Lab Report in Machine Learning

# Laboration 3

## TDDE01

David Grönberg
Davgr686

LINKÖPINGS UNIVERSITET

16-12-2019

# Contents

# List of Figures

# 1.  Assignments

## 1.1   Assignment 1

The task of this assignment is to predict hourly temperatures for a given date and place (defined by latitude and longitude coordinates) in Sweden. SMHI has provided us with temperature measurements for previous dates and places. This is done by the use of a kernel that is the sum of three Gaussian kernels:

- One that accounts for the euclidean distance between two stations. (denoted *distance_kernel*)

- One that accounts for the difference in days between two temperature measurements. (denoted *date_kernel*)

- One that accounts for the difference in time between two temperature measurements. (denoted *time_kernel*)

The chosen smoothing coefficients for each Gaussian kernel are presented below:

- *distance_kernel*: $h = 50,000$

- *date_kernel*: $h = 7$

- *time_kernel*: $h = 6$

To confirm we have chosen appropriate smoothing coefficients, the kernels width are plotted. Here, we use *2011-01-05* and *12:00* as the target date and time with coordinates representing a place in Stockholm.

Figure 1.1, 1.2, 1.3 shows a plot representing how sensible each kernel is. In figure 1.2 - where the x-axis is the day of the year and y-axis is the weight - the kernel gives points that are close to the target date higher values (i.e. days of the year that are close to 5). This smoothing coefficient seems reasonable to have because the weather can drastically change in a month and we would like to capture that change and take it into account in our model.

Figure 1.3 shows the width of the *time_kernel*. Here, the x-axis represents the hour of day, ranging from 0 to 23 and our target time is *12*. We can see that the kernel assigns higher values to times that are closer to our target time. As previously stated, we chose a smoothing coefficient of *6* here. This is because we want to catch the smaller change in temperature between day time and night time, but still believe all hours has some kind of relevance.
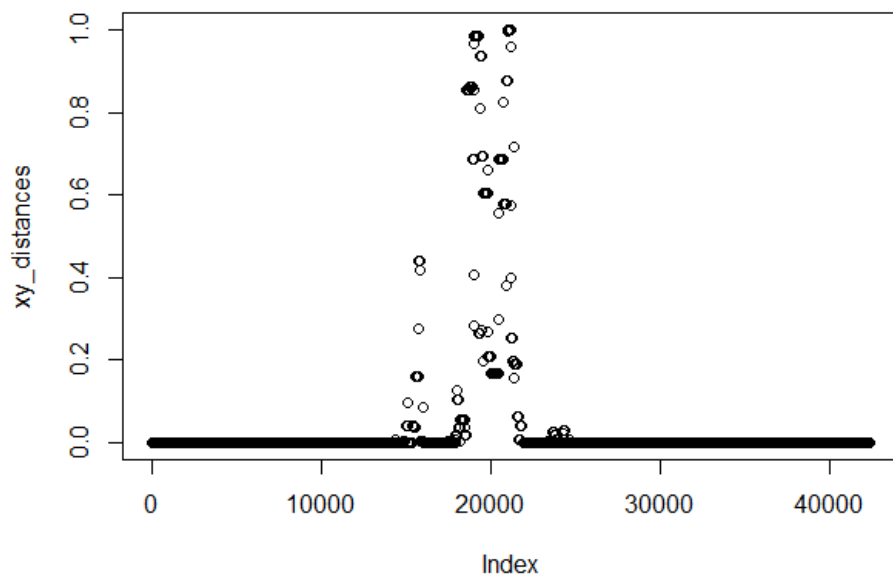
Figure 1.1: Width of *distance_kernel*
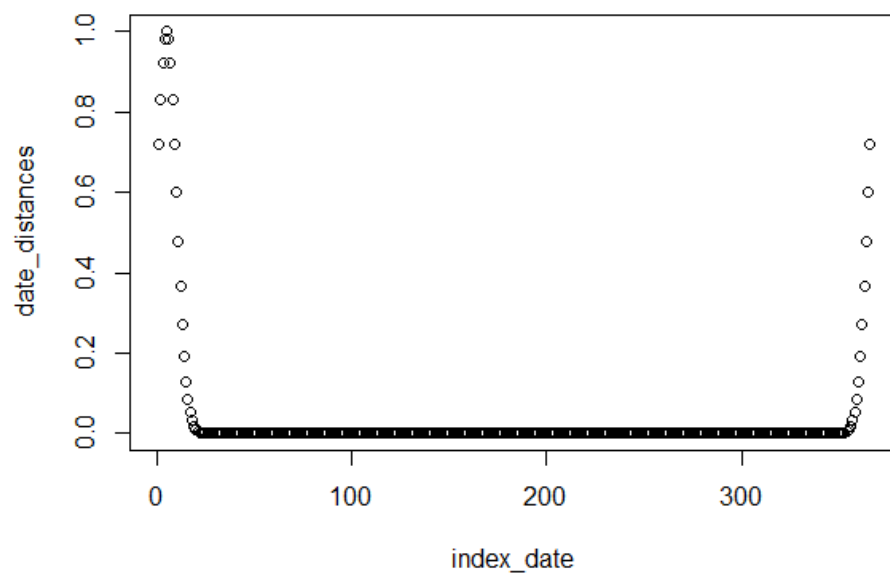
Figure 1.2: Width of *date_ kernel*
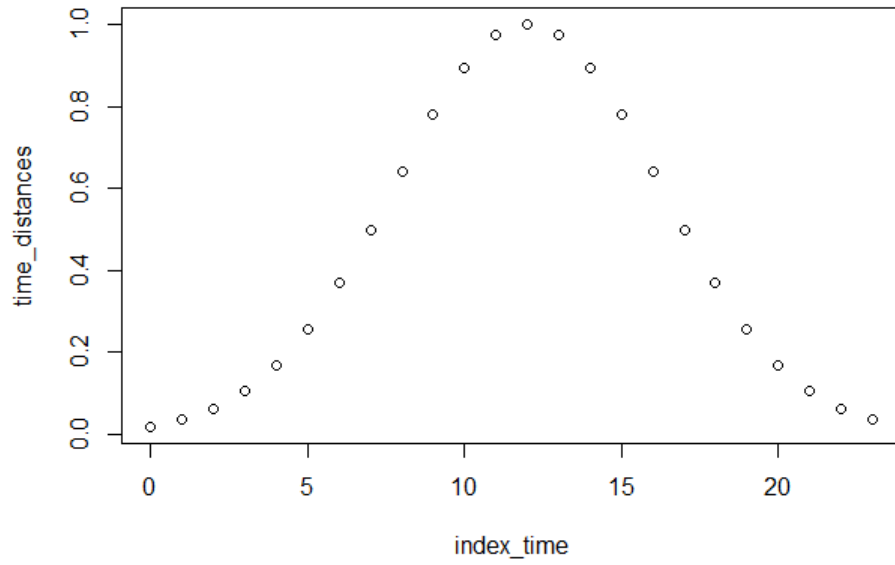
Figure 1.3: Width of *time_ kernel*

Figure 1.4 shows the predicted temperature for different times on the target date by summing the kernels. The model predicts temperatures between 3.5 and 5.0. These predictions can be plausible temperatures in the beginning of January, even though they seem rather high. We can see that the model predicts lower temperatures at night, compared to day time.
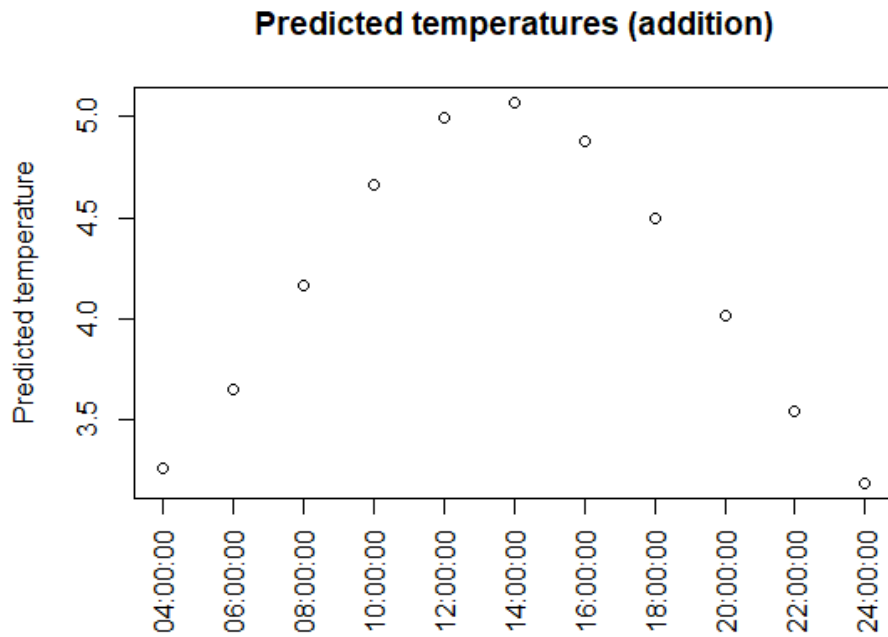
## Predicted temperatures (addition)



Figure 1.4: Temperature predictions using addition method

Next, instead of summing the kernels, we now multiply them and plot the results. Figure 1.5 shows the predicted temperatures for each time on the target date. Now we observe values ranging between -3 and -2, which seem more plausible for the target date. However, when inspecting the change in predicted temperature depending on time of day - the results seem rather strange. The model predicts that the lowest temperature will be observed around 08.00 and the highest temperature observed around 20.00 in the evening.

**Predicted temperatures (multiplication)**



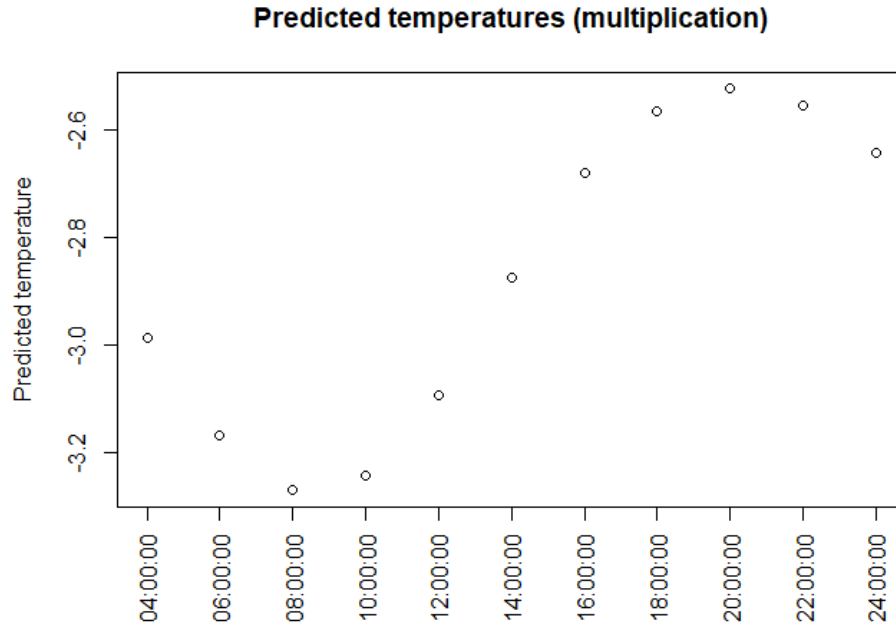Figure 1.5: Temperature predictions using multiplication method

To verify that the model takes the time of year into account, we make the model predict temperatures for every month of the year. Figure 1.6 shows that the model using addition predicts temperatures ranging only between 5-6, while figure 1.7 shows that the model using multiplication predicts temperatures with a larger range, yielding more plausible predictions.
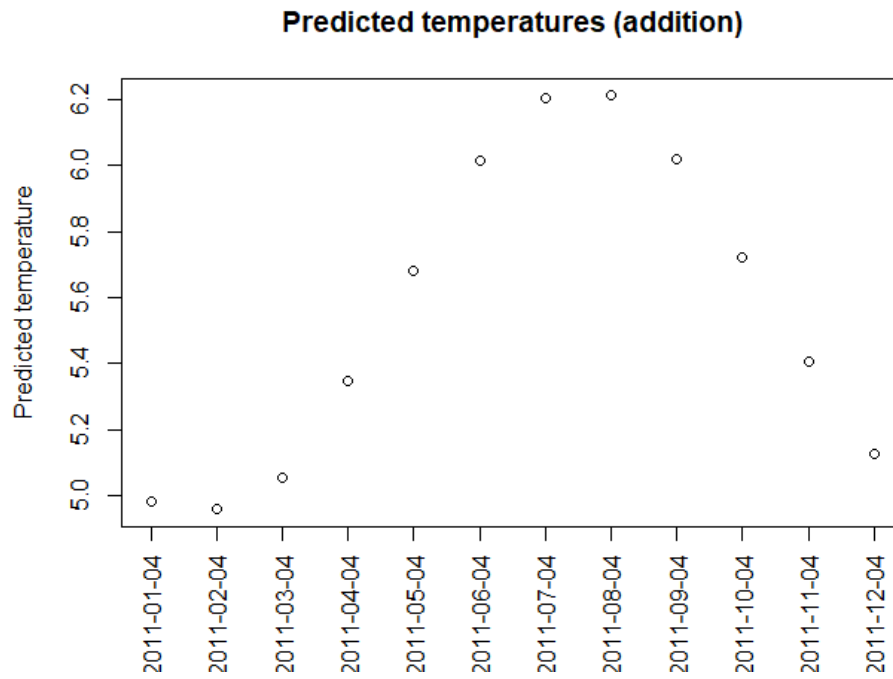
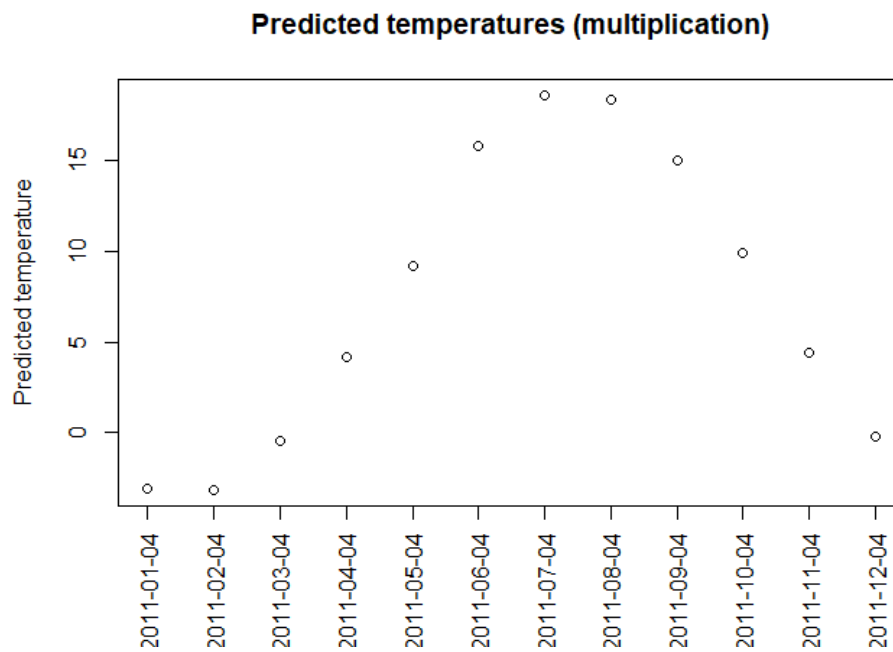Figure 1.6: Temperature predictions on dates using addition method



Figure 1.7: Temperature predictions on dates using multiplication method

7

## 1.2  Assignment 3

The task of this assignment is to train a neural network to learn the trigonometric sine function. First, a dataset containing mappings of variables and their corresponding sine output is divided into a train and validation set. Then, the *neuralnet* function is used to train the neural network, with a single hidden layer with 10 units.

To find the optimal threshold, we consider threshold values: $i/1000$ where $i = 1,...,10$ and then train a neural network with each threshold value using the train set and compute the Mean Square Error (MSE) on the validation set.

Figure 1.8 shows the Mean Squared Error on the validation set for each threshold value considered. The plot shows the lowest observed MSE when the threshold is 0.001, 0.003 and 0.008. After further investigation we find the lowest MSE when the threshold is 0.001.



Figure 1.8: MSE for different threshold values

Next, we fit a neural network with the optimal threshold and report its structure.

Figure 1.9: Optimal neural network

Lastly, the optimal neural network's predictions are plotted with the generated dataset in figure 1.10, where the predictions are coloured blue and the data is coloured red. As the plot shows, the majority of the points overlap, indicating that the neural network has high accuracy when predicting the sine function.

Figure 1.10: Predicted values and data

# 2. Code Appendix

## 2.1 Assignment 1

```r
1  library(geosphere)
2  library(lubridate)
3  require(stats)
4
5  set.seed(1234567890)
6  stations <- read.csv("stations.csv", fileEncoding = "latin1")
7  temps <- read.csv("temps50k.csv")
8  st <- merge(stations,temps,by="station_number")
9
10
11 h_distance <- 50000
12 h_date <- 7
13 h_time <- 6
14
15 target_point = c(18.0574, 59.3420) # sthlm
16 pred_date <- "2011-01-05" # The date to predict
17 times <- c("04:00:00", "06:00:00", "08:00:00", "10:00:00", "12:00:00", "14:00:00", "16:00:00",
       "18:00:00", "20:00:00", "22:00:00", "24:00:00")
18 temp <- vector(length=length(times))
19 # Students  code here
20 df_before <- subset(st, as.Date(date) < as.Date(pred_date))
21
22
23 get_date_diff <- function(day1, day2, h_date) {
24   diff <- abs(yday(day1) - yday(day2))
25   if (diff > 182) {
26     return (abs(diff - 365)/h_date)
27   }
28   return (diff/h_date)
29 }
30
31 get_time_diff <- function( time2, time1, h_time) {
32   time1.numeric <- as.numeric(substring(time1,1,2))
33   diff <- abs(time1.numeric - time2)
34   if (diff > 12) {
35     return (abs(diff - 24)/h_time)
36   }
37   return (diff/h_time)
38 }
39
40 get_xy_diff <- function(target, x, h_distance) {
41   x_lat <- as.numeric(x[4])
42   x_long <- as.numeric(x[5])
43   dist <- distHaversine(c(x_long, x_lat), target)
44   return (dist[1]/h_distance)
45 }
```
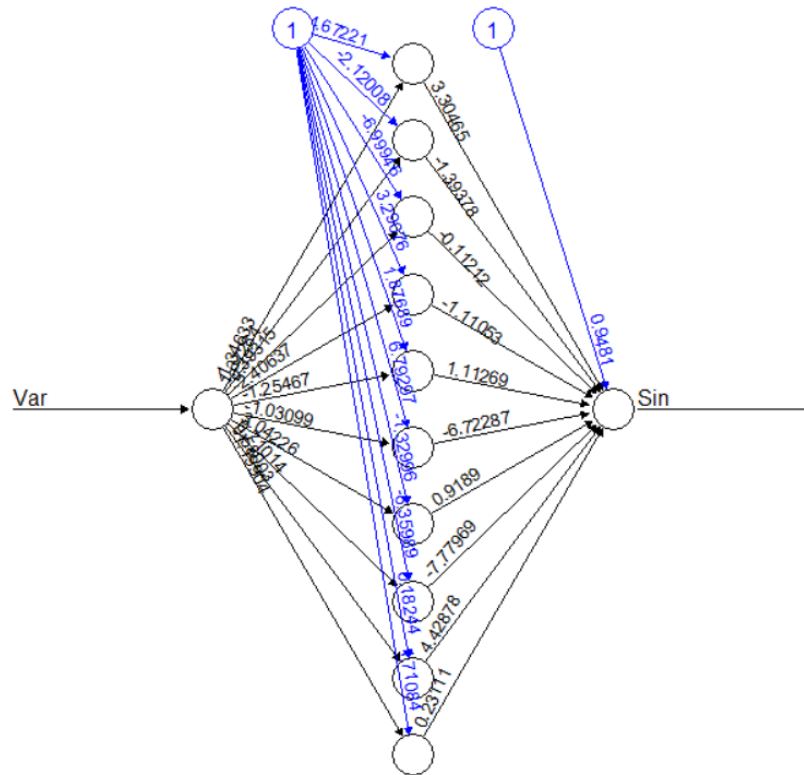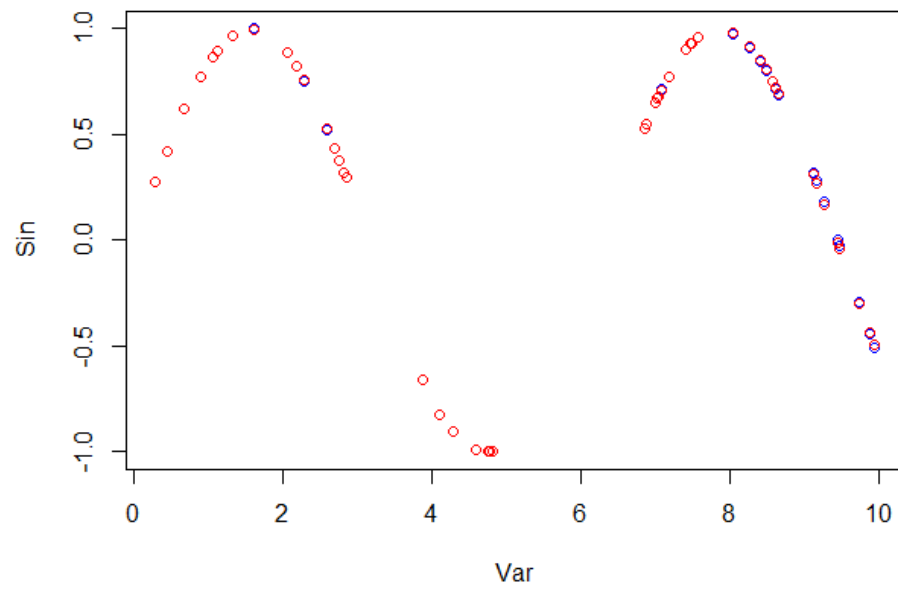
```r
46
47  euc_norm <- function(x) return (sqrt(sum(x^2)))
48
49
50  date_kernel <- function(X, target, h_date) {
51    u <- abs(sapply(X$date, function(x) get_date_diff(target, x, h_date)))
52    k <- exp(-u^2)
53    return (k)
54  }
55
56
57  time_kernel <- function(X, target, h_time) {
58    u <- abs(sapply(X$time, function(x) get_time_diff(target, x, h_time)))
59    k <- exp(-u^2)
60    return (k)
61  }
62
63  distance_kernel <- function(X, target, h_distance) {
64    u <- abs(apply(X, 1, function(x) get_xy_diff(target, x, h_distance)))
65    k <- exp(-u^2)
66    return (k)
67  }
68
69  xy_distances <- distance_kernel(df_before, target_point, h_distance)
70  plot(xy_distances)
71
72  index_date <- c(yday(df_before$date))
73  date_distances <- date_kernel(df_before, pred_date, h_date)
74  plot(y = date_distances, x = index_date)
75
76  index_time <- as.numeric(substring(df_before$time,1,2))
77  time_distances <- time_kernel(df_before, 12, h_time)
78  plot(y = time_distances, x = index_time)
79
80  denom = c()
81  nomi = c()
82
83
84  pred_dates <- c("2011-01-04", "2011-02-04", "2011-03-04", "2011-04-04", "2011-05-04", "
        2011-06-04", "2011-07-04", "2011-08-04", "2011-09-04", "2011-10-04","2011-11-04", "
        2011-12-04")
85
86  pred = c()
87  for (x in 1:length(pred_dates))
88  {
89    temp_date_distances <- date_kernel(df_before, pred_dates[x], h_date)
90    for (i in 1:length(df_before$date)) {
91      denom[i] <- (temp_date_distances[i] + time_distances[i] + xy_distances[[i]])
92      nomi[i] <- (denom[i] * df_before$air_temperature[i])
93    }
94    pred[x] <- sum(nomi) / sum(denom)
95
96  }
97
98
99  plot(pred, axes=FALSE, xlab="", ylab="Predicted temperature", main="Predicted temperatures (
        addition)")
100 axis(2)
101 axis(1, at=seq_along(pred),labels=as.character(pred_dates), las=2)
102 box()'
103
104
```

```
105
106  pred = c()
107  for (x in 1:length(times))
108  {
109    temp_time <- as.numeric(substring(times[x], 1, 2))
110    temp_time_distances <- time_kernel(df_before, temp_time, h_time)
111    for (i in 1:length(df_before$date)) {
112      denom[i] <- (date_distances[i] * temp_time_distances[i] * xy_distances[[i]])
113      nomi[i] <- (denom[i] * df_before$air_temperature[i])
114    }
115    pred[x] <- sum(nomi) / sum(denom)
116
117  }
118
119  plot(pred, axes=FALSE, xlab="", ylab="Predicted temperature", main="Predicted temperatures (
         multiplication)")
120  axis(2)
121  axis(1, at=seq_along(pred),labels=as.character(times), las=2)
122  box()
```

## 2.2   Assignment 3

```
 1  library(neuralnet)
 2  set.seed(1234567890)
 3  Var <- runif(50, 0, 10)
 4  trva <- data.frame(Var, Sin=sin(Var))
 5  tr <- trva[1:25,] # Training
 6  va <- trva[26:50,] # Validation
 7  #plot(tr)
 8  # Random initialization of the weights in the interval [-1, 1]
 9  winit <- runif(10, min=-1, max=1)# Your code here
10  MSE = c()
11  thresholds = c()
12
13
14  for(i in 1:10) {
15    nn <- neuralnet(Sin ~ Var, data=tr, startweights = winit, threshold = i/1000, hidden = c(10)
         )
16    pred <- predict(nn, va)
17    MSE[i] <- mean((pred - va$Sin)^2)
18    thresholds[i] <- (i/1000)
19  }
20
21  plot(y=MSE, x=thresholds, type = 'o', main="Validation MSE for different thresholds")
22  best_threshold <- thresholds[which.min(MSE)]
23
24  nn <- neuralnet(Sin ~ Var, data=trva, startweights = winit, threshold = best_threshold, hidden
         = c(10))
25  plot(nn)
26  # Plot of the predictions (black dots) and the data (red dots)
27  plot(prediction(nn)$rep1,  col = "blue", main="Predicions")
28  points(trva, col = "red")
```