875000	129300.132500	90138.902500																																											
50%	73051.080000	122699.795000	212716.240000	107978.190000																																											
75%	101602.800000	144842.180000	299469.085000	139765.977500																																											
max	165349.200000	182645.560000	471784.100000	192261.830000																																											
Inferential Stats 1. Ensure your data is Complete 2. Ensure your data is Strictly Numeric		<pre>data.head(5)</pre> <table><thead><tr><th></th><th>R&D Spend</th><th>Administration</th><th>Marketing Spend</th><th>State</th><th>Profit</th></tr></thead><tbody><tr><td>0</td><td>165349.20</td><td>136897.80</td><td>471784.10</td><td>New York</td><td>192261.83</td></tr><tr><td>1</td><td>162597.70</td><td>151377.59</td><td>443898.53</td><td>California</td><td>191792.06</td></tr><tr><td>2</td><td>153441.51</td><td>101145.55</td><td>407934.54</td><td>Florida</td><td>191050.39</td></tr><tr><td>3</td><td>144372.41</td><td>118671.85</td><td>383199.62</td><td>New York</td><td>182901.99</td></tr><tr><td>4</td><td>142107.34</td><td>91391.77</td><td>366168.42</td><td>Florida</td><td>166187.94</td></tr></tbody></table>		R&D Spend	Administration	Marketing Spend	State	Profit	0	165349.20	136897.80	471784.10	New York	192261.83	1	162597.70	151377.59	443898.53	California	191792.06	2	153441.51	101145.55	407934.54	Florida	191050.39	3	144372.41	118671.85	383199.62	New York	182901.99	4	142107.34	91391.77	366168.42	Florida	166187.94									
	R&D Spend	Administration	Marketing Spend	State	Profit																																										
0	165349.20	136897.80	471784.10	New York	192261.83																																										
1	162597.70	151377.59	443898.53	California	191792.06																																										
2	153441.51	101145.55	407934.54	Florida	191050.39																																										
3	144372.41	118671.85	383199.62	New York	182901.99																																										
4	142107.34	91391.77	366168.42	Florida	166187.94																																										
		<pre>finalData = pd.concat([pd.get_dummies(data['State'])], data.iloc[:,[0,1,2,4]] , axis = 1) finalData.head()</pre> <table><thead><tr><th></th><th>California</th><th>Florida</th><th>New York</th><th>R&D Spend</th><th>Administration</th><th>Marketing Spend</th><th>Profit</th></tr></thead><tbody><tr><td>0</td><td>0</td><td>0</td><td>1</td><td>165349.20</td><td>136897.80</td><td>471784.10</td><td>192261.83</td></tr></tbody></table>		California	Florida	New York	R&D Spend	Administration	Marketing Spend	Profit	0	0	0	1	165349.20	136897.80	471784.10	192261.83																													
	California	Florida	New York	R&D Spend	Administration	Marketing Spend	Profit																																								
0	0	0	1	165349.20	136897.80	471784.10	192261.83																																								

		1 1 0 0 162597.70 151377.59 443898.53 191792.06 2 0 1 0 153441.51 101145.55 407934.54 191050.39 3 0 0 1 144372.41 118671.85 383199.62 182901.99 4 0 1 0 142107.34 91391.77 366168.42 166187.94
	Inferential Stats 1. Seperate data as features and label and store the same in the form on Numpy Array 2. Split our data as Training Set and Testing Set where Training Set is used for training/convergence purpose whereas Testing set is used for Quality Check Purpose.	1. Seperate data as features and label and store the same in the form on Numpy Array <pre>features = finalData.iloc[:,0:6].values label = finalData.iloc[:,[6]].values</pre>
	# 2. Split our data as Training Set and Testing Set where Training Set Step1: Decide the Split Ratio (65%:35%)(80:20) Step2: Implement the same using Sci-kit Learn package X_train --> Training Features y_train --> Training Label X_test --> Testing Features y_test --> Testing Label For training -----> X_train,y_train For QA -----> X_test,y_test	<pre>from sklearn.model_selection import train_test_split X_train,X_test,y_train,y_test = train_test_split(features, label, test_size=0.2, random_state=10)</pre>
	Build the model Regression Algo -----> Linear Regression $y = mx + b$ (Slope Intercept Formula) profit = b0(California) + b1(Florida) + b2(New York)+ b3(R&D Spend) + b4(AdminSpend) + b5(MarketingSpend) + intercept Goal of this algo is to derive values of b0,b1,b2,b3,b4,b5,intercept based on the historical data !	<pre>from sklearn.linear_model import LinearRegression model = LinearRegression()</pre>
	Train the model Converging Training Set to the Algorithm # fit(features,label)	<pre>model.fit(X_train,y_train)</pre>
		<pre>model.intercept_ array([50001.73604086])</pre>
		<pre>model.coef_ array([[8.41023126e+01, 6.95447747e+02, -7.79550060e+02, 8.05859453e-01, -1.79706621e-02, 2.28153524e-02]])</pre>
	Profit = 8.41023126e+01(California) + 6.95447747e+02(Florida) - 7.79550060e+02(New York) + 8.05859453e-01 (RDSpend) - 1.79706621e-02(Admin) + 2.28153524e-02(MarkSpend) + 50001.73604086	

	<p>Quality Check (Guideline)</p> <p>1-Ensure your model that is converged is a Generalized Model</p> <ul style="list-style-type: none"> - Generalized model is a model that performs well with Known and Unknown data -Technique: Accuracy(Test Data) > Accuracy(Train Data) -----> Model is Generalized Model <p>2-Ensure the Model's accuracy score must be greater than or equal to CL score (*CL = 1 - SL)</p> <ul style="list-style-type: none"> -Technique: Accuracy(Test Data) >= CL 	
	<p>1. Ensure your model that is converged is a Generalized Model</p> <p>#In scikit learn you can get the accuracy score using score function</p>	<pre>testAccuracy = model.score(X_test,y_test) trainAccuracy = model.score(X_train,y_train) print("Test Score is {} and train Score is {}".format(testAccuracy,trainAccuracy))</pre> <p>Test Score is 0.9901105113397691 and train Score is 0.9385918220043519</p>
		As observed above, testScore > trainScore , thus the model is generalized !!!
	<p>SL = 0.05 CL = 0.95</p> <p>2. Ensure the Model's accuracy score must be greater than or equal to CL score (*CL = 1 - SL)</p> <p>Technique : Accuracy(Test Data) >= CL</p> <p>Since my testScore is greater than CL, model passed the Quality Check !</p>	
	<p>Deployment Test – method one</p>	<pre>rdSpend = float(input("Enter R&D Spend: ")) admSpend = float(input("Enter Administration Spend: ")) markSpend = float(input("Enter Marketing Spend: ")) state = input("Enter State: ") refState = ['California', 'Florida', 'New York'] if state in refState: if state == "California": stateDummy = np.array([[1,0,0]]) elif state == "Florida": stateDummy = np.array([[0,1,0]]) else: stateDummy = np.array([[0,0,1]]) finalFeatures = np.concatenate((stateDummy, np.array([[rdSpend,admSpend,markSpend]]), axis=1) profit = model.predict(finalFeatures)[0][0] print("Predicted profit is \$ {}".format(profit)) else: print("Model can't predict profit for the given {} state".format(state))</pre> <p>Enter R&D Spend: 234567 Enter Administration Spend: 45678 Enter Marketing Spend: 76543 Enter State: Florida Predicted profit is \$ 240650.70978691286</p>

<p>For Deployment Test, method 2 we need to do these before train/test split:</p> <p>-we handle Categorical Data with: Sci-kit Package/ OneHotEncoder.</p> <p>- do FeatureScaling with StandardScaler.</p> <p>1- Categorical Data Handling Using Sci-kit Package</p>	<pre>from sklearn.preprocessing import OneHotEncoder ohe = OneHotEncoder(sparse=False) fState = ohe.fit_transform(np.array(data['State']).reshape(-1,1)) features = np.concatenate([fState , np.array(data.iloc[:,[0,1,2]])] , axis = 1) features</pre>
<p>#FeatureScaling</p>	<pre>from sklearn.preprocessing import StandardScaler sc = StandardScaler() features = sc.fit_transform(features)</pre>
<p>Deployment Test / method 2</p>	<pre>rdSpend = float(input("Enter R&D Spend: ")) admSpend = float(input("Enter Administration Spend: ")) markSpend = float(input("Enter Marketing Spend: ")) state = input("Enter State: ") refState = ['California', 'Florida','New York'] if state in refState: stateDummy = ohe.transform(np.array([[state]])) finalFeatures = np.concatenate((stateDummy, np.array([[rdSpend,admSpend,markSpend]])) , axis=1) stdScaleFeatures = sc.transform(finalFeatures) profit = model.predict(stdScaleFeatures)[0][0] print("Predicted profit is \$ {}".format(profit)) else: print("Model can't predict profit for the given {} state".format(state))</pre>
	<pre>data['State'].unique()</pre>

		<code>array(['New York', 'California', 'Florida'], dtype=object)</code>
		<code>[0.000000e+00, 0.000000e+00, 1.000000e+00, 1.653492e+05, 1.368978e+05, 4.717841e+05]</code> <code>California Florida New York rdSpend admSpend markSpend</code>