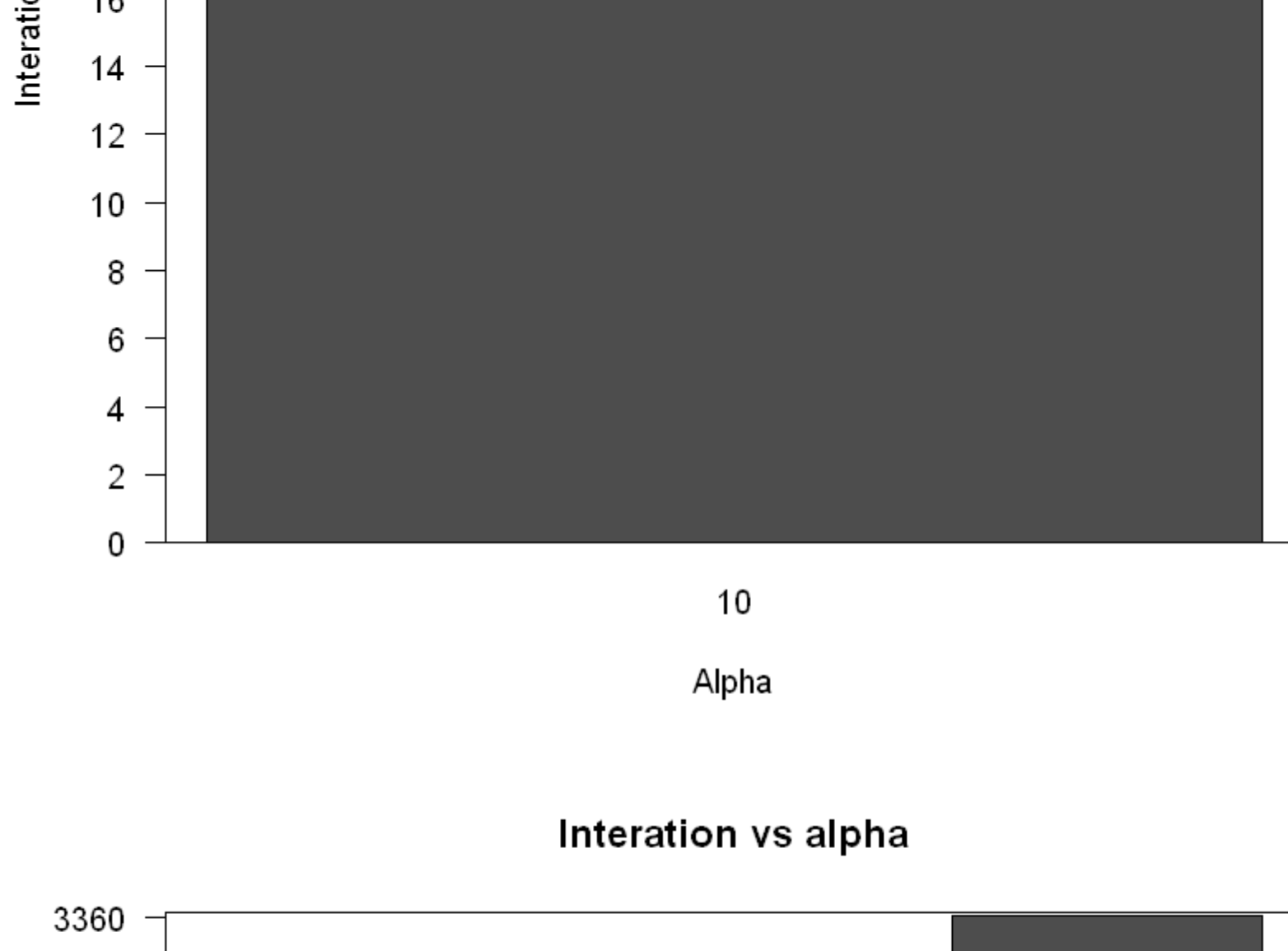
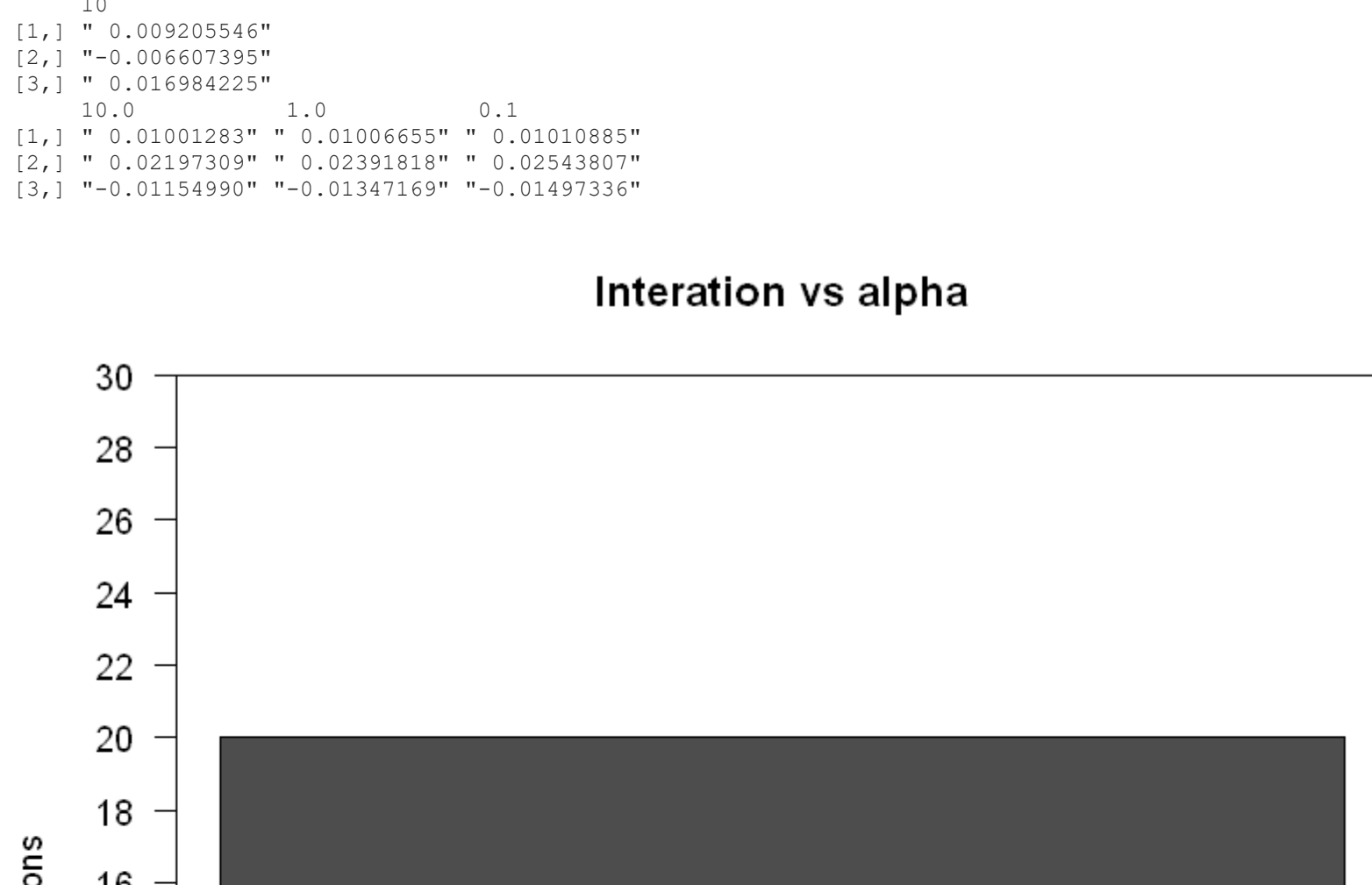






```
finalalpha <- do.call(cbind, finalalphs)
colnames(sx) <- successalphs
colnames(finalalpha) <- successalphs
as <- apply(sx,2,function(x) sol ~ x)
if <- format(as, scientific = FALSE)
print(as)
barplot(iterations, main="Iteration vs alpha", names.arg = successalphs,
ylab = "Iterations", ylim = c(0,max(iterations)+10), xlab = "Alpha", axes = FALSE)
ylabel <- seq(0, max(iterations)+10, by = max(iterations)/10)
axis(2, at = ylabel, las = 1)
box()
}
}
```

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



### 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 [35]:
stochasticGradientDes <- function(y1,A1,x1,advalpha,thrs, maxi, option){
  converged <- FALSE
  i=1
  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 <- deltaf(y1,A1,x1)
      if(is.infinite(deltaf)|| is.nan(norm(deltaf, type = "2"))){
        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, type = "2"))
      converged <- (norm(deltaf, type = "2") <= thrs)
      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
alphas <- c(1,0.1,0.01)
thrs <- 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,thrs,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
modifiedxs <- format(xsfull, scientific = FALSE)
modifiedits <- format(iterations, scientific = FALSE)
iterations
successalphs
```

	0.10	0.01
0.026916195	0.004764189	
0.538005611	0.521706534	
0.415810632	0.405792914	

	0.10	0.01
64	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
thrs <- 10**(2)
maxiter <- 1000
A2 <- A
y2 <- y
Siterations <- list()
SgradientTrend <- list()
Sxsfull <- list()
row = 1
x <- matrix(rep(0,ncol(A)),nrow = ncol(A), byrow = TRUE)
for (i in seq(1,5,1)) {
  answer <- stochasticGradientDes(y,A,x,alpha,thrs,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)
modifiedxs <- format(Sxsfull, scientific = FALSE)
modifiedits <- format(Siterations, scientific = FALSE)
Siterations
```

	1	2	3	4	5
-0.017651936	0.006640093	-0.075477447	-0.005380836	0.044466205	
0.810941528	0.486160311	0.707340578	0.650536353	0.556324112	
0.462324871	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),
col = "red", yaxp=c(1,1,10), ylab = "Iterations", xlab = "", axes = FALSE)
axis(2, at = ylabel, las = 1)
box()
```



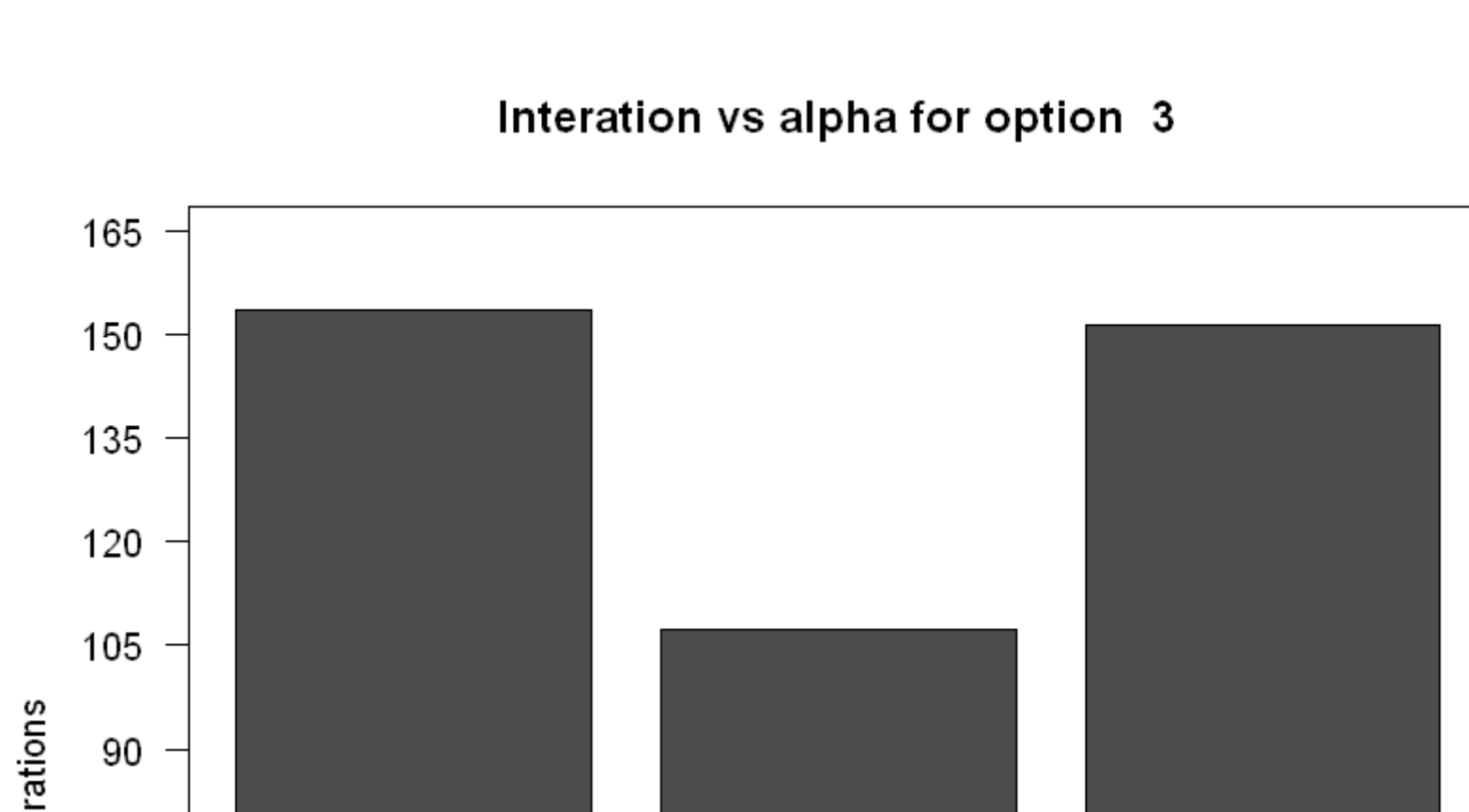
```
In [40]:
as <- apply(Sxsfull,2,function(x) sol ~ x)
as <- format(as, scientific = FALSE)
as
```

	1	2	3	4	5
0.017651936	-0.006640093	-0.075477447	-0.005380836	-0.044466205	
-0.140861953	0.189391626	-0.037261003	0.019543222	0.113755462	
-0.143060974	-0.06857222	-0.011444355	0.013938642	-0.029635588	

### SGD for Adaptive Step Size

```
In [53]:
options <- c(2,3)
for (opt in options){
  alphas <- c(1,0.1,0.01)
  thrs <- 10**(2)
  Siterations <- list()
  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,thrs,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) <- Ssuccessalphs
  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)
ylabel <- seq(0, max(Siterations)*max(Siterations)/10, by = max(Siterations)/10)
axis(2, at = ylabel, las = 1)
box()
}
```

### Iteration vs alpha for option 2



### Iteration 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]:
thrs <- c(10**(1),10**(2),10**(3))
for (thre in thrs){
  maxiter <- 1000
  alphas <- 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
  modifiedxs <- format(Sxsfull, scientific = FALSE)
  print(modifiedxs)
  barplot(Siterations, main=paste("Iteration vs alpha for threshold ", thre, sep = " "), names.arg = Ssuccessalphs,
col = "red", yaxp=c(1,1,10), ylab = "Iterations", xlab = "", axes = FALSE)
ylabel <- seq(0, max(Siterations)*max(Siterations)/10, by = max(Siterations)/10)
axis(2, at = ylabel, las = 1)
box()
}
```

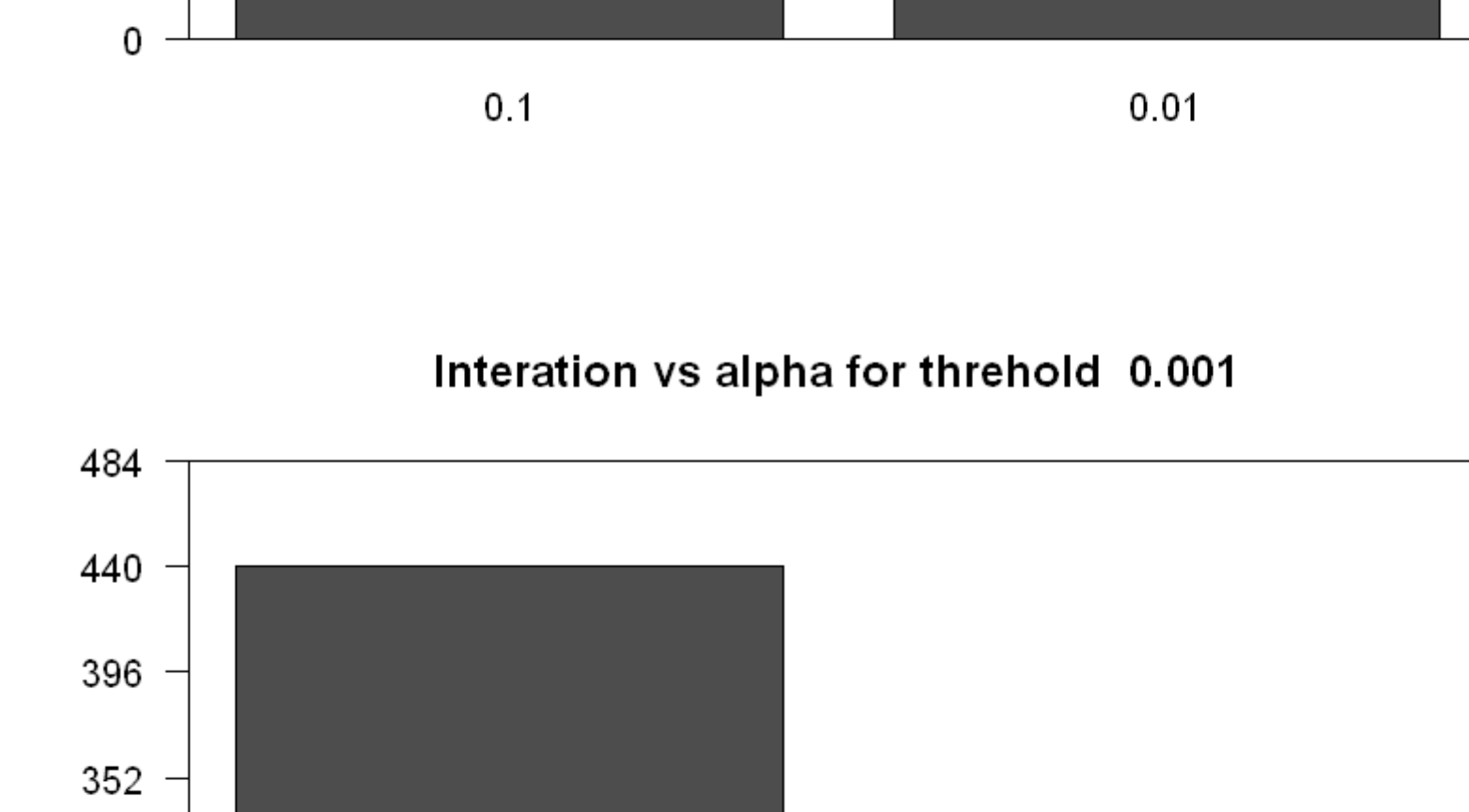
	0.1	0.01
[1,]	"0.09177923"	"0.01031580"
[2,]	"0.48941586"	"0.23232165"
[3,]	"0.54253265"	"0.26235312"

	0.1	0.01
[1,]	"0.09177923"	"0.01031580"
[2,]	"0.48941586"	"0.23232165"
[3,]	"0.54253265"	"0.26235312"

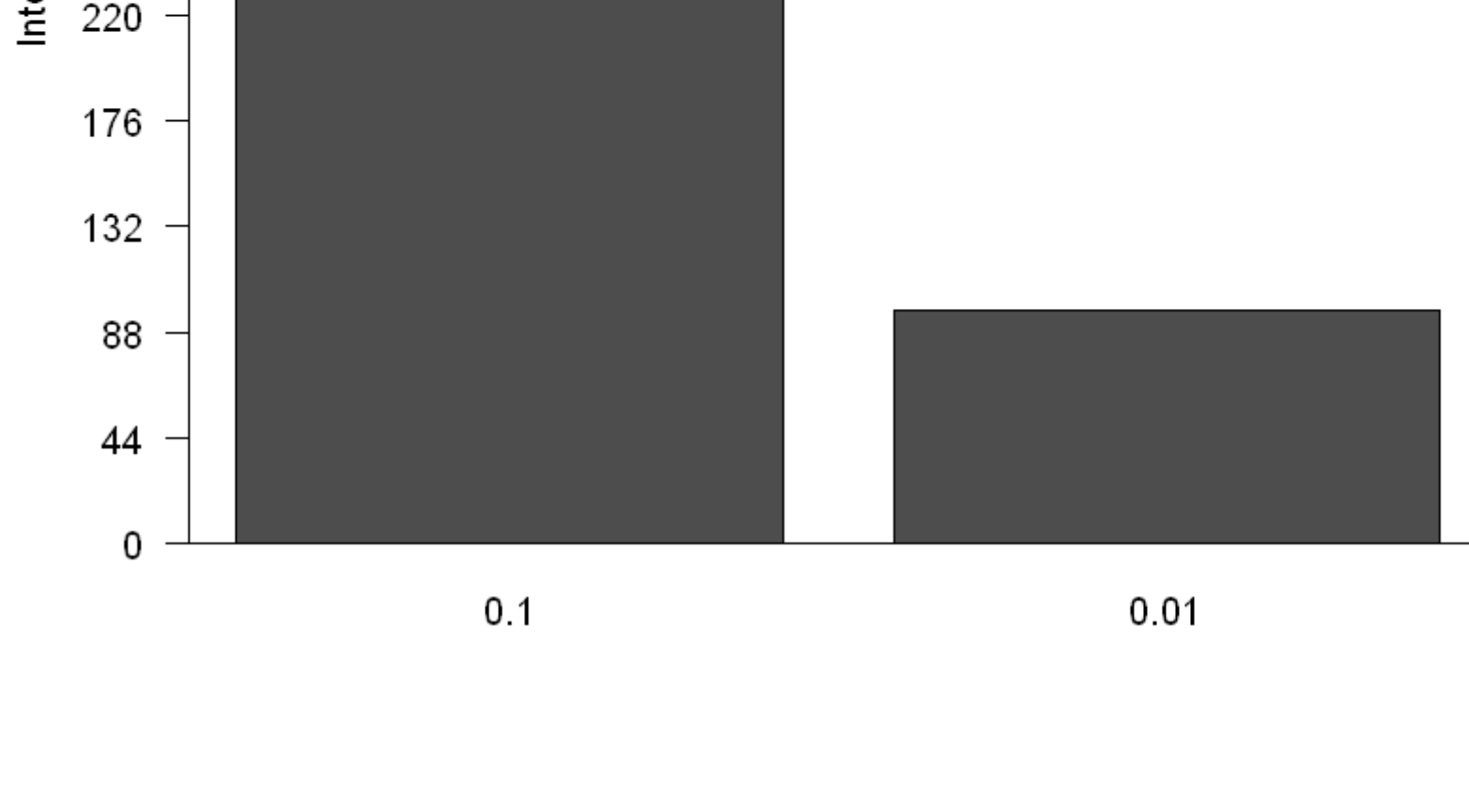
### Iteration vs alpha for threshold 0.1



### Iteration vs alpha for threshold 0.01



### Iteration vs alpha for threshold 0.001



### Observation on SGD

as threshold decreases no. of iteration increases which is reasonable as are aiming for more accurate solution

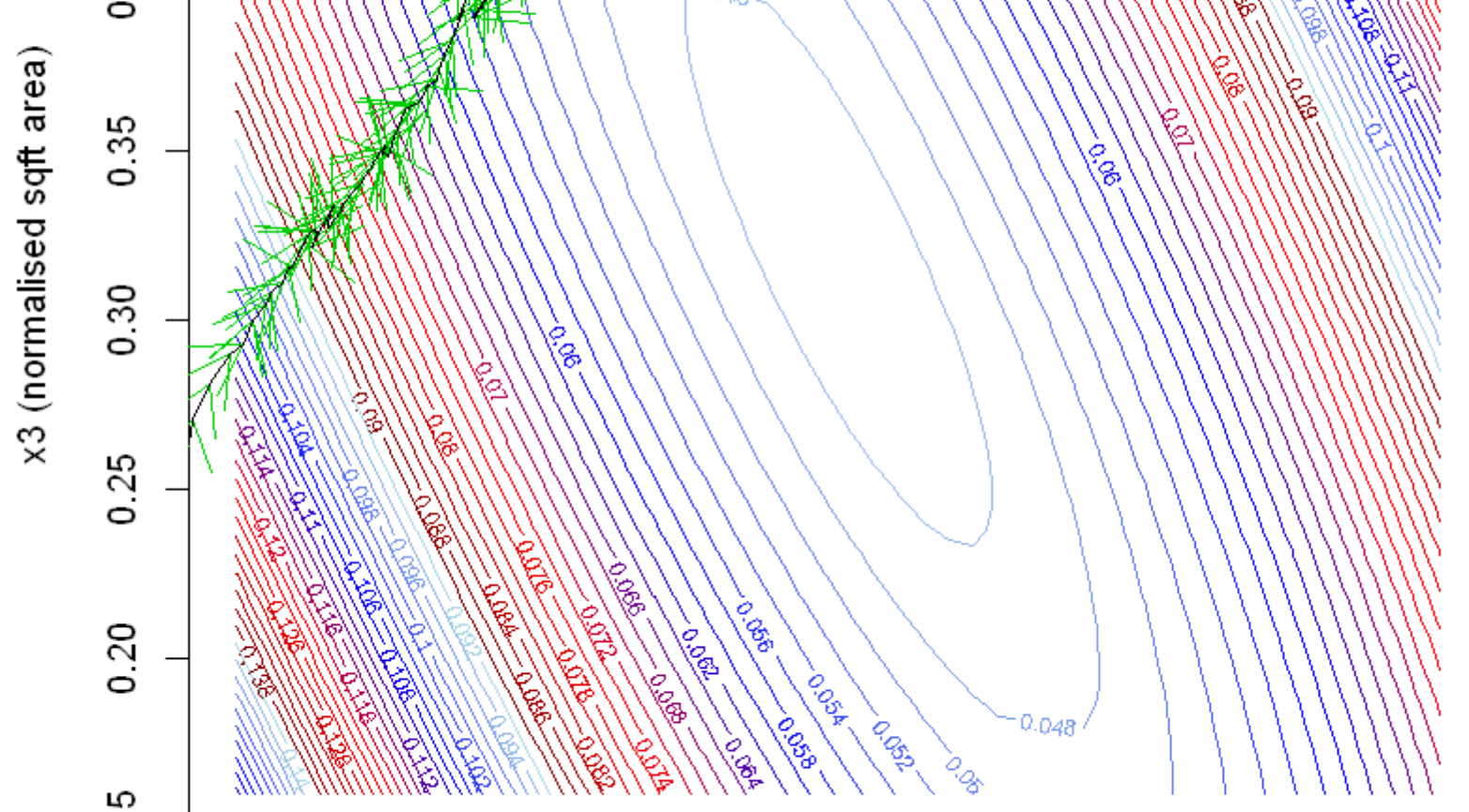
### Contour Plot for the SGD

```
In [129]:
# storing xs for plotting stochastic gradient flow
sgdX <- list()
thrs <- c(10**(2),10**(3))
converged <- FALSE
maxi <- 1000
alphas <- 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)
{
  iterationcount <- c()
  for (i in seq(1,5,1))
  {
    answer <- stochasticGradientDes(y,A,x,alpha,thrs,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
  modifiedxs <- format(Sxsfull, scientific = FALSE)
  print(modifiedxs)
  barplot(Siterations, main=paste("Iteration vs alpha for threshold ", thre, sep = " "), names.arg = Ssuccessalphs,
col = "red", yaxp=c(1,1,10), ylab = "Iterations", xlab = "", axes = FALSE)
ylabel <- seq(0, max(Siterations)*max(Siterations)/10, by = max(Siterations)/10)
axis(2, at = ylabel, las = 1)
box()
}
```

	0.1	0.01
[1,]	"0.09177923"	"0.01031580"
[2,]	"0.48941586"	"0.23232165"
[3,]	"0.54253265"	"0.26235312"

	0.1	0.01
[1,]	"0.09177923"	"0.01031580"
[2,]	"0.48941586"	"0.23232165"
[3,]	"0.54253265"	"0.26235312"

### Level sets for x1 = 1.016019e-16



### Observation on SGD

Here I have kept the contour plot for one of "nicer" run 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 [ ]:
# Github link for the code
https://github.com/jainanyam796/OptimizationForMachineLearning.git
```

```
In [ ]:
# For the contour plot formation from question 3 is used as there won't be any change in contour plot for prob1
```