

## Group report for Assignment 1

Group number: **XX**

Name of the group members with student numbers:

- s0000000, first name last name
- s0000001, first name last name
- s0000002, first name last name
- s0000003, first name last name

## General guidelines

Please add your answer to the places indicated. If there is no clear place indicated make sure that your answers are clear and structured according to the instructions in the assignment.

## Task 1: Pairwise correlation coefficients

Maximum obtainable points: **05**

1. The correlation coefficients are shown in Table 1. To obtain the correlation coefficients, the following code was used:

```

1 corr_mat = [0 0 0; 0 0 0; 0 0 0];
2
3 for r = 1:3
4     for c = 1:3
5         corr_mat(r, c) = coeff(lab1_1(:,r), lab1_1(:,
6             c));
7     end
8 end
9 function corr_coeff = coeff(f1, f2)
10     covar = sum((f1 - mean(f1)) .* (f2 - mean(f2)))/(
11         length(f1) - 1);
12     corr_coeff = covar/sqrt(var(f1) * var(f2));
13 end

```

2. Add the two scatter plots here.

Conclusions:

- Figure 1 shows the scatter plot of the height and body weight features. The figure displays an increase in height with an increase in body weight, indicating a positive linear correlation.

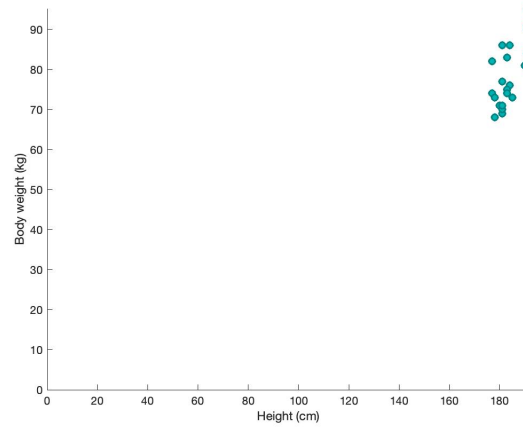


Figure 1: 2D scatter plot of Height (cm) and Body weight (kg) with  $\text{corr} = 0.7156$ .

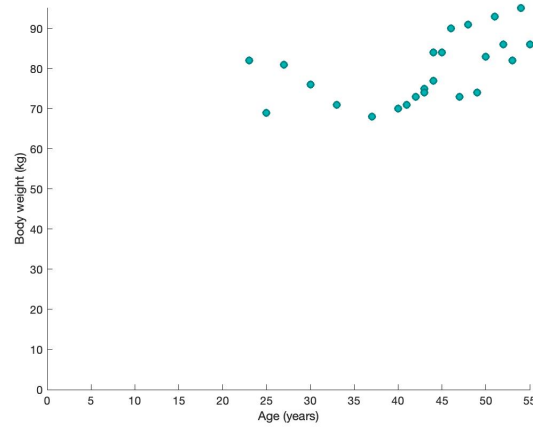


Figure 2: 2D scatter plot of age (years) and Body weight (kg) with  $\text{corr} = 0.5142$ .

	1	2	3
1	1	-0.0615	0.7156
2	-0.0615	1	0.5142
3	0.7156	0.5142	1

Table 1: Correlation coefficients.

- Figure 2 shows the scatter plot of the two features with the second largest correlation coefficient (age and body weight). The figure shows a decrease in body weight with an increase in age, after which it shows an increase in body weight with an increase in age. Therefore, the correlation between age and body weight features is non-linear.

## Task 2: Hamming distance

Maximum obtainable points: **15**

1. To sample from the files, a string array is created containing all file names. To create the S group, a file is randomly sampled from the string array, from which two rows are randomly sampled without replacement. To create the D group, two files are randomly sampled without replacement, and a row is selected. The Hamming distance for both the S and D groups is computed with the `pdist` function. The code used to obtain the groups S and D is shown below.

```

1 load(" Assignment1-files /Task_2/person01.mat");
2 path = " Assignment1-files /Task_2/";
3
4 %1
5 filenames = [];
6 for i = 1:20
7     h = sprintf('person%02d.mat',i);
8     filenames = [filenames; char(h)];
9 end
10
11 % a)
12 reps = 1000;
13 HD_S = zeros(1, reps);
14
15 for rep = 1:reps
16     file = datasample(filenames, 1);
17     load(path + file);
18     y = datasample(1:20,2,'Replace',false);
19     r1 = iricode(y(1), :);
20     r2 = iricode(y(2), :);
21     HD_S(1, rep) = pdist([r1;r2], 'hamming');
```

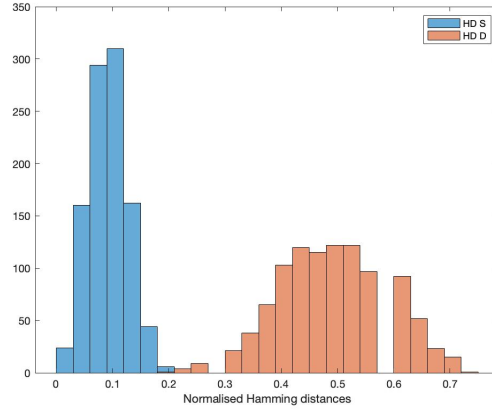


Figure 3: Histograms of S and D groups of Hamming distances.

```

22 end
23
24 % b)
25 HD.D = zeros(1, reps);
26
27 for rep = 1:reps
28     files = datasample(filenamees, 2, 'Replace', false
29     );
29     load(path + files(1, :));
30     r1 = iricode(datasample(1:20, 1), :);
31     load(path + files(2, :));
32     r2 = iricode(datasample(1:20, 1), :);
33     HD.D(1, rep) = pdist([r1;r2], 'hamming');
34 end

```

2. *How much do the two histograms overlap?*

The histogram is displayed in Figure 3, which shows that the two histograms slightly overlap at the normalized Hamming distance 0.2.

3. The means and variances of groups S and D are displayed in Table 2. They were computed using Matlab's `mean` and `var` functions.

	$S$	$D$
<i>Mean</i>	0.0861	0.4889
<i>Variance</i>	0.0015	0.0092

Table 2:

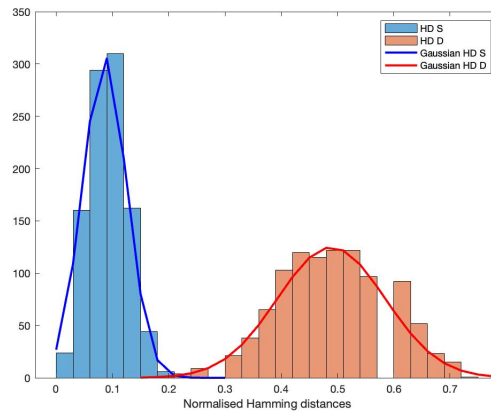


Figure 4: Gaussian distributions over histograms of S and D groups.

4. The histograms are shown in Figure 4. The Gaussian distributions were scaled by multiplying with the bin width and the number of samples so that the area under the curve would resemble that of the histograms. The Gaussian distributions were obtained through Matlab's `normpdf` function after calculating the standard deviations.
5. (a) Decision criterion = 0.19. The decision criterion was estimated by iterating over the values of the x-axis and obtaining the x-value for which the cumulative distribution function reached a value equal to or higher than the false acceptance error.  
 (b) False rejection rate = 0.0038. The false rejection rate was computed through subtracting the cdf value at the decision criterion of group S from the cdf at the maximum hamming distance of group S.
6. Person 5. By taking the mean Hamming distances of the test person's iriscodes to each other person's iriscodes using the mask, results in a minimum mean distance for person 5, making them the most likely to be the test person.

### Task 3: Covariance matrix

Maximum obtainable points: **05**

1. Means: [5.8 5.0 6.2]  
 Covariance matrix:

The code used to compute the covariance matrix and the means is listed below.

3.2000	0.2500	-0.4500
0.2500	2.5000	-3.7500
-0.4500	-3.7500	5.7000

```

1 feat_vec = [4 5 6; 6 3 9; 8 7 3; 7 4 8; 4 6 5];
2
3 cov_mat = get_cov(feat_vec);
4
5 means = [0 0 0];
6 for i = 1:3
7     means(i) = sum(feat_vec(:, i))/length(feat_vec(:,
8         i));
9
10 mvnpdf([5 5 6], means, cov_mat)
11 mvnpdf([3 5 7], means, cov_mat)
12 mvnpdf([4 6.5 1], means, cov_mat)
13
14 function cov_mat = get_cov(feat_vec)
15     cov_mat = [0 0 0; 0 0 0; 0 0 0];
16
17     for r = 1:3
18         f1 = feat_vec(:, r);
19         mean_f1 = sum(f1)/length(f1);
20         for c = 1:3
21             f2 = feat_vec(:, c);
22             mean_f2 = sum(f2)/length(f2);
23             cov_mat(r, c) = sum((f1 - mean_f1) .* (f2
24                 - mean_f2))/(length(f1) - 1);
25         end
26     end
end

```

2. Probability densities of:  $[5 \ 5 \ 6] = 0.0543$ ,  $[3 \ 5 \ 7] = 6.1287\text{e-}04$ ,  $[4 \ 6.5 \ 1] = 7.0300\text{e-}29$ .

## Task 4: 2D Gaussian

Maximum obtainable points: **05**

1. The Gaussian distribution is displayed in Figure 5, for which the code below was used.

```

1 x = -10:.1:10;

```

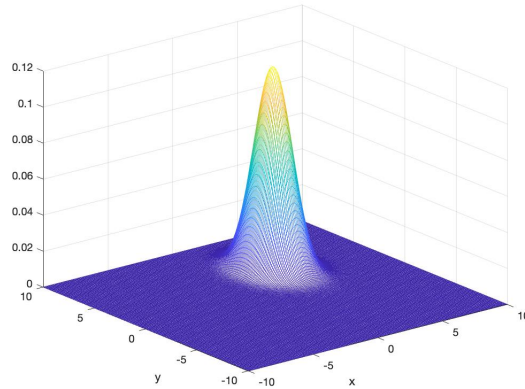


Figure 5: 3D plot of 2D Gaussian distribution.

```

2 [X1, X2] = meshgrid(x, x);
3 X = [X1(:) X2(:)];
4 sigma = [1 0; 0 2];
5 mean = [3 4];
6
7 gauss = mvnpdf(X, mean, sigma);
8 gauss = reshape(gauss, length(x), length(x));
9
10 mesh(-10:.1:10, -10:.1:10, gauss).

```

2. Please give the Mahalanobis distance between the points and the mean here. It is recommended to use code to calculate this, however you are also free to calculate it manually. Please include the code or the calculations as well.:  $[10 \ 10]' = 8.1854$ ,  $[0 \ 0]' = 4.1231$ ,  $[3 \ 4]' = 0$ ,  $[6 \ 8]' = 4.1231$

The code with which the distances were computed is listed below.

```

1 points = [10 10; 0 0; 3 4; 6 8];
2
3 for i = 1:length(points)
4     mahalanobis = mahal_dis(points(i,:), mean', sigma
5 );
6 end
7
8 function mahalanobis = mahal_dis(x, y, cov_mat)
9     mahalanobis = sqrt((x - y)' * inv(cov_mat) * (x -
10 y));
11 end

```

## Task 5: Naive Bayesian rule

Maximum obtainable points: **05**

1. (a) No Spam. The probability of  $P(\text{Spam}|\text{Customers}, \text{Watches}) = 1.350 \times 10^{-6}$ , and the probability of  $P(\text{NoSpam}|\text{Customers}, \text{Watches}) = 1.400 \times 10^{-5}$ . Therefore,  $P(\text{NoSpam}|\text{Customers}, \text{Watches}) > P(\text{Spam}|\text{Customers}, \text{Watches})$ .
- (b) Spam. The probability of  $P(\text{Spam}|\text{Fun}, \text{Vacation}) = 3.375 \times 10^{-8}$ , and the probability of  $P(\text{NoSpam}|\text{Fun}, \text{Vacation}) = 9.800 \times 10^{-9}$ . Therefore,  $P(\text{Spam}|\text{Fun}, \text{Vacation}) > P(\text{NoSpam}|\text{Fun}, \text{Vacation})$ .

The probabilities were computed with the code below.

```

1 spam_prior = 0.9;
2 nospam_prior = 0.1;
3
4 customers_spam = 0.005;
5 customers_nospam = 0.035;
6
7 watches_spam = 0.0003;
8 watches_nospam = 0.000004;
9
10 fun_spam = 0.00015;
11 fun_nospam = 0.0007;
12
13 vacation_spam = 0.00025;
14 vacation_nospam = 0.00014;
15
16 %a)
17 spam_prob = spam_prior * customers_spam *
    watches_spam
18 nospam_prob = nospam_prior * customers_nospam *
    watches_nospam
19
20 %b)
21 spam_prob = spam_prior * fun_spam * vacation_spam
22 nospam_prob = nospam_prior * fun_nospam *
    vacation_nospam

```

## Task 6: Decision Tree

Maximum obtainable points: **05**

The decision tree is shown in Figure 6.

For both features, size and colour, the gini impurities were computed in order to determine the root node.



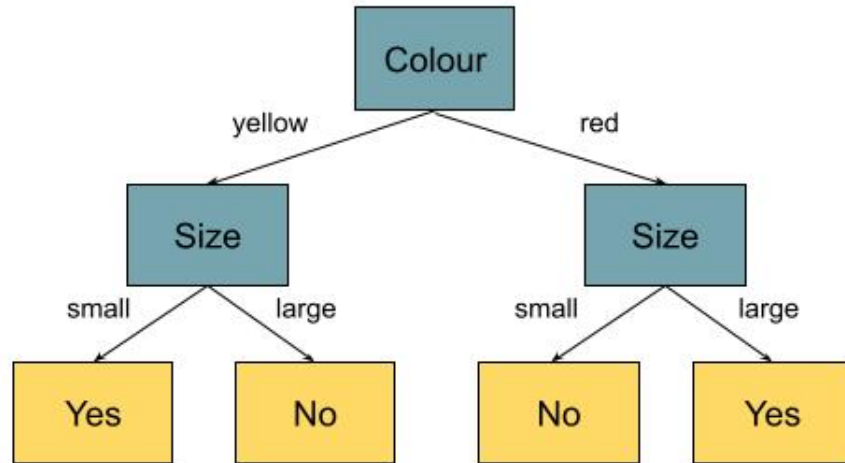


Figure 6: Decision tree best fitting the given data.

The size feature:

4 yes and 2 no corresponded to the small option, resulting in a gini impurity of  $1 - (\frac{4}{6})^2 - (\frac{2}{6})^2 = \frac{4}{9}$ .

1 yes and 1 no corresponded to the large option, resulting in a gini impurity of  $1 - (\frac{1}{2})^2 - (\frac{1}{2})^2 = \frac{1}{2}$ .

The total gini impurity is then:

$$\frac{6}{8} \cdot \frac{4}{9} + \frac{2}{8} \cdot \frac{1}{2} = \frac{11}{24}.$$

The colour feature:

4 yes and 1 no corresponded to the colour yellow, resulting in a gini impurity of  $1 - (\frac{4}{5})^2 - (\frac{1}{5})^2 = \frac{8}{25}$ .

1 yes and 2 no corresponded to the colour red, resulting in a gini impurity of  $1 - (\frac{1}{3})^2 - (\frac{2}{3})^2 = \frac{4}{9}$ .

The total gini impurity for the colour feature is then:

$$\frac{5}{8} \cdot \frac{8}{25} + \frac{3}{8} \cdot \frac{4}{9} = \frac{11}{30}.$$

Since  $\frac{11}{30} < \frac{11}{24}$ , the colour feature has the lowest total gini impurity, and thus is the root node of the decision tree.

In both the right and left branches are the nodes for size. For the left branch, the small option led to only yes, and the large option to only no, hence the nodes correspond to these findings in the data. Likewise the small decision for the size node in the right branch led to only no, and the large decision only yes.

## Task 7: Learning Vector Quantization

Maximum obtainable points: **10**

1. Please add scatter plots of the datasets and your choice and (brief) reasoning for the number of prototypes.
2.

Add your implementation of LVQ here.
3. Please include 1 plot containing the error curves for the four cases and four scatter plots with the different prototypes here. Note that the scatter plots should contain the labels assigned by LVQ to all points.

## Task 8: Cross-Validation

Maximum obtainable points: **05**

1. 

Add your code **for** 10-fold **cross** validation here.

2. Test error:

1 

Your code to compute the test **error** here (or  
alternatively include it in the code **for** the  
**cross** validation)

3. Please add the bar plot here.

## Task 9: ROC

Maximum obtainable points: **05**

1. Your answer to part 1 here.

2. You answer to part 2 here.

## Task 10: Transforms

Maximum obtainable points: **15**

Part 1:

Part 2:

Part 3:

Part 4:

## Task 11: K-Nearest Neighbor

Maximum obtainable points: **10**

1. 

Your implementation of k-nn here.

2. Please add scatter plots where the points are colored as the labels that K-nn (for  $k = 1, 3, 5, 7$ ) assigns them here.
3. Please add your answer here.
4. Please add your answer here.

## Task 12: K-means

Maximum obtainable points: **15**

1. Please add your answer to part 1.
2. Please add your answer to part 1.

## Individual contributions

- who did what?
- or, How the tasks were divided?
- or, How the load was distributed?