

# 1 Common libraries

## 1.1 dplyr

**dplyr**: A grammar of data manipulation, providing tools for filtering, selecting, mutating, and summarizing data efficiently. It is designed to work seamlessly with data frames and tibbles.

## 1.2 caret

**caret**: Provides tools for pre-processing, feature selection, and resampling.

## 1.3 kernlab

**kernlab**: Focused on kernel-based machine learning methods, this library supports algorithms like Support Vector Machines (SVM), Kernel Principal Component Analysis (KPCA), and Gaussian Processes.

## 1.4 neuralnet

**neuralnet**: A package for training and visualizing neural networks in R. It supports deep learning on numeric data and allows customization of activation functions and network architecture.

## 1.5 glmnet

**glmnet**: Implements regularized regression methods such as LASSO and Elastic Net, suitable for linear and logistic regression with high-dimensional data.

## 1.6 e1071

**e1071**: Provides tools for machine learning, including support for SVMs, Naïve Bayes classifiers, clustering algorithms, and statistical functions. It also includes implementations of fuzzy logic.

# 2 Models

## 2.1 Linear Regression

```
1 model = lm ( Fat ~ . , data = data_train)
2 #Or
3 model=glm(Fat~., data=data_train, family="gaussian")
```

## 2.2 Logistic Regression (classification)

```
1 glm(formula = as.factor(species)~., data = a, family = "binomial")
```

## 2.3 Poisson

```
1 glm(formula = Visitors~., data = data_b, family = "poisson"(link=log))
```

## 2.4 Lasso (L1) and Ridge (L2)

```
1 # Import glmnet package
2 library(glmnet)
3
4 # Create the lasso model by setting alpha = 1.
5 lasso <- glmnet(X_train_matrix, data_train$Fat, alpha = 1, family =
6   "gaussian")
7
8 # Cross validation can also be calculated for lasso and ridge models.
9 cv_lasso=cv.glmnet(X_train_matrix, data_train$Faaf, alpha=1
10   ,family="gaussian")
11
12 # Optimal lambda can be retrieved,
13 cv_lasso$lambda.min
14
15 # Predictions can be calculated by specifying optimal lambda
16 predictions = predict(cv_lasso, s = "lambda.min", newx =
17   X_test_matrix, type = "response")
18
19 # If known model can also be trained on the optimal lambda directly
20 lasso <- glmnet(X_train_matrix, data_train$Fat, alpha = 1, family =
21   "gaussian", lambda = optimal_lambda)
22
23 # Create a ridge model by setting alpha = 0.
24 ridge <- glmnet(X_train_matrix, data_train$Fat, alpha = 0, family =
25   "gaussian")
```

## 2.5 Principal Component Analysis

```
1 res <- princomp(dataSC)
2
3 # View summary of the PCA
4 summary(res)
5
6 print(res$loadings)
7 print(res$scores)
```

## 3 Lab 1 Task 1: K-Nearest Neighbor Classification

The data file `optdigits.csv` contains information about normalized bitmaps of handwritten digits from a preprinted form from a total of 43 people. The data were first derived as 32x32 bitmaps

which were then divided into nonoverlapping blocks of 4x4, and the number of on pixels was counted in each block. This generated a resulting image of size 8x8 where each element is an integer in the range 0..16. Accordingly, each row in the data file is a sequence corresponding to an 8x8 matrix, and the last element shows the actual digit from 0 to 9.

1. Import the data into R and divide it into training, validation, and test sets (50%/25%/25%) by using the partitioning principle specified in the lecture slides. Use the training data to fit a 30-nearest neighbor classifier with the function `kknn()` and `kernel="rectangular"` from the `kknn` package, and estimate:

- Confusion matrices for the training and test data (use `table()`).
- Misclassification errors for the training and test data.

Comment on the quality of predictions for different digits and on the overall prediction quality.

2. Find any 2 cases of digit "8" in the training data which were easiest to classify and 3 cases that were hardest to classify (i.e., having the highest and lowest probabilities of the correct class). Reshape features for each of these cases as an 8x8 matrix and visualize the corresponding digits (e.g., using the `heatmap()` function with parameters `Colv=NA` and `Rowv=NA`). Comment on whether these cases seem to be hard or easy to recognize visually.
3. Fit K-nearest neighbor classifiers to the training data for different values of  $K = 1, 2, \dots, 30$ , and plot the dependence of the training and validation misclassification errors on the value of  $K$  (in the same plot). How does the model complexity change when  $K$  increases, and how does it affect the training and validation errors? Report the optimal  $K$  according to this plot. Finally, estimate the test error for the model having the optimal  $K$ , compare it with the training and validation errors, and make necessary conclusions about the model quality.
4. Fit K-nearest neighbor classifiers to the training data for different values of  $K = 1, 2, \dots, 30$ , compute the error for the validation data as cross-entropy (when computing the log of probabilities, add a small constant within the log, e.g.,  $1e-15$ , to avoid numerical problems), and plot the dependence of the validation error on the value of  $K$ . What is the optimal  $K$  value here? Assuming that the response has a multinomial distribution, why might the cross-entropy be a more suitable choice of error function than the misclassification error for this problem?

```

1 ##### PART 1 #####
2
3 library(kknn)
4
5 data_frame = read.csv("../optdigits.csv", header = FALSE)
6
7 # Rename the last column to "label", represents the actual digit
8 colnames(data_frame)[65] = "label"
9
10 # Convert the label column to a factor for classification
11 data_frame$label = as.factor(data_frame$label) # Treat as values
    categories (0 - 9)
12
13 n = dim(data_frame)[1]
14 set.seed(12345)
15 id = sample(1:n, floor(n * 0.5))
16 data_train = data_frame[id, ]
17

```

```

18 id1 = setdiff(1:n, id)
19 set.seed(12345)
20 id2 = sample(id1, floor(n * 0.25))
21 data_validation = data_frame[id2, ]
22
23 id3 = setdiff(id1, id2)
24 data_test = data_frame[id3, ]
25
26 ##### PART 2 #####
27
28 # Performs k-nearest neighbor classification of a test set using a
  training set.
29 model_test = kknn(formula = label ~ ., train = data_train, test =
  data_test, k = 30, kernel = "rectangular")
30 model_train = kknn(formula = label ~ ., train = data_train, test =
  data_train, k = 30, kernel = "rectangular")
31
32 # Generate predictions for test and training data
33 predictions_test = predict(model_test)
34 predictions_train = predict(model_train)
35
36 # Generate confusion matrix (classification) comparing predictions
  with true data.
37 conf_matrix_test = table(predictions_test, data_test$label)
38 conf_matrix_train = table(predictions_train, data_train$label)
39 conf_matrix_test
40 conf_matrix_train
41
42 # Function to calculate the misclassification rate
43 # X = true labels, X1 = predictions
44 missclass=function(X,X1){
45   n=length(X)
46   return(1-sum(diag(table(X,X1)))/n)
47 }
48 missclass(predictions_test, data_test$label)
49 missclass(predictions_train, data_train$label)
50
51 # predictions_test    0    1    2    3    4    5    6    7    8    9
52 # 0   77    0    0    0    0    0    0    0    0    0
53 # 1    0   81    0    0    0    1    0    0    7    1
54 # 2    0    2   98    0    0    1    0    0    0    1
55 # 3    0    0    0  107    0    0    0    1    1    1
56 # 4    1    0    0    0   94    0    0    0    0    0
57 # 5    0    0    0    2    0   93    0    0    0    0
58 # 6    0    0    0    0    2    2   90    0    0    0
59 # 7    0    0    0    0    6    1    0  111    0    1
60 # 8    0    0    3    1    2    0    0    0   70    0
61 # 9    0    3    0    1    5    5    0    0    0   85
62
63
64 # predictions_train    0    1    2    3    4    5    6    7    8    9
65 # 0  202    0    0    0    1    0    0    0    0    1
66 # 1    0  179    1    0    3    0    2    3   10    3

```

```

67 # 2    0   11 190    0    0    0    0    0    0    0
68 # 3    0    0    0 185    0    1    0    0    2    5
69 # 4    0    0    0    0 159    0    0    0    0    2
70 # 5    0    0    0    1    0 171    0    0    0    0
71 # 6    0    0    0    0    0    0 190    0    2    0
72 # 7    0    1    1    1    7    1    0 178    0    3
73 # 8    0    1    0    0    1    0    0    1 188    3
74 # 9    0    3    0    1    4    8    0    0    2 183
75
76
77 # > missclass(predictions_test, data_test$label)
78 # [1] 0.05329154
79 # > missclass(predictions_train, data_train$label)
80 # [1] 0.04500262
81
82 # The overall quality of the models predictions is quite good, which
83 # can be seen
84 # above. For the test and training data set the misclassification
85 # rates are only
86 # 0.05329154 and 0.04500262 respectively. When checking the confusion
87 # matrix for
88 # the training dataset there are some values that stand out. These
89 # are the
90 # numbers 1, 4, 5, 8 and 9. These were wrongly classified the most by
91 # the model.
92 # This is probably because some handwritten numbers look the same
93 # when using a
94 # low resolution image.
95
96 ##### PART 3 #####
97
98 # Get the probability matrix from the trained model
99 prob_matrix = model_train$prob
100
101 # Identify rows with the label "8" in the training dataset
102 label_8_indices = which(data_train$label == "8")
103
104 # Extract probabilities of the true label "8" for these rows
105 label_8_prob = prob_matrix[label_8_indices, "8"]
106
107 # Sort the indices by descending probabilities
108 sorted_indices = order(label_8_prob, decreasing = TRUE)
109
110 # Select the indices of the two easiest and three hardest
111 # classifications
112 easiest_to_classify = label_8_indices[sorted_indices[1:2]]
113 hardest_to_classify =
114     label_8_indices[sorted_indices[(length(sorted_indices)-2):length(sorted_indices)]]
115
116 # Reshape to 8x8 matrices to visualize the digit with heat map
117 for (i in 1:2) {
118     easiest_case = matrix(as.numeric(data_train[easiest_to_classify[i],
119     1:64]), nrow = 8, ncol = 8, byrow=TRUE)

```

```

111 heatmap(easiest_case, Rowv = NA, Colv = NA, main = paste("Easiest
112 Case", i))
113 }
114 for (i in 1:3) {
115   hardest_case = matrix(as.numeric(data_train[hardest_to_classify[i],
116     1:64]), nrow = 8, ncol = 8, byrow=TRUE)
117   heatmap(hardest_case, Rowv = NA, Colv = NA, main = paste("Hardest
118 Case", i))
119 }
120 # The easiest two classifications are simple to recognize while the
121 # three hardest are almost impossible.
122 ##### PART 4 #####
123 # Decide how many data points from neighboring data you want to
124 # consider when making predictions.
125 min_k = 1
126 max_k = 30
127 k_values = min_k:max_k
128 missclass_vector = c()
129 missclass_vector_validation = c()
130 for (k_value in k_values) {
131   model_train = kknn(formula = label ~ ., train = data_train, test =
132     data_train, k = k_value, kernel = "rectangular")
133   model_validation = kknn(formula = label ~ ., train = data_train,
134     test = data_validation, k = k_value, kernel = "rectangular")
135   predict_train = predict(model_train)
136   predict_validation = predict(model_validation)
137   prob_validation = model_validation$prob
138   missclass_value = missclass(predict_train, data_train$label)
139   missclass_value_validation = missclass(predict_validation,
140     data_validation$label)
141   missclass_vector = append(missclass_vector, missclass_value)
142   missclass_vector_validation = append(missclass_vector_validation,
143     missclass_value_validation)
144 }
145 plot(k_values, missclass_vector, col = "red", xlab = "k (Number of
146   neighbors)", ylab = "Misclassification rate", main =
147   "misclassification rate vs k value", ylim = c(0, 0.05))
148 points(k_values, missclass_vector_validation, col = "blue")
149 missclass_vector_validation[which.min(missclass_vector_validation)]
150 legend("bottomright", legend = c("Training", "Validation"), col =
151   c("red", "blue"), lty = 1)

```

```

152 # When the K values are small the complexity of the model is high and
    there is a big risk for overfitting. While when the K values are
    big the complexity of the model decreases and there is risk for
    underfitting. By looking at the figure we can see that for low K
    values the red and blue model is better while for big K values
    both models is almost equally as bad.
153 #
154 # The most optimal K value for the blue model is K = 3 then the
    misclassification value is 0.02513089 and for the red model is K =
    1 and misclassification value is 0.00. This is because the red
    model is trained and tested on the same dataset. This makes the
    model look very good for low K values but in reality its
    overfitted. The blue model has worse misclassification value but
    its not trained and tested on the same dataset so it shows a more
    realistic rate for different K values.
155
156 ##### PART 5 #####
157
158 LOG_CONSTANT = 1e-15
159 cross_entropy = c()
160
161 for (k_value in k_values) {
162     model_train = kknn(formula = label ~ ., train = data_train, test =
        data_train, k = k_value, kernel = "rectangular")
163     model_validation = kknn(formula = label ~ ., train = data_train,
        test = data_validation, k = k_value, kernel = "rectangular")
164
165     prob_validation = model_validation$prob
166
167     # Cross entropy
168     cross_entropy_k = 0
169     for (i in 1:nrow(data_validation)) {
170         true_label = data_validation$label[i] # Gets label for each row
171         true_prob = prob_validation[i, as.numeric(true_label)] # Get prob
            for current row and true label
172         cross_entropy_k = cross_entropy_k - log(true_prob + LOG_CONSTANT)
            # calc cross entropy
173     }
174     cross_entropy = c(cross_entropy, cross_entropy_k /
        nrow(data_validation))
175 }
176
177 plot(k_values, cross_entropy, xlab = "k (Number of neighbors)", ylab
    = "Cross-entropy loss", main = "Cross-Entropy Loss vs. k", ylim =
    c(0, 1.0))
178 which.min(cross_entropy)
179
180 # From the figure we can see that the optimal K value is 6, because
    it has the lowest loss. The cross entropy is a more suitable
    choice for the error function because it evaluates the quality of
    the predicted probabilities. Compared to the misclassification
    which only cares if the predicted digit was correct or not, and
    not how high probability that the true digit actually has.

```

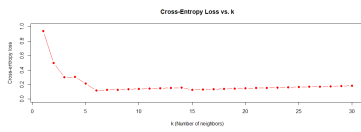


Figure 1: cross entropy vs k



Figure 2: easiest case 1

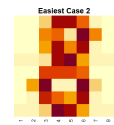


Figure 3: easiest case 2



Figure 4: hardest case 1



Figure 5: hardest case 2



Figure 6: hardest case 3

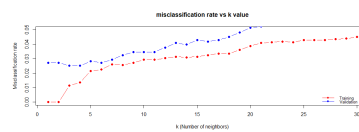


Figure 7: misclassification vs k

## 4 Lab 1 Task 2: Linear Regression and Ridge Regression

The data file `parkinson.csv` is composed of a range of biomedical voice measurements from 42 people with early-stage Parkinson's disease recruited to a six-month trial of a telemonitoring device for remote symptom progression monitoring. The purpose is to predict Parkinson's disease symptom score (motor UPDRS) from the following voice characteristics:

- **Jitter**: (%), Jitter(Abs), Jitter:RAP, Jitter:PPQ5, Jitter:DDP - Several measures of variation in fundamental frequency.
  - **Shimmer**: Shimmer(dB), Shimmer:APQ3, Shimmer:APQ5, Shimmer:APQ11, Shimmer:DDA - Several measures of variation in amplitude.
  - **NHR, HNR** - Two measures of the ratio of noise to tonal components in the voice.
  - **RPDE** - A nonlinear dynamical complexity measure.
  - **DFA** - Signal fractal scaling exponent.
  - **PPE** - A nonlinear measure of fundamental frequency variation.
1. Divide the data into training and test sets (60/40) and scale it appropriately. In the following steps, assume that `motor_UPDRS` is normally distributed and is a function of the voice characteristics. Since the data are scaled, no intercept is needed in the modeling.
  2. Compute a linear regression model from the training data, estimate training and test MSE, and comment on which variables contribute significantly to the model.
  3. Implement the following functions using basic R commands only (no external packages):
    - a) **LogLikelihood**: A function for a given parameter vector  $\theta$  and dispersion  $\sigma$  computes the log-likelihood function  $\log P(T|\theta, \sigma)$  for the stated model and the training data.
    - b) **Ridge**: A function that, for a given vector  $\theta$ , scalar  $\sigma$ , and scalar  $\lambda$ , adds a Ridge penalty  $\lambda \|\theta\|^2$  to the minus log-likelihood.



- c) **RidgeOpt**: A function that depends on the scalar  $\lambda$ , uses function 3b and function `optim()` with `method="BFGS"` to find the optimal  $\theta$  and  $\sigma$  for a given  $\lambda$ .
- d) **DF (Degrees of Freedom)**: function that for given scalar  $\lambda$  computes the degrees of freedom of the ridge model based on the training data.
4. Using the Ridge optimization function, compute the optimal  $\theta$  parameters for  $\lambda = 1$ ,  $\lambda = 100$ , and  $\lambda = 1000$ . Use the estimated parameters to predict `motor_UPDRS` values for training and test data and report the training and test MSE values. Determine which penalty parameter is most appropriate among the selected ones. Compute and compare the degrees of freedom of these models and make appropriate conclusions.

```

1 # ----- TASK 1 -----
2
3 data <- read.csv('parkinsons.csv')
4 labels <- colnames(data)
5
6 parameters<-c("Jitter...", "Jitter.Abs.", "Jitter.RAP",
7               "Jitter.PPQ5", "Jitter.DDP",
8               "Shimmer", "Shimmer.dB.", "Shimmer.APQ3",
9               "Shimmer.APQ5", "Shimmer.APQ11",
10              "Shimmer.DDA", "NHR", "HNR", "RPDE", "DFA", "PPE")
11 parameters_y<-c(parameters, "motor_UPDRS")
12
13 n=dim(data)[1]
14 set.seed(12345)
15 id=sample(1:n, floor(n*0.6))
16 train=data[id,]
17
18 valid=data[-id,]
19
20 library(caret)
21 scaler=preProcess(train)
22 trainS=predict(scaler,train)
23 validS=predict(scaler,valid)
24
25 trainS = trainS[, colnames(trainS) %in% parameters_y]
26 validS = validS[, colnames(validS) %in% parameters_y]
27
28 trainS_without_y_matrix = as.matrix(subset(trainS,
29                                           select=-motor_UPDRS))
30 validS_without_y_matrix = as.matrix(subset(validS,
31                                           select=-motor_UPDRS))
32
33 # ----- TASK 2 -----
34
35 y = trainS$motor_UPDRS
36 y_validS = validS$motor_UPDRS
37
38 fit=lm(y ~ . - motor_UPDRS, data=trainS)
39
40 coeff = coefficients(fit)[-1] # model coefficients, intercept not
41                               needed since the data is scaled
42 confint(fit, level=0.95) # CIs for model parameters

```

```

38 fitted(fit) # predicted values
39 residuals_trainS = residuals(fit) # residuals
40
41 # -----
42
43 predicted_validS = predict(fit, newdata=validS)
44 residuals_validS <- validS$motor_UPDRS - predicted_validS
45
46 # ANSWER:
47 mean(residuals_trainS^2) # MSE for trainS = 0.8785431
48 mean(residuals_validS^2) # MSE for validS = 0.9354477
49
50 summary(fit)
51
52 # ANSWER:
53 # summary(fit) will print some stuff. The three stars next to some
54 # variables mean that there is a high correlation
55 # between them and the target variable motor_UPDRS.
56 # Answer: Jitter.Abs., Shimmer.APQ5, Shimmer.APQ11, NHR, HNR, DFA and
57 # PPE
58
59 # ----- TASK 3 -----
60 Loglikelihood <- function(theta, sigma) # linear regression gives
61 # normal distribution
62 {
63   - nrow(trainS) / 2 * log(2 * pi) - nrow(trainS) / 2 * log(sigma^2)
64   - 1/(2*sigma^2) * sum((y - trainS_without_y_matrix %*% theta)^2)
65 }
66
67 Ridge <- function(theta, lambda, sigma)
68 {
69   -Loglikelihood(theta, sigma) + lambda * sum(theta^2)
70 }
71
72 RidgeOpt <- function(lambda)
73 {
74   optim(coeff, fn = Ridge, lambda = lambda, sigma = 1, method =
75         "BFGS")
76 }
77
78 DF <- function(lambda)
79 {
80   X = trainS_without_y_matrix
81   sum(diag(X %*% solve(t(X) %*% X + lambda * diag(ncol(X))) %*% t(X)))
82 }
83
84 # ----- TASK 4 -----
85
86 lambda_1 = RidgeOpt(1)$par
87 lambda_100 = RidgeOpt(100)$par
88 lambda_1000 = RidgeOpt(1000)$par
89
90 # -----

```

```

86 # Use the estimated parameters to predict the motor_UPDRS values for
    training data
87 predict_trainS_motor_UPDRS_lambda_1 = trainS_without_y_matrix %*%
    lambda_1
88 predict_trainS_motor_UPDRS_lambda_100 = trainS_without_y_matrix %*%
    lambda_100
89 predict_trainS_motor_UPDRS_lambda_1000 = trainS_without_y_matrix %*%
    lambda_1000
90
91 # calculate the MSE
92 MSE_trainS_lambda_1 = mean((predict_trainS_motor_UPDRS_lambda_1 -
    y)^2)
93 MSE_trainS_lambda_100 = mean((predict_trainS_motor_UPDRS_lambda_100 -
    y)^2)
94 MSE_trainS_lambda_1000 = mean((predict_trainS_motor_UPDRS_lambda_1000
    - y)^2)
95
96 # -----
97
98 # Use the estimated parameters to predict the motor_UPDRS values for
    test data
99 predict_validS_motor_UPDRS_lambda_1 = validS_without_y_matrix %*%
    lambda_1
100 predict_validS_motor_UPDRS_lambda_100 = validS_without_y_matrix %*%
    lambda_100
101 predict_validS_motor_UPDRS_lambda_1000 = validS_without_y_matrix %*%
    lambda_1000
102
103 # Calculate the MSE
104 MSE_validS_lambda_1 = mean((predict_validS_motor_UPDRS_lambda_1 -
    y_validS)^2)
105 MSE_validS_lambda_100 = mean((predict_validS_motor_UPDRS_lambda_100 -
    y_validS)^2)
106 MSE_validS_lambda_1000 = mean((predict_validS_motor_UPDRS_lambda_1000
    - y_validS)^2)
107
108 DF_lambda_1 = DF(1)
109 DF_lambda_100 = DF(100)
110 DF_lambda_1000 = DF(1000)
111
112 # -----
113
114 # ANSWER:
115 MSE_trainS_lambda_1 # 0.8786332
116 MSE_trainS_lambda_100 # 0.8850982
117 MSE_trainS_lambda_1000 # 0.9232605
118
119 MSE_validS_lambda_1 # 0.9349501
120 MSE_validS_lambda_100 # 0.9324709
121 MSE_validS_lambda_1000 # 0.9554836
122
123 DF_lambda_1 # 13.86074
124 DF_lambda_100 # 9.924887

```

```

125 DF_lambda_1000 # 5.643925
126
127 # We see that MSE_validS_lambda_100 is the lowest which
128 # means that the lambda = 100 is the most suitable choice. In this
    case it can be
129 # seen that a higher degree of freedom is not always the best choice
    since it
130 # can lead to overfitting.

```

## 5 Lab 1 Task 3: Logistic Regression and Basis Function Expansion

The data file `pima-indians-diabetes.csv` contains information about the onset of diabetes within 5 years in Pima Indians given medical details. The variables are (in the same order as in the dataset):

1. Number of times pregnant.
  2. Plasma glucose concentration 2 hours into an oral glucose tolerance test.
  3. Diastolic blood pressure (mm Hg).
  4. Triceps skinfold thickness (mm).
  5. 2-Hour serum insulin ( $\mu$  U/ml).
  6. Body mass index (weight in kg/(height in m)<sup>2</sup>).
  7. Diabetes pedigree function.
  8. Age (years).
  9. Diabetes (0 = no or 1 = yes).
1. Make a scatterplot showing Plasma glucose concentration on Age where observations are colored by Diabetes levels. Do you think that Diabetes is easy to classify by a standard logistic regression model that uses these two variables as features? Motivate your answer.
  2. Train a logistic regression model with  $y = \text{Diabetes}$  as target,  $x_1 = \text{Plasma glucose concentration}$ , and  $x_2 = \text{Age}$  as features, and make a prediction for all observations by using  $r = 0.5$  as the classification threshold.
    - Report the probabilistic equation of the estimated model (i.e., how the target depends on the features and the estimated model parameters probabilistically).
    - Compute the training misclassification error and make a scatterplot of the same kind as in step 1 but showing the predicted values of Diabetes as a color instead.
    - Comment on the quality of the classification by using these results.
  3. Use the model estimated in step 2 to:
    - a) Report the equation of the decision boundary between the two classes.
    - b) Add a curve showing this boundary to the scatterplot in step 2.

Comment whether the decision boundary seems to capture the data distribution well.
  4. Make the same kind of plots as in step 2 but use thresholds  $r = 0.2$  and  $r = 0.8$ . By using these plots, comment on what happens with the prediction when the  $r$  value changes.

5. Perform a basis function expansion trick by computing new features:

$$z_1 = x_1^4, \quad z_2 = x_1^3 x_2, \quad z_3 = x_1^2 x_2^2, \quad z_4 = x_1 x_2^3, \quad z_5 = x_2^4$$

Add them to the dataset and compute a logistic regression model with  $y$  as target and  $x_1, x_2, z_1, \dots, z_5$  as features.

- Create a scatterplot of the same kind as in step 2 for this model.
- Compute the training misclassification rate.
- What can you say about the quality of this model compared to the previous logistic regression model?
- How has the basis expansion trick affected the shape of the decision boundary and the prediction accuracy?

```
1 library(ggplot2)
2
3 # Set seed to 12345
4 set.seed(12345)
5
6 # Import data
7 data <- read.csv("../pima-indians-diabetes.csv", header = FALSE)
8
9 # Add headers to data
10 colnames(data) <- c("times_pregnant",
11                    "plasma",
12                    "blood_pressure",
13                    "skinfold_thickness",
14                    "insulin",
15                    "bmi",
16                    "diabetes_pedigree",
17                    "age",
18                    "diabetes")
19
20 # ----- TASK 1 -----
21
22 # Plot data where age is the x axis, plasma is the y axis and color
23 # is based
24 # on diabetes value
25 ggplot(data, aes(x = age, y = plasma, color = diabetes)) +
26   geom_point() + scale_color_gradient(low = "green", high = "red") +
27   labs(title = "Task 1",
28        x = "Age",
29        y = "Plasma Glucose Concentration",
30        color = "Diabetes") +
31   theme_minimal()
32
33 # By observing the result of the task, we deem that it is difficult
34 # to classify diabetes by a standard logistic regression model that
35 # uses these two variables as features. When looking at the graph it
36 # becomes quite clear that there isn't a strong enough connection
37 # between plasma levels and age to determine if a person has
38 # diabetes or not. There are for example several people between age
39 # 60 and 70 with a plasma level above 150 whom do not have diabetes.
```

Meanwhile there are also young people below the age of 30 with a plasma level below 100 whom have diabetes. In conclusion, we determine that diabetes is hard to classify using these two variable. However, we have observed that diabetes seems more common for older ages and higher plasma levels.

```

33
34 # ----- TASK 2 -----
35
36 model = glm(diabetes~age + plasma,data=data, family = binomial)
37 predictions <- predict(model, type = "response")
38
39 #Initialize all different R values to be used for future plotting.
40 r1 <- 0.5
41
42
43 #Calculate binary predictions based on all three r values.
44 binary_predictions_1 <- ifelse(predictions >= r1, 1, 0)
45
46 #Function to calculate missclassification error.
47 missclass = function(X,X1){
48   n=length(X)
49   return(1-sum(diag(table(X,X1)))/n)
50 }
51
52 #Calculate missclassification error.
53 misclassification_error = missclass(binary_predictions_1,
54   data$diabetes) # 0.2630208
55
56 # ANSWER:
57 # the family variable is set to binomial because the response
58 # variable follows a binomial distribution, since it is a binary
59 # variable.
60
61 #Plot binary predictions and the decision boundary where r = 0.5.
62 ggplot(data, aes(x = age, y = plasma, color = binary_predictions_1)) +
63   geom_point() + scale_color_gradient(low = "green", high = "red") +
64   labs(title = "Diabetes Predictions r = 0.5",
65     x = "Age",
66     y = "Plasma Glucose Concentration",
67     color = "Diabetes") +
68   theme_minimal()
69
70 # ANSWER:
71 # Probabilistic model and figure is under the code.
72
73 # ----- TASK 3 -----
74
75 #Obtain coefficients from model.
76 b0 <- coef(model)[1]
77 b1 <- coef(model)[2]
78 b2 <- coef(model)[3]

```

```

78 #Arrange X1 in order from smallest to largest.
79 x1 <- seq(min(data$age), max(data$age))
80
81 # Calculate x2 values for the decision boundary.
82 #  $b_0 + b_1 * x_1 + b_2 * x_2 = 0 \Rightarrow x_2$ 
83 x2 <- -(b0 + b1 * x1) / b2
84
85 #Create the decision boundary line based on x1 and x2 values
86 decision_boundary <- data.frame(x1 = x1, x2 = x2)
87
88 #Plot binary predictions and the decision boundary where r = 0.5.
89 ggplot(data, aes(x = age, y = plasma, color = binary_predictions_1)) +
90   geom_point() + scale_color_gradient(low = "green", high = "red") +
91   labs(title = "Diabetes Predictions r = 0.5 + decision boundary",
92        x = "Age",
93        y = "Plasma Glucose Concentration",
94        color = "Diabetes") +
95   geom_line(data = decision_boundary, aes(x = x1, y = x2), color =
96         "blue", linewidth = 1) + # Add decision boundary
97   theme_minimal()
98
99 # ANSWER:
100 # It seems to catch the data distribution well since people predicted
101 # to suffer from diabetes or above it or on it and people predicted
102 # not to suffer from diabetes are on it or below it. The decision
103 # boundary can be observed in figure.
104
105 # ----- TASK 4 -----
106
107 r2 <- 0.2
108 r3 <- 0.8
109
110 binary_predictions_2 <- ifelse(predictions >= r2, 1, 0)
111 binary_predictions_3 <- ifelse(predictions >= r3, 1, 0)
112
113 ggplot(data, aes(x = age, y = plasma, color = binary_predictions_2)) +
114   geom_point() + scale_color_gradient(low = "green", high = "red") +
115   labs(title = "Diabetes Predictions r = 0.2",
116        x = "Age",
117        y = "Plasma Glucose Concentration",
118        color = "Diabetes") +
119   theme_minimal()
120
121 ggplot(data, aes(x = age, y = plasma, color = binary_predictions_3)) +
122   geom_point() + scale_color_gradient(low = "green", high = "red") +
123   labs(title = "Diabetes Predictions, r = 0.8",
124        x = "Age",
125        y = "Plasma Glucose Concentration",
126        color = "Diabetes") +
127   theme_minimal()
128
129 # ANSWER:
130 # By observing these plots it becomes evident that r decides how

```

certain our model has to be in order to decide on whether a person suffers from diabetes or not. In other words, a person is predicted to suffer from diabetes if the predicted probability of the model for that person is greater than  $r$ . This results in a greater amount of people predicted to suffer from diabetes by the model for lower  $r$  values.

```

127
128 # ----- TASK 5 -----
129
130 data$z1 = data$plasma^4
131 data$z2 = data$plasma^3 * data$age
132 data$z3 = data$plasma^2 * data$age^2
133 data$z4 = data$plasma * data$age^3
134 data$z5 = data$age^4
135
136 model2 <- glm(diabetes ~ plasma + age + z1 + z2 + z3 + z4 + z5,
137               family = binomial, data = data)
138
139 predictions <- predict(model2, newdata = data, type = "response")
140 r <- 0.5
141 binary_predictions <- ifelse(predictions >= r, 1, 0)
142
143 ggplot(data, aes(x = age, y = plasma, color = binary_predictions)) +
144   geom_point() + scale_color_gradient(low = "green", high = "red") +
145   labs(title = "Diabetes Predictions r = 0.5",
146        x = "Age",
147        y = "Plasma Glucose Concentration",
148        color = "Diabetes") +
149   theme_minimal()
150
151 misclassification_error <- missclass(data$diabetes,
152                                     binary_predictions)
153 print(misclassification_error) # 0.2447917
154
155 # ANSWER:
156 # This is a small improvement compared to the misclassification
157   error observed in task 2, but one could argue that the two answers
158   are in a similar range and neither one is particularly good. even
159   though this new model has a slightly smaller misclassification
160   error than the previous one, it has one large issue which is that
161   the prediction rate for diabetes decreases for higher ages. This
162   seems like an unwanted behaviour since diabetes is usually more
163   common in higher ages, but this prediction model basically
164   predicts that people above a certain age (around 65) never have
165   diabetes.

```

Probabilistic equation of the model in task 2:

$$P(\text{diabetes} | \text{age, plasma}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 * \text{age} + \beta_2 * \text{plasma})}}$$



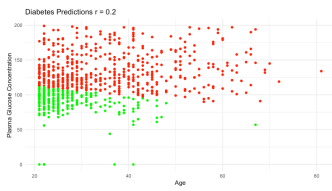


Figure 8: diabetes pred 0.2

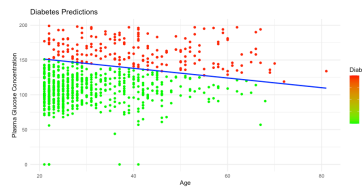


Figure 9: diabetes pred 0.5

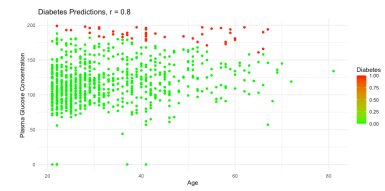


Figure 10: diabetes pred 0.8

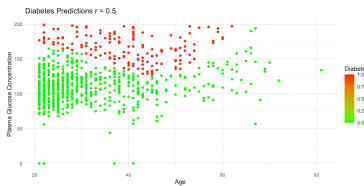


Figure 11: diabetes pred final

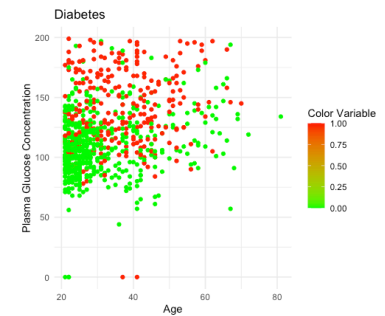


Figure 12: task 1 scatter plot

## 6 Lab 2 Task 1: Tecator Dataset Analysis

The file `tecator.csv` contains the results of a study aimed at investigating whether a near-infrared absorbance spectrum can predict the fat content of meat samples. For each meat sample, the data consists of a 100-channel spectrum of absorbance records and the levels of moisture (water), fat, and protein. The absorbance is  $-\log_{10}$  of the transmittance measured by the spectrometer. The moisture, fat, and protein levels are determined by analytic chemistry. Divide the data randomly into train and test sets (50/50) using the codes from the lectures.

1. Assume that Fat can be modeled as a linear regression in which absorbance characteristics (Channels) are used as features.
  - Report the underlying probabilistic model.
  - Fit the linear regression to the training data and estimate the training and test errors.
  - Comment on the quality of fit and prediction and therefore on the quality of the model.
2. Assume now that Fat can be modeled as a LASSO regression in which all Channels are used as features.
  - Report the cost function that should be optimized in this scenario.
3. Fit the LASSO regression model to the training data.
  - Present a plot illustrating how the regression coefficients depend on the log of the penalty factor ( $\log \lambda$ ).
  - Interpret this plot. What value of the penalty factor can be chosen if we want to select a model with only three features?
4. Repeat step 3 but fit Ridge regression instead of LASSO and compare the plots from steps 3 and 4.
  - Comment on the conclusions.
5. Use cross-validation with the default number of folds to compute the optimal LASSO model.

- Present a plot showing the dependence of the CV score on  $\log \lambda$  and comment on how the CV score changes with  $\log \lambda$ .
- Report the optimal  $\lambda$  and how many variables were chosen in this model.
- Does the information displayed in the plot suggest that the optimal  $\lambda$  value results in a statistically significantly better prediction than  $\log \lambda = -4$ ?
- Finally, create a scatter plot of the original test versus predicted test values for the model corresponding to the optimal  $\lambda$  and comment on whether the model predictions are good.

```

1 library(glmnet)
2
3 data=read.csv("../tecator.csv", header = TRUE)
4
5 # Set random seed so lab result is predictable
6 set.seed(12345)
7
8 # Create variable without protein and moisture columns.
9 channel_data <-data.frame(data[c(2:102)])
10
11 # Split the dataset into training (50%), and test (50%).
12 n = dim(channel_data)[1] # Total nr of rows in dataset
13 id = sample(1:n, floor(n * 0.5)) # Randomly selects 50 % rows for
   training
14
15 id1 = setdiff(1:n, id) # Remaining rows
16 data_train = channel_data[id, ] # Assign rows to training set
17 data_test = channel_data[id1, ] # Set remaining rows as test data
18
19 # Task 1
20 # Create a linear regression model where Fat is the target
21 # . indicates all other columns are prediction variables
22 model=lm(Fat~., data=data_train)
23 model
24
25 predict_train <- predict(model,data_train) # make predictions for
   training data
26 predict_test <- predict(model,data_test) # make predictions for test
   data
27
28 # calculate difference between predictions and actual values for both
   data sets
29 train_diff <- (data_train$Fat - predict_train)
30 test_diff <- (data_test$Fat - predict_test)
31
32 # Caclulate Squared Error for both data sets.
33 sum(train_diff^2)/dim(data_train)[1]
34 sum(test_diff^2)/dim(data_test)[1]
35
36 # Model is overfitted, probably due to 100 different variables.
37 # n is barley larger than p
38 # mse_test = 722.4294
39 # mse_train = 0.005709117

```

```

40
41 # Task 3
42 #Create a matrix without the fat column for the GLMNET model.
43 X_train_matrix <- as.matrix(data_train[c(1:100)])
44
45 # Create the lasso model by setting alpha = 1.
46 lasso <- glmnet(X_train_matrix, data_train$Fat, alpha = 1, family =
47   "gaussian")
48
49 # Look how many features the model has for the lambda = 0.818
50 coef(lasso, s = 0.818)
51 plot(lasso, xvar = "lambda")
52
53 # By looking at the graph, a model with only three coefficients is
54   approximately
55 # log lambda = -0,2 which results in lambda = 0.8187
56
57 # Task 4
58 #Create the ridge model by setting alpha = 0.
59 ridge <- glmnet(X_train_matrix, data_train$Fat, alpha = 0, family =
60   "gaussian")
61 plot(ridge, xvar = "lambda")
62
63 # task 5
64 # Train covariane lasso model
65 cv_lasso=cv.glmnet(X_train_matrix, data_train$Faas, alpha=1
66   ,family="gaussian")
67 cv_lasso
68 cv_lasso$lambda.min
69 plot(cv_lasso)
70
71 coef(cv_lasso, s = cv_lasso$lambda.min)
72 #8 coefficients are chosen in the model, intercept is NOT included
73
74 X_test_matrix <- as.matrix(data_test[c(1:100)])
75 predictions = predict(cv_lasso, s = "lambda.min", newx =
76   X_test_matrix, type = "response")
77
78 sum((predictions - mean(data_test$Fat))^2)/sum((data_test$Fat -
79   mean(data_test$Fat))^2)
80 sum((predictions - data_test$Fat)^2)
81
82 plot(data_test$Fat, col = "black", pch = 16, ylab = "Fat")
83 points(predictions, col="red", pch = 17)
84 legend("topleft",inset = c(0, -0.24),
85   xpd = TRUE, legend = c("Test data", "Predictions"),
86   col = c("black", "red"), pch = c(16, 17)
87 )

```

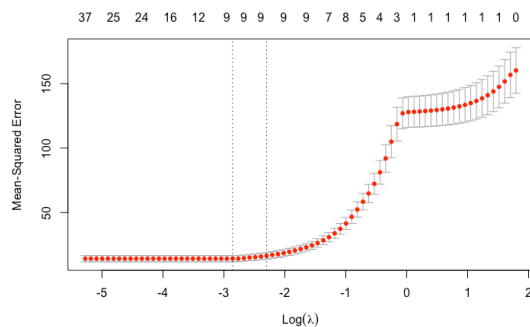


Figure 13: cv lasso

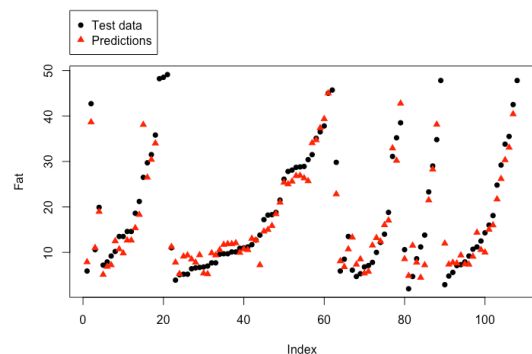


Figure 14: fat scatter

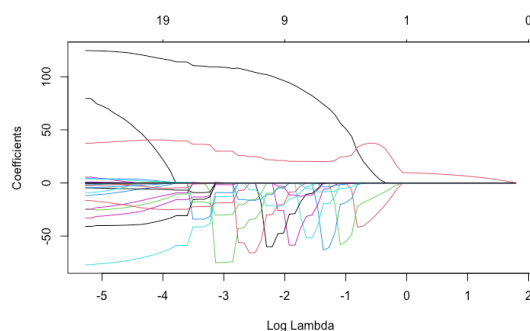


Figure 15: lasso regression

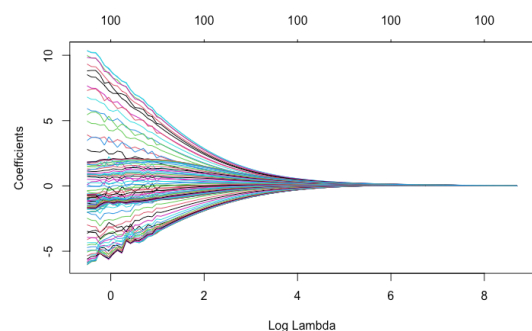


Figure 16: ridge regression plot

## 7 Lab 2 Task 2: Decision Trees and Logistic Regression for Bank Marketing

The data file `bank-full.csv` is related to direct marketing campaigns of a Portuguese banking institution. The marketing campaigns were based on phone calls. Often, more than one contact with the same client was required to assess if the product (bank term deposit) would be ('yes') or not ('no') subscribed.

### Input Variables

- **Bank client data:**

1. Age (numeric).
2. Job: type of job (categorical: 'admin.', 'blue-collar', 'entrepreneur', 'housemaid', 'management', 'retired', 'self-employed', 'services', 'student', 'technician', 'unemployed', 'unknown').
3. Marital: marital status (categorical: 'divorced', 'married', 'single', 'unknown'; note: 'divorced' means divorced or widowed).
4. Education (categorical: 'basic.4y', 'basic.6y', 'basic.9y', 'high.school', 'illiterate', 'professional.course', 'university.degree', 'unknown').
5. Default: has credit in default? (categorical: 'no', 'yes', 'unknown').

6. Housing: has housing loan? (categorical: 'no', 'yes', 'unknown').

7. Loan: has personal loan? (categorical: 'no', 'yes', 'unknown').

- **Related to the last contact of the current campaign:**

1. Contact: contact communication type (categorical: 'cellular', 'telephone').

2. Month: last contact month of the year (categorical: 'jan', 'feb', ..., 'nov', 'dec').

3. Day of week: last contact day of the week (categorical: 'mon', 'tue', 'wed', 'thu', 'fri').

4. Duration: last contact duration, in seconds (numeric). **Important note:** This attribute highly affects the output target (e.g., if `duration=0`, then `y='no'`). Yet, the duration is not known before a call is performed. Also, after the end of the call, `y` is obviously known. Thus, this input should only be included for benchmark purposes and should be discarded if the intention is to have a realistic predictive model.

- **Other attributes:**

1. Campaign: number of contacts performed during this campaign and for this client (numeric, includes the last contact).

2. Pdays: number of days that passed by after the client was last contacted from a previous campaign (numeric; 999 means client was not previously contacted).

3. Previous: number of contacts performed before this campaign and for this client (numeric).

4. Poutcome: outcome of the previous marketing campaign (categorical: 'failure', 'nonexistent', 'success').

## Output Variable (Target)

- `y`: Has the client subscribed a term deposit? (binary: 'yes', 'no').

## Tasks

1. Import the data to R, remove the variable `duration`, and divide it into training/validation/test sets as 40/30/30. Use data partitioning code specified in Lecture 2a.
2. Fit decision trees to the training data by changing the default settings one by one (i.e., not simultaneously):
  - a) Decision tree with default settings.
  - b) Decision tree with smallest allowed node size equal to 7000.
  - c) Decision tree with minimum deviance equal to 0.0005.

Report the misclassification rates for the training and validation data. Which model is the best among these three? Report how changing the deviance and node size affected the size of the trees and explain why.

3. Use the training and validation sets to choose the optimal tree depth in model 2c. Study the trees up to 50 leaves. Present a graph of the dependence of deviances for the training and validation data on the number of leaves. Interpret this graph in terms of the bias-variance tradeoff. Report the optimal number of leaves and identify which variables seem to be most important for decision-making in this tree. Interpret the information provided by the tree structure.

4. Estimate the confusion matrix, accuracy, and F1 score for the test data by using the optimal model from step 3. Comment on whether the model has good predictive power and which of the measures (accuracy or F1 score) should be preferred here.
5. Perform a decision tree classification of the test data using the following loss matrix:

$$L = \begin{bmatrix} 0 & 5 \\ 1 & 0 \end{bmatrix}$$

Report the confusion matrix for the test data. Compare the results with the results from step 4 and discuss how the rates have changed and why.

6. Use the optimal tree and a logistic regression model to classify the test data using the following principle:

$$Y = \begin{cases} y_{pos} & \text{if } P(Y = y_{pos}|X) > \pi, \\ y_{neg} & \text{otherwise.} \end{cases}$$

where  $\pi = 0.05, 0.1, 0.15, \dots, 0.9, 0.95$ . Compute the TPR and FPR values for the two models and plot the corresponding ROC curves. What are your conclusions? Why might the precision-recall curve be a better option here?

```

1 ##### Packages #####
2
3 library(tree)
4 library(rpart)
5
6 ##### PART 1 #####
7
8 # Import bank marketing campaign csv file with str as factors
9 data_frame = read.csv2("../bank-full.csv", stringsAsFactors = TRUE)
10
11 # Remove column "duration" using %in% operator
12 data_frame = data_frame[ , !names(data_frame) %in% c("duration")]
13
14 # Split the dataset into training (50%), validation (25%), and test
    (25%) sets
15 n = dim(data_frame)[1]          # Total nr of rows in dataset
16 set.seed(12345)
17 id = sample(1:n, floor(n * 0.4)) # Randomly selects 40 % rows for
    training
18 data_train = data_frame[id, ]    # Assign rows to training set
19
20 id1 = setdiff(1:n, id)           # Remaining rows
21 set.seed(12345)
22 id2 = sample(id1, floor(n * 0.30)) # Randomly selects 30 % rows for
    validation
23 data_validation = data_frame[id2, ] # Assign rows to validation set
24
25 id3 = setdiff(id1, id2)          # Remaining rows
26 data_test = data_frame[id3, ]    # Assign rest of 30 % rows to
    testing
27
28
29 ##### PART 2 #####

```

```

30
31 # Fit decision trees to training dataset
32 tree_a = tree(formula = y ~ ., data = data_train)
33 tree_b = tree(formula = y ~ ., data = data_train, control =
    tree.control(nrow(data_train), minsize = 7000))
34 tree_c = tree(formula = y ~ ., data = data_train, control =
    tree.control(nrow(data_train), mindev = 0.0005))
35
36 # Predictions on validation dataset
37 predictions_a = predict(tree_a, newdata = data_validation, type =
    "class")
38 predictions_b = predict(tree_b, newdata = data_validation, type =
    "class")
39 predictions_c = predict(tree_c, newdata = data_validation, type =
    "class")
40
41 # Misclassification rates for the training dataset
42 summary(tree_a) # Misclassification error rate: 0.1048
43 summary(tree_b) # Misclassification error rate: 0.1048
44 summary(tree_c) # Misclassification error rate: 0.09362
45
46 # Misclassification rates for the validation dataset
47 misclass = function(X, X1) {
48   n = length(X)
49   return(1 - sum(diag(table(X, X1))) / n)
50 }
51 misclass(predictions_a, data_validation$y) # Misclassification error
    rate: 0.1092679
52 misclass(predictions_b, data_validation$y) # Misclassification error
    rate: 0.1092679
53 misclass(predictions_c, data_validation$y) # Misclassification error
    rate: 0.1118484
54
55 #The best models are A and B, as they have low misclassification
    rates for both the training set and the validation set. The
    problem with Model C is that it has a good misclassification rate
    for the training set but a higher misclassification rate for the
    validation set, which is typical of overfitting.
56
57 #A larger minimum node size increases the model misclassification
    rate while a lower does not improve performance. Increasing the
    minsize value forces the tree to stop splitting when a node's
    number of observations falls below the threshold. This results in
    simpler and more general trees. Lowering the minsize can allow the
    tree to grow very deep, increasing the risk of overfitting. While
    high values the tree is constrained from splitting too much, which
    can lead to underfitting.
58
59 #A larger minimum deviance makes the model better for the validation
    set but worse for the training set. This is because a higher
    minimum deviance results in a smaller and simpler tree, reducing
    the risk of overfitting but increasing the risk of underfitting.
    Lowering the minimum deviance allows the tree to grow larger and

```

more complex, improving performance on the training set but increasing the risk of overfitting, which can lead to poor generalization and worse performance on the validation set.

```
##### PART 3 #####

# Init vector with length 50 to store score
train_score = rep(0, 50)
valid_score = rep(0, 50)

# Prune the C tree starting with at least 2 node leaves up to 50
# Calculate the deviance on both datasets and save to the score
  vector.
for(leaves in 2:50) {
  # Prune Tree C to current "leaves" number
  pruned_tree = prune.tree(tree_c, best = leaves)

  # Predictions on validation dataset using pruned tree
  predictions_valid = predict(pruned_tree, newdata = data_validation,
    type = "tree")

  train_score[leaves] = deviance(pruned_tree)
  valid_score[leaves] = deviance(predictions_valid)
}

# Graph of the dependence of deviance for the datasets on the number
  of leaves
plot(2:50, train_score[2:50], type = "b", col = "red", ylim =
  c(min(valid_score[2:50]), max(train_score[2:50])),
  main = "Optimal tree depth", ylab = "Deviance", xlab = "Number
    of leaves")
points(2:50, valid_score[2:50], type = "b", col = "blue")
legend("topright", c("train data", "validation data"), fill =
  c("red", "blue"))

# Optimal number of leaves that minimize training deviance.
optimal_train = which.min(train_score[2:50])
optimal_validation = which.min(valid_score[2:50])

optimal_tree = prune.tree(tree_c, best = optimal_validation)

# Display optimal tree
plot(optimal_tree)
text(optimal_tree, pretty=0)
optimal_tree
summary(optimal_tree)

# The optimal number of leaves is 47 for the training dataset and 21
  for the validation dataset. The bias-variance tradeoff explains
  how a model's complexity influences prediction accuracy and
  generalization to unseen data. A tree with many leaves (greater
  depth) has low bias but high variance, potentially leading to
  overfitting. A tree with fewer leaves has low variance but high
```



bias, which can result in underfitting. As model complexity increases, variance rises while bias decreases. In decision trees, the number of leaves directly impacts this balance as it defines the complexity of the tree. The key is to find the optimum where there is a balance between bias and variance so the error is minimized, see Figure \ref{fig:task\_3\_graph}.

```
# The variables which are the most important is poutcome which is the
# root node but also month, contact and pdays which appear
# frequently in splits in the tree, see Figure
# \ref{fig:task3_optimal_tree}.
```

```
##### PART 4 #####
```

```
# Predictions on the test dataset using optimal tree
```

```
predications_test = predict(optimal_tree, newdata = data_test, type =
"class")
```

```
# Creating confusion matrix
```

```
conf_matrix = table(data_test$y, predications_test)
```

```
conf_matrix
```

```
# Calculates the accuracy and f1 score from confusion matrix
```

```
accuracy_and_f1 = function(conf_matrix) {
```

```
  TP = conf_matrix["yes", "yes"] # True positives
```

```
  TN = conf_matrix["no", "no"] # True negatives
```

```
  FP = conf_matrix["no", "yes"] # False positives
```

```
  FN = conf_matrix["yes", "no"] # False negatives
```

```
  precision = TP / (TP + FP)
```

```
  recall = TP / (TP + FN)
```

```
  accuracy = (TN + TP) / (TN + FP + TP + FN)
```

```
  f1_score = 2 * (precision * recall) / (precision + recall)
```

```
  return(c(accuracy = accuracy, f1_score = f1_score))
```

```
}
```

```
# Print result
```

```
result = accuracy_and_f1(conf_matrix)
```

```
print(result["accuracy"])
```

```
print(result["f1_score"])
```

```
# The model has an accuracy of 0.8923 and an F1-score of 0.2849. The
# high accuracy might seem good, but 'its mainly due to the
# imbalance in the dataset. The F1-score, which considers both
# precision and recall, gives a better picture and shows that the
# model struggles to predict the less common class accurately. This
# shows why accuracy alone 'isnt a reliable metric for imbalanced
# datasets.
```

```
##### PART 5 #####
```

```

136 # Loss matrix defined in task
137 loss_matrix = matrix(c(0, 1, 5, 0), byrow = TRUE, nrow = 2,
138                       dimnames = list(c("no", "yes"), c("no", "yes")))
139
140 # Predictions on the test dataset using optimal tree
141 predications_test = predict(optimal_tree, newdata = data_test)
142
143 # Predictions multiplied with loss matrix
144 losses = predications_test %*% loss_matrix
145
146 # For each observation (row) it finds the class (column)
147 # with the smallest loss and returns the index of that class.
148 predicted_classes = apply(losses, 1, which.min)
149
150 # Map the correct column names (no, yes) for the predicted classes
151 predicted_classes = colnames(losses)[predicted_classes]
152
153 # Create confusion matrix
154 conf_matrix = table(data_test$y, predicted_classes)
155 conf_matrix
156
157 # Print result
158 result = accuracy_and_f1(conf_matrix)
159 print(result["accuracy"])
160 print(result["f1_score"])
161
162 # The model accuracy has decreased to 0.8732, while the F1-score has
    improved to 0.4826. This change is due to the introduction of a
    loss matrix that imposes a heavier penalty on false negatives
    compared to false positives. Correct predictions, represented by
    the diagonal entries of the matrix, are not penalized.
163
164 ##### PART 6 #####
165
166 tpr_and_fpr = function(conf_matrix) {
167   TP = conf_matrix["yes", "yes"] # True positives
168   TN = conf_matrix["no", "no"] # True negatives
169   FP = conf_matrix["no", "yes"] # False positives
170   FN = conf_matrix["yes", "no"] # False negatives
171
172   tpr = TP / (TP + FN) # true positive rate
173   fpr = FP / (FP + TN) # false positive rate
174
175   return(c(tpr = tpr, fpr = fpr))
176 }
177
178 # True positive rates and False positive rates
179 tpr_tree = c()
180 fpr_tree = c()
181 tpr_lrm = c()
182 fpr_lrm = c()
183
184 # Logistic regression model

```

```

185 lrm = glm(formula = y ~ ., data = data_train, family = "binomial")
186
187 # Predictions on test data with both models
188 predictions_test = predict(optimal_tree, newdata = data_test, type =
    "vector")
189 predictions_lrm = predict(lrm, newdata = data_test, type = "response")
190
191 for (pi in seq(from = 0.05, to = 0.95, by = 0.05)) {
192   # Optimal tree predictions
193   y_hat_optimal_tree = ifelse(predictions_test[, "yes"] > pi, "yes",
    "no")
194   conf_matrix = table(data_test$y, factor(y_hat_optimal_tree, levels
    = c("no", "yes")))
195   result = tpr_and_fpr(conf_matrix)
196   tpr_tree = append(tpr_tree, result["tpr"])
197   fpr_tree = append(fpr_tree, result["fpr"])
198
199   # Logistic regression model
200   y_hat_lrm = ifelse(predictions_lrm > pi, "yes", "no")
201   conf_matrix = table(data_test$y, y_hat_lrm)
202   result = tpr_and_fpr(conf_matrix)
203   tpr_lrm = append(tpr_lrm, result["tpr"])
204   fpr_lrm = append(fpr_lrm, result["fpr"])
205 }
206
207 # Plot ROC curves
208 plot(fpr_tree, tpr_tree, type = "b", col = "red", ylim =
    c(min(tpr_tree), max(tpr_tree))
209     ,xlim = c(min(fpr_tree), max(fpr_tree)),
210     main = "ROC curves", ylab = "True positive rate", xlab = "False
    positive rate")
211 points(fpr_lrm, tpr_lrm, type = "b", col = "blue")
212 legend("bottomright", c("Optimal tree", "Logistic regression model"),
    fill = c("red", "blue"))
213
214 # The ROC curve in Figure \ref{fig:roc_curves} shows that the optimal
    tree model performs slightly better than the logistic regression
    model, with higher True Positive Rate (TPR) across the False
    positive Rate (FPR) values. Since the test data has very
    imbalanced classes, the precision-recall curve would be a better
    choice here. Because it focuses on the minority class, which would
    be more informative for an imbalanced dataset such as this.

```

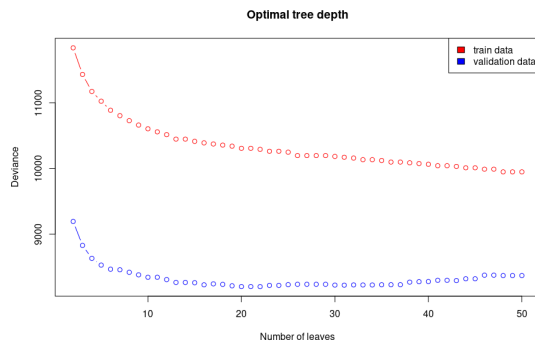


Figure 17: dependency deviation

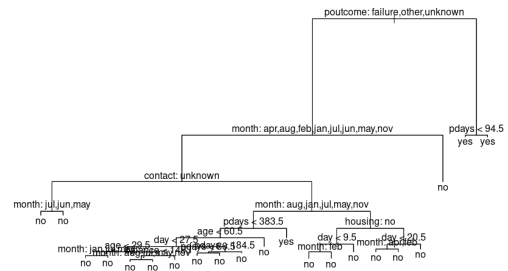


Figure 18: optimal tree

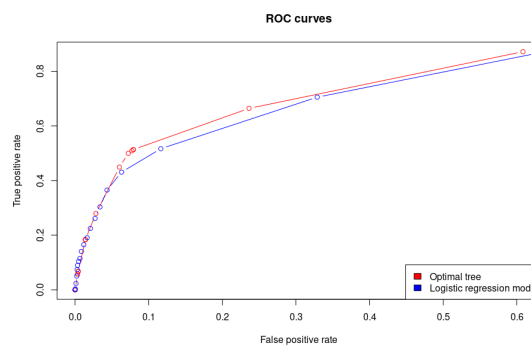


Figure 19: roc curves

## 8 Lab 2 Task 3: Principal Components and Implicit Regularization

The data file `communities.csv` contains the results of studies of the crime level in the United States based on various characteristics of a given location. The main variable studied is `ViolentCrimesPerPop`, which represents the total number of violent crimes per 100K population.

### Tasks

- Scale all variables except `ViolentCrimesPerPop` and implement PCA using the `eigen()` function.
  - Report how many components are needed to obtain at least 95% of the variance in the data.
  - What is the proportion of variation explained by each of the first two principal components?
- Repeat the PCA analysis using the `princomp()` function and make the trace plot of the first principal component.
  - Do many features have a notable contribution to this component?
  - Report which 5 features contribute the most (by absolute value) to the first principal component.

- Comment on whether these features have anything in common and whether they may have a logical relationship to the crime level.
  - Provide a plot of the PC scores in the coordinates (PC1, PC2) where the color of the points corresponds to `ViolentCrimesPerPop`. Analyze this plot (hint: use the `ggplot2` package).
3. Split the original data into training and test sets (50/50), scale both features and response appropriately, and estimate a linear regression model from the training data where `ViolentCrimesPerPop` is the target and all other data columns are features.
    - Compute the training and test errors for these data.
    - Comment on the quality of the model.
  4. Implement a function that depends on the parameter vector  $\theta$  and represents the cost function for linear regression without intercept on the training dataset.
    - Use the BFGS method (`optim()` function without the gradient specified) to optimize this cost with the starting point  $\theta_0 = 0$ .
    - Compute the training and test errors for every iteration number.
    - Present a plot showing the dependence of both errors on the iteration number.
    - Comment on which iteration number is optimal according to the early stopping criterion.
    - Compute the training and test errors in the optimal model, compare them with the results in step 3, and make conclusions.

## Hints

- **Hint 1:** Don't store parameters from each iteration (this will require a lot of memory). Instead, compute and store test errors directly.
- **Hint 2:** Discard some initial iterations (e.g., 500) in your plot to make the dependencies visible.

```

1 data <- read.csv('communities.csv')
2
3 # task 1
4
5 library(caret)
6 scaler<-preProcess(data)
7 dataS<-predict(scaler,data)
8 dataS$ViolentCrimesPerPop <- data$ViolentCrimesPerPop # no scaling on
  ViolentCrimesPerPop according to instructions
9 S <- (1/101) * t(as.matrix(dataS)) %*% as.matrix(dataS) # covariance
  matrix, how variables depends on each other
10
11 eig <- eigen(S)
12 lambda <- eig$values
13 lambda_perc <- lambda/sum(lambda)*100 # gives the variance of each
  'projection' (how much data the projection can represent)
14 sprintf("%.3f",lambda_perc)
15 # prints:
16 # [1] "25.025" "16.931" "9.298" "7.554" "5.660" "4.236" "3.226"
  "2.969" "2.067" "1.617" "1.572" "1.470"

```

```

17 # [13] "1.414" "1.030" "0.931" "0.890" "0.748" "0.705" "0.649"
    "0.638" "0.624" "0.568" "0.542" "0.519"
18 # [25] "0.504" "0.480" "0.466" "0.451" "0.431" "0.386" "0.365"
    "0.351" "0.337" "0.311" "0.287" "0.259"
19 # [37] "0.256" "0.245" "0.240" "0.222" "0.211" "0.205" "0.200"
    "0.190" "0.183" "0.165" "0.160" "0.142"
20 # [49] "0.138" "0.126" "0.112" "0.107" "0.104" "0.100" "0.090"
    "0.081" "0.077" "0.073" "0.069" "0.066"
21 # [61] "0.064" "0.062" "0.058" "0.051" "0.049" "0.046" "0.044"
    "0.042" "0.040" "0.038" "0.035" "0.034"
22 # [73] "0.032" "0.031" "0.028" "0.026" "0.024" "0.022" "0.021"
    "0.020" "0.018" "0.018" "0.016" "0.015"
23 # [85] "0.014" "0.013" "0.011" "0.009" "0.008" "0.006" "0.006"
    "0.005" "0.004" "0.004" "0.003" "0.002"
24 # [97] "0.002" "0.001" "0.001" "0.001" "0.001"
25 sum(lambda_perc[1:34]) # 94.9
26 sum(lambda_perc[1:35]) # 95.2 => we need 35 variables for a variance
    of 95 %
27
28 # task 2
29
30 res <- princomp(dataS)
31 lambda <- res$sdev^2 # variance = standard_deviation^2
32 sprintf("%2.3f", lambda/sum(lambda)*100) # prints as above
33 U1 <- (res$loadings)[,1]
34 plot(sort(abs(U1)), main="Traceplot, PC1")
35 top_5 <- tail(sort(abs(U1)), 5)
36 top_5
37 # PctPopUnderPov      pctWInvInc      PctKids2Par      medIncome
    medFamInc
38 # 0.1737183          0.1748076      0.1755406      0.1818171
    0.1831478
39 # It seems that poverty is the largest contributing factor
40
41 library(ggplot2)
42 S1 <- (res$scores)[,1]
43 S2 <- (res$scores)[,2]
44
45 ggplot() + geom_point(aes(x=S1, y=S2,
    color=dataS$ViolentCrimesPerPop)) +
46   labs(title = "PC1 and PC2 scores",
47         x = "PC1",
48         y = "PC2") +
49   theme_minimal()
50
51 # TASK 3
52
53 n <- dim(data)[1]
54 set.seed(12345)
55 id <- sample(1:n, floor(n*0.5))
56 train <- data[id,]
57 test <- data[-id,]
58

```

```

59 scaler<-preProcess(train)
60 trainS<-predict(scaler,train)
61 testS<-predict(scaler,test)
62
63 fit<-lm(ViolentCrimesPerPop ~ ., data=trainS)
64 summary <- summary(fit)
65
66 dim(summary$coefficients)
67
68 MSE_train <- mean(residuals(fit)^2) # MSE for train = 0.2752071
69 MSE_test <- mean((test$ViolentCrimesPerPop - predict(fit,
70   newdata=testS))^2) # MSE for test = 0.5408757
71
72 # We see that the MSE is quite good for train but quite high for test
73 # meaning it is alright but not more.
74 # Slightly overfitted
75 # through the summary(fit) we see that racepctblack, PctWorkMom,
76 # PctPersDenseHous are significant contributors (***)
77
78 # Task 4 #####
79
80 # cost function for linear regression w/o intercept
81 trainS_X <- as.matrix(subset(trainS, select=-ViolentCrimesPerPop))
82 testS_X <- as.matrix(subset(testS, select=-ViolentCrimesPerPop))
83 trainS_Y <- trainS$ViolentCrimesPerPop
84 testS_Y <- testS$ViolentCrimesPerPop
85
86 train_MSE_vec <- c()
87 test_MSE_vec <- c()
88
89 MSE <- function(theta)
90 {
91   MSE_train <- mean( (trainS$ViolentCrimesPerPop - trainS_X %*%
92     theta)^2 )
93   MSE_test <- mean( ( testS_Y - testS_X %*% theta )^2 )
94
95   train_MSE_vec <- c(train_MSE_vec, MSE_train)
96   test_MSE_vec <- c(test_MSE_vec, MSE_test)
97
98   return(MSE_train)
99 }
100
101 optim(rep(0, 100), fn = MSE, method = "BFGS")
102
103 interval <- 1:length(train_MSE_vec)
104 ggplot() +
105   geom_point(aes(interval, train_MSE_vec[interval]), color="red") +
106   geom_point(aes(interval, test_MSE_vec[interval]), color="blue") +
107   ylim(0,1) +
108   labs(title = "MSE train (red) and test (blue)", x = "Iteration", y
109     = "MSE Score")
110
111 # early stopping criterion says to stop at the lowest test MSE

```

```

107 optimal_MSE_train <- train_MSE_vec[which.min(test_MSE_vec)] # =
    0.3032999
108 optimal_MSE_test <- min(test_MSE_vec) # = 0.4002329
109 # We see that this is a slight elevation from the previous MSE of
    0.275 for test
110 # but a fairly good decrease from 0.541 for the test MSE.
111 # Overall this is probably a better model

```

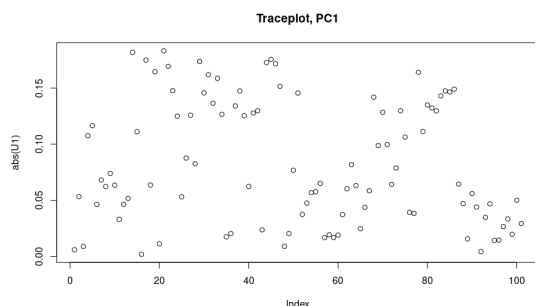


Figure 20: Traceplot PC1

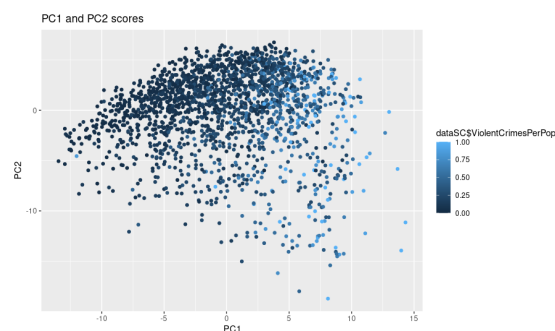


Figure 21: PC2 (y-axis) vs PC1 (x-axis)

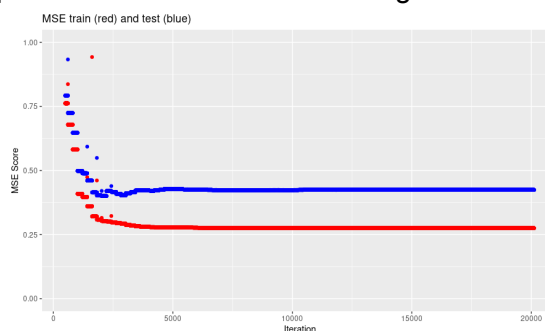


Figure 22: task 4: train = red, test = blue

## 9 Lab 3 Task 1: Kernel Method for Predicting Hourly Temperatures in Sweden

You are tasked with implementing a kernel method to predict hourly temperatures for a specific date and location in Sweden using the provided files `stations.csv` and `temps50k.csv`. These files contain data on weather stations and temperature measurements from the Swedish Meteorological and Hydrological Institute (SMHI).

### Objective

Provide a temperature forecast for a given date and location in Sweden. The forecast should include predicted temperatures from 4 am to midnight at 2-hour intervals.



## Methodology

Use a kernel composed of three Gaussian kernels:

1. **Physical Distance Kernel:** Accounts for the physical distance between a station and the point of interest. Use the `distHaversine` function from the `geosphere` package.
2. **Date Distance Kernel:** Accounts for the temporal distance between the measurement date and the prediction date.
3. **Time Distance Kernel:** Accounts for the temporal distance between the measurement time and the prediction time.

Choose appropriate smoothing coefficients ( $h_{\text{distance}}$ ,  $h_{\text{date}}$ ,  $h_{\text{time}}$ ) for each kernel manually. Show kernel values as a function of distance to demonstrate how closer points have higher values.

## Steps

1. Filter out temperature measurements that are posterior to the prediction date and time.
2. Compute kernel values for the physical distance, date distance, and time distance.
3. Combine the three kernels:
  - First, by summing them.
  - Then, by multiplying them.
4. Generate forecasts for temperatures at the specified times.
5. Compare results obtained using the summation and multiplication methods and explain differences.

## Template Code

Below is the R template code:

```
1 set.seed(1234567890)
2 library(geosphere)
3
4 # Load data
5 stations <- read.csv("stations.csv", fileEncoding = "latin1")
6 temps <- read.csv("temps50k.csv")
7 st <- merge(stations, temps, by = "station_number")
8
9 # Smoothing coefficients (to be determined manually)
10 h_distance <- # Physical distance smoothing coefficient
11 h_date <- # Date smoothing coefficient
12 h_time <- # Time smoothing coefficient
13
14 # Prediction parameters
15 a <- 58.4274 # Latitude of the point to predict
16 b <- 14.826 # Longitude of the point to predict
17 date <- "2013-11-04" # Date to predict
18 times <- c("04:00:00", "06:00:00", "08:00:00", ..., "24:00:00") #
19 Times
```

```

20 # Initialize temperature predictions
21 temp <- vector(length = length(times))
22
23 # Kernel implementation and prediction (students' code here)
24
25 # Visualization
26 plot(temp, type = "o", xlab = "Time", ylab = "Temperature (C)",
27       main = "Predicted Temperatures")

```

## Analysis

After implementing the above, repeat the exercise by combining the kernels through multiplication instead of summation. Discuss the differences in results and provide insights into why the two methods yield different outcomes.

```

1 # setwd(paste0(getwd(), "/tdde01-machine-learning/lab3/"))
2
3 set.seed(1234567890)
4 library(geosphere)
5
6 stations <- read.csv("stations.csv", fileEncoding = "latin1")
7 temps <- read.csv("temps50k.csv")
8 st <- merge(stations, temps, by="station_number")
9
10 h_pos <- c(58.41086, 15.62157) # Linköping's latitude and longitude
11 h_time <- "15:10:22"
12 h_date <- "2013-11-04"
13 times <- c("04:00:00", "06:00:00", "08:00:00", "10:00:00", "12:00:00",
14           "14:00:00", "16:00:00", "18:00:00", "20:00:00", "22:00:00",
15           "24:00:00")
16
17 ##### functions #####
18
19 ## finding the width of the first kernel
20 # max_lat_long <- c( max(st$latitude), max(st$longitude) )
21 # min_lat_long <- c( min(st$latitude), min(st$longitude) )
22 # max_dist = distHaversine(
23 #   max_lat_long,
24 #   min_lat_long
25 # )
26 # l <- 6378137 / 10 # earth radius / 10
27 # exp( - ( max_dist^2 ) / ( 2 * l^2 ) ) # gives practically zero
28   which is good
29
30 earth_distance <- c()
31 gaussian_kernel_earth <- function(x, data)
32 {
33   x_prim <- matrix( c(data$latitude, data$longitude), nrow =
34                     length(data$longitude), ncol=2)
35
36   l = 6378137 / 10 # standard deviation or in this case width = earth
37     radius / 10

```

```

35 distance <- c()
36 for (i in 1:dim(x_prim)[1])
37 {
38   distance <- c(distance, distHaversine(x, x_prim[i,]))
39 }
40 earth_distance <- distance
41 exp( - ( ( distance )^2 ) / ( 2 * 1^2 ) )
42 }
43
44 #
45 -----
46 days_distance <- c()
47 gaussian_kernel_days <- function(x, data)
48 {
49   x_prim <- data$date
50   l = sqrt(365 / 2) # standard deviation in gaussian kernel =
51     sqrt(variance)
52   x <- rep( x, length(x_prim) )
53   x_prim <- x_prim
54   days_diff <- as.numeric( difftime(x, x_prim) ) %% 365
55   days_distance <- days_diff
56
57   exp( - ( ( days_diff )^2 ) / ( 2 * 1^2 ) )
58 }
59
60 #
61 -----
62 hours_distance <- c()
63 gaussian_kernel_hours <- function(x, data)
64 {
65   x_prim <- as.POSIXct(data$time, format = "%H:%M:%S" )
66   x <- as.POSIXct(rep(x, length(data$time)) , format = "%H:%M:%S" )
67   l = sqrt(24 / 2) # standard deviation in gaussian kernel =
68     sqrt(variance)
69   time_diff <- abs( as.numeric( difftime(x,x_prim, units="hours") ) )
70     %% 24
71   hours_distance <- time_diff
72   exp( - ( ( time_diff^2 ) / ( 2 * 1^2 ) ) )
73 }
74
75 delete_posterior_dates <- function(current_date, data)
76 {
77   delete_index <- c()
78   for (i in 1:dim(data)[1])
79   {
80     if (difftime(current_date, data$date[i]) <= 0)
81     {
82       delete_index <- c(delete_index, i)
83     }
84   }
85 }

```

```

83 | data[-delete_index,]
84 | }
85 |
86 | # delete_posterior_dates("1960-04-20", st)
87 |
88 | # further tweaking of the width (l)
89 | # dist_kernel <- gaussian_kernel_earth( h_pos, st )
90 | # mean( dist_kernel )
91 | # max( dist_kernel )
92 | # min( dist_kernel )
93 | #
94 | # days_kernel <- gaussian_kernel_days(h_date, st)
95 | # mean(days_kernel)
96 | # max(days_kernel)
97 | # min(days_kernel)
98 | # #
99 | # # hours_kernel <- gaussian_kernel_hours(times[1], st)
100 | # mean(hours_kernel)
101 | # max(hours_kernel)
102 | # min(hours_kernel)
103 |
104 | total_kernel <- function(pos, date, times, data, mult_add)
105 | {
106 |   data <- delete_posterior_dates(date, data)
107 |
108 |   dist_kernel <- gaussian_kernel_earth( pos, data )
109 |   days_kernel <- gaussian_kernel_days( date, data )
110 |
111 |   total_kernel_vector <- c()
112 |   for (i in 1:length(times))
113 |   {
114 |     time_kernel <- gaussian_kernel_hours( times[i], data )
115 |     # print(time_kernel)
116 |     if (mult_add == "add")
117 |       total_kernel_vector <- c(total_kernel_vector, dist_kernel +
118 |         days_kernel + time_kernel)
119 |     else
120 |       total_kernel_vector <- c(total_kernel_vector, dist_kernel *
121 |         days_kernel * time_kernel)
122 |   }
123 |   total_matrix <- matrix(total_kernel_vector, nrow = dim(data)[1],
124 |     ncol = length(times))
125 |
126 |   temp_pred <- c()
127 |   air_temperatures <- data$air_temperature
128 |   for (i in 1:length(times))
129 |   {
130 |     # Nadaraya-Watson weighted average, eq 2.41, page 35, course book
131 |     temp_pred <- c(temp_pred, (air_temperatures %*% total_matrix[,i])
132 |       / sum(total_matrix[,i]) )
133 |   }
134 |   temp_pred
135 | }

```

```

132
133 #####
134
135 earth_kernel <- gaussian_kernel_earth( h_pos, st )
136 days_kernel <- gaussian_kernel_days( h_date, st )
137 hours_kernel <- gaussian_kernel_hours( times[1], st )
138
139 plot( earth_distance, earth_kernel, main = "earth" )
140 plot( days_distance, days_kernel, main = "days" )
141 plot( hours_distance, hours_kernel, main = "hours" )
142
143 time_of_day <- c()
144 for (i in 1:length(times))
145 {
146     time_of_day <- c(time_of_day, (strsplit(times[i], split =
147         ":")[[1]][1] )
148 }
149 time_of_day <- as.numeric(time_of_day)
150 times <- c("04:00:00", "06:00:00", "08:00:00", "10:00:00", "12:00:00",
151     "14:00:00", "16:00:00", "18:00:00", "20:00:00", "22:00:00",
152     "24:00:00")
153 pos <- c(58.41086,15.62157)
154
155 date <- "2010-06-21"
156 temp <- total_kernel(pos, date, times, st, "add")
157 plot(time_of_day, temp, type="o", main=date)
158 print(date)
159
160 date <- "2010-12-24"
161 temp <- total_kernel(pos, date, times, st, "add")
162 plot(time_of_day, temp, type="o", main=date)
163 print(date)
164
165 date <- "2010-06-21"
166 temp <- total_kernel(pos, date, times, st, "mult")
167 plot(time_of_day, temp, type="o", main=date)
168 print(date)
169
170 date <- "2010-12-24"
171 temp <- total_kernel(pos, date, times, st, "mult")
172 plot(time_of_day, temp, type="o", main=date)
173 print(date)
174
175 # We see that the sum kernel seems to always give values around 4
176 # degree celsius
177 # whereas the multiplication kernel seems to give more accurate
178 # values.
179 # This is thought to be because multiplying values between 0 and 1
180 # really accentuates
181 # the 'strong' values close to one and it really diminishes the
182 # values close to 0,
183 # whereas when summing the values it is possible that two values even
184 # each other out

```

```

179 # making it essentially only depend on the last value which might not
    # contribute majorly
180 # to the result.

```

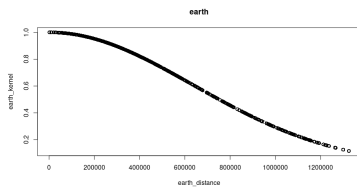


Figure 23: geo kernel

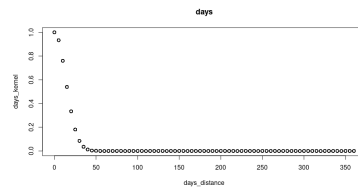


Figure 24: days kernel

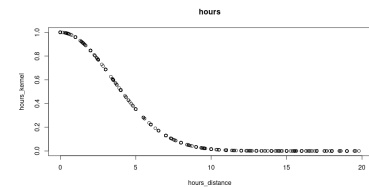


Figure 25: hours kernel

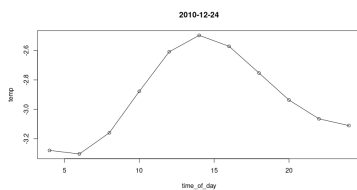


Figure 26: winter mult

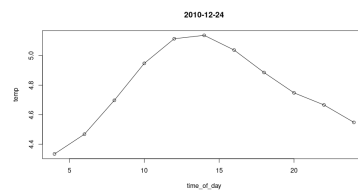


Figure 27: winter sum

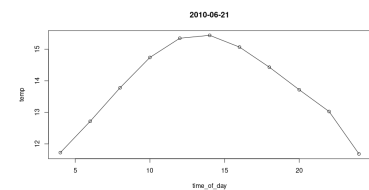


Figure 28: summer mult

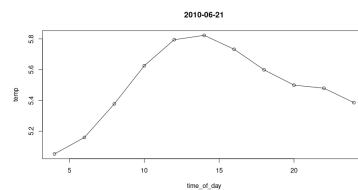


Figure 29: summer sum

## 10 Lab 3 Task 2 SVM Model Selection for Spam Dataset

The code in the file `Lab3Block1 2021 SVMs St.R` performs SVM model selection to classify the spam dataset. To do so, the code uses the function `ksvm` from the R package `kernlab`, which also includes the spam dataset. All the SVM models to select from use the radial basis function kernel (also known as Gaussian) with a width of 0.05. The  $C$  parameter varies between the models.

Run the code in the file `Lab3Block1 2021 SVMs St.R` and answer the following questions:

1. Which filter do you return to the user? `filter0`, `filter1`, `filter2` or `filter3`? Why?
2. What is the estimate of the generalization error of the filter returned to the user? `err0`, `err1`, `err2` or `err3`? Why?
3. Once an SVM has been fitted to the training data, a new point is essentially classified according to the sign of a linear combination of the kernel function values between the support vectors and the new point. You are asked to implement this linear combination for `filter3`.
  - You should make use of the functions `alphaindex`, `coef`, and `b` that return the indexes of the support vectors, the linear coefficients for the support vectors, and the negative intercept of the linear combination.

- See the help file of the kernlab package for more information.
- You can check if your results are correct by comparing them with the output of the function predict where you set type = "decision".
- Do so for the first 10 points in the spam dataset.
- Feel free to use the template provided in the file Lab3Block1 2021 SVMs St.R.

```

1 library(kernlab)
2 set.seed(1234567890)
3
4 data(spam)
5 foo <- sample(nrow(spam))
6 spam <- spam[foo,]
7 tr <- spam[1:3000, ]
8 va <- spam[3001:3800, ]
9 trva <- spam[1:3800, ]
10 te <- spam[3801:4601, ]
11 by <- 0.3
12 err_va <- NULL
13 for(i in seq(by,5,by)){
14   filter <- ksvm(type~.,data=tr,kernel="rbfdot",
15   kpar=list(sigma=0.05),C=i,scaled=FALSE)
16   mailtype <- predict(filter,va[,58])
17   t <- table(mailtype,va[,58])
18   err_va <-c(err_va,(t[1,2]+t[2,1])/sum(t))
19 }
20 filter0 <- ksvm(type~.,data=tr,kernel="rbfdot",
21 kpar=list(sigma=0.05),C=which.min(err_va)*by,scaled=FALSE)
22 mailtype <- predict(filter0,va[,58])
23 t <- table(mailtype,va[,58])
24 err0 <- (t[1,2]+t[2,1])/sum(t)
25 err0
26 filter1 <- ksvm(type~.,data=tr,kernel="rbfdot",
27 kpar=list(sigma=0.05),C=which.min(err_va)*by,scaled=FALSE)
28 mailtype <- predict(filter1,te[,58])
29 t <- table(mailtype,te[,58])
30 err1 <- (t[1,2]+t[2,1])/sum(t)
31 err1
32 filter2 <- ksvm(type~.,data=trva,kernel="rbfdot",
33 kpar=list(sigma=0.05),C=which.min(err_va)*by,scaled=FALSE)
34 mailtype <- predict(filter2,te[,58])
35 t <- table(mailtype,te[,58])
36 err2 <- (t[1,2]+t[2,1])/sum(t)
37 err2
38 filter3 <- ksvm(type~.,data=spam,kernel="rbfdot",
39 kpar=list(sigma=0.05),C=which.min(err_va)*by,scaled=FALSE)
40 mailtype <- predict(filter3,te[,58])
41 t <- table(mailtype,te[,58])
42 err3 <- (t[1,2]+t[2,1])/sum(t)
43 err3
44
45 # Questions
46

```

```

47 # 1. Which filter do we return to the user ? filter0, filter1,
    filter2 or filter3? Why?
48 # 2. What is the estimate of the generalization error of the filter
    returned to the user? err0, err1, err2 or err3? Why?
49 # 3. Implementation of SVM predictions.
50
51 sv<-alphaindex(filter3)[[1]]
52 co<-coef(filter3)[[1]]
53 inte<- - b(filter3)
54
55 k<-NULL
56 for(i in 1:10){ # We produce predictions for just the first 10 points
    in the dataset.
57     k2<-NULL
58     # Calculate k_values for each support vector.
59     for(j in 1:length(sv)){
60         # Calculate the difference between point and current support
            vector point
61         diff <- spam[i, -58] - spam[sv[j],-58]
62         # Calculate the value of the gaussian kernel for the current
            points
63         k_val<- exp(-0.05 * sum(diff^2))
64         # Add to vector of all k_values of the current point
65         k2 <- c(k2, k_val)
66     }
67     # Calculate decision function for the current point
68     k <- c(k, sum(co * k2) + inte)
69 }
70 print(k)
71 print(spam[1:10, 58])
72 predict(filter3,spam[1:10,-58], type = "decision")

```

Answers:

1.

- Filter0: trained on the training data, tested on validation data.
- Filter1: trained on the training data, tested on test data.
- Filter2: trained on training + validation data, tested on test data.
- Filter3: trained on the entire dataset, tested on test data.

A model performs better if it is trained on more data. As long as that data is accurate. Filter 0 and 1 are trained on only the training set, but we want our model to be trained on as much data as possible. Filter 2 is trained on training + validation data. However, filter3 is trained on the entirety of the spam data set. Therefore it's the most appropriate filter to return to someone who needs to use a model for making predictions on new data.

2.

- err0: obtained by filter0 which is trained on training data and validated on validation data. not appropriate because not trained on validation data. Validation data was used to choose best C. So it has already "seen" some of the data on which it was tested.
- err1: obtained by filter1 which is trained on training data and validated on tested data. Not



appropriate because not trained on validation data, I.E not the most amount of data possible which isn't test data.

- err2: obtained by filter2 which is trained on training + validation and test data. Thus an appropriate generalization error because the model is trained on the largest possible amount of data without having "seen" any of the test data.
- err3: obtained by filter3 which is trained on all data and validated on test data. Not appropriate because it has "seen" the test data.

Based on the reasoning above, err2 is the estimate of the generalization error of the filter returned to the user. The other errors are validated on the wrong dataset, are not trained on a satisfiable amount of data or have been trained on the data on which it is being validated.

3.

The decision function for the SVM for a point  $x_*$  is the following.

$$\hat{y}(x_*) = \text{sign}(\hat{\alpha}^\top K(X, x_*) + \beta)$$

Where  $\hat{\alpha}^\top$  is a vector where non-zero values represent the support vectors.  $K(X, x_*)$  is a vector created by applying the kernel on  $x_*$  and each of the training points in  $X$ .

$$K(X, x_*) = \exp(-\sigma \|X - x_*\|^2)$$

$\beta$  is the intercept of the SVM.

ANSWER TASK 3: Using the decision function, we can classify an email as spam accordingly.

$$class = \begin{cases} \text{nospam,} & \text{if } \text{sign}(\hat{y}(x_*)) = - \\ \text{spam,} & \text{if } \text{sign}(\hat{y}(x_*)) = + \end{cases}$$

The result of the SVM predictions on the first 10 points in the spam dataset can be seen in table 1

Data Point	Prediction (Class)	Decision Value ( $\hat{y}(x_*)$ )
1	Nospam	-1.0703
2	Spam	1.0003
3	Spam	0.9996
4	Nospam	-0.9999
5	Nospam	-0.9995
6	Spam	1.0001
7	Nospam	-0.8586
8	Nospam	-0.9997
9	Spam	0.9998
10	Nospam	-1.0001

Table 1: SVM Predictions for the First 10 Data Points in the spam dataset

## 11 Lab 3 Task 3: Neural Networks

This assignment is to be solved using the `neuralnet` package.

## Tasks

1. Train a neural network to learn the trigonometric sine function.
  - Sample 500 points uniformly at random in the interval  $[0, 10]$ .
  - Apply the sine function to each point. The resulting value pairs are the data points available to you.
  - Use 25 of the 500 points for training and the rest for testing.
  - Use one hidden layer with 10 hidden units. You do not need to apply early stopping.
  - Plot the training and test data, and the predictions of the learned neural network on the test data.
  - Comment on your results.
2. In question (1), you used the default logistic (a.k.a. sigmoid) activation function, i.e., `act.fct = "logistic"`.
  - Repeat question (1) with the following custom activation functions:

$$h_1(x) = x, \quad h_2(x) = \max\{0, x\}, \quad h_3(x) = \ln(1 + e^x)$$

(a.k.a. linear, ReLU, and softplus).

- See the help file of the `neuralnet` package to learn how to use custom activation functions.
  - Plot and comment on your results.
3. Sample 500 points uniformly at random in the interval  $[0, 50]$  and apply the sine function to each point.
    - Use the neural network learned in question (1) to predict the sine function value for these new 500 points.
    - You should get mixed results. Plot and comment on your results.
  4. In question (3), the predictions seem to converge to some value.
    - Explain why this happens. To answer this question, you may need to access the weights of the learned neural network.
    - You can do this by running `nn` or `nn$weights`, where `nn` is the learned neural network.
  5. Sample 500 points uniformly at random in the interval  $[0, 10]$  and apply the sine function to each point.
    - Use all these points as training points for learning a neural network that predicts  $x$  from  $\sin(x)$ , i.e., unlike before when the goal was to predict  $\sin(x)$  from  $x$ .
    - Use the learned neural network to predict the training data.
    - You should get bad results. Plot and comment on your results.
    - **Help:** Some people get a convergence error in this question. This can be solved by stopping the training before reaching convergence by setting `threshold = 0.1`.

```
1 ##### Packages #####
2
3 library(neuralnet)
4
5 ##### Part 1 #####
```

```

6
7 n = 500
8 set.seed(1234567890)
9 points = runif(n, min=0, max=10)
10 data <- data.frame(points, sin=sin(points))
11
12 # Split the dataset into 25 points for training and rest for testing
13 training_amount = 25
14 train_data = data[1:training_amount, ]
15 test_data = data[training_amount:n, ]
16
17 # Random initialization of the weights in the interval [-1, 1]
18 winit = runif(10, -1, 1)
19 nn = neuralnet(formula = sin ~ ., data = train_data, hidden = 10,
20               startweights = winit)
21
22 # Plot of the training data (black), test data (blue), and
23   predictions (red)
24 plot(train_data, cex=2, main = "Predictions of sin(x) from x")
25 points(test_data, col = "blue", cex=1)
26 points(test_data[, 1], predict(nn, test_data), col="red", cex=1)
27
28 legend("bottomright",
29       legend = c("Training Data", "Test Data", "Predictions"),
30       col = c("black", "blue", "red"),
31       pch = 1,
32       pt.cex = c(2, 1, 1))
33
34 # it can be seen that the model's predictions for the sine points
35   using the test data are highly accurate. The predicted values
36   closely overlap with the actual test data, indicating strong model
37   performance.
38
39 ##### Part 2 #####
40
41 h1 = function(x) x # Linear
42 h2 = function(x) ifelse(x >= 0, x, 0) # ReLU
43 h3 = function(x) log(1 + exp(x)) # Softplus
44 nn1 = neuralnet(formula = sin ~ ., data = train_data, hidden = 10,
45               act.fct = h1, startweights = winit)
46 nn2 = neuralnet(formula = sin ~ ., data = train_data, hidden = 10,
47               act.fct = h2, startweights = winit)
48 nn3 = neuralnet(formula = sin ~ ., data = train_data, hidden = 10,
49               act.fct = h3, startweights = winit)
50
51 plot(train_data, cex=2, main = "Predictions of sin(x) from x")
52 points(test_data, col = "blue", cex=1)
53 points(test_data[, 1], predict(nn1, test_data), col="red", cex=1)
54 points(test_data[, 1], predict(nn2, test_data), col="green", cex=1)
55 points(test_data[, 1], predict(nn3, test_data), col="orange", cex=1)
56
57 legend("bottomright",
58       legend = c("Training Data", "Test Data", "NN1 (Linear)", "NN2

```

```

    (ReLU)", "NN3 (Softplus)"),
51   col = c("black", "blue", "red", "green", "orange"),
52   pch = 1,
53   pt.cex = c(2, 1, 1, 1, 1))
54
55 # Figure shows that the Linear model performs poorly in predicting
    sin(x) from the data points, as a linear model cannot fit a
    sinusoidal curve. The ReLU activation function shows promising
    results for the initial data points but becomes worse after four
    points. This decline occurs because ReLU is a piecewise linear
    function, meaning it is linear in segments. Therefore, it
    struggles to capture the smooth, continuous oscillations of a
    sinusoidal function. On the other hand, the Softplus activation
    function allows the neural network to perfectly predict the test
    data. This is because Softplus is a smooth approximation of ReLU,
    making it well-suited for modeling non-linear functions.
56
57 ##### Part 3 #####
58
59 set.seed(1234567890)
60 points = runif(n, min=0, max=50)
61 data <- data.frame(points, sin=sin(points))
62
63 predictions = predict(nn, data)
64
65 plot(data, cex=2, ylim = c(min(predictions), max(predictions)), main
    = "Predictions of sin(x) from x")
66 points(data[, 1], predictions, col="red", cex=1)
67
68 legend("topright",
69       legend = c("New Data", "Predictions"),
70       col = c("black", "red"),
71       pch = 1,
72       pt.cex = c(2, 1))
73
74 # As shown in Figure the original NN model predicts the new data
    points perfectly until the after approximately x = 10 points. This
    is because the NN model is not trained for points larger than 10,
    leaving it unable to generalize the behavior of the sine function
    outside its training range.
75
76 ##### Part 4 #####
77
78 nn$weights
79
80 #[[1]]
81 #[[1]][[1]]
82 #      [,1]      [,2]      [,3]      [,4]      [,5]
83 #      [,6]      [,7]      [,8]      #[,9]      [,10]
84 # [1,]  3.3046893 -0.3271777  0.4043669 -0.7669631 11.955191
    -0.2057986 -6.472200  7.904868 #-0.5708964 0.04342472
    # [2,] -0.8033506 -0.8324709 -0.1499125 -0.8256431 -1.802597
    0.7993943  3.082136 -2.320714 #0.1543628 0.76288667

```

```

85 # [[1]][[2]]
86 #           [,1]
87 # [1,]  0.7993476
88 # [2,] -4.9649102
89 # [3,] -4.8295057
90 # [4,] 22.1703831
91 # [5,] -5.5590648
92 # [6,] -4.3747945
93 # [7,]  0.3489759
94 # [8,] -0.7224382
95 # [9,]  1.8612110
96 # [10,] -9.4390597
97 # [11,]  0.4400295
98
99
100 #When examining the NN weights, especially in the second layer, we
    can observe that most of them are #negative. This dominance of
    negative weights will influence the networks output, especially
    for #outputs outside the training range. When the NN encounters
    data it hasnt been trained on, it relies #heavily on the these
    learned weights to make predictions. The effect of this causes the
    network's #output to converge to a negative value for such unseen
    data.
101
102 ##### Part 5 #####
103
104 set.seed(1234567890)
105 points = runif(n, min=0, max=10)
106 data <- data.frame(points, sin=sin(points))
107
108 nn = neuralnet(formula = points ~ ., data = data, hidden = 10,
    startweights = winit, threshold = 0.1)
109
110 plot(x = data[, 2], y = data[, 1], cex=2, main = "Predictions of x
    from sin(x)", xlab = "sin", ylab = "points")
111 points(x = data[, 2], y = predict(nn, data), col = "red", cex = 1)
112
113 legend("bottomleft",
114       legend = c("Training Data", "Predictions"),
115       col = c("black", "red"),
116       pch = 1,
117       pt.cex = c(2, 1))
118
119 # The result of the NN trying to predict x from sin(x) can be seen in
    Figure \ref{fig:a3_t5}. This happens because its easier to predict
    sin(x) from x because the sine function provides a predictable
    output for each input x. The other way around is much harder, for
    the NN to learn the inverse of sine function, which is periodic,
    meaning that multiple values of x produce the same sin(x).

```

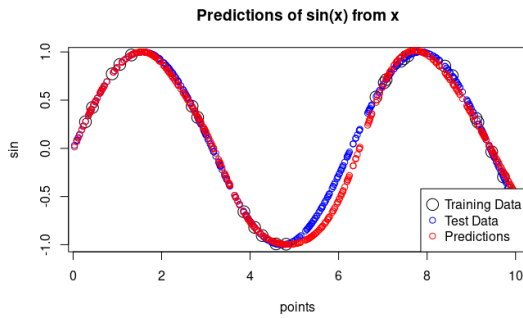


Figure 30: task 1

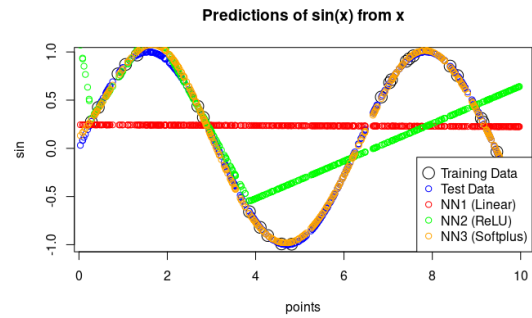


Figure 31: task 2

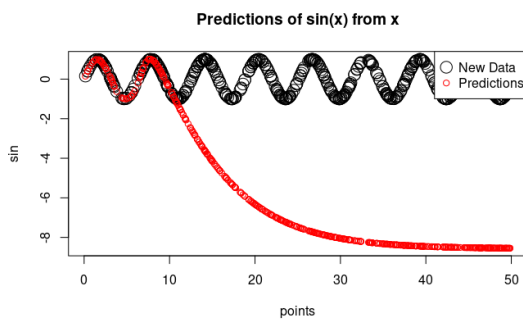


Figure 32: task 3

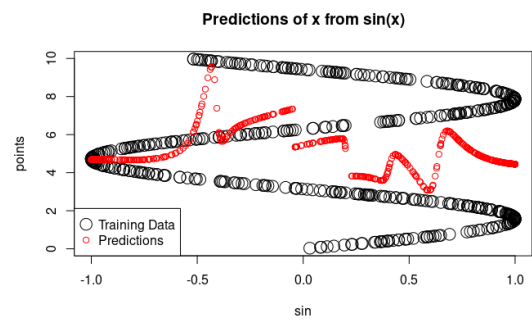


Figure 33: task 5

## 12 Exam 1

### 12.1 Assignment 1 (10p)

File `Bikes.csv` contains counts of public bicycles rented per hour in the Seoul Bike Sharing System, with corresponding weather data and holiday information.

1. Divide the data randomly into training and test data (70/30) and scale them appropriately. Compute a LASSO regression model by cross-validation in which `Rented Bikes` is the target variable and all remaining variables are features, and present a dependence of the cross-validation error on the penalty parameter. Which value of the penalty parameter is the optimal one? Interpret this model in terms of bias-variance tradeoff. Finally, report the equation showing how the predicted number of Rented bikes (scaled) depends on the features (scaled) in the estimated LASSO model corresponding to “`lambda.1se`” penalty value. (4p)
2. Consider `lambda.1se` LASSO model, use the training error MSE as an estimate of the target variance parameter and report the **95% prediction interval** for the first observation (first row) in the test data. Explain the model assumptions making these computations possible. (2p)
3. Consider the same partitioned and scaled data as in step 1. Assume now that Dew Point Temperature is related to the features `Humidity` and `Visibility` as a linear model without intercept, and assume that the loss and the error functions are given by the following

formula:

$$L(y, \hat{y}) = E(y, \hat{y}) = |y - \hat{y}|$$

Implement a code optimizing the cost function by the BFGS optimizer, and plot a dependence of the cost values and the test errors on the iteration number. Is early stopping needed? Report the optimal iteration number. Make 3 scatter plots of (Humidity, Visibility) where observations are colored by:

- original target values,
- predicted target values from the optimal model,
- target values from the model corresponding to 5 iterations, and compare the plots with respect to the complexity of the model and quality of prediction. **(4p)**

Hint: to make scatter plots, you may use this kind of code from package ggplot2:

```
1 df = data.frame(x = your_variable1, y = your_variable2, color = your_variable3)
2 ggplot(df, aes(x = x, y = y, color = color)) + geom_point()
```

## 12.2 Assignment 1 Answers

```
1 n=dim(data)[1]
2 set.seed(12345)
3 id=sample(1:n, floor(n*0.6))
4 train=data[id,]
5 test=data[-id,]
6
7 library(caret)
8
9 #Assignment 1 Part 1
10 scaler=preProcess(train)
11 trainS=predict(scaler,train)
12 testS=predict(scaler, test)
13
14 covariates=trainS[,-1]
15 response=trainS[,1]
16
17 library(glmnet)
18
19 set.seed(12345)
20 model=cv.glmnet(as.matrix(covariates),
21                 response, alpha=1,family="gaussian")
22 model$lambda.min
23
24 coef(model, s="lambda.1se")
25
26 ## 11 x 1 sparse Matrix of class "dgCMatrix"
27 ##              s1
28 ## (Intercept)  -8.742707e-17
29 ## Hour         2.264222e-01
30 ## Temperature  2.852890e-01
31 ## Humidity     -9.808210e-02
32 ## Wind.speed   .
```

```

33 ## Visibility .
34 ## Dew.point.temperature .
35 ## Solar.Radiation .
36 ## Rainfall -1.702395e-02
37 ## Snowfall -6.135284e-02
38 ## Holiday 1.086443e-01
39
40 #Assignment 1 Part 2
41 Ypred=predict(model, as.matrix(covariates), s="lambda.1se" )[,1]
42 MSE=mean((Ypred-response)^2)
43
44 testX=testS[1,-1]
45 Ytest=predict(model, as.matrix(testX), s="lambda.1se" )
46 CI=c(Ytest-1.96*sqrt(MSE), Ytest+1.96*sqrt(MSE))
47
48 print(CI)
49
50 ## [1] -2.031268 1.240571
51
52 #Assignment 1 part 3
53
54 costF<-function(Y, Yfit){
55   R=abs(Y-Yfit)
56   return(mean(R))
57 }
58
59 x=as.matrix(trainS[,c(4,6)])
60 xt=as.matrix(testS[,c(4,6)])
61
62 Fs=list()
63 k=0
64 TestE=list()
65 Theta=list()
66
67 myCost= function(theta){
68   f=costF(trainS$Dew.point.temperature, x%*%theta)
69   .GlobalEnv$k= .GlobalEnv$k+1
70   .GlobalEnv$Fs[[k]]=f
71   .GlobalEnv$TestE[[k]]=costF(testS$Dew.point.temperature, xt%*%theta)
72   .GlobalEnv$Theta[[k]]=theta
73   return(f)
74 }
75
76 res=optim(rep(0,2), fn=myCost, method="BFGS")
77
78
79 plot((as.numeric(Fs)), type="l", col="blue")
80 points((as.numeric(TestE)), type="l", col="red")
81
82 which.min(TestE)
83
84 ## [1] 49
85

```



```

86 Pred1=xt%%Theta[[which.min(TestE)]]
87 Pred2=xt%%Theta[[5]]
88
89
90 df=data.frame(x=xt[,1], y=xt[,2], color=testS$Dew.point.temperature)
91 ggplot(df,aes(x=x,y=y,color=color))+geom_point()
92
93 df=data.frame(x=xt[,1], y=xt[,2], color=Pred1)
94 ggplot(df,aes(x=x,y=y,color=Pred1))+geom_point()
95
96 df=data.frame(x=xt[,1], y=xt[,2], color=Pred2)
97 ggplot(df,aes(x=x,y=y,color=Pred2))+geom_point()
98
99 #Note If the student makes these plots with training data or full
    data, no point reduction should be done since assignment does not
    specify which data to use in plots
100
101 #The optimal prediction model captures the trend and the original
    scale well, while the model with 4 iterations captures the trend
    but not the scale, all predictive values are near zero implying
    that the model is too simple for this data (underfitted).

```

## 12.3 Exam 1: Assignment 2 (10p)

### 12.3.1 EXERCISE 1 – 5 POINTS

In January 2023, the students were asked to implement the backpropagation algorithm for training a neural network for regression as it appears in the course textbook and slides. The solution is available to you in the file TDD0211January2023.R. Now, you are asked to incorporate dropout into this solution. Recall that dropout is a regularization technique whose detailed description you can find in the course textbook and slides.

Run your implementation with a dropout rate  $r = 1$  equal to 0, 0.1, and 0.05, i.e.,  $r = 1, 0.99, 0.95$ . Comment on the results.

### 12.3.2 EXERCISE 2 – 5 POINTS

You are asked to implement the perceptron algorithm. This algorithm for binary classification is a predecessor of modern neural networks.

Consider a binary classification problem with class labels  $t \in \{-1, +1\}$ . Then, the class label assigned to a point  $\mathbf{z}$  is given by:

$$y(\mathbf{w}, \mathbf{z}) = \begin{cases} +1 & \text{if } \mathbf{w}^T \mathbf{z} \geq 0 \\ -1 & \text{if } \mathbf{w}^T \mathbf{z} < 0 \end{cases}$$

The values of  $\mathbf{w}$  are iteratively determined with the help of the learning data  $\{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$ . Specifically, the  $i$ -th weight is updated in each iteration as follows:

$$w_i^{(t+1)} = w_i^{(t)} + \alpha \sum_{n=1}^N \left( t_n - y(\mathbf{w}^{(t)}, \mathbf{z}_n) \right) z_{n,i}$$

where  $\alpha$  is the learning rate. Finally, you can assume that  $\mathbf{z}_n$  and  $\mathbf{w}$  are of dimension two, i.e.,  $\mathbf{z}_n = (x_{n,1}, x_{n,2})$  and  $\mathbf{w} = (w_1, w_2)$ .

You can stop the learning process after 100 iterations. You can use an alpha value of 0.0001. Plot the non-misclassification rate on the dataset below as a function of the number of iterations.

```
set.seed(1234)

x <- array(NA, dim = c(100, 2))

t <- array(NA, dim = c(100))

x[,1] <- runif(100,0,3)

x[,2] <-runif(100,0,9)

t <- ifelse(x[,2]<(x[,1]^2),-1,1)

plot(x[,1],x[,2],col=t+2)
```

Finally, explain when the perceptron algorithm works best and why it works in those cases.

## 12.4 Assignment 2 Answers

```
1 # exercise 1 - 5 p
2 # 3 p for implementation
3 # 1 p if it runs correctly
4 # 1 p for answering the question correctly
5 # lines corresponding to dropout implementation are marked with ###,
6   the rest of the code was provided in the exam
7
8
9 set.seed(1234)
10
11 # produce the training data in dat
12
13 x <- runif(500,-4,4)
14 y <- sin(x)
15 dat <- cbind(x,y)
16 plot(dat)
17
18 gamma <- 0.01
19 r <- 1 ### dropout rate 1-r
20
21 h <- function(z){
22   # activation function (sigmoid)
```

```

22 |
23 |   return(1/(1+exp(-z)))
24 | }
25 |
26 | hprime <- function(z){
27 |
28 |   # derivative of the activation function (sigmoid)
29 |
30 |   return(h(z) * (1 - h(z)))
31 | }
32 |
33 | yhat <- function(x){
34 |
35 |   # prediction for point x
36 |
37 |   q0 <- x
38 |   z1 <- w1 %*% q0 + b1
39 |   q1 <- as.matrix(apply(z1,1,h), nrow = 2, ncol = 1)
40 |   z2 <- w2 %*% q1 + b2
41 |   return(z2)
42 | }
43 |
44 | yhat2 <- function(x){
45 |
46 |   ### final prediction for point x (it involves r)
47 |
48 |   q0 <- x
49 |   z1 <- (r*w1) %*% q0 + b1
50 |   q1 <- as.matrix(apply(z1,1,h), nrow = 2, ncol = 1)
51 |   z2 <- (r*w2) %*% q1 + b2
52 |   return(z2)
53 | }
54 |
55 | MSE <- function(){
56 |
57 |   # mean squared error
58 |
59 |   res <- NULL
60 |   for(i in 1:nrow(dat)){
61 |     res <- c(res,(dat[i,2] - yhat(dat[i,1])) ^ 2)
62 |   }
63 |   return(mean(res))
64 | }
65 |
66 | # initialize parameters
67 |
68 | w1 <- matrix(runif(2,-.1,.1), nrow = 2, ncol = 1)
69 | b1 <- matrix(runif(2,-.1,.1), nrow = 2, ncol = 1)
70 |
71 | w2 <- matrix(runif(2,-.1,.1), nrow = 1, ncol = 2)
72 | b2 <- matrix(runif(1,-.1,.1), nrow = 1, ncol = 1)
73 |
74 | res <- NULL

```

```

75 for(i in 1:100000){
76   if(i %% 1000 == 0){
77     res <- c(res,MSE())
78   }
79
80   # forward propagation
81
82   j <- sample(1:nrow(dat),1)
83   q0 <- dat[j,1]
84   m0 <- sample(c(0,1),1,replace = TRUE, c(1-r,r)) ### drop input unit
85   z1 <- w1 %*% (q0*m0) + b1
86   q1 <- as.matrix(apply(z1,1,h), nrow = 2, ncol = 1)
87   m1 <- sample(c(0,1),2,replace = TRUE, c(1-r,r)) ### drop hidden
      units
88   z2 <- w2 %*% (q1*m1) + b2
89
90   # backward propagation
91
92   dz2 <- - 2 * (dat[j,2] - z2)
93   dq1 <- t(w2) %*% dz2
94   dz1 <- dq1 * hprime(z1)
95   dw2 <- dz2 %*% t(q1)
96   db2 <- dz2
97   dw1 <- dz1 %*% t(q0)
98   db1 <- dz1
99
100  # parameter updating
101
102  w2 <- w2 - (gamma * dw2 * m1) ### update non-dropped hidden units
103  b2 <- b2 - gamma * db2
104  w1 <- w1 - (gamma * dw1 * c(m0,m0) * m1) ### update non-dropped
      hidden units
105  b1 <- b1 - gamma * db1
106 }
107 plot(res, type = "l")
108
109 plot(dat)
110 points(dat[,1],lapply(dat[,1],yhat2),col="red") ### apply final
      prediction
111
112 # Lower r value implies larger regularization and, thus, worse fit of
      the training data. The ultimate goal of dropout is
113 # to get a better fit of the test data (i.e., low generalization
      error). However, this is too simple an example to observe it.
114
115 # exercise 2 - 5 p
116 # 3 p for implementation
117 # 1 p if it runs correctly
118 # 1 p for answering the question correctly
119
120 set.seed(1234)
121
122 x <- array(NA, dim = c(100,2))

```

```

123 t <- array(NA, dim = c(100))
124 x[,1] <- runif(100,0,3)
125 x[,2] <- runif(100,0,9)
126 t <- ifelse(x[,2] < (x[,1])^2, -1, 1)
127 plot(x[,1], x[,2], col=t+2)
128
129 w <- c(0,0)
130 alpha <- 0.0001
131
132 res <- NULL
133 for(i in 1:100){
134   res <- c(res, sum(t != (2 * ((w %*% t(x[,1:2])) > 0) - 1))) # note
135     that t is both class label and transpose operation :-(
136
137   w[1] <- w[1] + alpha * sum((t - w %*% t(x[,1:2])) * x[,1])
138   w[2] <- w[2] + alpha * sum((t - w %*% t(x[,1:2])) * x[,2])
139 }
140 plot(x[,1], x[,2], col=(2 * ((w %*% t(x[,1:2])) > 0) - 1)+2)
141 plot(res, type = "l")
142
143 # It works for linearly separable datasets. If the true label is
144   larger than the prediction, then the weights get
145   increased for positive input values and decreased for negative
146   input values. The opposite when the true label is smaller.

```

## 13 Exam 2

### 13.1 Assignment 1 (10p)

File `Rice.csv` contains a total of 3810 rice grain's images taken for the two species (Cammeo and Osmancik), which were processed and feature inferences were made. 7 morphological features were obtained for each grain of rice which are the variables in the data set.

1. Divide the data randomly into training and test data (70/30). Assume that columns `Area`, ..., `Extent` are denoted as  $x_1, \dots, x_p$ . Perform basis function expansion with basis functions  $\phi_1 = x_1, \dots, \phi_p = x_p, \phi_{p+1} = x_1^2, \dots, \phi_{2p} = x_p^2$  and fit two logistic regression models: one with features  $x_1, \dots, x_p$  and another with  $\phi_1, \dots, \phi_{2p}$ , both with target `Class` to the training data. Report the training and test misclassification errors for both models and report which model is better and why. Finally, report the estimated probabilistic model for the model with features  $x_1, \dots, x_p$ . **(4p)**
2. Use cross-validation to fit a decision tree model to the training data with target `Class` and all other columns ( $x_1, \dots, x_p$ ) as features and report the optimal number of leaves and training and test misclassification errors for the optimal tree. Comment on the prediction quality of the optimal tree. After this, write a loop that grows decision trees with parameter `mindev` = 0.001, 0.002, ... 0.01 (without cross-validation) and estimates training and test misclassification errors for these trees. Plot the dependence of these errors on `mindev`, comment on the trends observed in the plot and comment how the trends are expected to look in theory and why. **(4p)**

3. Use the model estimated by the cross-validation in step 2 and assume that a user wants to prune this tree. If we use misclassification error as impurity measure, which leaves must be pruned first, 14 and 15, or 4 and 5? Report necessary mathematical calculations to support your answer. (2p)

### 13.2 Assignment 1: Answers

```
1 #Part 1
2
3 df=read.csv("Rice.csv", stringsAsFactors = T)
4 library(dplyr)
5
6 n=dim(df)[1]
7 set.seed(12345)
8 id=sample(1:n, floor(n*0.7))
9 train=df[id,]
10 test=df[-id,]
11
12
13 train1=train%>%mutate_if(is.numeric, list(x2=function(x) x^2))
14 test1=test%>%mutate_if(is.numeric, list(x2=function(x) x^2))
15
16 m1=glm(Class~., data=train, family = binomial)
17
18 m2=glm(Class~., data=train1, family=binomial)
19
20 missclass=function(X,X1){
21   n=length(X)
22   return(1-sum(diag(table(X,X1)))/n)
23 }
24
25
26
27 missclass(train$Class, predict(m1, newdata = train,
28   type="response")>0.5)
29 # [1] 0.06861642
30
31 missclass(test$Class, predict(m1, newdata = test,
32   type="response")>0.5)
33 # [1] 0.07524059
34
35 missclass(train$Class, predict(m2, newdata = train1,
36   type="response")>0.5)
37 # [1] 0.06786652
38
39 missclass(test$Class, predict(m2, newdata = test1,
40   type="response")>0.5)
41 # [1] 0.07611549
```

```

42
43 #We can see that after basis function expansion the training error
    gets smaller and test error gets larger which indicates
    overfitting. The model with original features is best one
44
45 m1
46
47 ##
48 ## Call:  glm(formula = Class ~ ., family = binomial, data = train)
49 ##
50 ## Coefficients:
51 ##          (Intercept)                Area                Perimeter
    Major_Axis_Length
52 ##          -7.910612                0.006971                0.107500
    -0.112395
53 ## Minor_Axis_Length                Eccentricity                Convex_Area
    Extent
54 ##          0.553997                -3.747786                -0.011754
    0.315431
55 ##
56 ## Degrees of Freedom: 2666 Total (i.e. Null);  2659 Residual
57 ## Null Deviance:          3656
58 ## Residual Deviance: 920.9      AIC: 936.9
59
60 levels(train$Class)
61
62 ## [1] "Cammeo"    "Osmancik"
63
64 #

```

The probabilistic model is

$$P(\text{Class} = \text{Osmancik}) = \frac{1}{1 + \exp(z)}$$

where

$$z = 7.910612 - 0.006971 \cdot \text{Area} - 0.107500 \cdot \text{Perimeter} + 0.112395 \cdot \text{MajorAxisLength} - 0.553997 \cdot \text{MinorAxisLength} + 3.7$$

```

1 #Part 2
2 library(tree)
3 fit=tree(Class~., data=train)
4
5 set.seed(12345)
6 cv.res=cv.tree(fit)
7 finalTree0=prune.tree(fit, best=cv.res$size[which.min(cv.res$dev)])
8 Yfit1=predict(finalTree0, newdata=train,
9               type="class")
10 Yfit2=predict(finalTree0, newdata=test,
11               type="class")
12
13 cv.res$size[which.min(cv.res$dev)]
14

```

```

15 # [1] 5
16
17 missclass(train$Class,Yfit1)
18
19 # [1] 0.06861642
20
21 missclass(test$Class,Yfit2)
22
23 # [1] 0.07786527
24
25 #The tree has a good prediction quality since test error is only
    around 7 percent.
26
27 mind=seq(0.001, 0.01, 0.001)
28 lM=length(mind)
29 TrE=numeric(lM)
30 TeE=numeric(lM)
31
32 for (i in 1:lM){
33   fit=tree(Class~., data=train,
34             control=tree.control(nobs=nrow(train),mindev=mind[i]) )
35   TrE[i]=missclass(predict(fit, newdata=train, type="class"),
36                     train$Class)
37   TeE[i]=missclass(predict(fit, newdata=test, type="class"),
38                     test$Class)
39 }
40
41 plot(mind, TrE, col="blue", ylim=c(0.05,0.15))
42 points(mind, TeE, col="red")

```

```

1 #Part 3
2
3 print(finalTree0)
4
5 ## node), split, n, deviance, yval, (yprob)
6 ##      * denotes terminal node
7 ##
8 ## 1) root 2667 3656.0 Osmancik ( 0.437945 0.562055 )
9 ##    2) Major_Axis_Length < 191.075 1488 657.3 Osmancik ( 0.057796
    0.942204 )
10 ##      4) Major_Axis_Length < 181.225 1116 141.6 Osmancik (
    0.011649 0.988351 ) *
11 ##      5) Major_Axis_Length > 181.225 372 368.4 Osmancik ( 0.196237
    0.803763 ) *
12 ##      3) Major_Axis_Length > 191.075 1179 670.3 Cammeo ( 0.917727
    0.082273 )
13 ##      6) Perimeter < 467.639 198 258.4 Cammeo ( 0.641414 0.358586
    ) *
14 ##      7) Perimeter > 467.639 981 240.1 Cammeo ( 0.973496 0.026504 )
15 ##      14) Major_Axis_Length < 205.866 409 188.2 Cammeo ( 0.938875
    0.061125 ) *
16 ##      15) Major_Axis_Length > 205.866 572 14.7 Cammeo ( 0.998252

```



```

17         0.001748 ) *
18 ImpurityDiff0=0.057796*1488-(0.011649*1116 +0.196237*372)
19 ImpurityDiff0
20
21 ## [1] 1.421085e-14
22
23 ImpurityDiff1=0.026504*981 -(0.061125*409+0.001748*572 )
24 ImpurityDiff1
25
26 ## [1] 0.000443
27
28 #4 and 5 need to be pruned first since provide smallest impurity
    decrease

```

### 13.3 Assignment 2 (10p)

You are asked to implement the backpropagation algorithm for training a neural network for regression as it appears in the course textbook and slides. To do so, fill in the code provided below. The actual network has one hidden layer with two units.  $x$  denotes inputs,  $y$  denotes targets,  $\hat{y}$  denotes activation units, and  $z$  denotes hidden units. As the result of applying the activation function to  $z$ ,  $\sigma(z)$  denotes the squashed error, and  $J$  denotes the loss (mean squared error). The steps involved in this task are: (1) forward propagation, (2) backward propagation, (3) parameter update as marked in substeps 1-3B. Except the one indicated with a \*, that is element-wise product. (3B). Notice the use of matrix transposition in some steps (3B) in R. Comment your code and add results. The exercise will be graded as follows: Forward propagation 2p, backward propagation 4p, parameter updating 1p, and results 3p.

#### Steps:

- **Forward propagation.**

$$\begin{aligned}
 q^{(0)} &= x \\
 z^{(1)} &= W^{(1)}q^{(0)} + b^{(1)} \\
 q^{(1)} &= h(z^{(1)}) \\
 z^{(2)} &= W^{(2)}q^{(1)} + b^{(2)} \\
 J(\theta) &= (y - z^{(2)})^2
 \end{aligned}$$

- **Backward propagation.**

$$\begin{aligned}
 dz^{(2)} &= -2(y - z^{(2)}) \\
 dq^{(1)} &= W^{(2)T} dz^{(2)} \\
 dz^{(1)} &= dq^{(1)} \odot h'(z^{(1)}) \\
 dW^{(2)} &= dz^{(2)} q^{(1)T} \\
 db^{(2)} &= dz^{(2)} \\
 dW^{(1)} &= dz^{(1)} q^{(0)T} \\
 db^{(1)} &= dz^{(1)}
 \end{aligned}$$

- **Parameter updating.**

$$W_{t+1}^{(2)} = W_t^{(2)} - \gamma dW^{(2)}$$

$$b_{t+1}^{(2)} = b_t^{(2)} - \gamma db^{(2)}$$

$$W_{t+1}^{(1)} = W_t^{(1)} - \gamma dW^{(1)}$$

$$b_{t+1}^{(1)} = b_t^{(1)} - \gamma db^{(1)}$$

### Code Template:

You are requested to use the template below. Note that the algorithm performs 100000 iterations. In each iteration, one randomly selected training point is used to update the parameters (i.e., this essentially corresponds to stochastic gradient descent with a mini-batch of size 1).

```

1 # produce the training data in R
2 x = runif(500, -4, 4)
3 y = sin(x)
4 dat = cbind(x, y)
5 plot(dat)
6
7 gamma = 0.01
8
9 h = function(z){
10   # activation function (sigmoid)
11   return(1/(1+exp(-z)))
12 }
13
14 hprime = function(z){
15   # derivative of the activation function (sigmoid)
16   return(h(z) * (1 - h(z)))
17 }
18
19 yhat <- function(x){
20   # prediction for point x
21 }
22
23 MSE <- function(res){
24   # mean squared error
25 }
26
27 # initialize parameters
28 res <- NULL
29 for(it in 1:100000){
30   if(it %% 10000 == 0){
31     res <- c(res, MSE())
32   }
33
34   # forward propagation
35   j <- sample(1:nrow(dat), 1)
36   x0 <- dat[j, 1]
37
38   # backward propagation

```

```

39 # parameter updating
40 }
41
42 plot(res, type = "l")
43 plot(dat)
44 points(dat[,1], apply(dat[,1], 1, yhat))

```

## 13.4 Assignment 2: Answers

```

1  set.seed(1234)
2
3  # produce the training data in dat
4
5  x <- runif(500, -4, 4)
6  y <- sin(x)
7  dat <- cbind(x, y)
8  plot(dat)
9
10 gamma <- 0.01
11
12 h <- function(z){
13
14   # activation function (sigmoid)
15
16   return(1/(1+exp(-z)))
17 }
18
19 hprime <- function(z){
20
21   # derivative of the activation function (sigmoid)
22
23   return(h(z) * (1 - h(z)))
24 }
25
26 yhat <- function(x){
27
28   # prediction for point x
29
30   q0 <- x
31   z1 <- w1 %*% q0 + b1
32   q1 <- as.matrix(apply(z1, 1, h), nrow = 2, ncol = 1)
33   z2 <- w2 %*% q1 + b2
34   return(z2)
35 }
36
37 MSE <- function(){
38
39   # mean squared error
40
41   res <- NULL
42   for(i in 1:nrow(dat)){
43     res <- c(res, (dat[i,2] - yhat(dat[i,1])) ^ 2)

```

```

44     }
45     return(mean(res))
46 }
47
48 # initialize parameters
49
50 w1 <- matrix(runif(2,-.1,.1), nrow = 2, ncol = 1)
51 b1 <- matrix(runif(2,-.1,.1), nrow = 2, ncol = 1)
52
53 w2 <- matrix(runif(2,-.1,.1), nrow = 1, ncol = 2)
54 b2 <- matrix(runif(1,-.1,.1), nrow = 1, ncol = 1)
55
56 res <- NULL
57 for(i in 1:100000){
58   if(i %% 1000 == 0){
59     res <- c(res,MSE())
60   }
61
62   # forward propagation
63
64   j <- sample(1:nrow(dat),1)
65   q0 <- dat[j,1]
66   z1 <- w1 %*% q0 + b1
67   q1 <- as.matrix(apply(z1,1,h), nrow = 2, ncol = 1)
68   z2 <- w2 %*% q1 + b2
69
70   # backward propagation
71
72   dz2 <- - 2 * (dat[j,2] - z2)
73   dq1 <- t(w2) %*% dz2
74   dz1 <- dq1 * hprime(z1)
75   dw2 <- dz2 %*% t(q1)
76   db2 <- dz2
77   dw1 <- dz1 %*% t(q0)
78   db1 <- dz1
79
80   # parameter updating
81
82   w2 <- w2 - gamma * dw2
83   b2 <- b2 - gamma * db2
84   w1 <- w1 - gamma * dw1
85   b1 <- b1 - gamma * db1
86 }
87 plot(res, type = "l")
88
89 plot(dat)
90 points(dat[,1],lapply(dat[,1],yhat),col="red")

```

## 14 Theory

### 14.1 Linear Model

The general formula for a linear model is:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon$$

In matrix form, the same model can be expressed as:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

Where:

- $y$ : Dependent variable (response).  $\mathbf{X}$ : Design matrix of predictors, of size  $n \times (p + 1)$ , where the first column is typically ones (for the intercept).
- $\beta_0$ : Intercept.
- $\beta_1, \beta_2, \dots, \beta_p$ : Coefficients for the independent variables  $x_1, x_2, \dots, x_p$ .  $\boldsymbol{\beta}$ : Coefficient vector, including the intercept, of size  $(p + 1) \times 1$ .
- $\epsilon$ : Error term, capturing the variation not explained by the model.  $\boldsymbol{\epsilon}$ : Vector of error terms, of size  $n \times 1$ .
- $\mathbf{y}$ : Vector of responses, of size  $n \times 1$ .

### 14.2 Logistic Regression

The probability of the positive class is given by:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p)}}$$

In terms of the logit function:

$$\text{logit}(p) = \ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$$

In matrix form:

$$p = \frac{1}{1 + e^{-\mathbf{X}\boldsymbol{\beta}}}$$

### 14.3 Bias Variance tradeoff

#### 1. Bias:

- Bias refers to the error introduced by approximating a real-world problem, which may be highly complex, by a simplified model.
- High bias implies that the model is too simplistic and may not capture the underlying patterns in the data.
- This can lead to systematic errors, and the model may consistently underperform.

## 2. Variance:

- Variance refers to the model's sensitivity to small fluctuations or noise in the training data.
- High variance implies that the model is too complex and has learned to capture the noise in the training data rather than the underlying patterns.
- This can lead to overfitting, where the model performs well on the training data but poorly on new, unseen data.

The tradeoff arises because increasing the complexity of a model typically reduces bias but increases variance, and vice versa. The goal is to find the right balance that minimizes the overall prediction error on new, unseen data. Bayes variance: the more leaves, the higher variance but smaller bias. If our model is too simple and has very few parameters then it may have high bias and low variance. On the other hand if our model has large number of parameters then it's going to have high variance and low bias. So we need to find the right/good balance without overfitting and underfitting the data.

Bias is the difference between the average prediction of our model and the correct value which we are trying to predict. Model with high bias pays very little attention to the training data and oversimplifies the model. It always leads to high error on training and test data.

Variance is the variability of model prediction for a given data point or a value which tells us spread of our data. Model with high variance pays a lot of attention to training data and does not generalize on the data which it hasn't seen before. As a result, such models perform very well on training data but has high error rates on test data.

## 14.4 Precision and recall

The formulas for Precision and Recall are as follows:

Precision measures the proportion of positive predictions that were actually correct. It focuses on the accuracy of the positive predictions made by the model.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

Recall measures the proportion of actual positive instances that were correctly identified by the model. It focuses on the completeness of the positive predictions.

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

## 14.5 F1 score

Good measure of performance if the dataset is imbalanced. Measures how good the model is at guessing "positive". The formula for the F1 score is given by:

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$