

# Aprendizado de Máquina - MO444

## Exercício 1

Aluno: Paulo Ricardo Finardi. RA: 144809

### 0 Preliminares

Foi utilizado Python 3 na resolução dos exercícios. O código abaixo é o *header* de todos os itens.

---

```
1 import numpy as np
2 import pandas as pd
3 import csv
4 import matplotlib.pyplot as plt
5
6 # read dada1.csv in python with Pandas
7 data1 = pd.read_csv('data1.csv')
8
9 # delete the last column of data1
10 data1.set_index('clase').to_csv('data.csv', index=None)
11
12 # saving new file in csv format
13 data = pd.read_csv('data.csv')
14
15 # convert the datas in numpy-array
16 data1 = data1.as_matrix(columns=None)
17 data = data.as_matrix(columns=None)
```

---

### 1 PCA

**Enunciado:** Faça o *Principal Component Analysis* (PCA) dos dados sem a última coluna. Se você quiser que os dados transformados tenham 80% da variância original, quantas dimensões do

PCA você precisa manter?

**Resposta:** São necessárias 12 dimensões. A variância original é fornecida em cada componente do PCA. Somamos componente a componente até obter a variância de 80%. O código para realização dessa tarefa é fornecido a seguir.

---

```
1 # function PCA with sklearn
2 def doPCA(x): # x is the argument for the number of components in PCA
3     from sklearn.decomposition import PCA
4     pca = PCA(n_components=x)
5     pca.fit(data)
6     return pca
7
8 # this step is obsolete, but we'll repeat the PCA after we find the requested variance.
9 pca = doPCA(data.shape[1])
10
11 # the elements of the _variance contains the total variance.
12 _variance = pca.explained_variance_ratio_
13
14 # store the requested variance
15 requested_variance = 0
16
17 # the number of the variance requested
18 lim = 0.80
19
20 # loop to find the quantity of elements
21 for i in range (len(_variance)):
22     requested_variance += _variance[i]
23     if requested_variance >= lim:
24         element = i
25         requested_variance -= _variance[i] # just for adjust the index
26         break
27
28 print('With %d elements, we get the requested variance, that is: %.3f.' %
29       (element, requested_variance))
30
```

```

31 # repeating PCA
32 pca = doPCA(element)
33 variance_ = pca.explained_variance_ratio_
34
35 # prints
36 print('The %d principal components:' % (element))
37 np.set_printoptions(precision=3)
38 print(variance_)
39 print('\nThe sum of the components is: %.3f.' % np.sum(variance_))

```

---

Na linha 38, a função `print` fornece:

[0.312, 0.139, 0.076, 0.051, 0.049, 0.041, 0.032, 0.03, 0.02, 0.017, 0.015, 0.014],

e a soma dessas componentes vale 0.798 ou 79.8% (dada pela linha 39).

## 2 Regressão logística

**Enunciado:** Considere as primeiras 200 linhas dos dados como o conjunto de treino, e as 276 ultimas como o conjunto de dados. Treine uma regressão logística no conjunto de treino dos dados originais e nos dados transformados. Qual a taxa de acerto no conjunto de teste nas 2 condições (sem e com PCA)?

**Resposta:** Nessa etapa foi utilizado a matriz de confusão. A escolha da matriz de confusão vem do fato que ela oferece uma medida efetiva do modelo de classificação ao mostrar o número de classificações corretas versus as classificações preditas para cada classe, sobre os conjuntos de testes. As métricas utilizadas foram as seguintes (os termos em inglês foram mantidos):

- True Positive Rate (TPR) – fornece a sensibilidade do classificador;
- False Positive Rate (FPR) – fornece a quantidade de alarmes falsos do classificador;
- Accuracy – descrição de erros sistemáticos;
- Precision – descrição de erros aleatórios, variabilidade estatística.
- Recall – razão entre  $TP/(TP + FN)$ ;
- F1-Score – média harmônica entre Precision e Recall.

As Figuras (1) e (2) exibem as matrizes de confusão dos conjuntos de teste (*sem PCA*) e (*com PCA*) respectivamente.

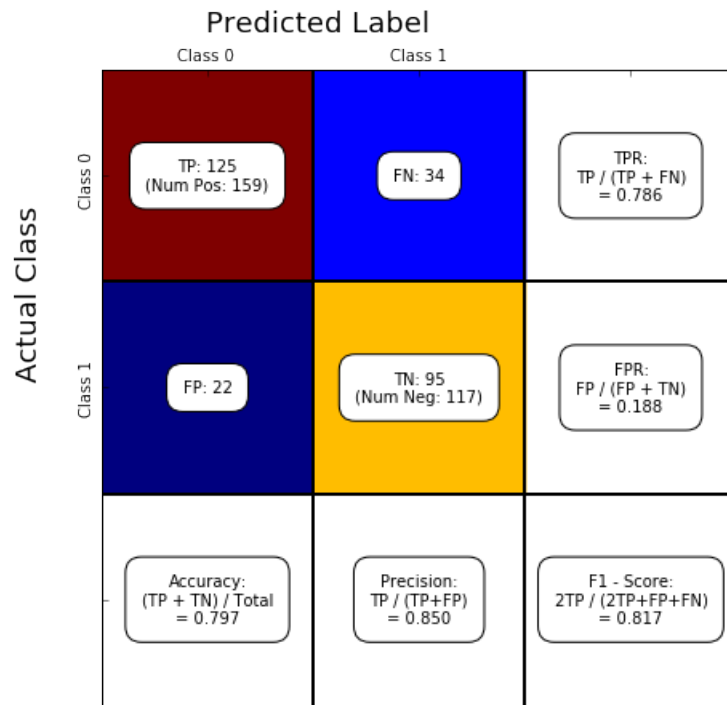


Figura 1: Matriz de confusão— dados *sem PCA*.

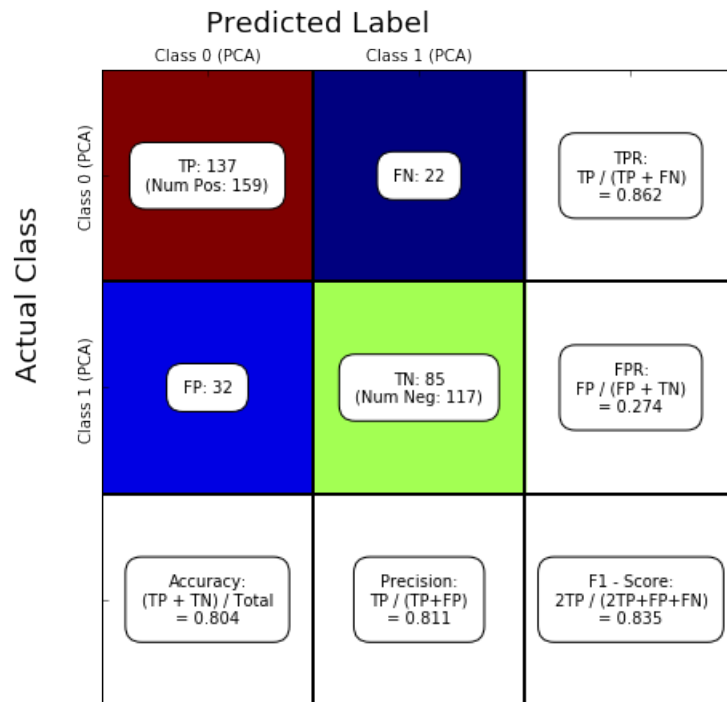


Figura 2: Matriz de confusão— dados *com PCA*.

Os dados com PCA obtiveram melhor resultado (exceto em Precision). A seguir o código:

---

```
1  # data with pca
2  transformed_data = pca.transform(data)
3
4  # y_total is the expected
5  y_total = data1[:, 166]
6
7  # split data in train and test sets
8  X_train_PCA, X_train, y_train = transformed_data[0:200, :], data[0:200,:], y_total[0:200]
9  X_test_PCA, X_test, y_test = transformed_data[200:,:], data[200:,:], y_total[200:]
10
11 # metrics for the confusion matrix
12 from sklearn import metrics
13
14 # function Linear Regression (LR) with sklearn
15 def doLR(X, y):
16     from sklearn.linear_model \
17     import LogisticRegression as LR
18     lr = LR()
19     lr.fit(X, y)
20     return lr
21
22 # LR in PCA data
23 lrPCA = doLR(X_train_PCA, y_train)
24
25 # LR in data
26 lr = doLR(X_train, y_train)
27
28 # predicted sets
29 predictedPCA = lrPCA.predict(X_test_PCA)
30 predicted = lr.predict(X_test)
31
32 # confusion matrix
33 m_confusion_PCA = metrics.confusion_matrix(y_test, predictedPCA)
34 m_confusion = metrics.confusion_matrix(y_test, predicted)
```

```

35
36 # function for confusion matrix
37 def show_confusion_matrix(C,class_labels=['0','1']):
38
39     # true negative, false positive, etc...
40     tp = C[0,0]; fn = C[0,1]; fp = C[1,0]; tn = C[1,1];
41
42     NP = fn+tp # num positive examples
43     NN = tn+fp # num negative examples
44     N  = NP+NN
45
46     fig = plt.figure(figsize=(6.5,6.5))
47     ax  = fig.add_subplot(111)
48     ax.imshow(C, interpolation='nearest', cmap='jet')
49
50     # draw the grid boxes
51     ax.set_xlim(-0.5,2.5)
52     ax.set_ylim(2.5,-0.5)
53     ax.plot([-0.5,2.5],[0.5,0.5], '-k', lw=2)
54     ax.plot([-0.5,2.5],[1.5,1.5], '-k', lw=2)
55     ax.plot([0.5,0.5],[-0.5,2.5], '-k', lw=2)
56     ax.plot([1.5,1.5],[-0.5,2.5], '-k', lw=2)
57
58     # set xlabels
59     ax.set_xlabel('Predicted Label', fontsize=18)
60     ax.set_xticks([0,1,2])
61     ax.set_xticklabels(class_labels + [''])
62     ax.xaxis.set_label_position('top')
63     ax.xaxis.tick_top()
64
65     # these coordinate might require some tinkering.
66     ax.xaxis.set_label_coords(0.34,1.06)
67
68     # set ylabels
69     ax.set_ylabel('Actual Class', fontsize=18, rotation=90)
70     ax.set_yticklabels(class_labels + [''],rotation=90)

```

```

71 ax.set_yticks([0,1,2])
72 ax.yaxis.set_label_coords(-0.09,0.65)
73
74 # fill in initial metrics: tp, tn, etc...
75 ax.text(1,1,
76         'TN: %d\n(Num Neg: %d)'%(tn, NN),
77         va='center',
78         ha='center',
79         bbox=dict(fc='w',boxstyle='round,pad=1'))
80
81 ax.text(1,0,
82         'FN: %d'%fn,
83         va='center',
84         ha='center',
85         bbox=dict(fc='w',boxstyle='round,pad=1'))
86
87 ax.text(0,1,
88         'FP: %d'%fp,
89         va='center',
90         ha='center',
91         bbox=dict(fc='w',boxstyle='round,pad=1'))
92
93 ax.text(0,0,
94         'TP: %d\n(Num Pos: %d)'%(tp, NP),
95         va='center',
96         ha='center',
97         bbox=dict(fc='w',boxstyle='round,pad=1'))
98
99 # secondary metrics: accuracy, true pos rate, etc...
100 ax.text(2,1,
101         'FPR:\nFP / (FP + TN)\n= %.3f'%(fp / (fp+tn)),
102         va='center',
103         ha='center',
104         bbox=dict(fc='w',boxstyle='round,pad=1'))
105
106 ax.text(2,0,

```

```

107         'TPR:\nTP / (TP + FN)\n= %.3f'%(tp / (tp+fn)),
108         va='center',
109         ha='center',
110         bbox=dict(fc='w',boxstyle='round,pad=1'))
111
112     ax.text(0,2,
113             'Accuracy:\n(TP + TN) / Total\n= %.3f'%((tp+tn)/N),
114             va='center',
115             ha='center',
116             bbox=dict(fc='w',boxstyle='round,pad=1'))
117
118     ax.text(2,2,
119             'F1 - Score:\n2TP / (2TP+FP+FN)\n= %.3f'%(2*tp/(2*tp+fp+fn)),
120             va='center',
121             ha='center',
122             bbox=dict(fc='w',boxstyle='round,pad=1'))
123
124     ax.text(1,2,
125             'Precision:\nTP / (TP+FP)\n= %.3f'%(tp/(tp+fp)),
126             va='center',
127             ha='center',
128             bbox=dict(fc='w',boxstyle='round,pad=1'))
129
130     plt.tight_layout()
131     plt.show()
132
133     # prints
134     show_confusion_matrix(m_confusion, ['Class 0 ', 'Class 1'])
135     show_confusion_matrix(m_confusion_PCA, ['Class 0 (PCA)', 'Class 1 (PCA)'])

```

---



### 3 LDA

**Enunciado:** Treine o *Linear Discriminant Analysis* (LDA) nos conjuntos de treino com e sem PCA e teste nos respectivos conjuntos de testes. Qual a acurácia nas 2 condições?

**Resposta:** Adotamos a mesma estratégia do item anterior. Treinamos o LDA e geramos novas matrizes de confusão para o classificador LDA. Figuras (3) e (4).

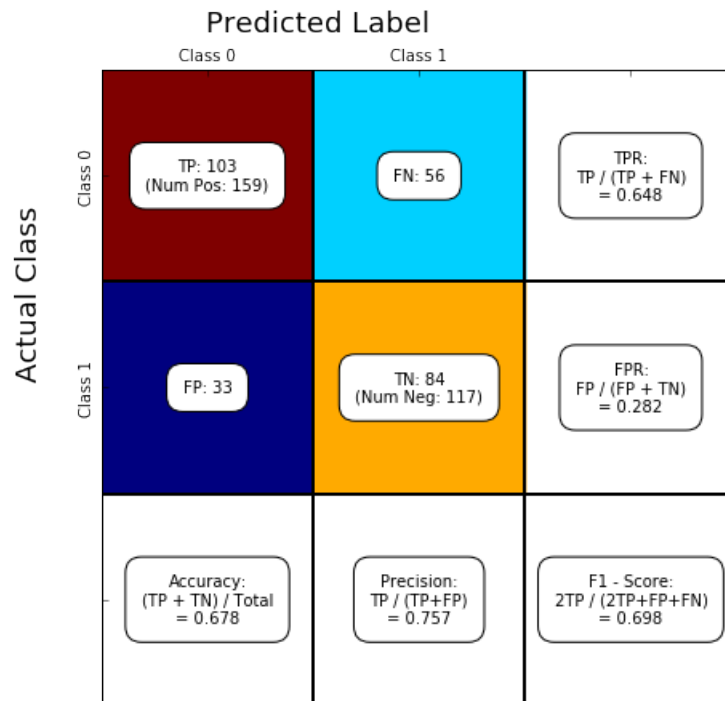


Figura 3: Matriz de confusão— LDA dados *sem PCA*.

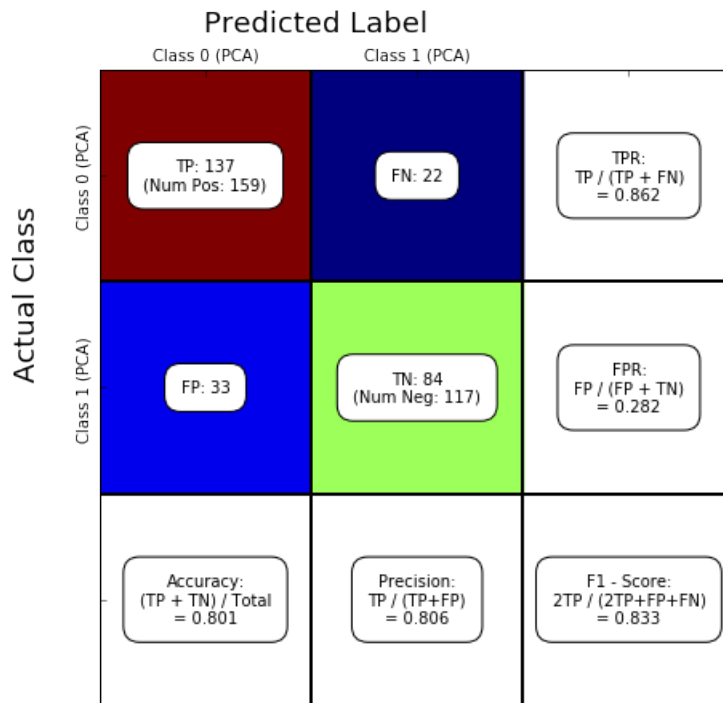


Figura 4: Matriz de confusão— LDA dados *com PCA*.

Conclusão: Diferente do item da regressão logística (2), a variação da acurácia nos dois conjuntos foi maior. No conjunto de dados *sem PCA* a acurácia foi de 0.678 ou 67,8%, já no conjunto de dados *com PCA* a acurácia foi de 0.801 ou 80,1%. Na sequência o código.

---

```

1  # LDA function
2  def doLDA(X, y, n):
3      # X is the train set, y is the expected, and n is the number of components
4      from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
5      lda = LDA(n_components=n)
6      lda.fit(X, y)
7      return lda
8
9  # lda in PCA data
10 ldaPCA = doLDA(X_train_PCA, y_train, element)
11
12 # lda in data
13 lda = doLDA(X_train, y_train, element)
14

```

```
15 # predicted sets
16 predicted_ldaPCA = ldaPCA.predict(X_test_PCA)
17 predicted_lda    = lda.predict(X_test)
18
19 # confusion matrix
20 lda_confusionPCA = metrics.confusion_matrix(y_test, predicted_ldaPCA)
21 lda_confusion    = metrics.confusion_matrix(y_test, predicted_lda)
22
23 # prints
24 show_confusion_matrix(lda_confusionPCA, ['Class 0 (PCA)', 'Class 1 (PCA)'])
25 show_confusion_matrix(lda_confusion, ['Class 0', 'Class 1'])
```

---

## 4 Melhor combinação

**Enunciado:** Qual a melhor combinação de classificador e PCA ou não?

**Resposta:** Considerando as matrizes de confusão das Figuras (1), (2), (3) e (4) a que apresenta os melhores resultados é a Figura (2), ou seja, **o melhor resultado** é obtido nos dados *com PCA e regressão logística*. Esse resultado é bem sutil; ele obteve 32 FP contra 33 FP no classificador LDA. Com essa combinação, também temos a melhor métrica **F1-Score** (dois centésimos maior do que o PCA-LDA).