

HOMEWORK1

Read the data

We will dataset is about the house prices in a region in Europe together with the square footage and number of bedroom information denoted by the column (variable) with names sqft and bd respectively. There are n=10 million data points.

```
In [1]: options(warn=1)
houseData <- read.csv("house.csv")
```

Inspect the data read

R read the data as dataframe by default so there is no need to convert the "houseData" to dataframe.

```
In [2]: print(paste("Is houseData a dataframe", is.data.frame(houseData), sep = " ", quote=FALSE))
head(houseData)

[1] Is houseData a dataframe TRUE
  X bd sqft price
1 1 1370 451000
2 4 2060 899000
3 4 1738 944000
4 3 1588 715000
5 4 2092 897000
6 3 1456 703000
```

```
In [3]: cat("Size of data frame",dim.data.frame(houseData))

Size of data frame 10000000 4
```

Question 1 - Rewrite the objective in the form

```
In [66]: y <- as.matrix(houseData$price)
A <- as.matrix(houseData[,2:3])
ones = matrix(rep(1,nrow(y)), nrow = nrow(y), ncol = 1)
print("Price")
head(y)
print("Features")
head(A)
[1] "Price"
451000
899000
944000
715000
897000
703000
[1] "Features"
bd sqft
1 1370
4 2060
4 1738
3 1588
4 2092
3 1456
```

Normalise Data

```
In [67]: y <- (y - mean(y))/sd(y)
A[,1] = (A[,1] - mean(A[,1]))/sd(A[,1])
A[,2] = (A[,2] - mean(A[,2]))/sd(A[,2])
print("After Normalisation")
head(y)
print("Price")
head(A)
[1] "After Normalization"
[1] "Price"
-1.2490641
1.4362784
1.7060115
0.3333699
1.4242903
0.2614411
[1] "Features"
bd sqft
1 -1.3411019 -0.6753744
1.3424687 1.9160003
1.3424687 0.7066921
0.4479452 0.1434948
1.3424687 2.0361800
0.4479452 0.1433496
1.3424687 2.0361800
0.4479452 -0.3523915
```

Complete Data matrix A with ones for intercept

```
In [68]: A = cbind(ones,A)
colnames(A)[1] <- "intercept"
head(A)
ncol(A)

intercept bd sqft
1 -1.3411019 -0.6753744
1 1.3424687 1.9160003
1 1.3424687 0.7066921
1 0.4479452 0.1434948
1 1.3424687 2.0361800
1 0.4479452 -0.3523915
3
```

Objective Function in the Matric form with square of L2 Norm

```
In [69]: fx <- function(x,y,A){
  return ((norm((y-(A**x))), type = "2"))**2/(2*nrow(A))
}
##This is sample x row matrix [x1,x2,x3] initialize as zero
x = matrix(rep(0,ncol(A)),nrow = ncol(A), byrow = TRUE)
dim(x)
```

```
1.3
2.1
```

Question 2 - Gradient descent

Definition of function for gradient descent and gradient calculation

```
In [70]: deltafx <- function(y1,A1,x1){
  yhat <- A1*x1
  diff <- y1 - yhat
  At <- t(A1)
  dF <- -(At*x1+diff)/nrow(A1)
  return(diff)
}

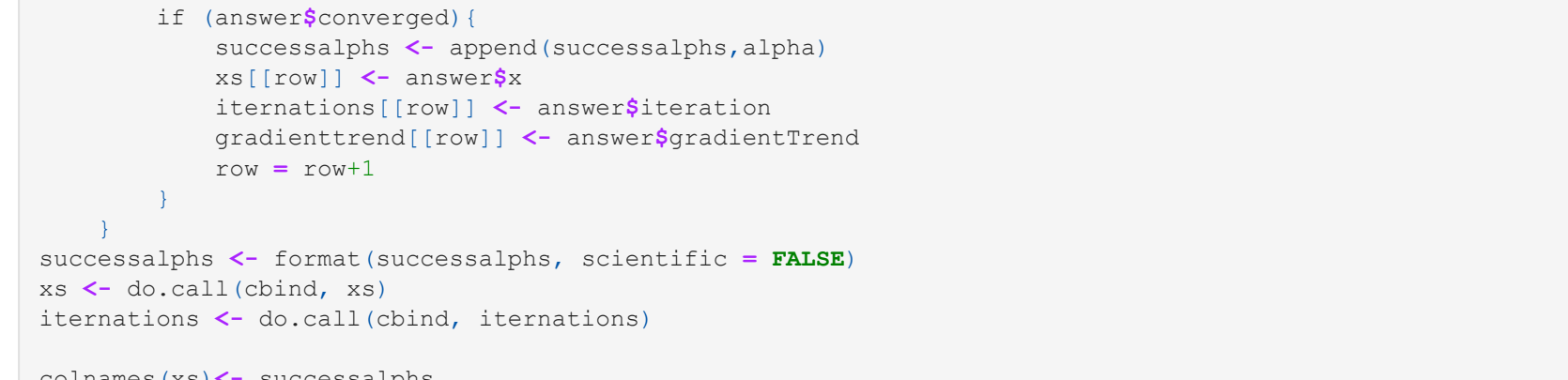
gradientDescent <- function(y1,A1,x1,alpha,thrs, maxi){
  converged <- FALSE
  i=1
  while(!converged && i <= maxi){
    {
      deltafx <- deltafx(y1,A1,x1)
      if (is.nan(norm(deltafx))){
        break
      }
      x1 <- x1 + (alpha*deltafx)
      gradientTrend <- append(gradientTrend,norm(deltafx))
      converged <- (norm(deltafx) <= thrs)
      i <- i+1
    }
  }
  return (list("x1"= x1,"gradientTrend" =gradientTrend, "iteration" = i-1,"converged" = converged))
}
```

```
In [71]: #running gradient descent on complete data
alphas <- c(1,0.1)
thrs <- 10**(=-2)
maxiter <- 1000
maxiter <- 1000
A2 <- A
y2 <- y
successalphs <- c()
iterations <- list()
gradienttrend <- list()
xsfull <- list()
x <- matrix(rep(0,ncol(A)),nrow = ncol(A), byrow = TRUE)
row = 1
for (alpha in alphas) {
  answer <- gradientDescent(y,A,x,alpha,thrs,maxiter)
  if (answer$converged){
    successalphs <- append(successalphs,alpha)
    xsfull[[row]] <- answer$x
    iterations[[row]] <- answer$iteration
    gradienttrend[[row]] <- answer$gradientTrend
    row = row+1
  }
}
successalphs <- format(successalphs, scientific = FALSE)
xsfull <- do.call(cbind, xsfull)
iterations <- do.call(cbind, iterations)
colnames(xsfull) <- successalphs
colnames(iterations) <- successalphs
modifiedxs <- format(xsfull, scientific = FALSE)
modifiedxs
iterations
successalphs
```

	1.0	0.1
intercept	0.0000000000000007827117	0.0000000000000008478084
bd	0.6714973007165752072671	0.63946203261536271877929
sqft	0.32223947593443569292759	0.34986143845254865427208

1.0	0.1
31	108

```
In [72]: barplot(iterations, main="Iteration vs alpha", names.arg = successalphs,
  ylab = "Iterations", ylim = c(0,max(iterations)+20), xlab = "Alpha", axes = FALSE)
ylabel <- seq(0, max(iterations)+20, by = max(iterations)/10)
axis(2, at = ylabel, las = 1)
box()
```



Running gradient descent for various values of alpha with less no. number of rows

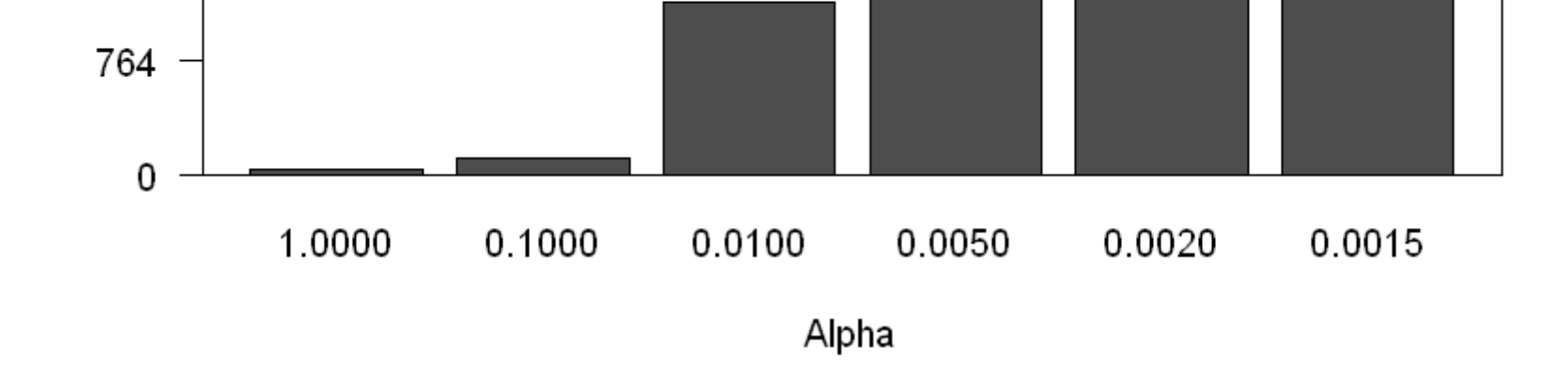
```
In [73]: alphas <- c(10,1,0.1,0.01,0.005,0.002,0.0015,0.0010,0.0008,0.0006,0.0004,0.0002,0.0001)
thrs <- 10**(=-2)
maxiter <- 1000
#considering first 1000 rows
A2 <- A[1:1000,]
y2 <- matrix(y[1:1000,],nrow=1000,byrow=TRUE)

successalphs <- c()
iterations <- list()
gradienttrend <- list()
xs <- list()
row = 1
for (alpha in alphas)
{
  x <- matrix(rep(0,ncol(A)),nrow = ncol(A), byrow = TRUE)
  answer <- gradientDescent(y2,A2,x,alpha,thrs,maxiter)
  if (answer$converged){
    successalphs <- append(successalphs,alpha)
    xs[[row]] <- answer$x
    iterations[[row]] <- answer$iteration
    gradienttrend[[row]] <- answer$gradientTrend
    row = row+1
  }
}
successalphs <- format(successalphs, scientific = FALSE)
xs <- do.call(cbind, xs)
iterations <- do.call(cbind, iterations)
colnames(xs) <- successalphs
colnames(iterations) <- successalphs
xs
iterations
successalphs
```

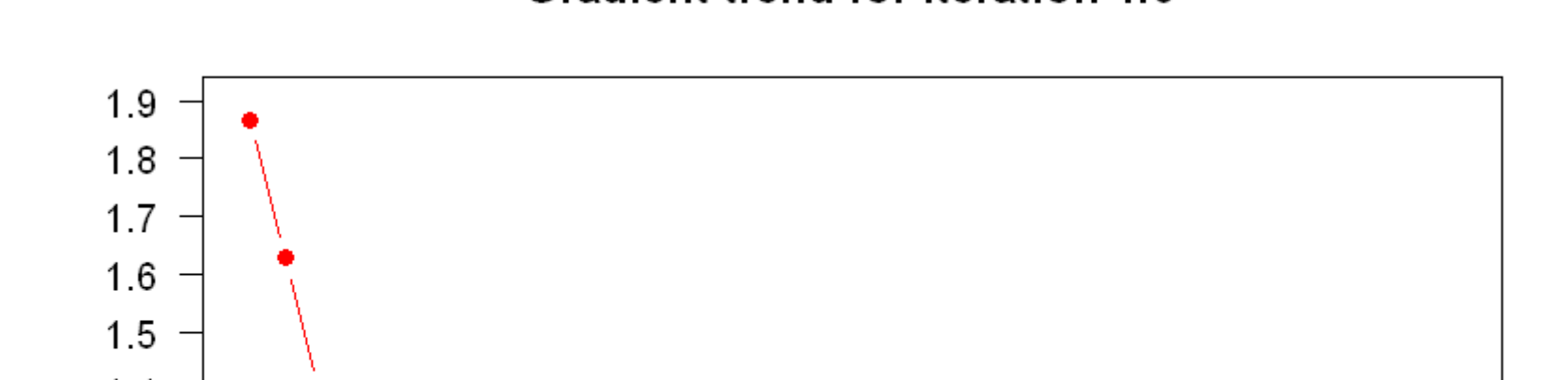
	1.0000	0.1000	0.0100	0.0050	0.0020	0.0015
intercept	-0.009055142	-0.01009188	-0.01010718	-0.01010849	-0.01010900	-0.01010900
bd	0.6738787865	0.64525086	0.64470713	0.64465470	0.64463654	0.6446367
sqft	0.301272612	0.33363521	0.33417774	0.33422421	0.3342423	0.3342420

1.0000	0.1000	0.0100	0.0050	0.0020	0.0015
36	115	1147	2293	5732	7643

```
In [74]: barplot(iterations, main="Iteration vs alpha", names.arg = successalphs,
  ylab = "Iterations", ylim = c(0,max(iterations)+100), xlab = "Alpha", axes = FALSE)
ylabel <- seq(0, max(iterations)+100, by = max(iterations)/10)
axis(2, at = ylabel, las = 1)
box()
```



```
In [75]: ## Plot of gradient trend for best alpha
plot(gradienttrend[which("1.0000",names(successalphs))], type = "b", main="Gradient trend for iteration 1.0",
  col = "red", xlab = "Iterations", ylab = "Gradient Value", axes = FALSE)
ylabel <- seq(0, 2, by = 0.1)
xlabel <- seq(0, 70, by = 5)
axis(1, at = xlabel, las = 1)
axis(2, at = ylabel, las = 1)
box()
```



```
In [76]: sol <- solve(t(A**A))**t(A)**y
modifiedsol <- format(sol, scientific = FALSE)
modifiedsol
```

	1.0	0.1
intercept	0.0000000000000000016019	
bd	0.6700795745981675821668	
sqft	0.319263896737376973850	

Calculating difference between actual solution and solutions obtained for various value of alphas for full data

```
In [77]: as <- apply(xsfull[,2,function(x) sol-x)
as <- format(as, scientific = FALSE)
as
```

	1.0	0.1
0.000000000000000002333068		0.000000000000000000168201
-0.0014772611840762507995		-0.03061754198280486338746
-0.0029757917669799554261		-0.0306175416948105668870

Conclusion

For complete dataset algorithm converge for values of alpha = 1 with 31 iterations and alpha = 0.1 with 108 iterations

Multiple Alphas Gradient descent was run for alphas = [1.0,1.0,0.01,0.005,0.002,0.0015,0.0008,0.0006,0.0004,0.0002,0.0001] There was no convergence for alpha greater than 1. Algorithm converged for alphas [1.0,1.0,0.01,0.005,0.002,0.0015] with in 10,000 iterations, with minimum iterations 36 for alpha = 1.

Also there is negligible difference between solutions using gradient descent and actual solutions.

Question 3 - Plot Contour plot with movement of gradient descent and cost function for two features with 3rd kept as constant

Generating random values for x2 and x3 while keeping x1 at optimum value

```
In [78]: #random values are generated and arranged in a matrix to feed to cost function
xrange2 <- seq(0,2,by=0.5),to = sol[2]*(1.5), length.out = 50
xrange3 <- seq(0,2,by=0.5),to = sol[3]*(1.5), length.out = 50
x2 = rep(xrange2,length(xrange2))
x3 = c()
for(i in xrange3){
  list = rep(i,length(xrange3))
  x1 = append(xr3,list)
}
xrandom = cbind(xr1,xr2,xr3)
head(xrandom)
```

xr1	xr2	xr3
1.016019e-16	0.3350398	0.1596319
1.016019e-16	0.3487149	0.1596319
1.016019e-16	0.3623900	0.1596319
1.016019e-16	0.3760651	0.1596319
1.016019e-16	0.3897402	0.1596319
1.016019e-16	0.4034153	0.1596319

Solving fx (cost function) for xrandom

```
In [81]: fxValues <- c()
for (i in seq(1,nrow(xrandom))) {
  x1 <- as.matrix(xrandom[i,], nrow = 3)
  fxValues[i] <- fx(x1,A)
}
fxValues <- matrix(fxValues,nrow = length(xrange2),ncol = length(xrange3) )
dim(fxValues)
head(fxValues)
```

	1.50	2.50
0.1586738	0.158219	0.1530125
0.1525226	0.1495756	0.1468410
0.1462184	0.1435162	0.1408564
0.1402713	0.1376438	0.1355089
0.1345111	0.1319650	0.1294484
0.1289319	0.1264602	0.1244081
0.1234911	0.1212480	0.1216319
0.1181943	0.1162987	0.1192815
0.1130343	0.1115790	0.1173444
0.1080162	0.1070443	0.1157290
0.1031662	0.1026788	0.1142850
0.0984858	0.0984657	0.1130343
0.0939584	0.0943984	0.1119281
0.0895719	0.0904719	0.1109303
0.0853184	0.0866719	0.1100343
0.0811849	0.0829884	0.1092423
0.0771584	0.0794149	0.1085568
0.0732369	0.0759474	0.1079663
0.0694164	0.0725799	0.1074663
0.0656929	0.0693024	0.1070528
0.0620614	0.0661149	0.1067213
0.0585189	0.0630124	0.1064687
0.0550614	0.0600000	0.1062878
0.0516849	0.0570725	0.1061662
0.0483874	0.0542350	0.1060963
0.0451589	0.0514875	0.1060713
0.0419954	0.0488250	0.1060849
0.0388929	0.0462425	0.1061284
0.0358474	0.0437350	0.1061969
0.0328569	0.0413075	0.1062844
0.0299174	0.0389550	0.1063959
0.0270269	0.0366825	0.1065264
0.0241824	0.0344850	0.1066719
0.0213819	0.0323575	0.1068294
0.0186224	0.0302950	0.1069969
0.0159029	0.0282925	0.1071714
0.0132234	0.0263450	0.1073509
0.0105839	0.0244525	0.1075324
0.0080794	0.0226100	0.1077149
0.0056949	0.0208175	0.1078964
0.0034254	0.0190750	0.1080769
0.0012709	0.0173775	0.1082554

Gradient descent is for with x1 kept at its optimum values.

```
In [22]: #storing x for plotting
gradx <- list()

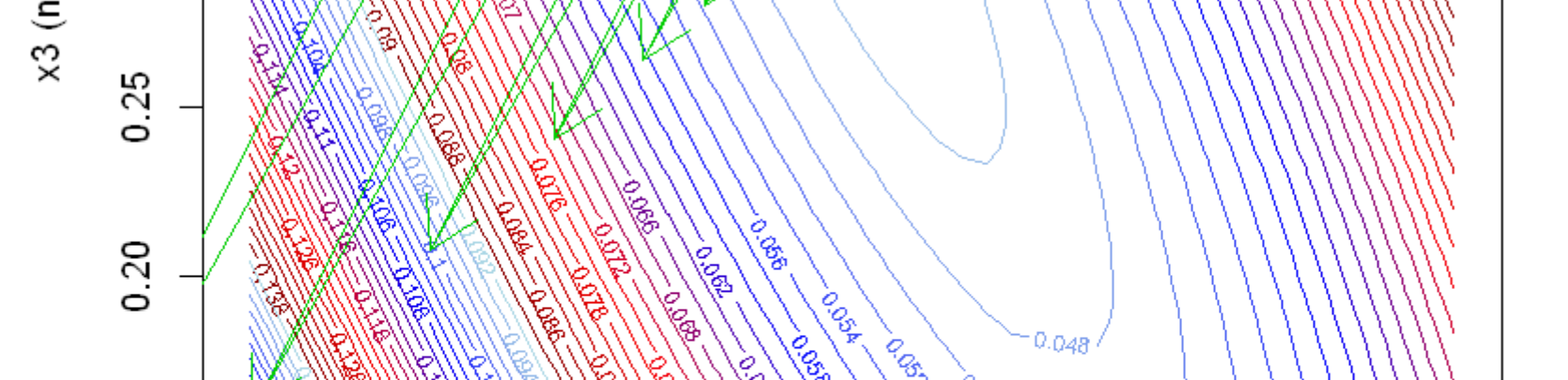
x2 <- matrix(c(sol[1],0,0),nrow = ncol(A), byrow = TRUE)
thrs <- 10**(=-2)
maxiter <- 1000
converged <- FALSE
maxi <- 100
gradX[[1]] <- x2
i=2
alpha = 1
while(!converged && i <= maxi){
  {
    delta <- deltafx(y,A,x2)
    if (is.nan(norm(delta))){
      break
    }
    x2 <- x2 + (alpha*delta)
    # xi (inside column matrix x2) is changed to optimum value
    x2[i] = sol[i]
    gradX[[i]] <- x2
    gradientTrend <- append(gradientTrend,norm(deltafx))
    converged <- (norm(deltafx) <= thrs)
    i <- i+1
  }
}
#creating matrix of solution of xs to plot direction of gradient
mgradx <- do.call(cbind, gradx)
tmgradx <- t(mgradx)
head(tmgradx)

xe2 <- tmgradx[,2]
xe3 <- tmgradx[,3]
```

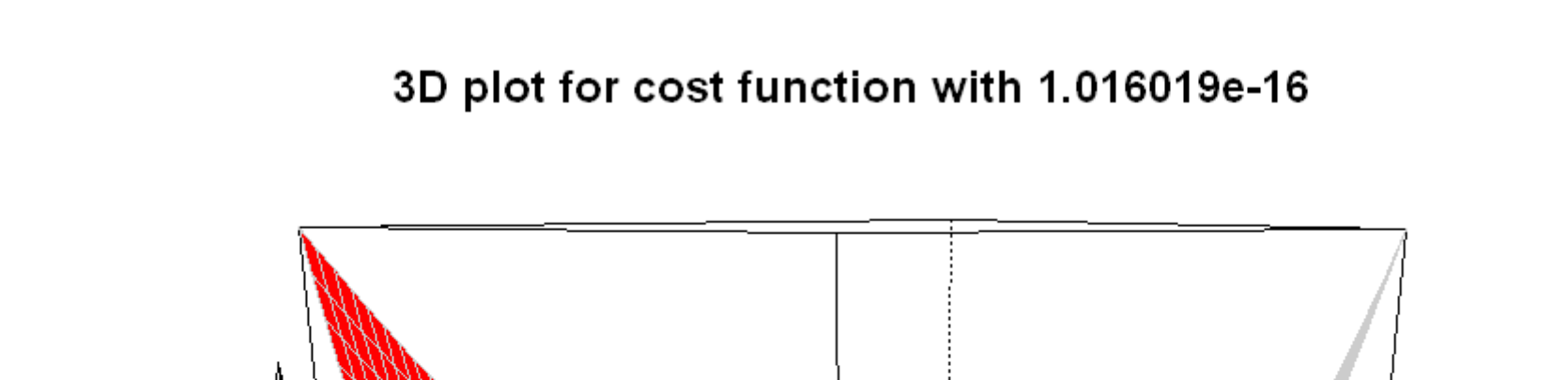
intercept	bd	sqft
1.016019e-16	0.0000000	0.000000000
1.016019e-16	0.9381574	0.88191285
1.016019e-16	0.1976374	0.09416576
1.016019e-16	0.8590888	0.15961676
1.016019e-16	0.3369824	0.05955765
1.016019e-16	0.8033412	0.59895719

Contour plot with Gradient Flow

```
In [23]: cols = rev(colorRampPalette(c("darkred","red","blue","lightblue"))(24))
contour(xrange2,xrange3,z = fxValues,nlevels = 80, main = "Level sets for x1 = 1.016019e-16", xlab = "x2 (n
  x = seq(length(xe2)-1)
arrows(xe2[x], xe3[x], xe2[x+1], xe3[x+1], col = 3)
```



```
In [24]: persp(xrange2,xrange3,fxValues, theta = 40,main = "3D plot for cost function with 1.016019e-16", xlab = "x2
  col="red", border="grey30")
```



As graph is not perfect cereal bowl shape to inspect the possible reason corelation is calculated b/w "bd" and "sqft". This same thing will be present in their solution and range around it

```
In [23]: cor(c(A[,2]), c(A[,3]), use="complete.obs", method="pearson")
0.839674958214801
```

Conclusion

Level set is plotted for a range of x2 and x3 with x1 kept at optimum it is observed that gradient to flow perpendicular to level sets and as visible in the plot gradient jumped around a lot before converging to minima. For even though not shown in code 3 iteration were done one with x = [0,0,0], one with x started with = [x1,0,0] and final one shown in arc where x1 is kept x1 for complete gradient process

3D plot is not perfect cereal bowl shape reason for this high corelation between two features of the dataset as indicated indicated by correlation factor 0.839674958214801

Question 4-Repeat the previous question with different choice of the stepsize alpha in your implementation

Conclusion for part a)

Expanding on the analysis done for Question 2, starting with alpha = 10 we moved to alpha = 0.0001 and for higher alpha there was no convergence. Convergence started with alpha = 1 with minimum no. of iterations for alpha = 1.

For this problem any alpha greater than 1 will lead to divergence.

```
In [38]: barplot(iterations, main="Iteration vs alpha", names.arg = successalphs,
  ylab = "Iterations", ylim = c(0,max(iterations)*max(iterations)/10), xlab = "Alpha", axes = FALSE)
ylabel <- seq(0, max(iterations)*max(iterations)/10, by = max(iterations)/10)
axis(2, at = ylabel, las = 1)
box()
```



```
In [39]: gradientDescent2 <- function(y1,A1,x1,advalpha,thrs, maxi, option){
  converged <- FALSE
  i=1
  while(!converged && i <= maxi){
    {
      deltafx <- deltafx(y1,A1,x1)
      if (is.nan(norm(deltafx))){
        break
      }
      #option for alpha
      if (option == 1)
      {
        alpha <- advalpha/i
      } else if (option == 2)
      {
        alpha <- advalpha/sqrt(i)
      }
      x1 <- x1 + (alpha*deltafx)
      gradientTrend <- append(gradientTrend,norm(deltafx))
      converged <- (norm(deltafx) <= thrs)
      i <- i+1
    }
  }
  return (list("x1"= x1,"gradientTrend" =gradientTrend, "iteration" = i-1,"converged" = converged, "alpha" = alpha))
}
```

```
In [43]: gdwadpstep <- function(alphas,opts){
  thrs <- 10**(=-2)
  maxiter <- 1000
  #considering first 1000 rows as below data set takes lot of time this many iterations
  A3 <- A[1:1000,]
  y3 <- matrix(y[1:1000,],nrow=1000,byrow=TRUE)
  successalphs <- c()
  finalalphs <- list()
  iterations <- list()
  gradienttrend <- list()
  xs <- list()
  row = 1
  for (alpha in alphas)
  {
    x <- matrix(rep(0,ncol(A)),nrow = ncol(A), byrow = TRUE)
    answer <- gradientDescent2(y3,A3,x,alpha,thrs,maxiter,opt)
    if (answer$converged){
      successalphs <- append(successalphs,alpha)
      finalalphs[[row]] <- answer$alpha
      xs[[row]] <- answer$x
      iterations[[row]] <- answer$iteration
      gradienttrend[[row]] <- answer$gradientTrend
      row = row+1
    }
  }
  if (is.null(successalphs))
  {
    print("Convergence not achieved")
  } else {
    successalphs <- format(successalphs, scientific = FALSE)
    iterations <- do.call(cbind, iterations)
    finalalphs <- do.call(cbind, finalalphs)
  }
}
```

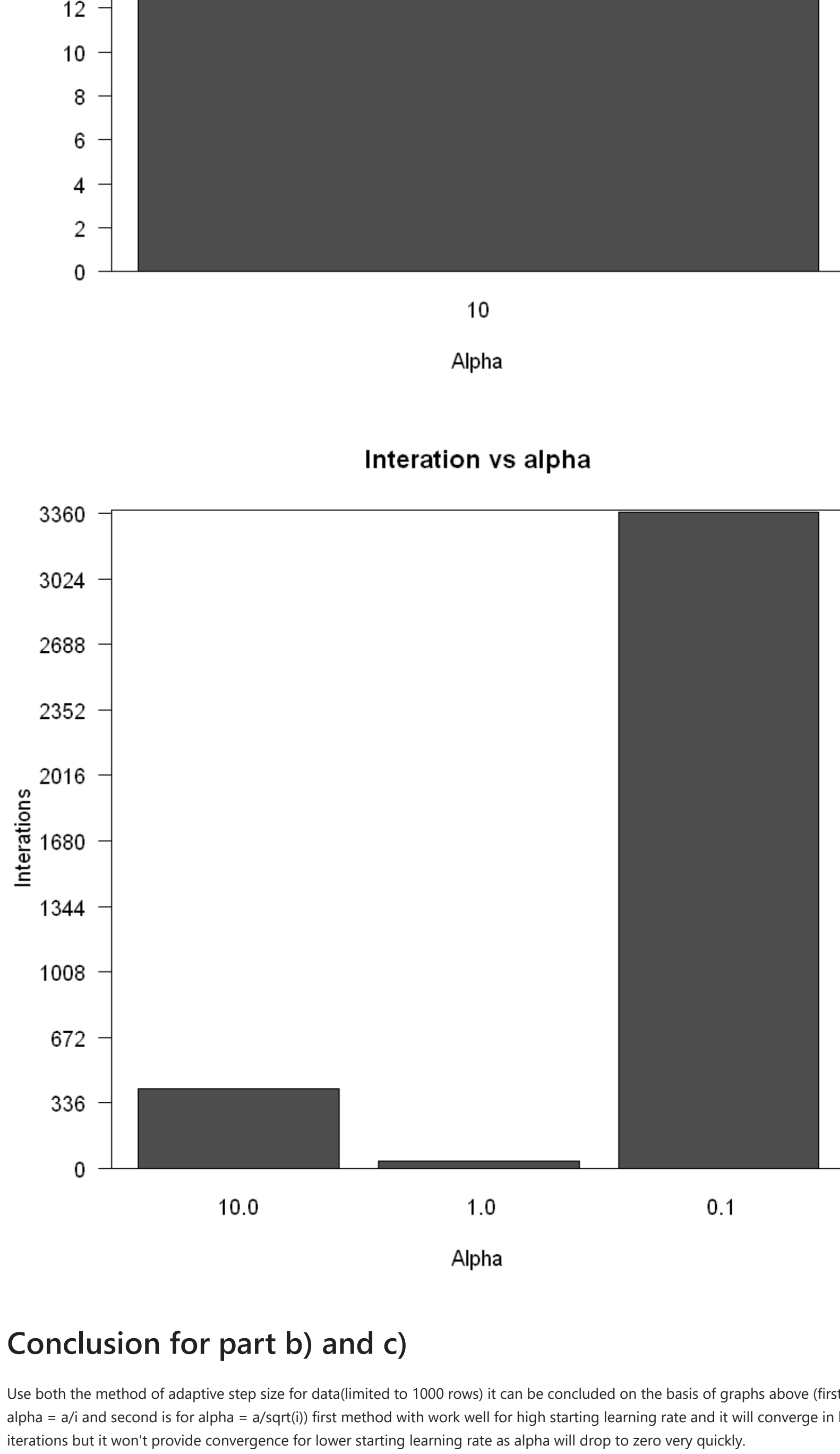


```
colnames(xs)<- successalphs
colnames(iterations)<- successalphs
as <- apply(xs,2,function(x) sol ~x)
as <- format(as, scientific = FALSE)
print(as)
barplot(iterations, main="Iteration vs alpha", names.arg = successalphs,
ylab = "Iterations", ylim = c(0,max(iterations)*10), xlab = "Alpha", axes = FALSE)
ylab1 <- seq(0, max(iterations)*10, by = max(iterations)/10)
axis(2, at = ylab1, las = 1)
box(2)
```

```
In [45]: alpha <- c(10,1,0.1,0.01,0.001)
options <- c(1,2)
for (opt in options)
{
  gdwthadpstep(alpha,opt)
}
```

```
10
[1,] " 0.00950546"
[2,] " -0.006607395"
[3,] " 0.016984225"
1.0
[1,] " 0.01001283" " 0.01006655" " 0.01010885"
[2,] " 0.02197309" " 0.02391818" " 0.02543807"
[3,] " -0.01154890" " -0.01347169" " -0.01497335"
```

Interaction vs alpha



Interaction vs alpha



Conclusion for part b) and c)

Use both the method of adaptive step size for data(limited to 1000 rows) it can be concluded on the basis of graphs above (first is for alpha = a/i and second is for alpha = a/sqrt(i)) first method with work well for high starting learning rate and it will converge in less iterations but it won't provide convergence for lower starting learning rate as alpha will drop to zero very quickly.

Second method provide convergence for smaller learning rates giving chance to better solutions . Though both the approaches have same similar difference from solution by solving linear equations.

Question 5 - Stochastic Gradient Method

```
In [16]: stochasticGradientDes <- function(y1,A1,x1,advalpha,thres, maxi, option){
  i=1
  x1<-x1
  gradientTrend <- c()
  alpha <- advalpha
  while(!converged && i <= maxi){
    {
      index <- sample(nrow(A1),1)
      y1 <- matrix(y1[index],nrow = 1,byrow = TRUE)
      A1 <- matrix(A1[index],nrow = 1,byrow = TRUE)
      deltaf <- deltafx(y1,A1,x1)
      if(is.na(norm(deltaf))){
        break
      }
    }
    #option for alpha
    if (option == 1)
      alpha <- advalpha
    } else if (option == 2)
      alpha <- advalpha/i
    } else if (option == 3)
      alpha <- advalpha/sqrt(i)
    }
    x1 <- x1 - (alpha*deltaf)
    gradientTrend <- append(gradientTrend,norm(deltaf))
    converged <- (norm(deltaf) < thres)
    i <- i+1
  }
  return (list("x1"= x1,"gradientTrend" =gradientTrend, "iteration" = i-1,"converged" = converged, "alpha"
```

```
In [23]: # Stochastic Gradient on complete data for alpha = 1 and 0.1
alpha <- c(1,0.1,0.0,0.1)
thres <- 10**(2)
maxiter <- 1000
A2 <- A
y2 <- y
successalphs <- c()
iterations <- list()
gradienttrend <- list()
xsfull <- list()
x <- matrix(rep(0,ncol(A)),nrow =ncol(A), byrow = TRUE)
row = 1
for (alpha in alphas) {
  answer <- stochasticGradientDes(y,A,x,alpha,thres,maxiter,1)
  if (answer$converged){
    successalphs <- append(successalphs,alpha)
    xsfull[[row]] <- answer$x
    iterations[[row]] <- answer$iteration
    gradienttrend[[row]] <- answer$gradientTrend
    row = row+1
  }
}
```

```
successalphs <- format(successalphs, scientific = FALSE)
xsfull <- do.call(cbind, xsfull)
iterations <- do.call(cbind, iterations)

colnames(xsfull)<- successalphs
colnames(iterations)<- successalphs
modifiedx <- format(xsfull, scientific = FALSE)
modifiedxs
iterations
successalphs
```

	0.10	0.01
0.026916195	0.004764189	
0.538005611	0.521706534	
0.415810632	0.425792914	
0.10	0.01	
62	217	
1	0.10	
2	0.01	

Observation on SGD

It is fast compare to normal GD, but it doesn't guarantee same solution in each iteration or as a matter of fact best solution. Above I ran SGD for three alphas out of which it converged for two but the iteration and solution is different in each run, this will be shown in the next section.

```
In [38]: alpha <- 0.1
thres <- 10**(2)
maxiter <- 1000
A2 <- A
y2 <- y
Siterations <- list()
Sgradienttrend <- list()
Sxsfull <- list()
x <- matrix(rep(0,ncol(A)),nrow =ncol(A), byrow = TRUE)
row = 1
for (i in seq(1,5,1)) {
  answer <- stochasticGradientDes(y,A,x,alpha,thres,maxiter,1)
  if (answer$converged){
    Sxsfull[[row]] <- answer$x
    Siterations[[row]] <- answer$iteration
    Sgradienttrend[[row]] <- answer$gradienttrend
    row = row+1
  }
}

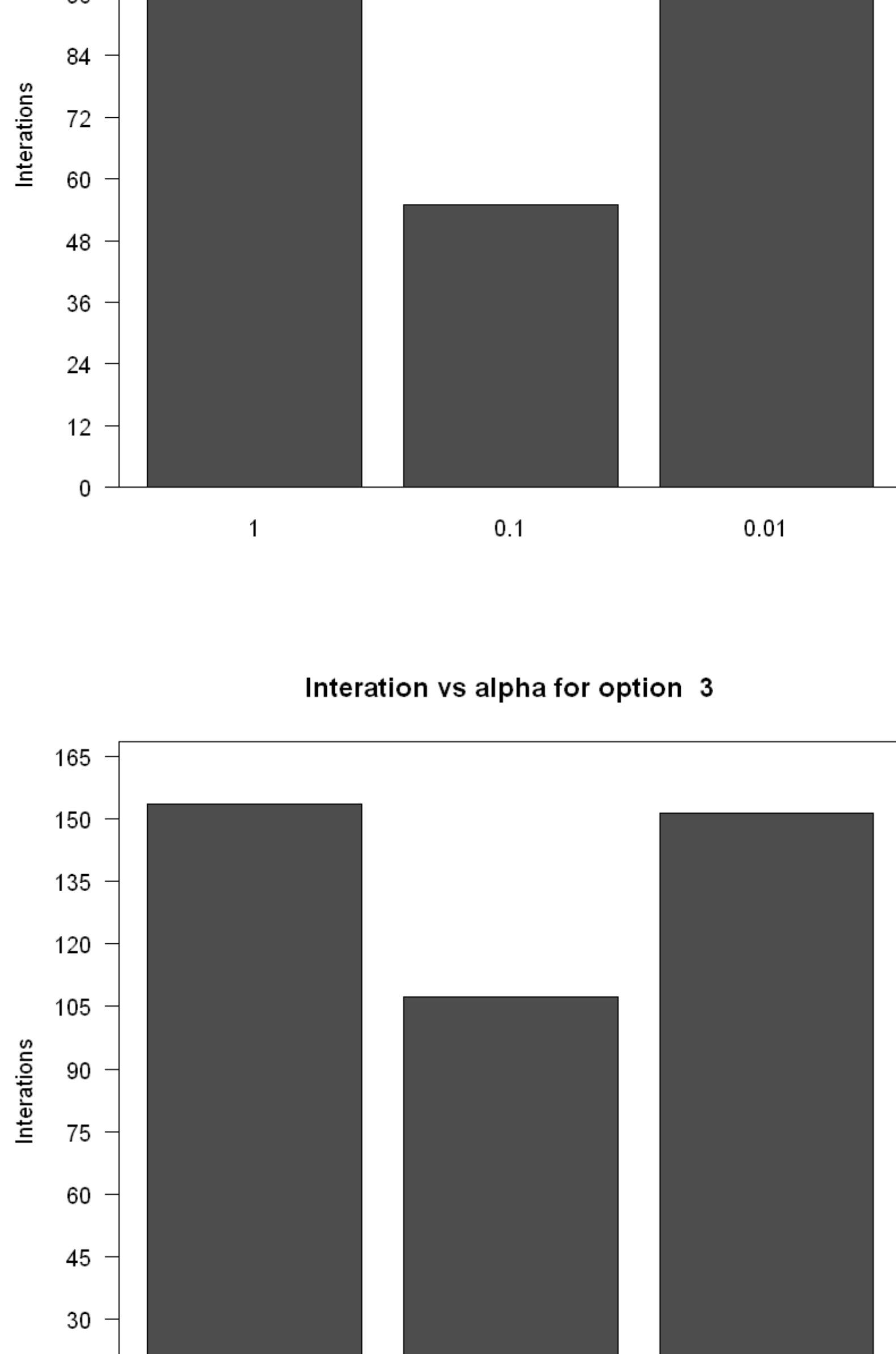
Sxsfull <- do.call(cbind, Sxsfull)
Siterations <- do.call(cbind, Siterations)

colnames(Sxsfull)<- seq(1,5,1)
colnames(Siterations)<- seq(1,5,1)
modifiedx <- format(Sxsfull, scientific = FALSE)
modifiedxs
Siterations
```

	1	2	3	4	5
-0.017651936	-0.006640093	-0.075477447	-0.005380836	-0.044466205	
0.810941528	0.4860311	0.707340578	0.650536353	0.55632412	
0.462324671	0.38821118	0.330708252	0.305325255	0.348899485	
1	2	3	4	5	
110	31	172	145	56	

```
In [39]: barplot(Siterations, main="Iteration for same alpha in different run", names.arg = seq(1,5,1),
ylab = "Iterations", ylim = c(0,max(Siterations)*max(Siterations)/10), xlab = "", axes = FALSE)
contour(xrange2,xrange3,z = $a$Values,levels = 80, main = "Level sets for x1 = 1.016019e-16", xlab = "", axes = FALSE)
axis(2, at = ylab1, las = 1)
box(1)
```

Interaction for same alpha in different run



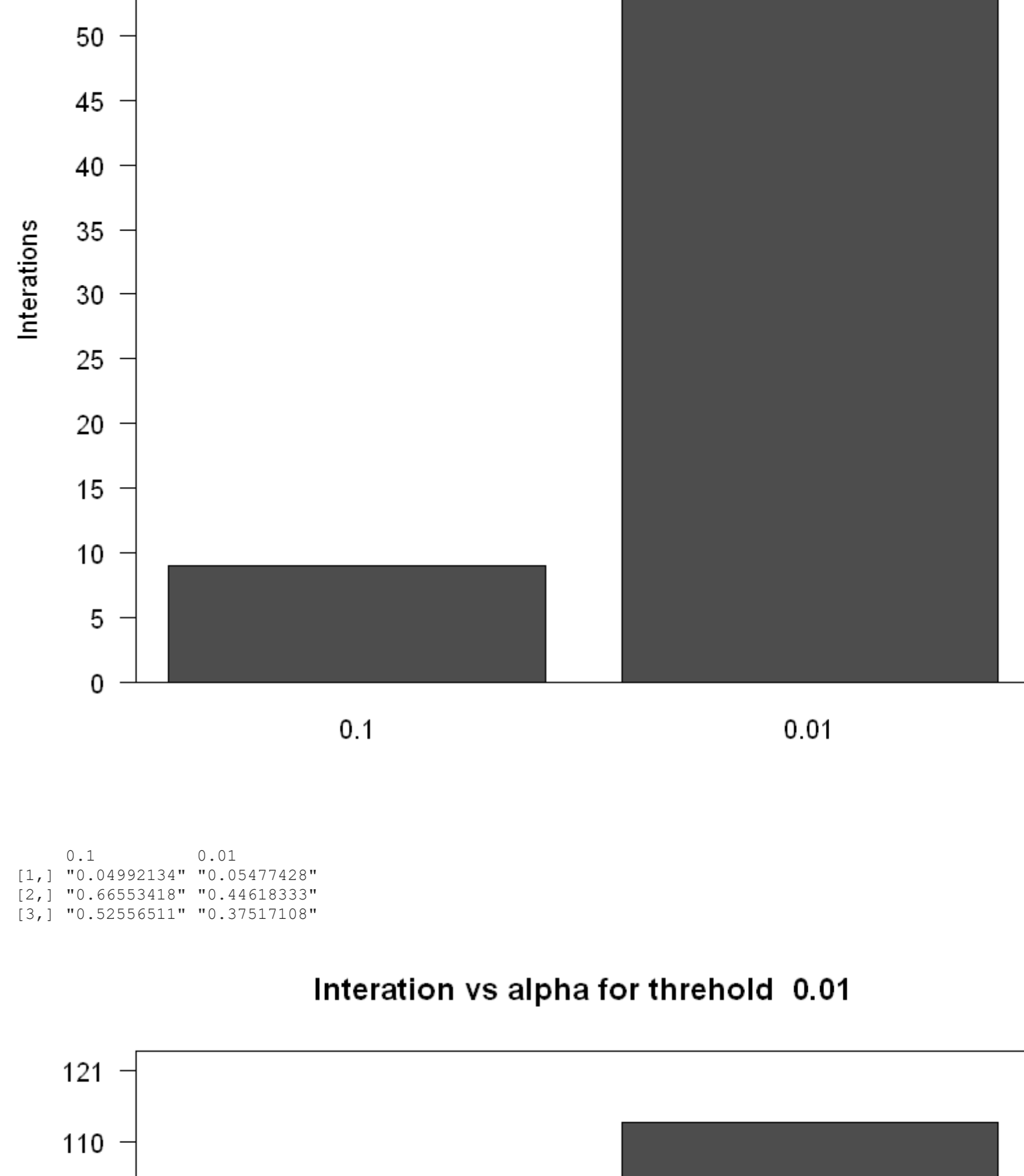
```
In [40]: as <- apply(Sxsfull,2,function(x) sol ~x)
as <- format(as, scientific = FALSE)
as
```

	1	2	3	4	5
0.017651936	-0.006640093	-0.075477447	-0.005380836	-0.044466205	
-0.140861953	0.181919264	-0.037261003	0.019543222	0.113755462	
-0.143060974	-0.069557222	-0.011444355	0.013938642	-0.029635588	

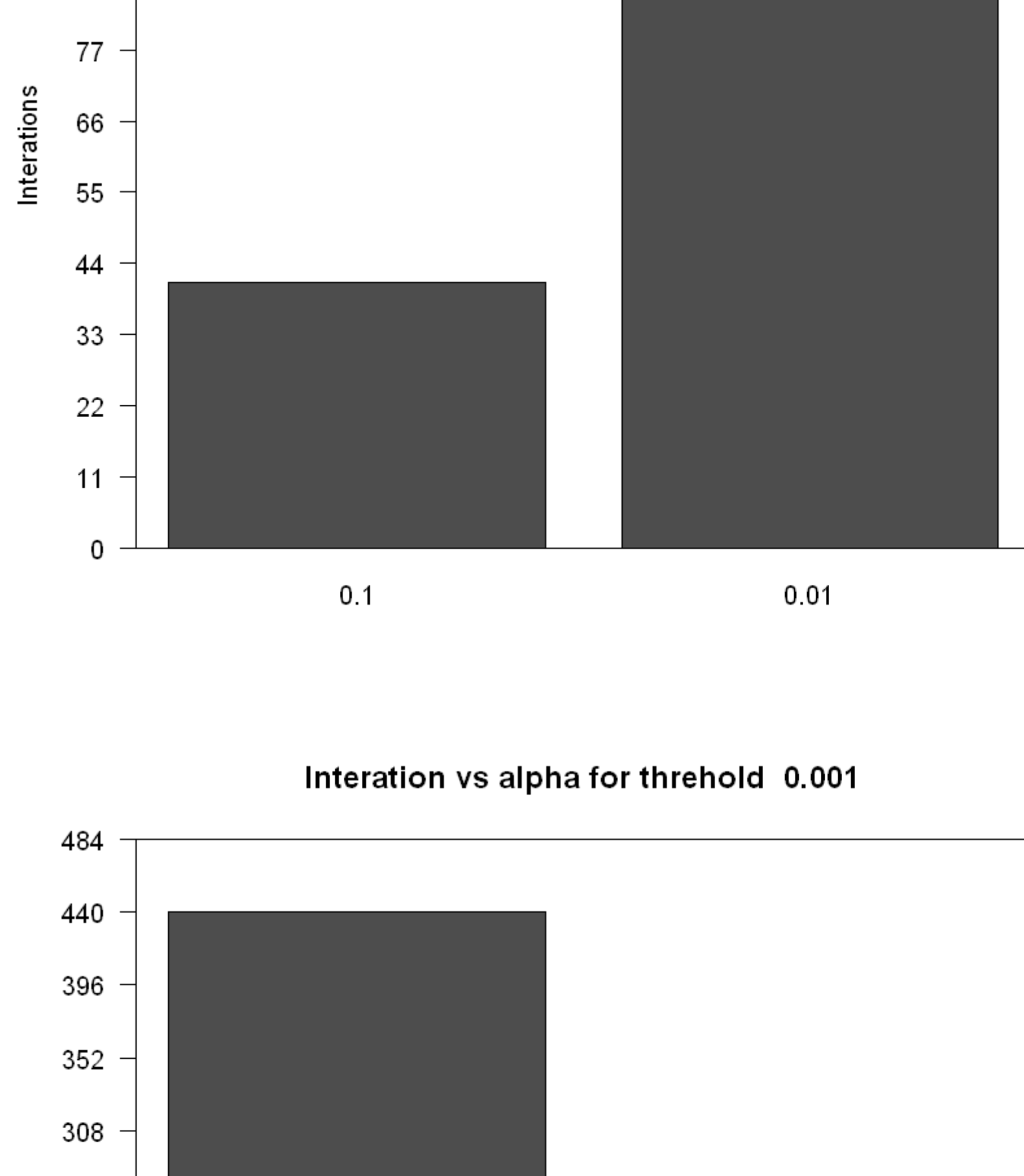
SGD for Adaptive Step Size

```
In [53]: options <- c(2,3)
for (opt in options){
  alpha <- c(1,0.1,0.01)
  thres <- 10**(2)
  Siterations <- list()
  Ssuccessalphs <- c()
  Sgradienttrend <- list()
  Sxsfull <- list()
  x <- matrix(rep(0,ncol(A)),nrow =ncol(A), byrow = TRUE)
  row = 1
  for (alpha in alphas) {
    iterationcount <- c()
    for (i in seq(1,5,1))
      {
        answer <- stochasticGradientDes(y,A,x,alpha,thres,maxiter,opt)
        if (answer$converged){
          iterationcount = append(iterationcount,answer$iteration)
        }
      }
    Ssuccessalphs <- append(Ssuccessalphs,alpha)
    Siterations[[row]] <- mean(iterationcount)
    row = row+1
  }
  Siterations <- do.call(cbind, Siterations)
  colnames(Siterations)<- successalphs
  barplot(Siterations, main=paste("Iteration vs alpha for option ", opt,sep = " "), names.arg = Ssuccessalphs)
  ylab = "Iterations", ylim = c(0,max(Siterations)*max(Siterations)/10), xlab = "", axes = FALSE)
  contour(xrange2,xrange3,z = $a$Values,levels = 80, main = "Level sets for x1 = 1.016019e-16", xlab = "", axes = FALSE)
  axis(2, at = ylab1, las = 1)
  box(1)
}
```

Interaction vs alpha for option 2



Interaction vs alpha for option 3



Observation on SGD

It is quite difficult to draw any conclusion for adaptive alpha for SGD as each iteration gives different result. So I ran 5 iteration for each and capture mean no. of iteration for each. It is observable that option 2 (alpha = a/i) gives answer in less number of iteration.

SGD for different threshold

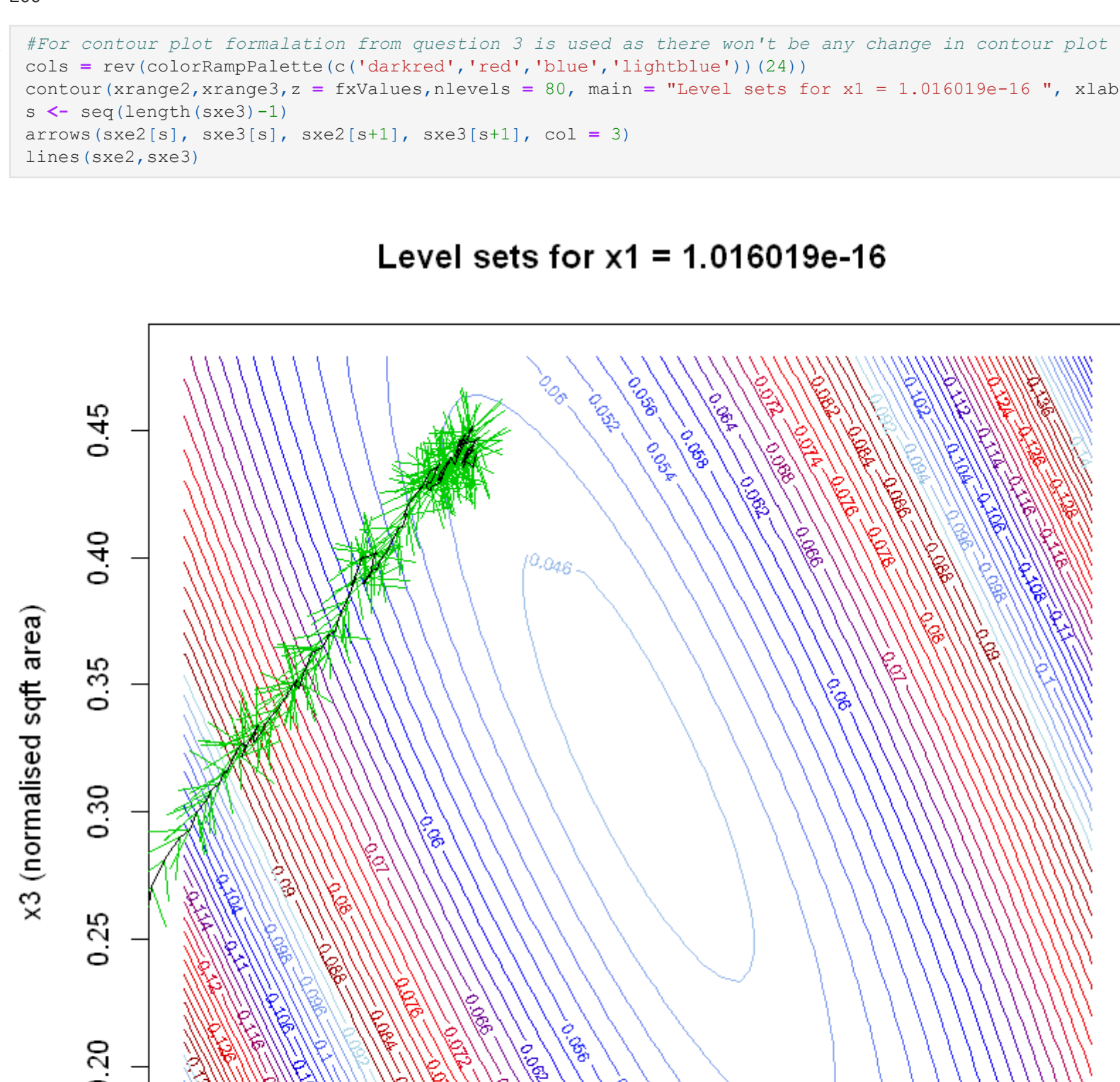
```
In [65]: thres <- c(10**(1),10**(2),10**(3))
for (thre in thres){
  maxiter <- 1000
  alpha <- c(0.1,0.01)
  A2 <- A
  y2 <- y
  Siterations <- list()
  Ssuccessalphs <- c()
  Sgradienttrend <- list()
  Sxsfull <- list()
  x <- matrix(rep(0,ncol(A)),nrow =ncol(A), byrow = TRUE)
  row = 1
  for (alpha in alphas) {
    answer <- stochasticGradientDes(y,A,x,alpha,thre,maxiter,1)
    if (answer$converged){
      Ssuccessalphs <- append(Ssuccessalphs,alpha)
      Sxsfull[[row]] <- answer$x
      Siterations[[row]] <- answer$iteration
      Sgradienttrend[[row]] <- answer$gradienttrend
      row = row+1
    }
  }
  Sxsfull <- do.call(cbind, Sxsfull)
  Siterations <- do.call(cbind, Siterations)

  colnames(Sxsfull)<- Ssuccessalphs
  colnames(Siterations)<- Ssuccessalphs
  modifiedx <- format(Sxsfull, scientific = FALSE)
  print(modifiedx)

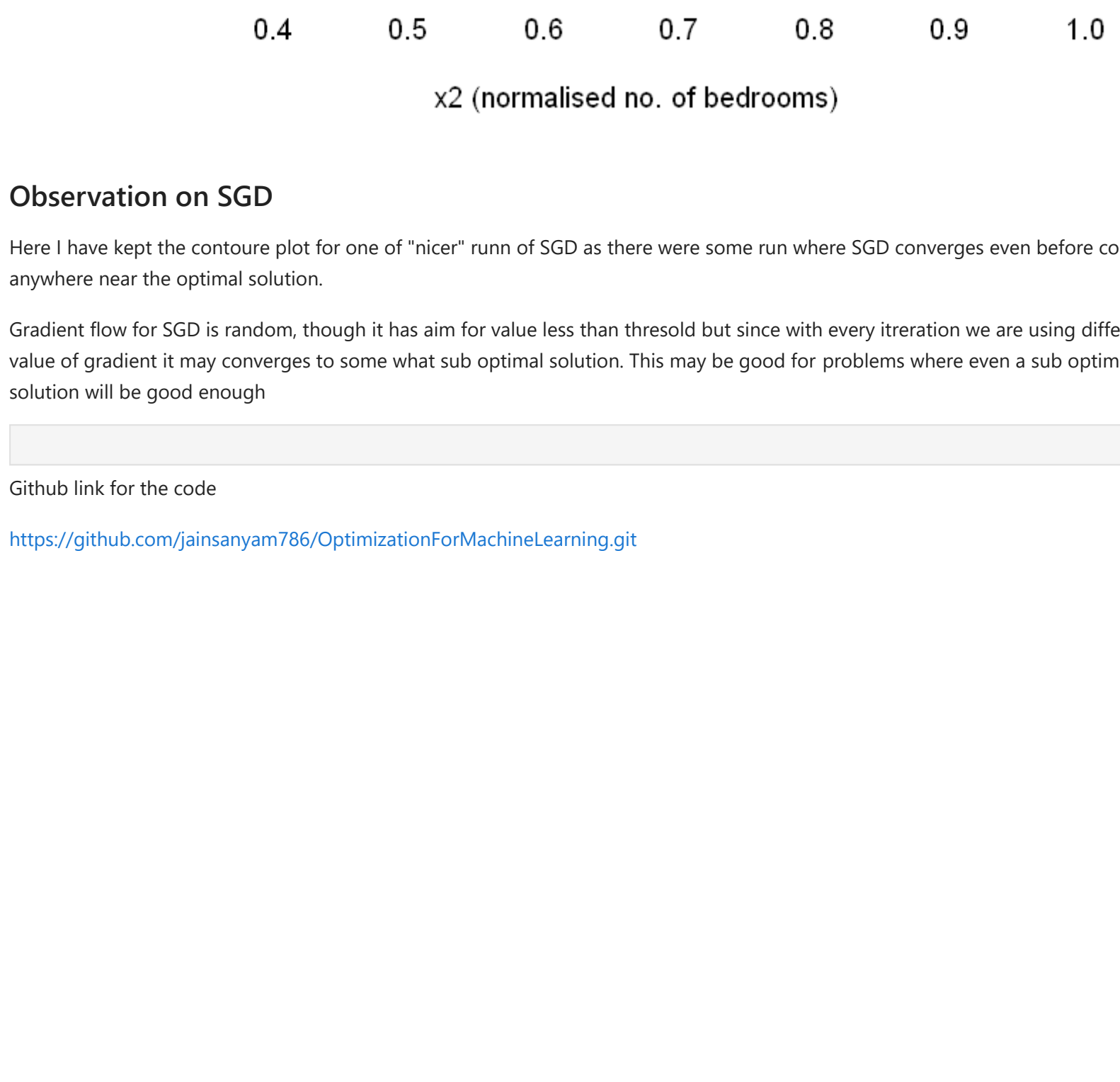
  barplot(Siterations, main=paste("Iteration vs alpha for threshold ", thre,sep = " "), names.arg = Ssuccessalphs)
  ylab = "Iterations", ylim = c(0,max(Siterations)*max(Siterations)/10), xlab = "", axes = FALSE)
  contour(xrange2,xrange3,z = $a$Values,levels = 80, main = "Level sets for x1 = 1.016019e-16", xlab = "", axes = FALSE)
  axis(2, at = ylab1, las = 1)
  box(1)
}
```

	0.1	0.01
[1,]	" 0.0917923" " -0.01031580"	
[2,]	" 0.48341586" " 0.28323165"	
[3,]	" 0.54253265" " 0.26233532"	
0.1	0.01	
[1,]	" -0.09717073" " -0.01137831"	
[2,]	" 0.19742232" " 0.47050362"	
[3,]	" 0.47486920" " 0.40409319"	

Interaction vs alpha for threshold 0.1



Interaction vs alpha for threshold 0.01



Interaction vs alpha for threshold 0.001



Observation on SGD

Here I have kept the contour plot for one of "nicer" runn of SGD as there were some run where SGD converges even before coming anywhere near the optimal solution.

Gradient flow for SGD is random, though it has aim for value less than threshold but since with every iteration we are using different value of gradient it may converges to some what sub optimal solution. This may be good for problems where even a sub optimal solution will be good enough

```
In [13]:
```

GitHub link for the code

<https://github.com/gajansanyam786/OptimizationForMachineLearning.git>

Level sets for x1 = 1.016019e-16

Observation on SGD

Here I have kept the contour plot for one of "nicer" runn of SGD as there were some run where SGD converges even before coming anywhere near the optimal solution.

Gradient flow for SGD is random, though it has aim for value less than threshold but since with every iteration we are using different value of gradient it may converges to some what sub optimal solution. This may be good for problems where even a sub optimal solution will be good enough

```
In [13]:
```