

Math 5365

Data Mining 1

Homework 14

Mary Barker

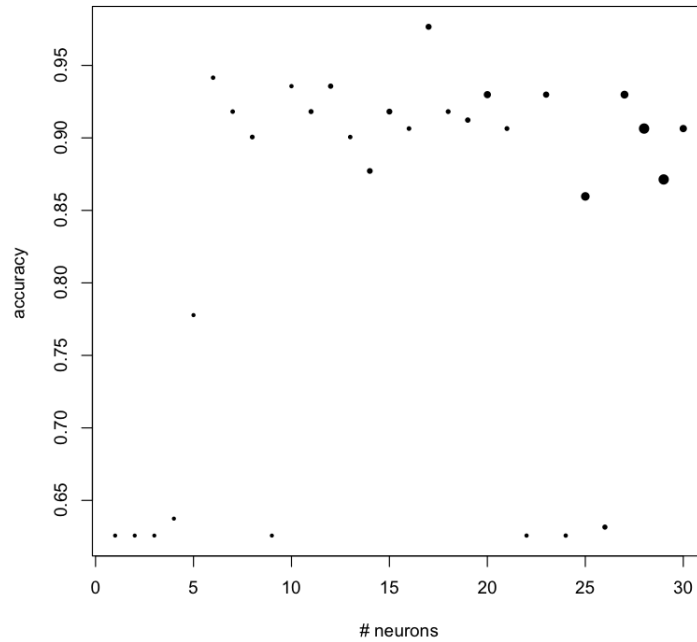
1. Use gradient descent to minimize the function  $f(x, y) = x^2 + y^2 + y$  on the open disk  $x^2 + y^2 < 1$ . This is essentially the example on p. 7 of the Lagrange multipliers slides, except that we are only interested in interior points on this problem.

(Hints: Start by randomly selecting a vector  $(x, y)$  in the disk, and then apply gradient descent. If the vector escapes the disk during the algorithm, randomly select a new vector. Store vectors that converge to a minimum in a matrix, and repeat the entire process a large number of times to ensure that there is only one local minimum inside the disk. You may need to take some care with the learning rate to prevent vectors from escaping the disk too often.)

Using 300 initial guesses, the gradient descent function accurately computed the minimum point and gradient to within  $10e-10$  accuracy.

2. Compare the merits of using a 2-layer neural network (with 1 hidden layer) and a multi-layer neural network to the wdbc data set. Let's say you only have 10 minutes to train a network with nnet or mlp. Which approach produces the model with the higher classification accuracy? It may be interesting to compare the two packages for a variety of time frames and to consider different network topologies for mlp, e.g., `size=c(2, 4)` vs. `size=c(2, 2, 2)`.

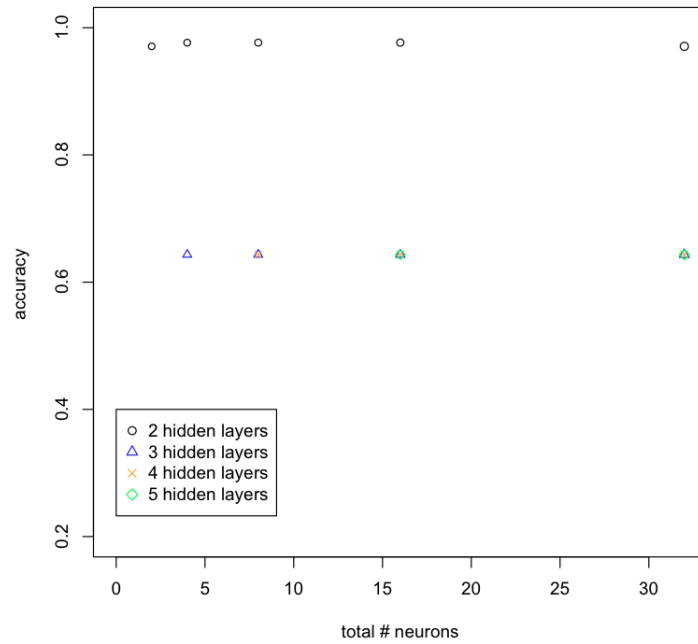
The first thing I considered was the effect of increasing the size of the hidden layer. The results, together with total time taken for each corresponding model are shown below. Note that the size of the data point corresponds to the length of time taken to produce the model.



As can be seen, the overall runtime increases with the size of the hidden layer. However, the accuracy also increases dramatically between 1 and 5 neurons and then remains largely at the same level (with variation). The highest accuracy occurred when using 17 neurons. With this model, the accuracy was approximately 98%. The total time taken for this model is slightly above average at a value of 0.254 seconds where the maximum time taken was 0.799 seconds, the minimum time taken was 0.008, and the mean time was 0.1976 seconds.

For the multi-layer neural networks, I compared results using 2 through 5 hidden layers. In each case, the number of neurons in each layer was equal. Thus for a network with 4 hidden layers and 16 total neurons, each layer contains 4 neurons.

The results are shown in the graph below. As in the single hidden layer case above, the size of the point indicates the length of time taken to create the model.



The average time taken was 0.8603333 seconds, with a minimum of 0.818 seconds using a neural network with 4 hidden layers and 1 neuron in each layer, and a maximum of 1.022 seconds for a neural network with 2 hidden layers and 8 neurons in each layer. The variations in time with these cases are so small that time taken has little meaning except when compared with the case of a single hidden network.

The maximum accuracy using multi-layer neural networks was roughly 98%, the same as for 1 hidden layer. This level of accuracy was achieved using 4 neurons in each case. Thus it can be seen that for this problem, the most accurate models are not the most complex ones, but rather the single- and double-hidden layer models. In addition, the most time-efficient model is the single-layer model. Therefore the best model for this problem is a network with a single hidden layer and 4 neurons.

```

1
2 library(doMC)
3 registerDoMC()

```

```

4
5 source('~Dropbox/Tarleton/data_mining/generic_functions/dataset_ops.R')
6
7 #1 Use gradient descent to minimize the function  $f(x, y) = x^2 + y^2 + y$ 
8 # on the open disk  $x^2 + y^2 < 1$ . This is essentially the example on
9 # p. 7 of the Lagrange multipliers slides, except that we are only
10 # interested in interior points on this problem.
11
12 # (Hints: Start by randomly selecting a vector  $(x, y)$  in the disk, and
13 # then apply gradient descent. If the vector escapes the disk during the
14 # algorithm, randomly select a new vector. Store vectors that converge
15 # to a minimum in a matrix, and repeat the entire process a large number
16 # of times to ensure that there is only one local minimum inside the disk.
17 # You may need to take some care with the learning rate to prevent vectors
18 # from escaping the disk too often.)
19
20 f <- function(x){
21   return( sum(x^2) + x[2])
22 }
23
24 grad_f <- function(x){
25   return(c(2*x[1], 2*x[2] + 1))
26 }
27
28 params <- function(x){
29   if( sum(x^2) < 1){
30     return(TRUE)

```

```

31   }else{
32       return(FALSE)
33   }
34 }
35
36 initial_guess <- function(){
37     x <- runif(2, -1, 1)
38     return(x)
39 }
40
41 grad_descent <- function(f, grad_f, params, initial_guess, x0, eps, tol, niter){
42     i = 1; NOT_ZERO = TRUE
43     x = x0
44     while((i < niter) && (NOT_ZERO)){
45         i = i + 1
46         x0 = x
47         x = x0 - eps * grad_f(x)
48         if(f(x) < f(x0)){
49             eps = eps * 1.1
50         }else{
51             eps = eps * 0.5
52         }
53         if(params(x)){
54             if(abs(sum(grad_f(x)^2)) < tol){
55                 NOT_ZERO = FALSE
56             }
57         }else{

```

```

58     x <- initial_guess()
59     eps = 0.75 * eps
60     i = 1
61   }
62 }
63 return(list(grad = grad_f(x), x = x, dx = eps,
64           nit = i, converged = !NOT_ZERO))
65 }
66
67 #####
68 # Initialize matrix of initial guesses and compute gradient descent #
69 #####
70 n_guesses = 300
71 x_m <- matrix(runif(n_guesses * 2, -1, 1), nrow = n_guesses, ncol = 2)
72 vals <- matrix(,nrow=n_guesses,ncol=2)
73
74 vals <- foreach(i=1:n_guesses, .combine=rbind) %dopar% {
75   grad_descent(f, grad_f, params, initial_guess, x_m[i,], 0.9, 1.0e-17, 10000)$x
76 }
77
78 dx <- paste('TOTAL VARIATION IN COMPUTED CRITICAL POINT',
79           '\n variation in x-values -> ',
80           toString(max(vals[,1]) - min(vals[,1])),
81           '\n','variation in y-values -> ',
82           toString(max(vals[,2]) - min(vals[,2])))
83 writeLines(dx)
84 print(apply(vals, 1, grad_f))

```

```

85
86 #2 Compare the merits of using a 2-layer neural network (with 1 hidden layer)
87 # and a multi-layer neural network to the wdbc data set. Let's say you only
88 # have 10 minutes to train a network with nnet or mlp. Which approach produces
89 # the model with the higher classification accuracy? It may be interesting to
90 # compare the two packages for a variety of time frames and to consider
91 # different network topologies for mlp, e.g., size=c(2, 4) vs. size=c(2, 2, 2).
92
93 library(nnet)
94 library(RSNNS)
95
96 source('~/.Dropbox/Tarleton/data_mining/class_notes/extras.R')
97 wdbc <- read.csv('~/.Dropbox/Tarleton/data_mining/dfiles/wdbc.data',
98                 header=F, sep=',')
99 wdbc <- wdbc[,-1]
100 splitset <- splitdata(wdbc,0.7,F)
101 train <- splitset$train
102
103 ## Consider first merely increasing size
104 len = 30
105
106
107 # 1    2    4    8    16
108 #    11   22   44   88
109 #      1111 2222 4444
110 #          11111111 22222222
111 #              1111111111111111

```

```

112
113 # accnnet <- rep(-1,len)
114 # for(i in 1:len){
115   # model <- nnet(V2~., wdbc[train,], size=i,trace=F)
116   # predvals <- predict(model,wdbc[-train,],type='class')
117   # accnnet[i] <- confmatrix(wdbc$V2[-train], predvals)$accuracy
118 # }
119
120 normwdbc <- standardize(wdbc, 2:ncol(wdbc))
121 train_vals <- normwdbc[train,-1]
122 train_targ <- decodeClassLabels(normwdbc[train,1])
123 test_vals <- normwdbc[-train,-1]
124 test_targ <- decodeClassLabels(normwdbc[-train,1])
125
126 num = 4
127 accnnet <- rep(-1, num)
128 accmlpnet <- matrix(-1, num, num)
129 nnet_times <- rep(-1, num)
130 mlp_times <- matrix(-1, num, num)
131
132 for(i in 1:num){
133   t_0 <- proc.time()
134   model <- nnet(V2~., wdbc[train,], size=i, linout=F, trace=F)
135   t_1 <- proc.time()
136   nnet_times[i] <- t_1[3] - t_0[3]
137   predvals <- predict(model, wdbc[-train,],type='class')
138   accnnet[i] <- confmatrix(wdbc$V2[-train], predvals)$accuracy

```



```

139
140   for(j in 1:num){
141       if(j >= i){
142           reps <- 2^(i - 1)
143           vec <- rep(2^(j - i), reps)
144           t_0 <- proc.time()
145           model <- mlp(train_vals,
146                       train_targ,
147                       size=vec,
148                       maxit = 50,
149                       learnFuncParams = c(0.1),
150                       inputsTest = test_vals,
151                       targetsTest = test_targ,
152                       linout = TRUE)
153           t_1 <- proc.time()
154           mlp_times[i,j] <- t_1[3] - t_0[3]
155           idx <- (model$fittedTestValues[,1] >= 0.5) * 1
156           predvals <- idx
157           predvals[idx == 1] <- 'B'
158           predvals[idx == 0] <- 'M'
159           accmlpnet[i,j] <- confmatrix(wdbc$V2[-train], predvals)$accuracy
160       }
161   }
162 }
163 mlp_times

```