

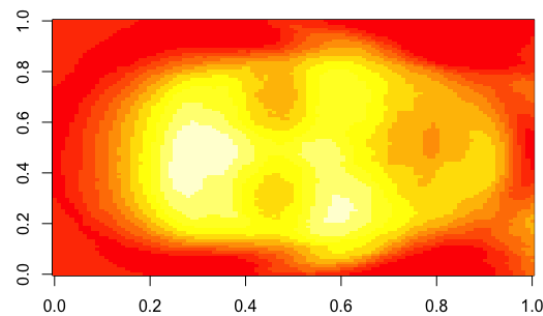
Machine Learning - Coursework 1

- Calculate and plot the average face of the training set, then write a function to find a PCA basis of size M, where the inputs will be M and X, the matrix containing the training set. Clearly describe all aspects of your function, then use it to plot the first 5 eigenfaces of the training set.

To calculate the average face of the training set, I simply calculate the average pixel values for each pixel across the trainings set.

```
1 library(rARPACK)
2 library(philentropy)
3
4
5 # I load the raw csv's
6 faces.train.inputs <- read.csv("./2018_ML_Assessed_Coursework_1_Data/
7                               Faces_Train_Inputs.csv", head=FALSE)
8 faces.train.label <- read.csv("./2018_ML_Assessed_Coursework_1_Data/
9                               Faces_Train_Labels.csv", head=FALSE)
10 faces.test.inputs <- read.csv("./2018_ML_Assessed_Coursework_1_Data/
11                               Faces_Test_Inputs.csv", head=FALSE)
12 faces.test.label <- read.csv("./2018_ML_Assessed_Coursework_1_Data/
13                               Faces_Test_Labels.csv", head=FALSE)
14
15
16 # I turn the input values into a list of 320 matrices, each matrix a 112 x 92 value
17 # of pixels corresponding to each image .. I need to use lapply again on the
18 #result because apply gives the matrices in a weird form
19 faces.train.inputs.cleaned <- lapply(apply(X=faces.train.inputs,
20                                           MARGIN=1,
21                                           function(x) list(matrix(as.numeric(x),
22                                                                    nrow = 112))), "[", 1)
23
24 # Here I calculate the average face
25 avg.face <- Reduce('+', faces.train.inputs.cleaned) /
26   length(faces.train.inputs.cleaned)
27 image(avg.face)
```

Figure 1: Average face of the training set



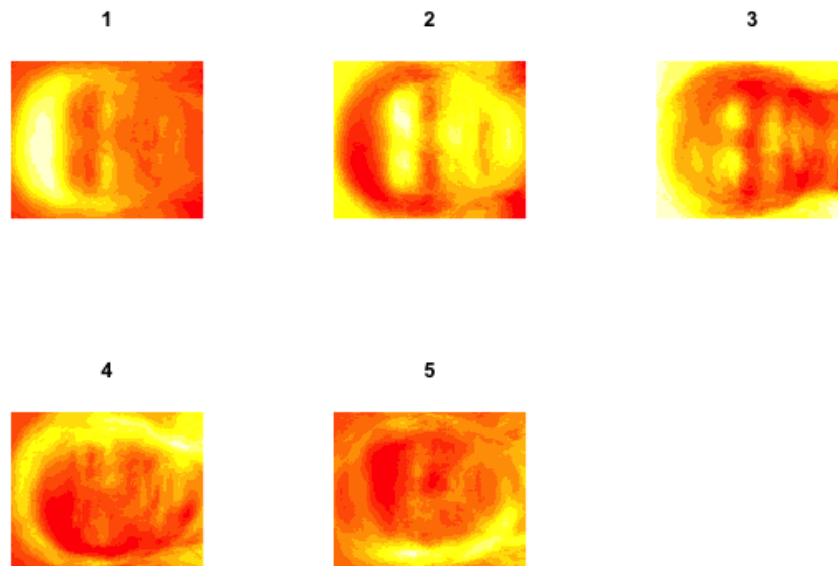
The following function returns the PCA basis of size M as specified. I added a default parameter which also allows access to the eigenvalues associated with each vector of the eigenbasis.

The function works exactly as the PCA described in lectures. Calculate the centralised data matrix, X , and then calculate the first M eigenvectors for the covariance matrix $\frac{XX^T}{n}$.

```

1 find.pca.basis <- function(M,X, return.full.results = FALSE){
2
3   n <- dim(X)[1] # The number of images
4
5   # Turn the input data into a matrix and transpose it
6   X.data.matrix <- data.matrix(t(X))
7
8   # Centralise the data matrix
9   means <- rowMeans(X.data.matrix) # calculate row means
10  data.matrix.centralised <- X.data.matrix - means %*% t(rep(1,n)) # and subtract
11
12  # Calculate the covariance matrix as defined in lectures
13  covariance.matrix <- (data.matrix.centralised %*% t(data.matrix.centralised)) / n
14
15  # Now I need to compute the first M eigenvectors/ eigenvalues using the
16  # R package rARPACK
17  results <- eigs_sym(covariance.matrix,k=M,which="LM")
18
19  if (return.full.results){
20    return(results)
21  } else{
22    return(results$vectors)
23  }
24
25 }
```

Figure 2: First 5 Eigenfaces of the training set



The code below then plots the first 5 eigenfaces of the training set. This means simply calculating the eigenbasis for $M = 5$, and then plotting the vectors in the resulting eigenbasis.

```
1 eigenbasis <- find.pca.basis(5, faces.train.inputs)
2
3 par(mfrow=c(2,3))
4 for (i in 1:5){
5   # eigenbasis[,i] corresponds to the i'th eigenvector.
6   image(matrix(eigenbasis[,i], nrow = 112), useRaster=TRUE, axes=FALSE, main=i)
7 }
8 par(mfrow=c(1,1))
```

The results can be seen in figure 2.

- Choose a single face and project it into a PCA basis for dimension $M = 5, 10, 50$, then plot the results.

Here is the code to project the first image of the training set onto the PCA basis for dimensions 5, 10 and 50.

```
1 # initialise the dimensions and face used.
2 dimensions <- c(5,10,50)
3 single.face <- 1
4 means <- as.vector(avg.face) # convert the mean face back into a vector
5
6 # iterate through the dimensions considered
7 for (i in dimensions){
8
9   # compute the eigenbasis using the function created
10  eigenbasis <- find.pca.basis(i,faces.train.inputs)
11  # calculate the projection values in this pca basis
12  projection.vals <- t(as.numeric(faces.train.inputs[single.face,]) -
13                      means) %*% eigenbasis
14  # calculate the actual face in this basis
15  projection.vector <- eigenbasis %*% as.numeric(as.list(projection.vals))
16
17  # carry out the plot
18  image(matrix(projection.vector, nrow = 112),useRaster=TRUE,
19        axes=FALSE,main=paste("dimension",i,sep=" "))
20 }
21
22 # and here's the original image
23 image(matrix(as.numeric(faces.train.inputs[single.face,]), nrow = 112),
24        useRaster=TRUE, axes=FALSE, main="Original Image")
```

Please see Figure 3 for the results of this code.

- Plot a graph of the mean squared error of each lower dimensional approximation of this chosen face, with the dimensionality plotted along the x-axis. Is there a clear point at which we can choose a good approximation? Discuss how we should choose the appropriate dimensionality of the approximation.

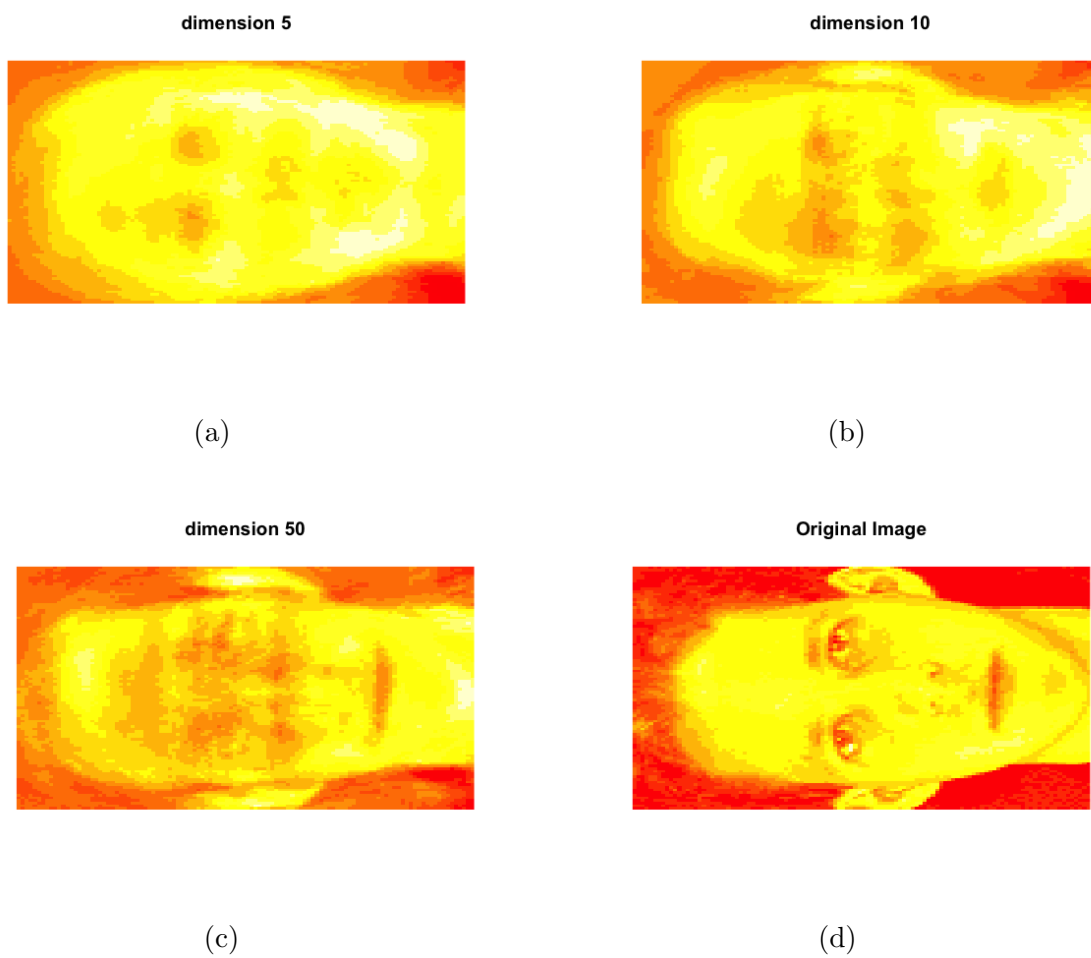


Figure 3: The projection of training image 1 onto PCA bases of different dimensions