# Project 2: Fisher's Linear Discriminant Analysis

## 1 Fisher's LDA via a toy problem

### 1.1 Introduction

The aim of the Fisher's Linear Discriminant Analysis is to find a separation vector $v$ to classify data points. This vector can be found by solving the maximization problem:

$$\max \frac{(v^T m_A - v^T m_B)}{v^T(\Sigma_A + \Sigma_B)v} \tag{1}$$

where $\Sigma_C$ is the covariance matrix and $m_C$ are the sampled means for $C \in \{A, B\}$. We can rewrite (1) as

$$R = \frac{v^T S v}{v^T M v} \tag{2}$$

where $S = (m_A - m_B)(m_A - m_B)^T$ and $M = \Sigma_A + \Sigma_B$. From the exercises we know that the vector $v$ maximizes (2) if and only if $Sv = \lambda M v$ is satisfied. Since $S$ has rank one, we know $Sv$ has the direction $m_A - m_B$. This implies $Mv$ has to be in the direction to satisfy $Sv = \lambda M v$.

Together we can conclude $v = M^{-1}(m_A - m_B) = (\Sigma_A + \Sigma_B)^{-1}(m_A - m_B)$ maximizes (1). The separation vector $v$ will then be used to classify data points $x$ based on its projection. We classify $x$ to class A, if $v^T x > c$ for a constant $c$, and to class B, if $v^T x \leq c$.

### 1.2 Building a classifier for the toy problem

For the given toy problem in the task sheet, we aim to build a classifier. First, for given datasets *DataA* and *DataB* we split the data points into training and testing data for each dataset, where each column represents a sample.

We will now create the mean value vector $m_C$ and the covariance matrix $\Sigma_C$ for $C \in \{A, B\}$ to create the separation vector $v$.

```
1  % sample mean & covariance
2  mA = mean(TrainA')';
3  mB = mean(TrainB')';
4  sA = cov(TrainA');
5  sB = cov(TrainB');
```

```matlab
6 % separation vector
7 v = (sA+sB)\(mA-mB);
8 v = v/norm(v);
```

Now we can classify the data points, i.e the column vectors, for the data set *DataA* and *DataB*, while counting the false classifications, as described before:

```matlab
1  c = v'*(mA+mB)/2;
2  classifyA = 0;
3  classifyB = 0;
4  missA = 0;
5  missB = 0;
6  for i=1:size(TestA,2)
7      if v'*TestA(:,i) > c
8          classifyA = classifyA+1;
9
10     else
11         missA = missA + 1;
12     end
13 end
14
15 for i=1:size(TestB,2)
16     if v'*TestB(:,i) <= c
17         classifyB = classifyB+1;
18
19     else
20         missB = missB + 1;
21     end
22 end
```

From the vectors *MissA* and *MissB* we can compute the miss rate and the success rate of the classifier.

```matlab
1  missRate = (missA+missB)/(size(TestA,2)+size(TestB,2));
2  successRate = 1 - missRate
```

In this case we get the *successRate* = 1 for our toy problem.

# 2 UCI Benchmark Problems

## 2.1 Building a classifier for different data sets

We will now build the classifier for the data set *sonar.mat* provided by the UCI Machine Learning Repository. The classifier for *ionosphere.mat* works analogously (apart from the singularity problem we will discuss later).

The dataset *sonar.mat* consists of the $208 \times 60$-matrix *sonar_ data*, where 208 data points are stored in the rows, each with 60 features stored in the columns. Further there is the vector *sonar_ label* with 0-1 entries representing the classes of each row of the data matrix (0 means class A and 1 means class B).

First we will separate *sonar_ data* into the vector *SonarA* and *SonarB*, where the data points with label 0 are stored in *SonarA* and with label 1 in *SonarB*.

```
1  load sonar.mat
2  sizeOfSonarA = sum(sonar_label == 0);
3  sizeOfSonarB = sum(sonar_label == 1);
4
5  sonarA = zeros(sizeOfSonarA, 60);
6  sonarB = zeros(sizeOfSonarB, 60);
7
8  %Store all sonar_data == 0 in SonarA
9
10 for i = 1:size(sonar_label)
11     if sonar_label(i) == 0
12         sonarA(i,:) = sonar_data(i,:);
13     end
14 end
15
16 %Deleting the 0 rows
17 ind = find(sum(sonarA,2)==0) ;
18 sonarA(ind,:) = [] ;
19
20 %Same for SonarB (sonar_data == 1)
21
22 for i = 1:size(sonar_label)
23     if sonar_label(i) == 1
```

```
24            sonarB(i,:) = sonar_data(i,:);
25        end
26  end
27
28
29  ind = find(sum(sonarB,2)==0) ;
30  sonarB(ind,:) = [] ;
```

This will give us 2 matrices *SonarA* and *SonarB* with data points in class A and B. We now can split these matrices into training and testing data:

```
1  %Create test and train data of SonarA and SonarB
2  m = ceil(size(sonarA,1)*(7/10));
3
4  trainSonarA = sonarA(1:m,:);
5  testSonarA = sonarA(m+1:end,:);
6
7  m = ceil(size(sonarB,1)*(7/10));
8
9  trainSonarB = sonarB(1:m, :);
10 testSonarB = sonarB(m+1:end,:);
```

From the test vectors *testSonarA* and *testSonarB* we can create the mean value vector and the covariance matrix to create the separation vector *vSonar* and the constant *cSonar*.

```
1  %Create test and train data of SonarA and SonarB
2  mSonarA = mean(trainSonarA)';
3  mSonarB = mean(trainSonarB)';
4  sSonarA = cov(trainSonarA);
5  sSonarB = cov(trainSonarB);
6
7  %Create the separation vector and the constant for the
       classification
8  vSonar = (sSonarA+sSonarB)\(mSonarA-mSonarB);
9  vSonar = vSonar/norm(vSonar);
10
11 cSonar = vSonar'*(mSonarA + mSonarB)/2;
```

Like in the toy problem from the first part, we have all the tools to test and report the success rate for classification by the linear discriminant analysis.

```
1  classifySonarA = 0;
2  classifySonarB = 0;
3  missSonarA = 0;
4  missSonarB = 0;
5
6  for i=1:size(testSonarA,1)
7      if vSonar'*testSonarA(i,:)' > cSonar
8          classifySonarA = classifySonarA+1;
9
10     else
11         missSonarA = missSonarA + 1;
12     end
13 end
14
15 for i=1:size(testSonarB,1)
16     if vSonar'*testSonarB(i,:)' <= cSonar
17         classifySonarB = classifySonarB+1;
18
19     else
20         missSonarB = missSonarB + 1;
21     end
22 end
23
24 missRateSonar = (missSonarA+missSonarB)/(size(testSonarA,1)+
       size(testSonarB,1));
25 successRateSonar = 1 - missRateSonar
```

This gives us following output:

```
successRateSonar =

    0.7742
```

For the data set *ionosphere.data* we follow the same procedure. One thing we notice is that for the covariance matrices $\Sigma_A$ and $\Sigma_B$ of *trainIonoA* and *trainIonoB* the sum $\Sigma_A + \Sigma_B$

is singular. Most likely one feature has such small values that it leads to rounding errors and the values are stored as 0, instead of very small numbers. This leads to a 0 column, which means the $det(\Sigma_A + \Sigma_B) = 0$, i.e the matrix is singular.

If we disregard that $\Sigma_A + \Sigma_B$ is singular, we get following success rate for ionosphere classification:

```
successRateIono =

     0
```

To avoid the singularity problem, instead of the inverse, we will now compute the pseudo inverse of $\Sigma_A + \Sigma_B$ to compute the separation vector *vIono*:

```
1  vIono = pinv(sIonoA+sIonoB)*(mIonoA-mIonoB);
2  vIono = vIono/norm(vIono);
```

This leads to following success rate:

```
successRateIono =

     0.8544
```

## 2.2   Trying different thresholds for constant c

In this part we are trying out different constants between $v^T m_A$ and $v^T m_B$ to see how sensitive the classification rate is to the choice of the constant $c$.

Computing the boundaries of the *sonar.mat* data set we get:

$$v^T m_A = -0,0150 \text{ and } v^T m_B = -0,0353$$

For $c = -0.02$ we get following output:

```
successRateSonar =

     0.7581
```

and for $c = -0.033$ we get:

```
successRateSonar =

    0.7097
```

After trying out several different constants we can conclude that the closer the constant $c$ gets to the boundaries $v^T m_A$ and $v^T m_B$, the worse the LDA classification, i.e the success rate, gets.
The same phenomenon we can observe while changing the constant $c$ in the $ionosphere.data$ data set. Here we notice, the further the boundaries $v^T m_A$ and $v^T m_B$ are apart, the worse the LDA classification gets, when $c$ is close to the boundaries.