

Exercises of LDA classifier

Eva Cernadas

FMLCV Course

CITIUS: Centro de Tecnoloxías Intelixentes da USC

Universidade de Santiago de Compostela

1 de diciembre de 2021

1. Programs in Matlab

To practice the Linear Discriminant Analysis (LDA) classifier, i provide the datasets `wine.data` (with 3 classes) and `hepatitis.data` (with 3 classes), which were downloaded from the UCI repository (<https://archive.ics.uci.edu/ml/datasets.php>). I also share the following Matlab code of the LDA classifier:

```
1 %lda: implements the lda classifier
2 %output: y the predicted output
3 %inputs: x matrix with the training patterns (each pattern one
      row)
4 %          c vector with the desired output in training set
5 %          xtest matrix with the test patterns
6 function y=lda(x, c, xtest)
7 [N, I]=size(x);
8 mx=mean(x); stdx=std(x);
9 %preprocessing: mean 0, desviation 1
10 %x=bsxfun(@rdivide, bsxfun(@minus, x, mx), stdx);
11 x=(x-mx)./stdx;
12
13 cl=unique(c); C=numel(cl);
14 nc=zeros(C,1); %number of patterns per class
15 mc=zeros(C,I); %mean of each class
16 S=zeros(I); %total covariance
17 w=zeros(C, I+1); %coeficients of LDA
18
19 for i=1:C
20     j=(c==cl(i)); nc(i)=sum(j);
21     u=x(j, :); mc(i, :)=mean(u);
```

```

22         S=S+(nc(i)-1)*cov(u)/(N-C);
23     end
24     pr=nc/N;    % probabilities
25     for i=1:C
26         u=mc(i,:); t=u/S;
27         w(i,1)=log(pr(i))-t*u'/2;    % offset
28         w(i,2:end)=t;    % linear term
29     end
30     % standardized xtest
31     % preprocessing: mean 0, desviation 1
32     %xtest=bsxfun(@rdivide,bsxfun(@minus,xtest,mx),stdx);
33     xtest=(xtest-mx)./stdx;
34
35     L=[ones(size(xtest,1),1) xtest] * w';    % linear scores
36     %implement softmax function
37     P=exp(L)./repmat(sum(exp(L),2),[1 C]);    % class probabilities
38     [~,y]=max(P,[],2);    % class predicted by LDA

```

Firstly, we are going to use the `lda.m` function to train and test the classifier using the whole dataset. The Matlab program could be:

```

1  clear all;
2  % 3 classes
3  % dataset='wine'; x=load('wine.data');
4  % 2 classes
5  dataset='hepatitis'; x=load('hepatitis.data');
6  c=x(:,1); x(:,1) = []; [N,I]=size(x);
7  cl=unique(c); C=numel(cl);
8  y=lda(x, c, x);
9  [kappa, acc, cm]=evaluate(c, y, C);
10 disp('Confusion matrix='); disp(cm);
11 fprintf('dataset %s: accuracy=%0.2f%%\n',dataset, acc)
12 fprintf('dataset %s: kappa=%0.2f%%\n',dataset, kappa)

```

which use the following function `evaluate()` to calculate the confusion matrix, accuracy and Cohen kappa:

```

1  % Return: kappa, accuracy and confusion matrix
2  % Inputs: tc (true class), pc (predicted class) and C (number of
           classes)
3  function [kappa, acc, cm]=evaluate(tc, pc, C)
4       cm=zeros(C); np=length(tc);
5       for i=1:np
6           j=tc(i); k=pc(i); cm(j,k)=cm(j,k)+1;
7       end

```

```

8         s=sum(sum(cm)); pa=trace(cm); acc=100*pa/s; pe=0;
9         for k=1:C
10             pe=pe+sum(cm(k,:))*sum(cm(:,k))/s;
11         end
12         kappa=100*(pa-pe)/(s-pe);
13 end

```

Secondly, we will apply the LDA classifier using cross-validation with two dataset (training and testing sets). The validation set is not necessary because there is no hyperparameter to tune. I also share the function code to do this operation:

```

1 %createFolds: create the folds for cross-validation
2 %Inputs: x (matrix of patterns), x (desired output) and K (
   number of folds)
3 %Outputs: tx matrix with training patterns(rows)
4 %          tc vector with the desired output for training
   patterns
5 %          vx, vc: idem to validation set
6 %          sx, sc: idem to test set
7 function [tx,tc,vx,vc,sx,sc]=createFolds(x, c, K)
8 rand('seed',0);
9 [N,n]=size(x); %Number of patterns and features
10 val=unique(c); %output values
11 Q=numel(val); %number of classes
12
13 for j=1:Q
14     fprintf(' class %d: %d patterns\n',j,sum(c==j))
15 end
16
17 ntf=K-2; %Number of training folds
18 nvf=1; %Number of validation folds: the number of test folds is
   K-ntf-nvf
19 %creation of folds
20 npc=zeros(1,Q); %No. Patterns per class
21 %ntp/nvp/nsp=no. train/valid/test patterns of each class;
22 %npf=no. patterns of each class per fold
23 ntp=zeros(1,Q); nvp=zeros(1,Q);
24 nsp=zeros(1,Q); npf=zeros(1,Q);
25 tx=cell(1,K); tc=cell(1,K);
26 vx=cell(1,K); vc=cell(1,K);
27 sx=cell(1,K); sc=cell(1,K);
28 for i=1:Q
29     t=find(c==i); j=numel(t); npc(i)=j; k=randperm(j);
30     ind=t(k); %ind=indices of patterns of each class

```

```

31  npf(i)=floor(j/K); ntp(i)=ntf*npf(i);
32  nvp(i)=nvf*npf(i); nsp(i)=j-ntp(i)-nvp(i);
33  start=1;
34  for k=1:K
35      p=start; u=[];
36      for l=1:ntp(i) %indices of train patterns
37          u=[u ind(p)]; p=p+1;
38          if p>npc(i); p=1; end
39      end
40      tx{k}=[tx{k}; x(u,:)]; tc{k}=[tc{k}; c(u)]; u=[];
41      for l=1:nvp(i) %indices of validation patterns
42          u=[u ind(p)]; p=p+1;
43          if p>npc(i); p=1; end
44      end
45      vx{k}=[vx{k}; x(u,:)]; vc{k}=[vc{k}; c(u)]; u=[];
46      for l=1:nsp(i) %indices of test patterns
47          u=[u ind(p)]; p=p+1;
48          if p>npc(i); p=1; end
49      end
50      sx{k}=[sx{k}; x(u,:)]; sc{k}=[sc{k}; c(u)];
51      start=start+npf(i);
52  end
53 end

```

The Matlab code to use to use the function `createFolds()` could be:

```

1  clear all;
2  %3 classes
3  dataset='wine'; x=load('wine.data'); %first column is the output
4  %2 classes
5  %dataset='hepatitis'; x=load('hepatitis.data'); %first column is
   the output
6  c=x(:,1); x(:,1) = []; [N,I]=size(x);
7  cl=unique(c); C=numel(cl);
8  K=4 %number of folds
9  [tx,tc,vx,vc,sx,sc]=createFolds(x, c, K);
10 cmt=zeros(C); %confusion matrix
11 kappa=zeros(1,K); acc=zeros(1,K);
12 for i=1:K
13     ti=[tx{i}; vx{i}]; %join training and validation sets for
        training
14     ci=[tc{i}; vc{i}]; %idem for desired output
15     y=lda(ti, ci, sx{i});
16     [kappa(i), acc(i), cm]=evaluate(sc{i}, y, C);

```

```

17     fprintf('Confusion matrix fold %d=\n', i); disp(cm);
18     fprintf('fold %a: kappa=%.1f%%accuracy=%.1f%%\n', i, kappa(i),
        acc(i))
19     cmt = cmt + cm;
20 end
21 kappa_mean=mean(kappa); acc_mean=mean(acc); cmt=cmt/K;
22 disp('Final confusion matrix='); disp(cmt);
23 fprintf('dataset %s: kappa=%.1f%%accuracy=%.1f%%\n', dataset,
        kappa_mean, acc_mean)

```

2. Programs in Python

The LDA classifier can be executed using the object `sklearn.linear_discriminant.LinearDiscriminantAnalysis` object. The training and test on the whole dataset can be executed using the following program:

```

1 from numpy import *
2 from sklearn.discriminant_analysis import *
3 from sklearn.metrics import *
4 from sklearn.model_selection import *
5
6 dataset='hepatitis'; # hepatitis (2 clases), wine (3 clases)
7 nf='../data/%s.data'%dataset;x=loadtxt(nf)
8 y=x[:,0];x=delete(x,0,1)
9 # preprocessing: mean 0, desviation 1
10 x=(x-mean(x,0))/std(x,0)
11 print('LDA dataset %s:'%dataset)
12 #-----
13 # training and test on the whole dataset
14 #-----
15 lda=LinearDiscriminantAnalysis().fit(x,y)
16 z=lda.predict(x)
17 kappa=cohen_kappa_score(y,z); acc=accuracy_score(y,z)
18 print('Train+Test: kappa=%.1f%% accuracy=%.1f%%'\
19       %(100*kappa,100*acc))
20 cf=confusion_matrix(y,z)
21 print('confusion matrix:'); print(cf)
22 #-----
23 # 4-fold cross-validation using cross_val_predict sklearn
   function
24 #-----
25 lda=LinearDiscriminantAnalysis()
26 K=4;z=cross_val_predict(lda,x,y,cv=K)
27 kappa=cohen_kappa_score(y,z); acc=accuracy_score(y,z)
28 print('%i-fold CV: kappa=%.1f%% accuracy=%.1f%%'\
29       %(K,100*kappa,100*acc))
30 cf=confusion_matrix(y,z)
31 print('confusion matrix:'); print(cf)

```

```

32
33 C=len(unique(y))
34 if C==2:
35     pre=precision_score(y,z)
36     re=recall_score(y,z)
37     f1=f1_score(y,z)
38     print('precision=%.1f%% recall=%.1f%% f1=%.1f%%'\
39           %(100*pre,100*re,100*f1))

```

In order to perform 4-fold cross-validation, the following program uses the corresponding function createFolds() for splitting data into train, validation and test sets:

```

1 # LDA sintonizando o no. V de vecinhos con validacion cruzada
2 # K-fold e particions de entrenamiento, validacion e teste
3 from numpy import *
4 from sklearn.discriminant_analysis import *
5 from sklearn.metrics import *
6
7 #dataset='wine'
8 dataset='hepatitis'
9 nf='%s.data'%dataset;x=loadtxt(nf)
10 y=x[:,0]-1;x=delete(x,0,1);C=len(unique(y))
11 print('LDA dataset %s'%dataset)
12
13 def createFolds(x,y,K):
14     from numpy.random import shuffle,seed
15     seed(100)
16     [N,n]=x.shape;C=len(unique(y));ntf=K-2;nvf=1
17     tx=[];ty=[];vx=[];vy=[];sx=[];sy=[]
18     for i in range(K):
19         tx.append(zeros([1,n]));ty.append(array([], 'int'))
20         vx.append(zeros([1,n]));vy.append(array([], 'int'))
21         sx.append(zeros([1,n]));sy.append(array([], 'int'))
22     for i in range(C):
23         t=where(y==i)[0];npc=len(t);shuffle(t)
24         npf=int(npc/K);ntp=npf*ntf
25         nvp=npf*nvf;nsp=npc-ntp-nvp;start=0
26         for k in range(K):
27             p=start;u=array([], 'int')
28             for l in range(ntp):
29                 m=t[p];u=append(u,m);p=(p+1)%npc
30                 tx[k]=vstack((tx[k],x[u]))
31                 ty[k]=append(ty[k],y[u]);u=array([], 'int')
32             for l in range(nvp):
33                 m=t[p];u=append(u,m);p=(p+1)%npc
34                 vx[k]=vstack((vx[k],x[u]))
35                 vy[k]=append(vy[k],y[u]);u=array([], 'int')
36             for l in range(nsp):
37                 m=t[p];u=append(u,m);p=(p+1)%npc
38                 sx[k]=vstack((sx[k],x[u]))
39                 sy[k]=append(sy[k],y[u]);start=start+npf

```

```

40     for k in range(K):
41         tx[k]=delete(tx[k],0,0);vx[k]=delete(vx[k],0,0)
42         sx[k]=delete(sx[k],0,0)
43     return [tx,ty,vx,vy,sx,sy]
44
45 K=4;
46 tx,ty,vx,vy,sx,sy=createFolds(x,y,K)
47
48 # preprocessing: mean 0, deviation 1
49 for k in range(K):
50     med=mean(tx[k],0);dev=std(tx[k],0)
51     tx[k]=(tx[k]-med)/dev
52     vx[k]=(vx[k]-med)/dev
53     sx[k]=(sx[k]-med)/dev
54 kappa=zeros(K);acc=zeros(K);cm=zeros([C,C])
55 print(' %10s %10s %10s'%( 'k', 'kappa(%)', 'acc(%)'),end='')
56 if C==2:
57     pre=zeros(K);re=zeros(K);f1=zeros(K)
58     print(' %15s %10s %10s'%( 'Precision(%)', 'Recall(%)', 'F1(%)'),
59           end='')
60 print(' ')
61 for k in range(K):
62     x=vstack((tx[k],vx[k]));y=concatenate((ty[k],vy[k]))
63     modelo=LinearDiscriminantAnalysis().fit(x,y)
64     z=modelo.predict(sx[k]);y=sy[k]
65     kappa[k]=100*cohen_kappa_score(y,z);acc[k]=100*
66         accuracy_score(y,z);cm+=confusion_matrix(y,z)
67     print(' %10i %10.2f %10.2f'%(k+1,kappa[k],acc[k]),end='')
68     if C==2:
69         pre[k]=100*precision_score(y,z)
70         re[k]=100*recall_score(y,z)
71         f1[k]=100*f1_score(y,z)
72         print(' %15.2f %10.2f %10.2f'%(pre[k],re[k],f1[k]),end='')
73     print(' ')
74 kappa_mean=mean(kappa);acc_mean=mean(acc);cm/=K
75 print(' kappa_mean=%.2f%% acc_mean=%.2f%%'%(kappa_mean,acc_mean
76       ),end='')
77 if C==2:
78     pre_mean=mean(pre);re_mean=mean(re);f1_mean=mean(f1)
79     print(' precision_mean=%.2f%% recall_mean=%.2f%% F1_mean=%.2f
80           %%%'%(pre_mean,re_mean,f1_mean))
81 else:
82     print(' ')

```

3. Exercises to do by the students

The lab work for the students is:

1. Download the datasets `wine.data` and `hepatitis.data` from the TEAMS.

2. Calculate the accuracy, Cohen kappa and confusion matrix for both datasets using the LDA classifier using the whole dataset as training and test set.
3. Repeat the process using cross-validation with 4 folds.
4. Implement the cross-validation using the leave-one-pattern-out approach and provide the results. In this case, the process training-test is repeated N times, each one excluding a pattern.
5. Use the LDA classifier for the classification of the textures in the exercise of previous unit 4.

Submit before 8 January by TEAMS the results and difficulties founded. It can be done individually or by groups.