

### Key Points of the Project

**1. Training a model via linear regression**

Linear Regression means data is modelled using straight line. Line of Best Fit is predicted.

**2. Using the trained model to do the classification**

Classification is supervised learning, in which response is categorical. Here, we classify the iris.data based on its three different target (Setosa-0, Virginica-1, Versicolor-2)

**3. Using cross-validation.**

I use K-fold cross validation here. First, split the dataset into K equal folds. Use fold1 as the testing set and the union of the other folds as training set. Then, calculate testing accuracy.

### Code Explanation:

- The code is designed using Python and contains a function for classification problem . First, the input file "iris.data" is read as a dataframe using numpy library.
- Then the three iris species names are changed to '0','1','2' labels in float format and the entire data is manually split into training set and testing set (80/20) .
- With training data as X and training labels as Y we find Beta cap  $\hat{\beta}$  using linear regression model formula,

$$\hat{\beta} = (A^T A)^{-1} A^T Y$$

- Now, using the estimator values we classify the model and also predict label values over the testing data,

$$f(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots w_d x_d = w_0 + \sum_{j=1}^d w_j x_j$$

- In order to cross validate we fit the data again into predefined regression model and Compute the predicted values from N-fold cross validation and pass it in regression score function to get the accuracy of predicted output which comes out to be 1.0 mostly.
- Since there is no discrepancy in no. of folds, the N value is chosen to be **10** which is most widely used as the more folds we have, we will be reducing the error due to the bias but increasing the error due to variance which nullifies each other.

# RESULT SCREENSHOT

Spyder (Python 3.7)

File Edit Search Source Run Debug Consoles Projects Tools View Help

Editor - C:\Users\Shreyas Mohan\Documents\Fall 19\Machine Learning\Project1\_sxm9806\code.py

```
19 ls = list(test[i,:])
20 for w, x in zip(est_beta, ls):
21     func_x = func_x + (w * x)
22     if abs(func_x) > abs(func_x-1):
23         y_pred = 1
24         if abs(func_x - 1) > abs(func_x - 2):
25             y_pred = 2
26 else:
27     y_pred = 0
28     if abs(func_x) > abs(func_x - 2):
29         y_pred = 2
30 acc = testLabel[i]/y_pred
31 accuracy[i] = acc
32 i+=1
33 avg_acc = np.mean(list(accuracy.values()))
34 print('Classification score: ', avg_acc)
35
36
37 #Data preprocessing
38 df = np.loadtxt("iris.data", dtype="str", delimiter=",")
39 rows = df.shape[0]
40 split_train = (rows*80)//100
41 train = df[0:split_train, :-1].astype(np.float)
42 test = df[split_train:, :-1].astype(np.float)
43 train = np.hstack((train,np.ones((train.shape[0], 1))))
44 test = np.hstack((test,np.ones((test.shape[0], 1))))
45 #print(train)
46 #print(test)
47 #Assigning integer values from 0 to 2 for data Labels
48 trainLabel = df[0:train.shape[0],-1]
49 trainLabel[trainLabel == "Iris-setosa"] = "0"
50 trainLabel[trainLabel == "Iris-virginica"] = "1"
51 trainLabel[trainLabel == "Iris-versicolor"] = "2"
52 #print(trainLabel)
53 testLabel = df[train.shape[0], :-1]
54 testLabel[testLabel == "Iris-setosa"] = "0"
55 testLabel[testLabel == "Iris-virginica"] = "1"
56 testLabel[testLabel == "Iris-versicolor"] = "2"
57 #print(testLabel)
58 #Linear regression
59 #Beta_cap = (A^T.A)^-1*Y
60 est_beta = None
61 matA = train
```

Variable explorer

Name	Type	Size	Value
matY	float64	(120, 1)	[[0.]
predictions	float64	(120,)	[-2.22044605e-16 -2.22044605e-16 -2.22044605e-16 ... 1.00000000e+00 ...]
psize	int	1	30
rows	int	1	150
s	int	1	0
set_size	int	1	30
split_train	int	1	120
t1	float64	(5, 5)	[[ 0.09616686 -0.06265953 -0.07386963 0.0698607 -0.17524193]
t2	float64	(5, 120)	[[ 0.00645537 0.01855176 -0.00582655 ... -0.01409629 0.05330729
test	float64	(30, 5)	[[6.9 3.2 5.7 2.3 1. ]
testLabel	float64	(30,)	[1. 1. 1. ... 1. 1. 1.]
test_A	DataFrame	(30, 4)	Column names: 0, 1, 2, 3

Python console

```
In [10]: runfile('C:/Users/Shreyas Mohan/Documents/Fall 19/Machine Learning/Project1_sxm9806/code.py', wdir='C:/Users/Shreyas Mohan/Documents/Fall 19/Machine Learning/Project1_sxm9806')
[-0.04558273 -0.63329621 0.52108015 -0.45601488 1.91854987]
2 -fold Cross-validation scores: 1.0
3 -fold Cross-validation scores: 1.0
4 -fold Cross-validation scores: 1.0
5 -fold Cross-validation scores: 1.0
6 -fold Cross-validation scores: 1.0
7 -fold Cross-validation scores: 1.0
8 -fold Cross-validation scores: 1.0
9 -fold Cross-validation scores: 1.0

In [11]:
```

Permissions: RW End-of-lines: CRLF Encoding: ASCII Line: 26 Column: 8 Memory: 56 %

9:30 PM 9/10/2019