

Machine Learning: Assignment 2

Karsten Standnes - STNKAR012

October 2017

1 Introduction

In this assignment I'll be looking at the performance of different models when fitting a data created by a unknown mode. In "Task 1" two linear models with an offset is fitted to dataset and the relative performance is measured. "Task 2" looks at making legendre based model and fit them to a dataset made from a sinus function. Different values for λ , the regularization parameter is used. In "Task 3" a set of facial images are processed and used to produce eigenfaces. This task evolves around eigenfaces and how they can be used to reconstruct a picture.

Contents

1	Introduction	1
2	Task 1	2
2.1	Task 1 i	2
2.2	Task 1 ii	3
3	Task 2	4
3.1	Task 2 i	4
3.2	Task 2 ii	5
3.3	Task 2 iii	6
4	Task 3	7
4.1	Task 3 i	7
4.2	Task 3 ii	8
4.3	Task 3 iii	8
5	Appendix	10
5.1	R-Code: Question 1	10
5.2	R-Code: Question 2	16
5.3	R-Code: Question 3	23

2 Task 1

In this task we look at the linear function of the form:

$$y_i = 0.8x_i + \epsilon_i; \quad -1 \leq x \leq 1, i = 1, \dots, N$$

(1)

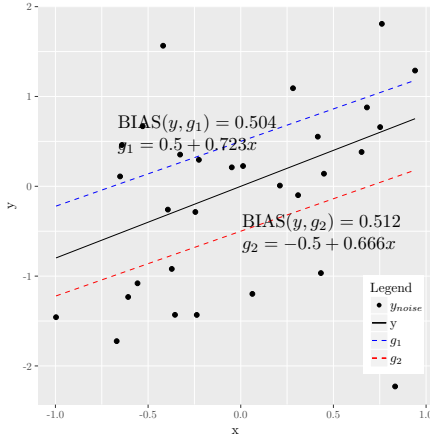
where: $\epsilon_i \sim \text{Normal}(0, 1)$

2.1 Task 1 i

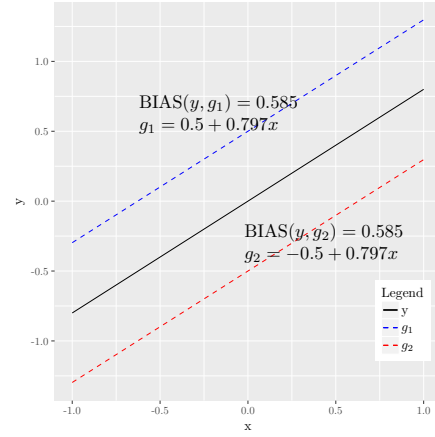
First a dataset of size $N = 30$ is created from using equation 1 using a set of x values generated from the $\text{Uniform}(-1, 1)$. Then two models are attempted fit two the dataset. The models fitted are:

$$g_1(x) = 0.5 + b_1x$$

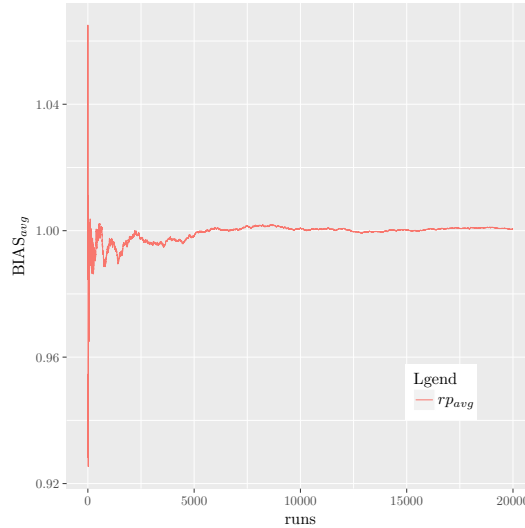
$$g_2(x) = -0.5 + b_2x$$



(a) g_1 and g_2 fitted to the noisy data.



(b) Average fitted g_1 and g_2 over 20000 data sets.



(c) Relative bias $\frac{\text{BIAS}(y, g_1)}{\text{BIAS}(y, g_1)}$ for 1 to 20000 runs.

Figure 1: .

In figure 1a we see the underlining function together with a dataset created from equation 1 and the linear functions $g_1(x)$ and $g_2(x)$ fitted to the data. Both functions in this case performs fairly well with pretty similar slope, with $g_1(x)$ performing slightly better. This being the result from only one dataset means that it can't be viewed as representative for the performance of the functions. To get a better picture, 20000 dataset are created and averaged over. Figure 1b shows $g_1(x)$ and $g_2(x)$ after 20000 runs. We can see that the bias and slope is the same for both giving a relative performance (rp) equal to 1. In figure 1c below the relative performance is plotted as a function of average relative bias after 1 to 20000 runs. After around 5000 runs the relative bias starts to converge towards 1, meaning the models perform equally well. The functions have a constant ± 0.5 making it impossible to perfectly fit the underlying function. Intuitively the best models will have the slope with a offset of ± 0.5 . This is confirmed when the slope of the functions converge to $b_1, b_2 \approx 0.8$, the slope of the original underlying function.

2.2 Task 1 ii

By using the model in equation 1, 10000 datasets of size $N = 30$ are simulated. Each set is split into 21 combinations of validation and training sets, where the validation set is of the size i and the training set of the size $30 - i$, with $i \in [5, 25]$. For each of the set the models $g_1(x)$ and $g_2(x)$ are attempted fit to the training data and validated on the validation set. For each dataset and i the best of the two models is selected. For each model selected the error measures E_{out} and E_{val} is calculated and the average over all the datasets is taken to give the errors for each combination of training and validation set. Below the averaged best for each i is plotted.

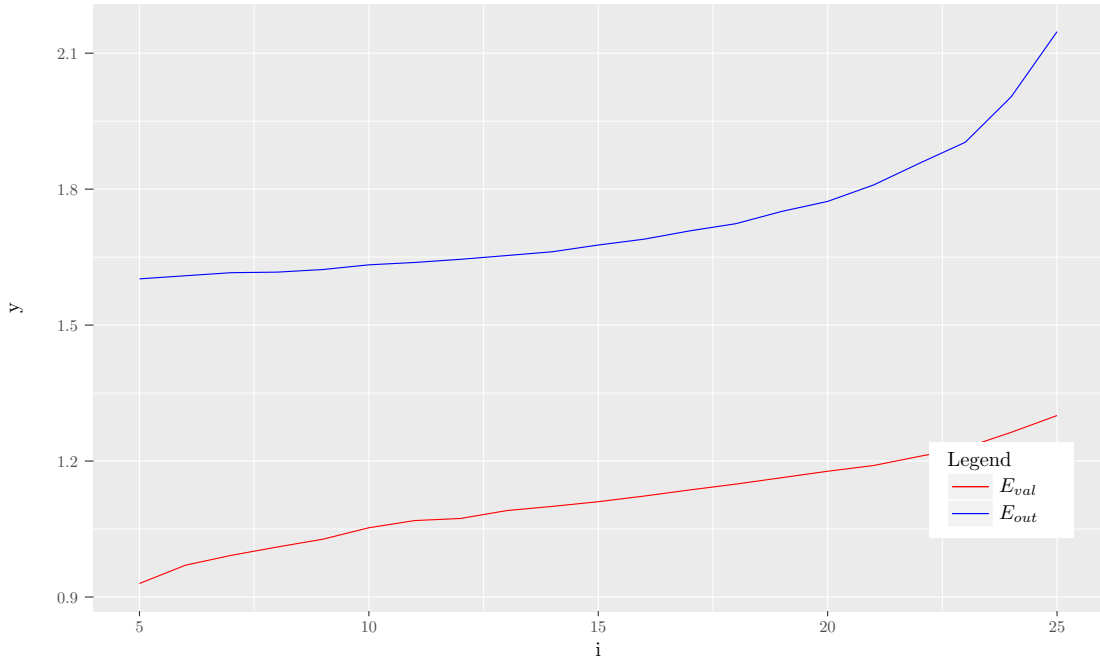


Figure 2: Error measures E_{out} and E_{val} are plotted for different sizes of training and validation sets.

In figure 2 it is a clear difference between the value of E_{val} (use MSE) and E_{out} (use BIAS) which is pretty stable for different sizes of the training and validation set. We can see that both the errors for increases for increasing i . E_{out} starts with incrementing slowly, then more rapidly around $i \sim 15 - 20$, while E_{val} has a more steady increase with a faster increase in the start ($i < 10$). We can deduce from this that E_{val} will consequently give a lower value than E_{out} and that decreasing the training set in advantage for increasing the validation set give a worse performance, at least for datasets of size $N = 30$. This makes sense since the less data the models have available during

training, the more likely it is to fit the noise of the data. Then increasing the validation set does not help other than giving a more precise measure of the performance. That being said, having a large enough validation set is also important to keep up the precision of the measure of error. Here the conclusions of it being a disadvantage to decrease the training set can be drawn using the average of 10000 sets of data.

3 Task 2

In "Task 2" we look at the function of the form:

$$y_i = \sin(\pi x_i) + \epsilon_i; \quad -1 \leq x \leq 1, i = 1, \dots, N$$

(2)

where: $\epsilon_i \sim \text{Normal}(0, 1)$

For each of the subtasks a dataset/datasets of size $N = 50$ with $x \sim \text{Uniform}(-1, 1)$.

3.1 Task 2 i

Simulate a dataset using equation ?? and plot it together with the underlying model:

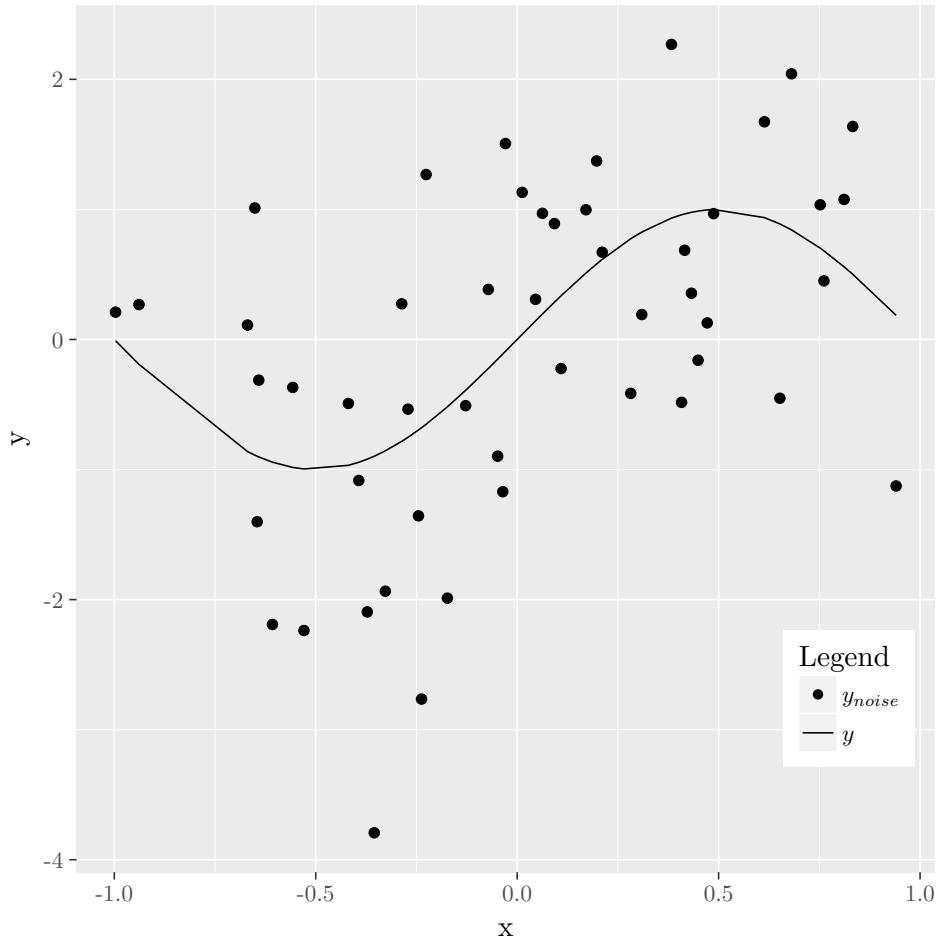


Figure 3: Underlying function with simulated dataset using equation 2

3.2 Task 2 ii

In this subtask two models are fitted to the noisy data, the models are based on

$$y_i = \sum_{q=0}^{Q_f=10} \beta_q L_q(x) \quad (3)$$

, $L_q(x)$ being the Legendre polynomial of q -th order. The Legendre polynomial is given by:

$$L(x) = 2^q \sum_{k=0}^q x^k \binom{q}{k} \binom{\frac{q+k-1}{2}}{q} \quad (4)$$

. The two models have a regularization parameter λ equal to 0 and 5 for the models. Using linearization to fit the models to the noisy data:

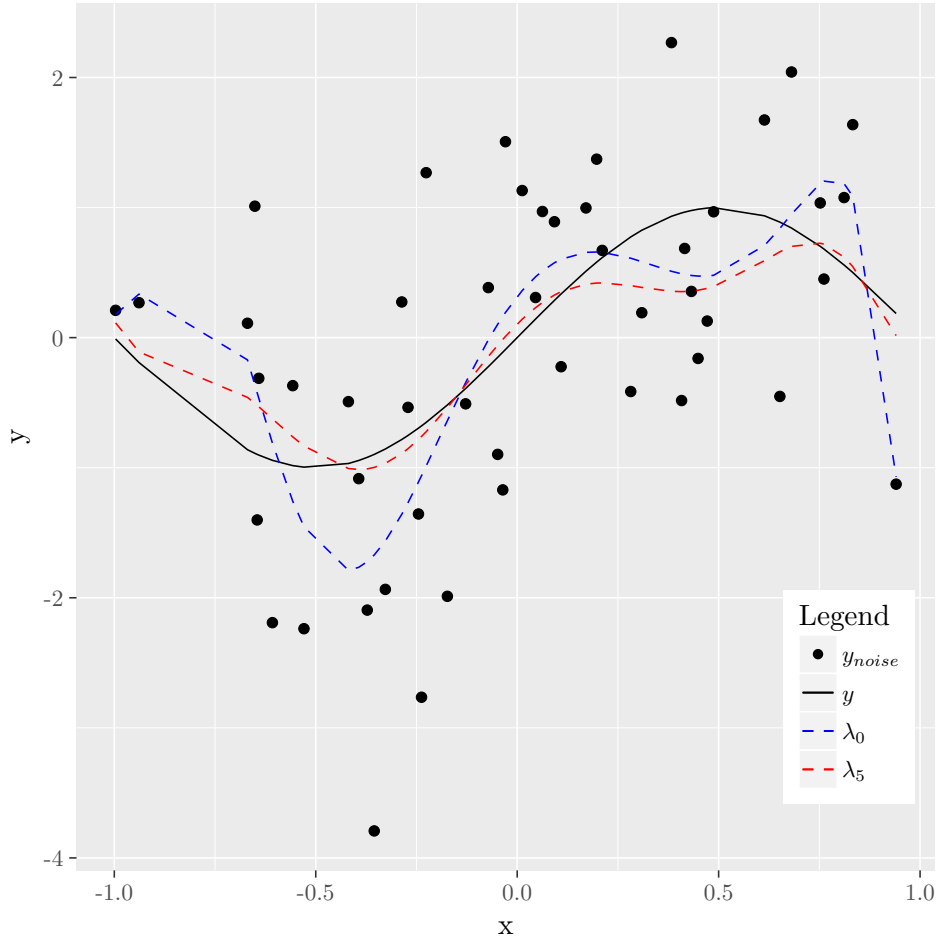


Figure 4: Two legendre models are fitted to the dataset with different values for λ .

From the figure 4 we can see that the model with a higher regularization parameter give a result that is much more similar to the underlying function. The model with $\lambda = 0$ has much more rapid change in direction, leading to it quickly missing the underlying function. This indicate that regularization increase the precision of the model compared to no regularization, at least for regularization parameter $\lambda = 5$.

3.3 Task 2 iii

To try to find the best regularization parameter for the legendre based model we can use cross-validation. In this task 10-fold cross-validation is used, meaning the dataset is split into 10 folds where each fold is used as a validation set with the other nine as the training set. This is done for all folds, then the mean-squared-error is calculated for the results of all the folds. Below are two plots, figure 5a giving the best λ based on the CV-error for 1000 dataset and figure 5b showing a model fitted with the best found λ on the data from figure 4.

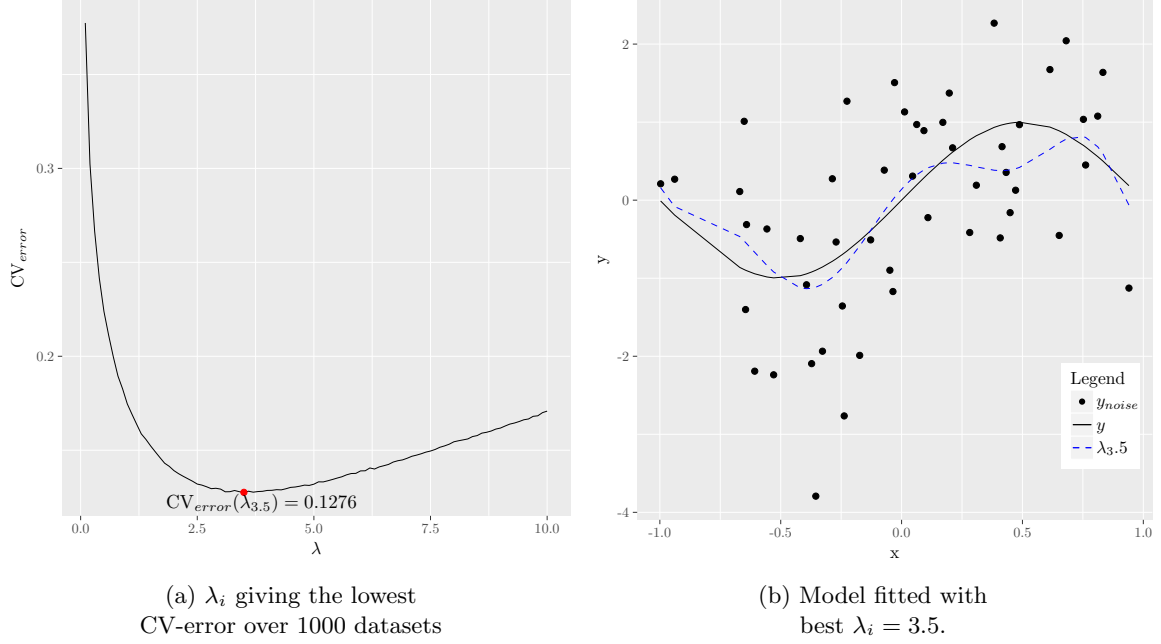


Figure 5: Task 2 iii

From the results above we can clearly see that the CV-error drastically decrease in the start when increasing the value of λ , then it has stabilize while curving slightly around $\lambda \in [2.5, 5.0]$ before it gradually increases in a seemingly linear fashion. When plotting the function with the obtained best regularization parameter $\lambda = 3.5$, it performs quite similar as the model with $\lambda = 5$ in figure 4. Looking at the CV-error plot this seems very reasonable, giving that both values give a very similar error. It's hard to say, but it might also seem like the model with $\lambda = 5$ performs better on this dataset. This does not destroy the result that $\lambda = 3.5$ is the best in figure 5a since an *average best* does not guarantee a *always best*.

4 Task 3

Using a image set of 400 pictures of 40 people with 10 of each person, we look at eigenfaces and reconstructing faces from them using principal component analysis PCA.

4.1 Task 3 i



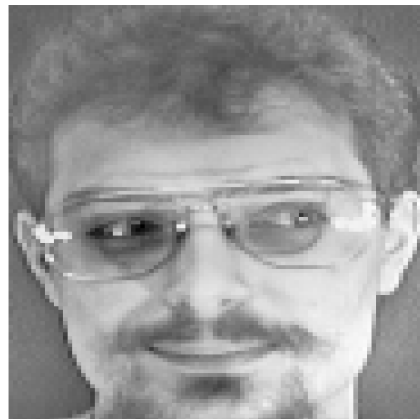
(a) Mean image/face



(b) Standard deviation image/face



(c) Original image/face



(d) Scaled image/face

Figure 6: Task 2 i - SD and mean image/face of the set together with the scaled and original image/face of *i* 168.pgm

4.2 Task 3 ii

Deriving the first ten eigenfaces by using the whole set of images:

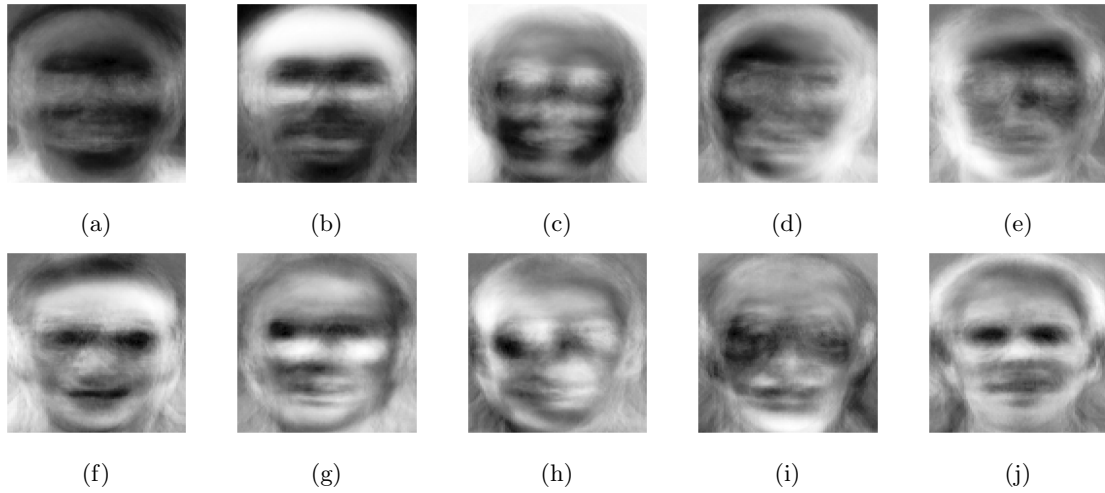


Figure 7: Ten first eigenfaces based on image set.

4.3 Task 3 iii

In this subtask we will look at how it is possible to reconstruct a image using eigenfaces. Below the image "115.pgm" together with reconstructed images based on 5, 50 and 200 eigenfaces.

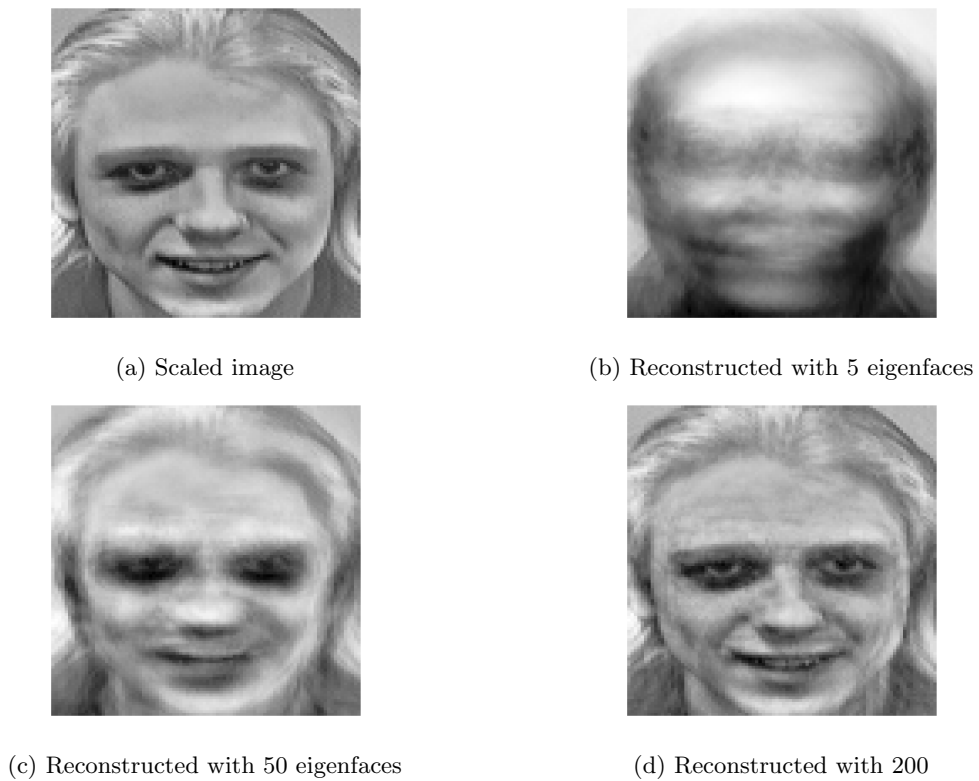


Figure 8: Scaled version of "115.pgm" together with reconstructed versions.

The purpose of using eigenfaces is to efficiently store and encode images. Eigenfaces are derived eigenvalues of the co-variance matrix of the set of images. By using eigenfaces it is possible to reduce the dimensionality by using a smaller set of images to represent the original set. The eigenvalues store the variation in the set of face images and can be combined to reconstruct a specific image. As seen in figure 8 that as the number of eigenfaces increase the face shown starts to look more and more similar to the scaled version of the real face. Another interesting observation is that as the number of eigenfaces used increase each image has a lesser impact in improving the quality. The difference in increasing from 5 to 50 eigenfaces is much more dramatic than from 50 to 200 eigenfaces. This indicates that to get a recognizable image a small amount of eigenfaces is needed and increasing the number after that will only improve the detail.

5 Appendix

5.1 R-Code: Question 1

```
1 ## Libraries and seed
2 library(ggplot2) #Ggplot library used for plotting
3 library(tcltk) #Used for loading bar
4 library(tikzDevice) #Using tikz to store ggplot as tex
5 set.seed(420)
6
7 ## Help functions
8
9 # calculated function from "quote" expression
10 make_data_set <- function(
11   expression,
12   x
13 )
14 {
15   y_dataset <- sapply(x, function(x) eval(expression))
16 }
17
18 # fit linear model to data with offset
19 fit_function <- function(
20   b_intercept,
21   x,
22   y,
23   N
24 )
25 {
26   g_i <- lm(y ~ 0+x, offset = rep(b_intercept,N))
27   slope <- unname(coef(g_i))
28 }
29
30 ## Main functions
31
32 # Task 1, fit two linear functions to y with noise
33 # if i >= 1, performs the average of the runs
34 # set i = 1 to get only one run
35 # plot the two estimated functions with the original y - function
36 # also plot relative performance for between the two models
37 # when averaging many runs, recommended i > 100 for average run
38 task1i <- function(i)
39 {
40   N = 30
41
42   sigma = 1
43   b_y <- 0.8
44
45   bias_g_1 <- 0
46   bias_g_2 <- 0
47
48   b_g_1_avg <- 0
49   b_g_2_avg <- 0
50
51   avg_relative_performance <- rep(0, i)
52
53   y_dataset_avg <- rep(0, N)
54
55   y_wthout_noise <- quote(0.8 * x)
```

```

56 y_wth_noise <- quote(0.8 * x + rnorm(1,0,1))
57
58 for(k in 1:i){
59
60     x <- runif(n = N, min = -1, max = 1)
61
62     y_dataset <- make_data_set(y_wth_noise, x)
63     y_dataset_avg <- y_dataset_avg + y_dataset
64
65     b_g_1 <- fit_function(0.5, x, y_dataset, N)
66     b_g_2 <- fit_function(-0.5, x, y_dataset, N)
67
68     b_g_1_avg <- b_g_1_avg + b_g_1
69     b_g_2_avg <- b_g_2_avg + b_g_2
70
71     bias_g_1 <- bias_g_1 + integrate(function(x) (0.5 + (b_g_1 - b_y
72         ) * x)^2, -1, 1)$value
73     bias_g_2 <- bias_g_2 + integrate(function(x) (-0.5 + (b_g_2 - b_
74         y) * x)^2, -1, 1)$value
75
76     avg_relative_performance[k] <- bias_g_1 / bias_g_2
77 }
78
79 # Task average of all values
80 bias_g_1 <- bias_g_1 / i
81 bias_g_2 <- bias_g_2 / i
82 b_g_1_avg <- b_g_1_avg / i
83 b_g_2_avg <- b_g_2_avg / i
84 y_dataset <- y_dataset_avg / i
85
86 if(i > 1){
87     x = seq(-1,1, length.out = N)
88 }
89
90 y_values <- eval(y_wthout_noise)
91
92 g_1 <- quote(0.5 + b_g_1_avg * x)
93 g_2 <- quote(-0.5 + b_g_1_avg * x)
94
95 g_1_values <- eval(g_1)
96 g_2_values <- eval(g_2)
97
98 print(bias_g_1)
99 print(bias_g_2)
100
101 ggplot_df <- data.frame(x, y_dataset, y_values, g_1_values, g_2_
102     values)
103 avg_relative_performance_df <- data.frame(x = 1:i, avg_relative_
104     performance)
105
106 # plot fitted function together with original function y
107 tikz(file = paste0("Pictures/Task1/tasklibias", i, ".tex"), width =
108     5, height = 5)
109 ggplot1 <- ggplot(data = ggplot_df, aes(x = x)) +
110     geom_line(aes(y = y_values, colour = "y")) +
111     geom_line(aes(y = g_1_values, colour = "$g_1$"), linetype
112         = "dashed") +
113     geom_line(aes(y = g_2_values, colour = "$g_2$"), linetype
114         = "dashed") +

```

```

108     xlab("x") +
109     ylab("y") +
110     annotate("text", x = -0.1, y = 0.7,
111             label = paste0("$\\mathrm{BIAS}(y,g_1) = ",
112                             round(bias_g_1,3),
113                             "$ \\\\",
114                             "$g_1 = 0.5 + ",
115                             round(b_g_1_avg,3),
116                             " x$")) +
117     annotate("text", x = 0.5, y = -0.4,
118             label = paste0("$\\mathrm{BIAS}(y,g_2) = ",
119                             round(bias_g_2,3),
120                             "$ \\\\",
121                             "$g_2 = -0.5 + ",
122                             round(b_g_2_avg,3),
123                             " x$"))
124 #Only plot datapoints when running one dataset
125 if(i == 1){
126     ggplot1 <- ggplot1 +
127     geom_point(aes(y = y_dataset, colour = "$y_{noise}$")) +
128     scale_colour_manual("Legend",
129                         breaks = c("$y_{noise}$", "y", "$g_1$", "$g_2$"),
130                         values = c("blue", "red", "black", "black"),
131                         guide = guide_legend(override.aes = list(
132                             linetype = c("blank", "solid", "dashed", "
133                                     dashed"),
134                             shape = c(16, NA, NA, NA)
135                         ))) +
136     theme(legend.position = c(0.9, 0.2))
137 }else{
138     ggplot1 <- ggplot1 +
139     scale_colour_manual("Legend",
140                         breaks = c("y", "$g_1$", "$g_2$"),
141                         values = c("blue", "red", "black"),
142                         guide = guide_legend(override.aes = list(
143                             linetype = c("solid", "dashed", "dashed")
144                         ),
145                         shape = c(NA, NA, NA)
146                     ))) +
147     theme(legend.position = c(0.9, 0.2))
148 }
149
150 ggsave(paste0("Pictures/Task1/task1i", i, ".png"))
151 print(ggplot1) #need to print plot to write to .tex file
152 dev.off()
153
154 # plot average relative performance rp and shows value after
155 # averaging 1 to i times
156 tikz(file = paste0("Pictures/Task1/task1libias_avg", i, ".tex"),
157      width = 5, height = 5)
158 ggplot2 <- ggplot(data = avg_relative_performance_df, aes(x = x)) +
159     geom_line(aes(y = avg_relative_performance, colour = "$rp_{avg}$")) +
160     xlab("runs") +
161     ylab("$\\mathrm{BIAS}_{avg}$") +
162     theme(legend.position = c(0.8, 0.2))
163 ggsave(paste0("Pictures/Task1/task1i_avg", i, ".png"))
164 print(ggplot2)

```

```

161     dev.off()
162   }
163
164   # finding the val. and out. error for different sizes of training and
165   # validation
166   # set out of a set with total size N = 30
167   taskl1i <- function(
168     n_training_sets = 10000
169   )
170   {
171     N = 30
172     indexes <- 1:N
173
174     # progress bar
175     pb <- tkProgressBar(title = paste0("Running ", n_training_sets, "
176       data sets"),
177       min = 0, max = n_training_sets, width = 300)
178
179     length_df <- N - 2*5 + 1
180     error_df <- data.frame(g_1_error_val = rep(0,length_df),
181       g_2_error_val = rep(0,length_df),
182       g_1_error_out = rep(0,length_df),
183       g_2_error_out = rep(0,length_df),
184       collection_val = rep(0,length_df),
185       collection_out = rep(0,length_df))
186
187     for(k in 1:n_training_sets){
188       setTkProgressBar(pb, k, label = paste(round(k/n_training_sets*
189         100, 0), "% done"))
190       x <- runif(n = N, min = -1, max = 1)
191       b_y <- 0.8
192       y_wth_noise <- quote(0.8 * x + rnorm(1,0,1))
193       y_dataset <- make_data_set(y_wth_noise, x)
194       for(i in 5:(N-5)){
195         train_indexes <- sample(indexes, N - i)
196         test_indexes <- subset(indexes, !(indexes %in% train_indexes
197           ))
198
199         x_train <- x[train_indexes]
200         x_test <- x[test_indexes]
201
202         y_train <- y_dataset[train_indexes]
203         y_test <- y_dataset[test_indexes]
204
205         b_g_1 <- fit_function(0.5, x_train, y_train, N - i)
206         b_g_2 <- fit_function(-0.5, x_train, y_train, N - i)
207
208         g_1 <- quote(0.5 + b_g_1 * x_test) # TO DO: pr ve
209           utf re fit_function og f ut verde f r settes i quote
210         g_2 <- quote(-0.5 + b_g_2 * x_test)
211
212         g_1_values <- eval(g_1)
213         g_2_values <- eval(g_2)
214
215         g_1_mse <- round(mean((y_test - g_1_values)^2),5)
216         g_2_mse <- round(mean((y_test - g_2_values)^2),5)
217
218         g_1_bias <- integrate(function(x) (0.5 + (b_g_1 - b_y) * x)
219           ^2, -1, 1)$value + 1

```

```

214         g_2_bias <- integrate(function(x) (-0.5 + (b_g_2 - b_y) * x)
215           ^2, -1, 1)$value +1
216     error_df$g_1_error_val[i-4] <- (error_df$g_1_error_val[i-4]
217       + g_1_mse)
218     error_df$g_2_error_val[i-4] <- (error_df$g_2_error_val[i-4]
219       + g_2_mse)
220     error_df$g_1_error_out[i-4] <- (error_df$g_1_error_out[i-4]
221       + g_1_bias)
222     error_df$g_2_error_out[i-4] <- (error_df$g_2_error_out[i-4]
223       + g_2_bias)
224     error_df$collection_val[i-4] <- (error_df$collection_val[i
225       -4] + min(g_1_mse, g_2_mse))
226     error_df$collection_out[i-4] <- (error_df$collection_out[i
227       -4] + min(g_2_bias, g_2_bias))
228   }
229 }
230 close(pb)
231 error_df$test_set_amount <- 5:(N-5)
232
233 error_df[-7] <- error_df[-7] / n_training_sets
234
235 # Individual errors for g_1 and g_2
236 ggplot1 <- ggplot(data = error_df, aes(x = test_set_amount)) +
237   geom_line(aes(y = g_1_error_val, colour = "g_1")) +
238   geom_line(aes(y = g_2_error_val, colour = "g_2")) +
239   geom_line(aes(y = g_1_error_out, colour = "g_1_2")) +
240   geom_line(aes(y = g_2_error_out, colour = "g_2_2"))
241 ggsave("Pictures/Task2/task2ii1.png")
242
243 # Best error chosen at each i
244 tikz(file = paste0("Pictures/Task1/task1ii.tex"), width = 5, height
245   = 5)
246 ggplot2 <- ggplot(data = error_df, aes(x = test_set_amount)) +
247   geom_line(aes(y = collection_val, colour = "$E_{val}$"))
248   +
249   geom_line(aes(y = collection_out, colour = "$E_{out}$"))
250   +
251   xlab("i") +
252   ylab("y") +
253   scale_colour_manual("Legend",
254     breaks = c("$E_{val}$", "$E_{out}$"),
255     values = c("blue", "red")) +
256   theme(legend.position = c(0.9, 0.2))
257 ggsave("Pictures/Task2/task2ii2.png")
258 print(ggplot2)
259 dev.off()
260 }
261
262 ## Run
263
264 main <- function()
265 {
266   #for(i in 1:10){
267   #   task1i(i)
268   #}
269   #task1i(20000)

```

```
263         task1i(1)
264         task1ii()
265     }
266
267 main()
268
269 ## Plotting against the machine
```

code/Task1.R

5.2 R-Code: Question 2

```
1 ## Libraries and seed
2 library(ggplot2)
3 library(tcltk)
4 library(tikzDevice)
5 set.seed(420)
6 ## Help functions
7
8 options(tikzMetricPackages = c("\\usepackage[utf8]{inputenc}", "\\usepackage[
  T1]{fontenc}", "\\usetikzlibrary{calc}", "\\usepackage{amssymb}"))
9
10 # Take a expression and compute the values based on input z
11 make_data_set <- function(
12   expression,
13   x
14 )
15 {
16   y_dataset <- sapply(x, function(x) eval(expression))
17   return(y_dataset)
18 }
19
20 # legendre expression within sum formula
21 l <- function(
22   x,q,k
23 ){
24   legendre <- x**k * choose(q,k) * choose((q+k-1)/2,q)
25   return(legendre)
26 }
27
28 # produces legendre polynomial
29 legendres <- function(
30   x,q
31 )
32 {
33   func <- rep(0,length(x))
34   for(k in 0:q){
35     func <- func + l(x,q,k)
36   }
37   #print(func)
38   legendre_polynomial <- 2**q * func
39   return(legendre_polynomial)
40 }
41
42 # use regularization paramter lambda to train weights use in the model
43 regularEstimation <- function(
44   data,
45   lambda,
46   Q_order
47 )
48 {
49   dimentions <- dim(data)
50   y = matrix(data[,2], nrow = dimentions[1])
51   Z = matrix(rep(0,(dimentions[1] * (Q_order + 1))), nrow = dimentions
     [1])
52   for(i in 1:dimentions[1]){
53     Z[i,1] <- 1
54     for(j in 2:(Q_order+1)){
55       Z[i,j] = legendres(data[i, "x"], j-1)
```



```

56         }
57     }
58     weights = solve((t(Z)%*%Z) + (lambda * diag(Q_order + 1)))%*%(t(Z)%*
59         %y)
60     return(weights)
61 }
62 # use weights from regularEstimation to regularize the legendre based model
63 regularLinearization <- function(
64     x,
65     data,
66     lambda,
67     Q_order
68 ){
69     legendreMatrix <- regularEstimation(data, lambda, Q_order)
70     pol <- legendreMatrix[1]
71     for(i in 2:(Q_order + 1)){
72         pol <- pol + (legendreMatrix[i] * (legendres(x, i-1)))
73     }
74     return(pol)
75 }
76 # creates n_folds of N indexes from 1 to N
77 create_cv_indexes <- function(N, n_folds){
78     indexes_per_fold <- ceiling(N/n_folds)
79     index_matrix <- matrix(0L, nrow = n_folds, ncol = indexes_per_fold)
80     index_available <- 1:50
81     for(i in 1:n_folds){
82         selected_indexes <- sample(index_available, indexes_per_fold)
83         index_available <- index_available[! index_available %in% selected_
84             indexes]
85         index_matrix[i, ] <- selected_indexes
86     }
87     return(index_matrix)
88 }
89
90 # calculate the cv error
91 cv_error <- function(
92     N,
93     lambda,
94     data,
95     n_folds,
96     cv_indexes,
97     Q_order
98 ){
99     {
100         indexes <- 1:N
101         cv_indexes <- create_cv_indexes(N, n_folds)
102         cv_error <- c()
103         for(i in 1:n_folds){
104             test_indexes <- cv_indexes[i, ]
105             train_indexes <- subset(indexes, !(indexes %in% test_indexes))
106
107             x_train <- data$x[train_indexes]
108             x_test <- data$x[test_indexes]
109
110             y_train <- data$y[train_indexes]
111             y_test <- data$y_real[test_indexes]
112

```

```

113         temp_data <- data.frame(x = x_train, y = y_train)
114
115         lambdaX <- regularLinearization(x_test, temp_data, lambda, 10)
116
117         cv_error <- c(cv_error, ((y_test - lambdaX)^2))
118     }
119     cv_error <- mean(cv_error)
120 }
121
122 ## Main functions
123
124 # plotting underlying function with data set
125 task2i <- function()
126 {
127     N = 50
128     x = runif(n = N, min = -1, max = 1)
129
130     y <- quote(sin(pi*x))
131     y_wth_noise <- quote(sin(pi*x) + rnorm(1,0,1))
132
133     y_dataset <- make_data_set(y_wth_noise, x)
134
135     ggplot_df <- data.frame(x, y_dataset)
136
137     tikz(file = "Pictures/Task2/task2i.tex", width = 5, height = 5)
138     ggplot1 <- ggplot(data = ggplot_df, aes(x = x)) +
139         geom_point(aes(y = y_dataset, colour = "$y_{noise}$")) +
140         geom_line(aes(y = sin(pi*x), colour = "$y$")) +
141         xlab("x") +
142         ylab("y") +
143         scale_colour_manual("Legend",
144                             breaks = c( "$y_{noise}$", "$y$"),
145                             values = c("black", "black"),
146                             guide = guide_legend(override.aes =
147                                 list(
148                                     linetype = c( "blank", "solid"),
149                                     shape = c( 16, NA)
150                                 ))) +
151         theme(legend.position = c(0.9, 0.2))
152     ggsave("Pictures/Task2/task2i.png")
153     print(ggplot1)
154     dev.off()
155 }
156
157 # Fit two legendre polynomial with different regularization paramters
158 task2ii <- function()
159 {
160     N = 50
161     sigma = 1
162     x <- runif(n = N, min = -1, max = 1)
163     y <- sin(pi*x) + rnorm(N,0,sigma^2)
164
165     data <- data.frame(x,y)
166
167     lambda0 <- regularLinearization(x, data, 0, 10)
168     lambda5 <- regularLinearization(x, data, 5, 10)
169
170     ggplot_df <- data.frame(x,y,lambda0,lambda5)

```

```

171     tikz(file = "Pictures/Task2/task2ii.tex", width = 5, height = 5)
172     ggplot1 <- ggplot(data = ggplot_df, aes(x = x)) +
173       geom_point(aes(y = y, colour = "$y_{noise}$")) +
174       geom_line(aes(y = sin(pi*x), colour = "$y$")) +
175       geom_line(aes(y = lambda0, colour = "$\\lambda_0$",
176         linetype = "dashed")) +
177       geom_line(aes(y = lambda5, colour = "$\\lambda_5$",
178         linetype = "dashed")) +
179       xlab("x") +
180       ylab("y") +
181       scale_colour_manual("Legend",
182         breaks = c( "$y_{noise}$", "$y$", "$\\lambda_0$", "$\\lambda_5$",
183           values = c("blue", "red", "black", "black"),
184           guide = guide_legend(override.aes =
185             list(
186               linetype = c( "blank", "solid", "dashed", "dashed"),
187               shape = c( 16, NA, NA, NA)
188             ))) +
189       theme(legend.position = c(0.9, 0.2))
190     ggsave("Pictures/Task2/task2ii.png")
191     print(ggplot1)
192     dev.off()
193   }
194
195 # plot result from task 2 iii when finding the best regularization parameter
196 plot_task2iii <- function
197 (
198   ggplot_df
199 )
200 {
201   lowest_error = ggplot_df[which.min(ggplot_df[,2]),]
202   str(lowest_error)
203   ?tikzTest
204   tikz(file = "Pictures/Task2/task2iii.tex", width = 5, height = 5)
205   ggplot1 <- ggplot(data = ggplot_df, aes(x = lambdas, y = error_
206     vector)) +
207     geom_line() +
208     geom_point(data = lowest_error,
209       aes(x = lambdas, y = error_vector),
210       color = "red") +
211     geom_text(data = lowest_error,
212       aes(label = paste0("$\\mathrm{CV}_{error}(\\"
213         lambda_{",
214         lambdas, "}) = ",
215         round(error_vector,4), "$"
216       )),
217       hjust = 0.4, vjust = 1.3) +
218     labs(x = "$\\lambda$", y = "$\\mathrm{CV}_{error}$")
219   # title = paste("CV error- regularisation $\\lambda$
220     between",
221     ggplot_df[1,"lambdas"], "and",
222     ggplot_df[nrow(ggplot_df), "lambdas"]
223   "]]))
224   print(ggplot1)
225   dev.off()

```

```

219     ggsave("Pictures/Task2/avg_2iiii1000runs.png")
220 }
221
222 # Using CV-error to fund the best regularization parameter for dataset made
223   by
224 # this model
225 task2iii <- function(
226 )
227 {
228     N = 50
229     sigma = 1
230     x <- runif(n = N, min = -1, max = 1)
231     y <- sin(pi*x) + rnorm(N, 0, sigma)
232     y_real <- sin(pi*x)
233
234     data <- data.frame(x, y, y_real)
235
236     lambda = 0.1
237     lambdas <- seq(from = 0.1, 10, by = 0.1)
238     n_folds = 10
239     error_vector <- integer(length(lambdas))
240     cv_indexes <- create_cv_indexes(N, n_folds)
241
242     pb <- tkProgressBar(title = paste0("Running ", length(lambdas), "
243       lambdas"),
244       min = 0, max = length(lambdas), width = 300)
245     for(i in 1:length(lambdas)){
246       setTkProgressBar(pb, i, label = paste(round(i/length(lambdas)*
247         100, 0), "% done"))
248       lambda = lambdas[i]
249       print(lambda)
250       error_vector[i] <- cv_error(N = N, lambda = lambdas[i], data =
251         data,
252         n_folds = 10, cv_indexes = cv_
253           indexes,
254           Q_order = 10)
255     }
256     close(pb)
257     write.csv(error_vector, "error_01")
258     print(error_vector)
259     ggplot_df <- data.frame(lambdas, error_vector)
260     write.csv(ggplot_df, "ggplot_01.csv", row.names = FALSE)
261     plot_task2iii(ggplot_df)
262 }
263
264 # averaging the CV-error for the different lambdas over n_runs dataset
265 # takes some hours to run on large amounts of dataset,
266 # if using this more parallelization would be a good step to improve
267   performance
268 task2iii_avg_runs <- function(
269   n_runs
270 )
271 {
272     N = 50
273     sigma = 1
274     n_folds = 10
275     lambdas <- seq(from = 0.1, 10, by = 0.1)
276     error_vector <- integer(length(lambdas))

```

```

272     cv_indexes <- create_cv_indexes(N, n_folds)
273     error_result <- error_vector
274
275     pb <- txtProgressBar(min = 0, max = n_runs, style = 3)
276     for(r in 1:n_runs){
277         setTxtProgressBar(pb, r)
278         x <- runif(n = N, min = -1, max = 1)
279         y <- sin(pi*x) + rnorm(N, 0, sigma)
280         y_real <- sin(pi*x)
281
282         data <- data.frame(x, y, y_real)
283
284         for(i in 1:length(lambdas)){
285             error_vector[i] <- cv_error(N = N, lambda = lambdas[i], data
                = data,
286                                     n_folds = 10, cv_indexes = cv_
                indexes,
287                                     Q_order = 10)
288         }
289         error_result <- error_result + error_vector
290     }
291     close(pb)
292     error_result <- error_result / n_runs
293     png("test_avg_2ii.png")
294     plot(error_result, type = "b")
295     dev.off()
296     avg_2iii <- data.frame(lambdas, error_result)
297     write.csv(avg_2iii, "avg_2iii.csv", row.names = FALSE)
298 }
299
300 # plot the best lambda found, a simplification of the plot for ii
301 task2iii_plot_best_lambda <- function(
302     best_lambda
303 )
304 {
305     N = 50
306     sigma = 1
307     x <- runif(n = N, min = -1, max = 1)
308     y <- sin(pi*x) + rnorm(N,0,sigma^2)
309
310     data <- data.frame(x,y)
311
312     lambda_best <- regularLinearization(x, data, best_lambda, 10)
313
314     ggplot_df <- data.frame(x,y,lambda_best)
315
316     tikz(file = "Pictures/Task2/task2iii_best_lambda.tex", width = 5,
        height = 5)
317     ggplot1 <- ggplot(data = ggplot_df, aes(x = x)) +
318         geom_point(aes(y = y, colour = "$y_{noise}$")) +
319         geom_line(aes(y = sin(pi*x), colour = "$y$")) +
320         geom_line(aes(y = lambda_best,
321                     colour = paste0("$\\lambda_", best_lambda,
322                                     "$")),
323                 linetype = "dashed") +
324         xlab("x") +
325         ylab("y") +
326         scale_colour_manual("Legend",
            breaks = c( "$y_{noise}$", "$y$",

```

```

327         paste0("$\\lambda_", best_lambda, "
328         $")),
329         values = c("blue", "black", "black"),
330         guide = guide_legend(override.aes =
331         list(
332             linetype = c("blank", "solid", "
333             dashed"),
334             shape = c(16, NA, NA)
335             ))) +
336         theme(legend.position = c(0.9, 0.2))
337     ggsave("Pictures/Task2/task2iii_best_lambda.png")
338     print(ggplot1)
339     dev.off()
340 }
341 ## Run
342
343 main <- function()
344 {
345     task2i()
346     task2ii()
347     #task2iii()
348
349     task2iii_avg_runs(1000)
350
351     #big rung for iii, stored to save computation if replot is nessesary
352     #ggplot_df <- read.csv("avg_2iii1000runs.csv")
353     #colnames(ggplot_df)[2] <- "error_vector"
354     #plot_task2iii(ggplot_df)
355     task2iii_plot_best_lambda(3.5)
356 }
357 main()

```

code/Task2.R

5.3 R-Code: Question 3

```
1 ## Libraries and seed
2 library(pixmap)
3 library(reshape2)
4 library(ggplot2)
5 library(RColorBrewer)
6 library(gtools)
7 ## Help functions
8
9 #Load all faces from "Faces" folder
10 load_faces <- function(
11   pattern,
12   directory
13 )
14 {
15   temp = paste0("Faces/", list.files(path = "Faces", pattern = "*.pgm"
16   ))
17   #Files are read in alpha order, need to use mixedsort to get
18   #right order
19   temp <- mixedsort(temp)
20   myfiles = lapply(temp, read.pnm)
21   return(myfiles)
22 }
23
24 #Get and transform the different needed data
25 get_image_data <- function(
26   faces_list,
27   n_pictures,
28   height,
29   width
30 )
31 {
32   image_data <- array(0, dim = c(n_pictures, height, width))
33   image_transpose <- array(0, dim = c(n_pictures, width, height))
34   image_up <- image_transpose
35   image_mean <- matrix(0, nrow = width, ncol = height)
36   image_standard_deviation <- image_mean
37   image_scale <- image_transpose
38   for(i in 1:length(faces_list)){
39     image <- faces_list[[i]]@grey
40     image_data[i, , ] <- image
41     image_transpose[i, , ] <- t(image)
42   }
43   for(i in 1:n_pictures){
44     for(j in 1:height){
45       image_up[i, , (height - j + 1)] <- image_transpose[i, , j]
46     }
47   }
48   for(k in 1:width){
49     for(l in 1: height){
50       image_up_value <- image_up[, k, l]
51       image_mean[k,l] <- mean(image_up_value)
52       image_standard_deviation[k,l] <- sd(image_up_value)
53     }
54   }
55 }
56
```

```

57     for(m in 1:n_pictures){
58         image_scale[m, , ] <- (image_up[m , , ] - image_mean) / image_
            standard_deviation
59     }
60
61 }
62
63 #Plot a face to pdf in two different ways
64 plot_face_pdf <- function(
65     gg_melt,
66     shades_of_grey = 256
67 )
68 {
69     gg <- ggplot(gg_melt, aes(x = Var1, y = Var2)) +
70         geom_tile(aes(fill = value)) +
71         theme_classic()
72     gg2 <- ggplot(gg_melt, aes(x = Var1, y = Var2)) +
73         geom_tile(aes(fill = value)) +
74         #scale_fill_gradientn(colours = brewer.pal(shades_of_grey,"Greys
            ")) +
75         scale_colour_brewer(palette = brewer.pal(shades_of_grey, "Greys"
            ), direction = -1)
76         theme_classic()
77     #print(gg)
78     print(gg2)
79 }
80
81 #Plot a face to the wanted folder as png
82 plot_face_png <-function(
83     gg_melt,
84     shades_of_grey = 256,
85     name
86 )
87 {
88     name <- paste0("Pictures/Task3/", name)
89     gg <- ggplot(gg_melt, aes(x = Var1, y = Var2)) +
90         geom_tile(aes(fill = value)) +
91         scale_fill_gradientn(colours = rev(brewer.pal(shades_of_grey,"
            Greys")))) +
92         theme_classic()
93     #Remove all axis and legends from the plot to only view the image
94     gg <- gg + theme(axis.line=element_blank(),axis.text.x=element_blank
        (),
95         axis.text.y=element_blank(),axis.ticks=element_blank(),
96         axis.title.x=element_blank(),
97         axis.title.y=element_blank(),legend.position="none",
98         panel.background=element_blank(),panel.border=element_blank(),
99         panel.grid.major=element_blank(),
100         panel.grid.minor=element_blank(),plot.background=element_blank())
101     ggsave(name)
102 }
103
104 #Convert data to input_matrix
105 convert_to_input_matrix <- function(
106     n_faces,
107     width,
108     height
109 )
110 {

```



```

111     X <- matrix(0, nrow = n_faces, ncol = width * height)
112     for(i in 1: n_faces){
113         picture <- image_scale[i, 1, ]
114         for(j in 2: width){
115             picture <- cbind(picture, image_scale[i, j, ])
116         }
117         X[i, ] <- picture
118     }
119     pca <- prcomp(X)
120     A <- X %%% t(X)
121     store <- prcomp(A)
122     image_eigen <- t(X) %%% store$rotation
123 }
124
125 #make a eigen face image
126 eigen_image <- function(
127     face_vector,
128     width,
129     height
130 )
131 {
132     eigen_face <- matrix(NA, ncol = height, nrow = width)
133     for(i in 1: width){
134         eigen_face[i, ] <- face_vector[(height * (i - 1) + 1):(height *
135             i)]
136     }
137     return(eigen_face)
138 }
139
140 #recreate face from principal component analysis PCA
141 recreate_face <- function(
142     image_nr,
143     n_faces
144 )
145 {
146     recreate <- tcrossprod(pca$x[, 1:n_faces], pca$rotation[, 1:n_faces
147         ])
148     recreate_matrix <- t(matrix(data = rev(recreate[image_nr, ]),
149         nrow = 112, ncol = 92))
150 }
151
152 ## Main functions
153
154 # plot mean, sd, original and scaled version og image_nr
155 task3i <- function(
156     image_nr
157 )
158 {
159     gg_mean <- melt(image_mean)
160     gg_standard_deviation <- melt(image_standard_deviation)
161
162     gg_up <- melt(image_up[image_nr, , ])
163     gg_scale <- melt(image_scale[image_nr, , ])
164
165     plot_face_png(gg_mean, name = "i/mean.png")
166     plot_face_png(gg_standard_deviation, name = "i/sd.png")
167     plot_face_png(gg_up, name = "i/org.png")
168     plot_face_png(gg_scale, name = "i/scale.png")
169 }

```

```

168
169 # plot the n_eigen_faces first eigenfaces
170 task3ii <- function(
171     n_eigen_faces,
172     width,
173     height
174 )
175 {
176     for(i in 1:n_eigen_faces){
177         eigen <- image_eigen[, i]
178         plot_face_png(melt(image_eigen(eigen, width, height)),
179                       name = paste0("ii/eigen_face", i, ".png"))
180     }
181 }
182
183 # Recreate face from different amount of eigenfaces.
184 task3iii <- function(
185     image_nr
186 )
187 {
188     rotate_matrix <- function(x) t(apply(x, 2, rev))
189     eigen_n_vector <- c(1, 5, 50, 200)
190     plot_face_png(melt(image_scale[image_nr, ]),
191                   name = paste0("iii/image",
192                                 image_nr, ".png"))
193     for(K in eigen_n_vector){
194         recreate <- recreate_face(image_nr, K)
195         #rotate matrix the right way, rotating 90 degrees twice
196         recreate <- rotate_matrix(rotate_matrix(recreate))
197         plot_face_png(melt(recreate), name = paste0("iii/image",
198                                                     image_nr,
199                                                     "eigen_faces",
200                                                     K, ".png"))
201     }
202 }
203
204 ## Run
205
206 main <- function()
207 {
208     # matrices and df
209     # image_data
210     # image_transpose
211     # image_up
212     # image_mean
213     # image_standard_deviation
214     # image_scale
215     # X
216     # pca
217
218     n_faces = 400
219     height = 112
220     width = 92
221     shades_of_grey = 256
222     faces <- load_faces()
223     get_image_data(faces, n_faces, height, width)
224     convert_to_input_matrix(n_faces, width, height)
225
226     task3i(image_nr = 168)

```

```
227     task3ii(n_eigen_faces = 10, width = width, height = height)
228     task3iii(image_nr = 115)
229
230 }
231
232 main()
```

code/Task3.R