

機器學習概論 HW1 Report

Part1

Split the train and test data

USING PANDAS TO FINISH CSV FILES EDITING

```
1 def read_and_split(filename):
2     file = pd.read_csv(filename, header=None)
3     file.columns = ['Label', 'Alcohol', 'Malic acid', 'Ash', 'Alcalinity of ash', 'Ma
4
5     type1_filt = (file['Label'] == 1)
6     type2_filt = (file['Label'] == 2)
7     type3_filt = (file['Label'] == 3)
8
9     type1 = file.loc[type1_filt]
10    type2 = file.loc[type2_filt]
11    type3 = file.loc[type3_filt]
12
13    type1_test = type1.sample(n=18, frac=None, random_state=200)
14    type1_train = type1.drop(type1_test.index)
15    type2_test = type2.sample(n=18, frac=None, random_state=200)
16    type2_train = type2.drop(type2_test.index)
17    type3_test = type3.sample(n=18, frac=None, random_state=200)
18    type3_train = type3.drop(type3_test.index)
19
20    training_set = pd.concat([type1_train, type2_train, type3_train])
21    testing_set = pd.concat([type1_test, type2_test, type3_test])
22
23    training_set.to_csv('training_set.csv')
24    testing_set.to_csv('testing_set.csv')
25
26    return training_set, testing_set
```

Part2

Seperate each feature into a group with Label

SPLIT INTO X AND Y

```
1 x_dataset = dataset[:, 1:] #features
2 y_dataset = dataset[:, 0] # Label
```

CLASSIFY INTO LABELS

```
1 for i in range(len(y_dataset)):
2     if(y_dataset[i] == 1):
3         type1.append(x_dataset[i])
4
5     elif(y_dataset[i] == 2):
6         type2.append(x_dataset[i])
7
8     elif(y_dataset[i] == 3):
9         type3.append(x_dataset[i])
```

FOR EACH FEATURE, IT HAS INDEX FOR ITS LABEL

```

1 type1_ft = np.zeros(shape=[13, len(type1)])
2 type2_ft = np.zeros(shape=[13, len(type2)])
3 type3_ft = np.zeros(shape=[13, len(type3)])
4 for idx, ele in enumerate(x_dataset[0:len(type1)]):
5     for i in range(13):
6         type1_ft[i][idx] = x_dataset[idx][i]
7 for idx, ele in enumerate(x_dataset[len(type1):len(type1)+len(type2)]):
8     for i in range(13):
9         type2_ft[i][idx] = x_dataset[len(type1)+idx][i]
10 for idx, ele in enumerate(x_dataset[len(type1)+len(type2):len(y_dataset)]):
11     for i in range(13):
12         type3_ft[i][idx] = x_dataset[len(type1)+len(type2)+idx][i]

```

Calculate mean, standard and normal distribution

GET 13 FEATURES MEAN, STANDARD AND NORMAL DISTRIBUTION WITH NUMPY AND SYMSTATS

$$f(x) = \frac{\exp(-x^2/2)}{\sqrt{2\pi}}$$

```

1 for i in range(13):
2     type1_mean.append(np.mean(type1[i]))
3     type2_mean.append(np.mean(type2[i]))
4     type3_mean.append(np.mean(type3[i]))

```

```

1 for i in range(13):
2     type1_std.append(np.std(type1[i]))
3     type2_std.append(np.std(type2[i]))
4     type3_std.append(np.std(type3[i]))

```

```

1 for i in range(13):
2     type1_norm.append(st.norm(type1_mean[i], type1_std[i]))
3     type2_norm.append(st.norm(type2_mean[i], type2_std[i]))
4     type3_norm.append(st.norm(type3_mean[i], type3_std[i]))

```

Prior

THE DISTRIBUTION OF LABELS WILL BE THE PRIOR DISTRIBUTION

```

prior[0] = len(type1) / len(y_dataset) #0.33
prior[1] = len(type2) / len(y_dataset) #0.42
prior[2] = len(type3) / len(y_dataset) #0.24

```

Likelihood

$$\operatorname{argmax}_{\theta} \frac{1}{h} \int_{x_j}^{x_j+h} f(x \mid \theta) dx,$$

```

1 likelihood = sc.integrate.quad(dis[label][i].pdf, data[i+1], data[i+1]+delta )[0]

```

Posterior

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Likelihood
Class Prior Probability

Posterior Probability
Predictor Prior Probability

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

```
1 post[label] = 1. * prior[label]
```

```
1 post[label] = post[label]*likelihood
```

MAP

```
1 predict = np.argmax(post)
2 total+=1
3 if predict == (data[0] -1):
4     correct += 1
5 else:
6     pass
7 print('accuracy: ', correct/total)
```

Part3

Dimension 13 features into 2 & 3 features

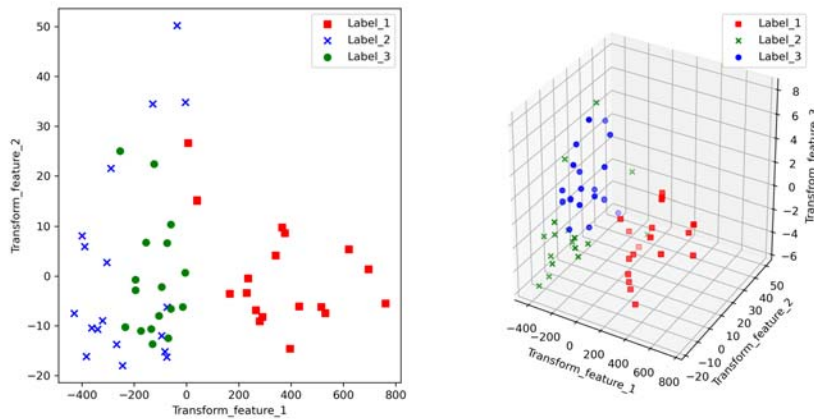
```
1 X_p = pca2.fit(X).transform(X)
```

```
1 X_p = pca3.fit(X).transform(X)
```

Visualization

```
1 for c, i, target_name, m in zip('rbg', labels, types, markers):
2     plt2D.scatter(X_p[y==i, 0], X_p[y==i, 1], c=c, label=target_name, marker=m)
```

```
1 for c, i, target_name, m in zip('rgb', labels, types, markers):
2     plt3D.scatter(X_p[y==i, 0], X_p[y==i, 1], X_p[y==i, 2], c=c, label=target_
```



Part4

Discussion with different Prior

FIRST, I JUST TRY SEVERAL CASES TO OBSERVE THE DIFFERENCE WITH EACH OF THEM.

THE MORE LABELS APPROACHES TO 0 , THE WORSE ACCURACY WILL BE.

IN THE OTHER HAND, BE MORE CLOSE TO THE TRAINING DISTRIBUTION OF 3 LABELS, ACCURACY WILL BE MORE HIGHER.

```

1  #acc = 0.96
2  ""
3  prior[0] = 0.05
4  prior[1] = 0.05
5  prior[2] = 0.9
6  ""
7  #acc = 0.87
8  ""
9  prior[0] = 0.9
10 prior[1] = 0.09999999999999999
11 prior[2] = 0.00000000000000001
12 ""
13 #acc = 0.57
14 ""
15 prior[0] = 0.00000000000000005
16 prior[1] = 0.99999999999999999
17 prior[2] = 0.00000000000000005
18 ""
19 #acc = 0.33
20 ""
21 prior[0] = 0.0
22 prior[1] = 0.1
23 prior[2] = 0.0
24 ""

```

Discover most important factors

THEN I CHECK CRITICAL ROLES IN 13 FEATURES, PROLINE & MAGNESIUM ARE 2 HIGHEST METRIC AFFECT RESULTS.

```
1 pd.DataFrame(pca2.components_,columns=X.columns,index = ['Transform_feature_1','Tr
```

	Alcohol	Malic acid
Transform_feature_1	0.001404923002591525	-0.00094
Transform_feature_2	-0.002664669843727335	-0.00675

PROLINE ,MAGNESIUM, AND ALCALINITY OF ASH ARE 3 HIGHEST METRIC

```
1 pd.DataFrame(pca3.components_,columns=X.columns,index = ['Transform_feature_1','Tr
```

	Alcohol	Malic acid
Transform_feature_1	0.001404923002591525	-0.00094
Transform_feature_2	-0.002664669843727335	-0.00675
Transform_feature_3	0.04599422668581956	0.138707

JUST USE THREE MOST IMPORTANT FEATURES TO CALCULATE MAP

ACCURACY = 0.81

```
1 for i in range(3):
2     if i == 0:
3         # 13th feature
4         i = 13-1
5     elif i == 1:
6         # 5th feature
7         i = 5-1
8     else:
9         # 4th feature
10        i = 4-1
11    likelihood = sc.integrate.quad(dis[label][i].pdf, data[i+1], data[i+1]+delta )
12    post[label] = post[label]*likelihood
```

WE STILL HAVE OVER 80% ACCURACY