

HW7_Wei_Yanran

Problem 2: Sums of Squares

Four methods are used in calculating sums of squares. These methods are loop, vector operations, dopar and parSapply methods. The time for running these methods are displayed in the table below. The dataset is $1e+06$ data selected from normal distribution with mean of 1 and standard deviation of 1. I shrink the dataset from $1e+7$ to $1e+6$ because my computer will get stuck when running parSapply method.

Table 1: Time for Different Methods

	user.self	sys.self	elapsed	Sum of Squares
Loop_Method	0.15	0	0.16	1002631
Vector_Method	0.08	0	0.07	1002631
Dopar_Method	0.11	0.05	0.48	1002631
ParSapply_Method	1.72	0.51	3.97	1002631

From the table, we can see that vector operations method perform the best, then dopar method and loop method. Parsapply method perform worst. Sum of squares value is same for all four methods. The cluster I set in Dopar and ParSapply method is 2.

Problem 3

Table 2: Running time for GD and Linear Methods

	user.self	sys.self	elapsed	user.child	sys.child
GD_time	0.28	0	1.09	NA	NA
Linear_Time	0	0	0	NA	NA

Table 3: Coefficients of GD Method

	Tolerance	theta1	theta2
result.1	1	1	2
result.2	0.1	1	2
result.3	0.01	1	2
result.4	0.001	1	2
result.5	1e-04	0.9999	2
result.6	1e-05	0.9993	1.998
result.7	1e-06	0.9743	2.001
result.8	1e-07	0.97	2.001
result.9	1e-08	0.9696	2.002
result.10	1e-09	0.9696	2.002

Table 4: Coefficients of Linear Regression Methods

	X1	X2
lm_h.coefficients	0.9696	2.002

In the parallelize method, cluster 2 was used. The parameter I made adjustment on is tolerance. As shown in the table *Coefficients of GD Method*, as tolerance becoming smaller, the coefficients becomes more accurate. When tolerance is less or equal than 1e-08, we get the same coefficients with that calculated by linear regression. But gradient descent method seems take more time than linear regression method.

problem 4

a

Dopar function was used to bootstrap 10000 samples from Y and X. Each sample contains 10 pairwise data and is chosen from the 200 pairwise dataset of X and Y. Linear regression model was built to test the relationship between x and y and to get the coefficients of the regression line. Cluster 2 was utilized.

b

Table 5: Overview of Beta

	x1	x2	x3	x4
result.1	-0.3042	1.915	3.021	0.006305
result.2	2.921	1.841	3.141	0.09711
result.3	1.292	2.001	2.871	0.01526
result.4	2.651	1.935	3.056	0.08267
result.5	1.99	2.064	3.268	-0.04152
result.6	1.575	2.006	3.231	1.055

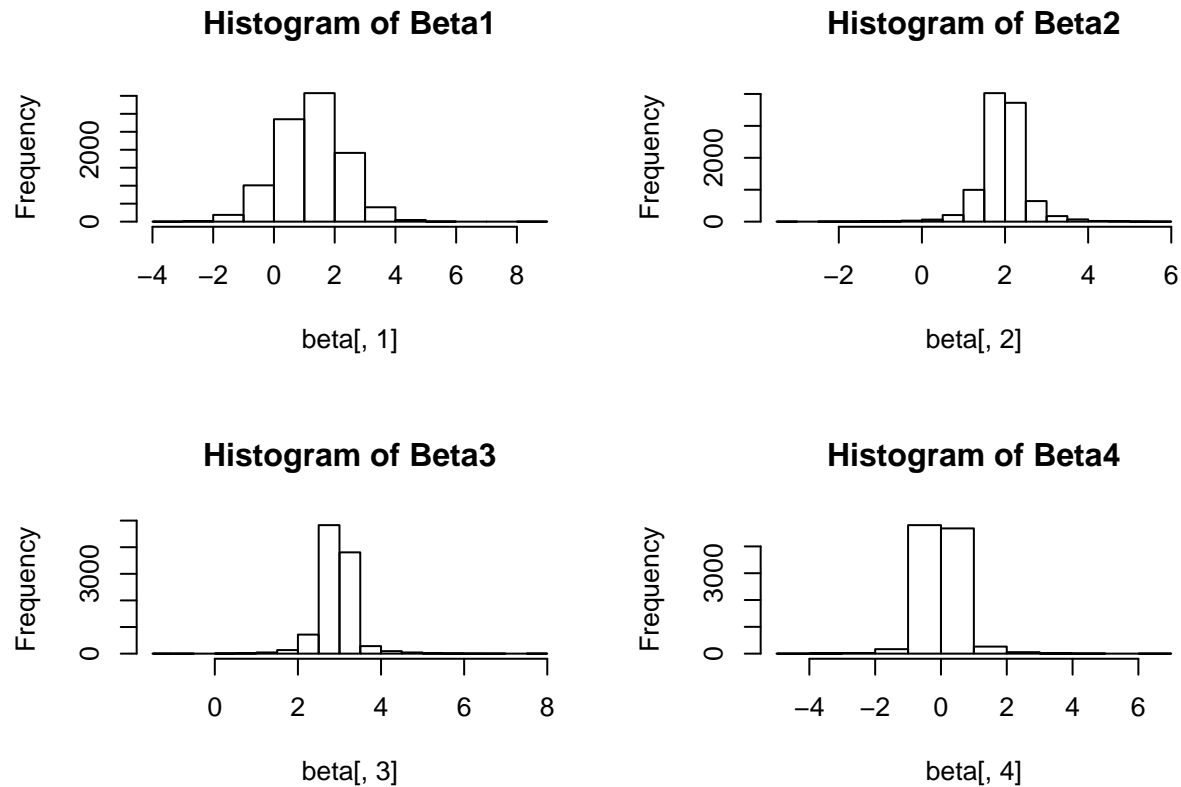
Table 6: Summary of Beta

x1	x2	x3	x4
Min. :-3.1588	Min. :-3.106	Min. :-1.347	Min. :-4.960764
1st Qu.: 0.5407	1st Qu.: 1.729	1st Qu.: 2.783	1st Qu.: -0.164395
Median : 1.2577	Median : 1.978	Median : 2.973	Median : 0.000068
Mean : 1.2406	Mean : 1.950	Mean : 2.931	Mean : 0.012636
3rd Qu.: 1.9607	3rd Qu.: 2.167	3rd Qu.: 3.055	3rd Qu.: 0.107264
Max. : 8.0750	Max. : 5.514	Max. : 7.877	Max. : 6.476911

As shown in the table *Overview of Beta*, we get a matrix of beta with 10000 rows of results. Each row contains the coefficients (betas) value of the regression between x and y.

In the table *Summary of Beta*, we can see the statistics summary of these betas. Take X1 as example. The column X1 shows the summary for beta1. The minimum value of beta1 is -3.973 while the maximum value is 13.715. The median value is 1.25 and the mean value is 1.239.

c



As shown in the histogram, the values of Beta1 is mainly less than 5. Beta2 is mainly located between 0 and 5. Beta3 is mainly located between 2 and 4. Beta4 is mainly located between -1 and 1.

```
set.seed(12345)
y <- rnorm(n = 1e+06, mean = 1, sd = 1)
library(pander)

##### Loop Method#####
t_lo <-
  system.time({
    mean_y <- mean(y)
    sum_difa <- 0
    for (i in 1:length(y)){
      sum_difa <- sum_difa + (y[i] - mean_y)^2
    }
  })
Loop_Method <- c(t_lo[1:3], sum_difa)

##### Vector Method#####
t_vec<-
  system.time({
    mean_y <- mean(y)
    diff_y <- y - mean_y
    sum_difb <- as.numeric(t(diff_y) %*% diff_y)
  })
```

```

Vector_Method <- c(t_vec[1:3], sum_difb)

##### Dopar Method#####
library(foreach)
library(doParallel)
cl <- makeCluster(2)
registerDoParallel(cl)
t_dop <-
  system.time({
    mean_y <- mean(y)
    sum_difc <- foreach(1) %dopar% sum((y - mean_y)^2)
  })
stopCluster(cl)
Dopar_Method <- c(t_dop[1:3], sum_difc)

##### ParSapply Method#####
mean_y <- mean(y)
t_parsap <-
  system.time({
    c2 <- makeCluster(2)
    clusterExport(c2, "mean_y")
    sum_difd <- sum(parSapply(c2, y, function(y) (y-mean_y)^2))
  })
stopCluster(c2)
ParSapply_Method <- c(t_parsap[1:3], sum_difd)

##### Summary#####
time_1 <- rbind(Loop_Method, Vector_Method, Dopar_Method, ParSapply_Method)
colnames(time_1)[4] <- c("Sum of Squares")
pander(time_1, caption = "Time for Different Methods")

##### Create data point#####
set.seed(1256)
theta <- as.matrix(c(1,2),nrow=2)
X <- cbind(1,rep(1:10,10))
h <- X%*%theta+rnorm(100,0,0.2)
m <- length(h)
cl <- makeCluster(2)
registerDoParallel(cl)

##### Gradient descent#####
GD_time <-
system.time({
dif1 <- 1
dif2 <- 1
alpha <- 0.001
ttheta <-
  foreach (i = 1:10, .combine = "cbind") %dopar% {
    tolerance <- 10^(1-i)
    while(dif1 > tolerance | dif2 > tolerance){
      grad_0 <- (1/m) * sum(((X%*%theta) - h))
      grad_1 <- (1/m) * t(X[, 2]) %*% ((X%*%theta) - h)
      theta_0 <- theta[1] - alpha * grad_0

```

```

theta_1 <- theta[2] - alpha * grad_1
dif1 <- abs(theta_0 - theta[1])
dif2 <- abs(theta_1 - theta[2])
theta[1] <- theta_0
theta[2] <- theta_1
}

return(c(tolerance, theta[1], theta[2]))
}
})
theta <- t(ttheta)
colnames(theta) <- c("Tolerance", "theta1", "theta2")
stopCluster(cl)

##### Linear regression#####
Linear_Time <-
system.time({
  lm_h <- lm(h~0+X)
Linear_coef <- t(data.frame(lm_h$coefficients))
})

##### Summary#####
pander(rbind(GD_time, Linear_Time), caption = "Running time for GD and Linear Methods")
pander(theta, caption = "Coefficients of GD Method")
pander(Linear_coef, caption = "Coefficients of Linear Regression Methods")
set.seed(1267)
n <- 200
X <- 1/cbind(1, rt(n, df = 1), rt(n, df = 1), rt(n, df = 1))
beta <- c(1, 2, 3, 0)
Y <- X %*% beta + rnorm(100, sd = 3)
data <- cbind(X, Y)
Y <- data[, 5]
X1 <- data[, 1]
X2 <- data[, 2]
X3 <- data[, 3]
X4 <- data[, 4]
cl <- makeCluster(2)
registerDoParallel(cl)
beta <- matrix(NA, nrow = 10000, ncol = 4)
colnames(beta) <- c("Intercept", "X2", "X3", "X4")
tbeta <-
  foreach (b = 1:10000, .combine = data.frame) %dopar% {
    i <- sample(1:200, size = 10, replace = TRUE)
    y <- Y[i]
    x1 <- X1[i]
    x2 <- X2[i]
    x3 <- X3[i]
    x4 <- X4[i]
    lm(y~0+x1+x2+x3+x4)$coef
  }
beta <- t(tbeta)
stopCluster(cl)
library(pander)

```

```
pander(head(beta), caption = "Overview of Beta")
pander(summary(beta), caption = "Summary of Beta")
par(mfrow = c(2, 2))
hist(beta[, 1], main = "Histogram of Beta1")
hist(beta[, 2], main = "Histogram of Beta2")
hist(beta[, 3], main = "Histogram of Beta3")
hist(beta[, 4], main = "Histogram of Beta4")
```