

Question # 1

Below is the R code that applied K-means algorithm on given image for $K = 3, 5$, and 10 . Figure. 1, Figure. 2, and Figure. 3 shows the results of k-means for foregoing values of K .

```
1 fileName <- "DM2016.txt"
2 conn <- file(fileName, open="r")
3 linn <- readLines(conn)
4
5 rgb = matrix(NA, 1018 * 1554, 5)
6 index = 1
7
8 for (i in c(2:1019))
9 {
10   y = 1
11   row = strsplit(linn[i], " ")
12   row = unlist(row)
13
14   for (pcol in seq(2, 4663, 3))
15   {
16     rgb[index, 1] = i - 1
17     rgb[index, 2] = y
18     rgb[index, 3] = as.numeric(row[pcol]) / 255
19     rgb[index, 4] = as.numeric(row[pcol + 1]) / 255
20     rgb[index, 5] = as.numeric(row[pcol + 2]) / 255
21     y = y + 1
22     index = index + 1
23   }
24 }
25
26 close(conn)
27
28 rgb_data = data.frame(x = as.vector(rgb[, 1]), y = as.vector(rgb[, 2]), r.value = as.
29   vector(rgb[, 3]), g.value = as.vector(rgb[, 4]), b.value = as.vector(rgb[, 5]))
30 #head(rgb_data)
31
32 kColors <- 10
33 kMeans <- kmeans(rgb_data[, c(3, 4, 5)],
34   centers = kColors)
35
36 clusterColour <- rgb(kMeans$centers[kMeans$cluster, ])
37
38 plot(y ~ x, data = rgb_data, main = "k-mean cluster analysis",
39   col = clusterColour, asp = 1, pch = ".",
40   axes = FALSE, ylab = "",
41   xlab = "k-means cluster analysis of 10 colours")
```

Conclusion : K-means clusters the colours in an image that's why it can't recover the shuffled image.

k-mean cluster analysis



k-means cluster analysis of 3 colours

Figure 1: K-mean clustering with $k=3$

k-mean cluster analysis



k-means cluster analysis of 5 colours

Figure 2: K-mean clustering with $k=5$

k-mean cluster analysis



k-means cluster analysis of 10 colours

Figure 3: K-mean clustering with k=10

Question # 2

Following code applied dimensionality reduction technique "PCA" on given image using function "prcomp" in R.

```
1 pca = prcomp(img[,1:ncol(img)], scale. = FALSE, retx = TRUE, center = FALSE)
```

Below is the R code that shows the commulative summary of variance captured by principal components in Figure. 4. We can see, the first 3 principal components captures more than 90% of the varaince.

```

1 plot((cumsum(pca$sdev^2)/sum(pca$sdev^2))*100,type="b",pch=4,col="blue",ylab="
    cumulative variance captured",xlab="principal components",main="Cumulative
    Variance (%) Captured \n by Principal Components")

```

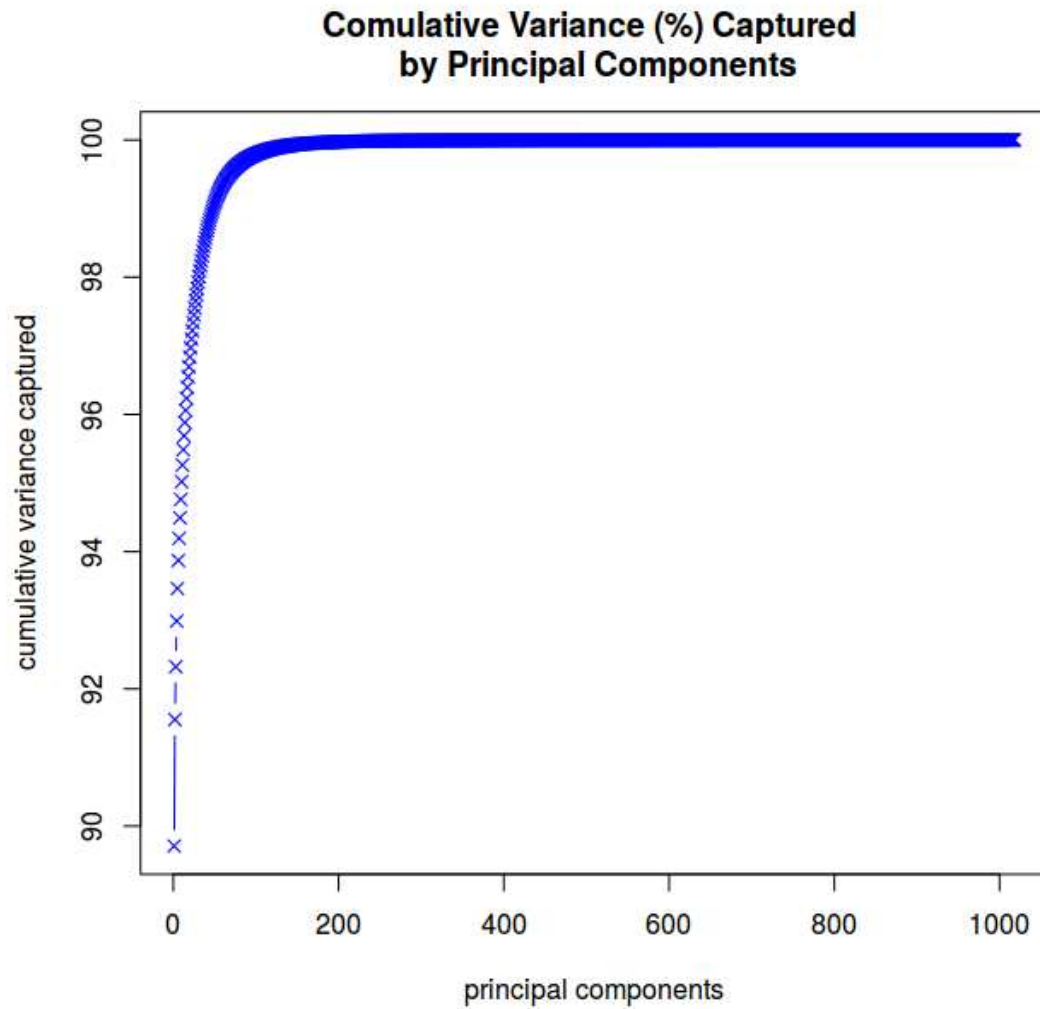


Figure 4: Cumulative Variance (%) captured by principal components

Figure. 5 illustrates the 3D plot of first 3 principal components.

```

1 library(scatterplot3d)
2 with(mtcars, {
3   scatterplot3d(pca$x[,1],
4                 pca$x[,2],

```

```

5     pca$x[,3],
6     color = "blue",
7     xlab = "pca 1",
8     ylab = "pca 2",
9     zlab = "pca 3",
10    main="3-D Scatterplot of PCA 1, PCA 2, PCA 3")
11 })

```

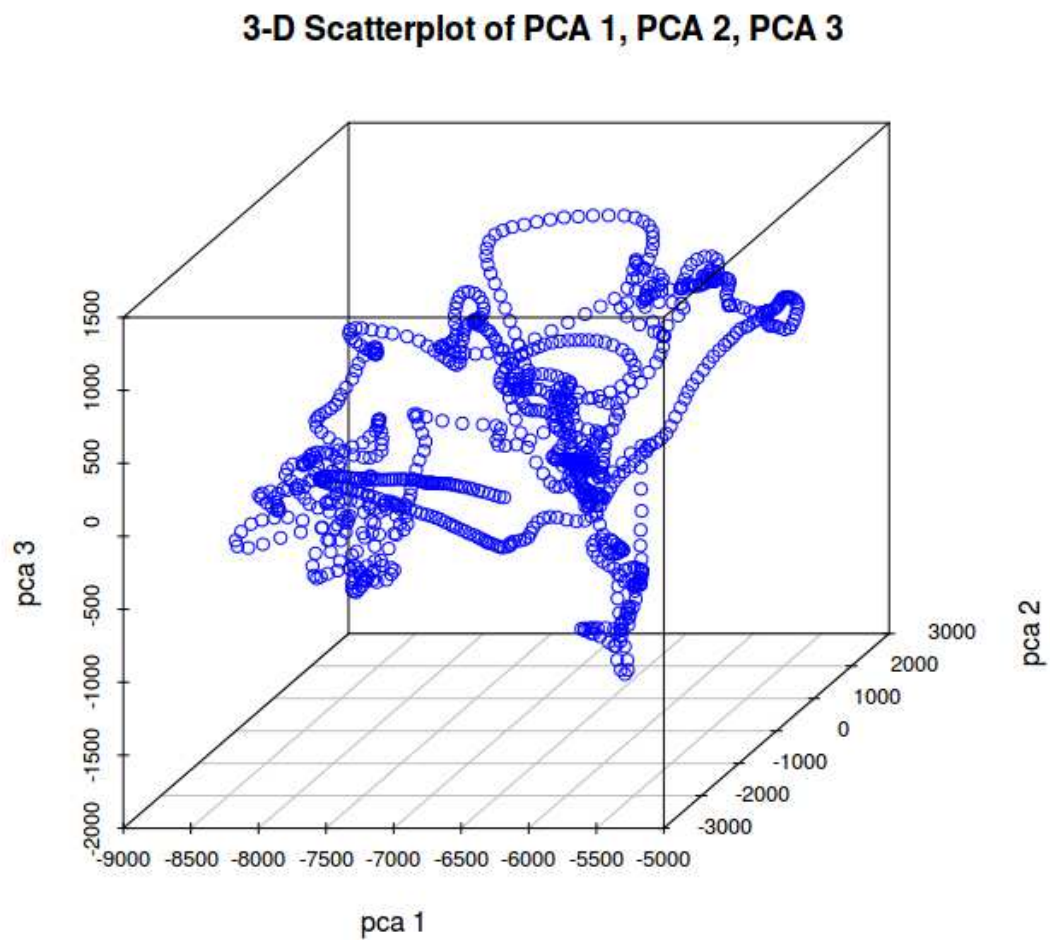


Figure 5: 3D scatter plot of PCA-1, PCA-2, and PCA-3

Question # 3

Following is the R code to compute pivot table summary using “dplyr” package in R.

```
1 library(dplyr)
2 library(RColorBrewer)
3 library(gplots)
4
5 data = read.csv("extract_medium.csv", sep = ";")
6
7 data$Marriage = factor(data$Marriage, level= c(1,2,3,4,5), labels=c("Married", "
  Widowed", "Divorced", "Seperated", "Never married (includes under 15 years)"))
8 education_labels = c("Not in universe (Under 2 years)", "No schooling completed", "
  Nursery school to 4th grade", "5th grade or 6th grade", "7th grade or 8th grade", "
  9th grade", "10th grade", "11th grade", "12th grade, no diploma")
9 education_labels = c(education_labels, c("High school graduate", "Some college, but
  less than 1 year", "One or more years of college, no degree", "Associate degree", "
  Bachelor's degree"))
10 education_labels = c(education_labels, c("Master's degree", "Professional degree", "
  Doctorate degree"))
11
12 data$Education = factor(data$Education, levels=c
  (0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16), labels = education_labels)
13 data$Sex= factor(data$Sex, level=c(1,2), labels = c("Male", "Female"))
14
15 group= group_by(data, Sex, Education)
16 sum_earnings = summarise(group,
17   avg_earnings = mean(Earnings))
18 summ_t = dcast(sum_earnings, Education ~ Sex, value.var = "avg_earnings")
19 row.names(summ_t) <- summ_t$Education
20 summ_t <- summ_t[, 2:3]
21
22 dt_matrix <- data.matrix(summ_t)
23
24 heatmap(dt_matrix, Rowv=NA, Colv=NA, col = brewer.pal(8, "Blues"), scale="column",
  margins=c(5,10), cexRow = .75, cexCol = .90)
```

Below is the pivot table summary of people's earning based on gender and education level.

Education	Male	Female
Not in universe(Under 2 years)	0.0000	0.0000
No schooling completed	3504.7847	947.7903
Nursery school to 4th grade	690.8592	387.3430
5th grade or 6th grade	5083.6720	2771.0552
7th grade or 8th grade	4073.9616	2256.1374
9th grade	9596.4985	3326.6576
10th grade	11185.8485	4281.3388
11th grade	11167.6389	5003.1646
12th grade,no diploma	19404.3564	9200.8858
High school graduate	24012.2758	11515.8326
Some college, but less than 1 year	28201.2106	16305.0455
One or more years of college no degree	28488.3473	16949.3175
Associate degree	35081.0018	21991.4803
Bachelor's degree	53294.2643	33193.7135
Master's degree	70755.1736	44412.2318
Professional degree	94245.2049	46821.9613
Doctorate degree	61467.6768	41476.0656

Heat of the above table is shown in Figure. 6

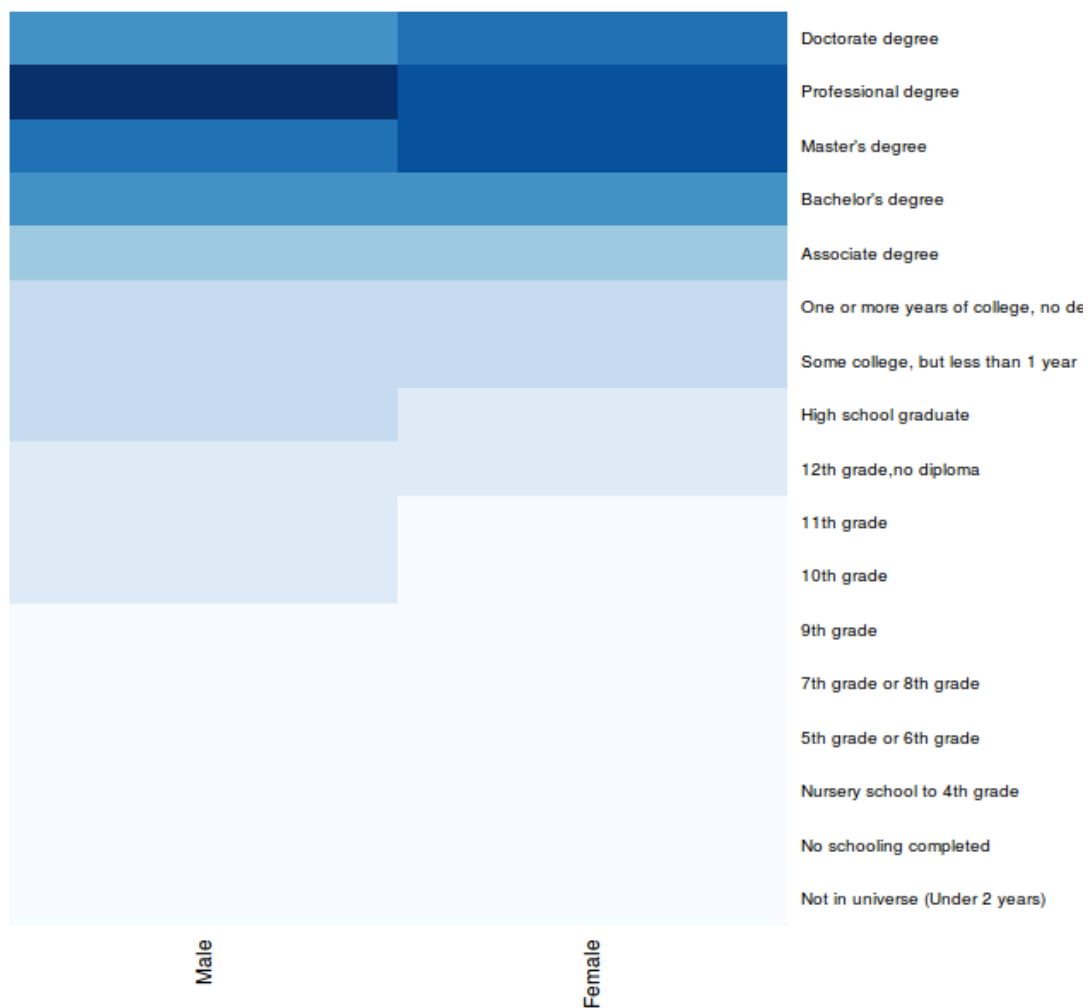


Figure 6: peoples' Earnings based on gender and education level

The message is clear and upsetting :- (, professional and master's degree holders are earning better than people with doctoral degree :-)

Question # 4

Figure. 7 explored the earnings based on marital status and education level. Prosperity formula, earn a profession degree and stay seperated. Another advice is never go for a doctoral degree :-)

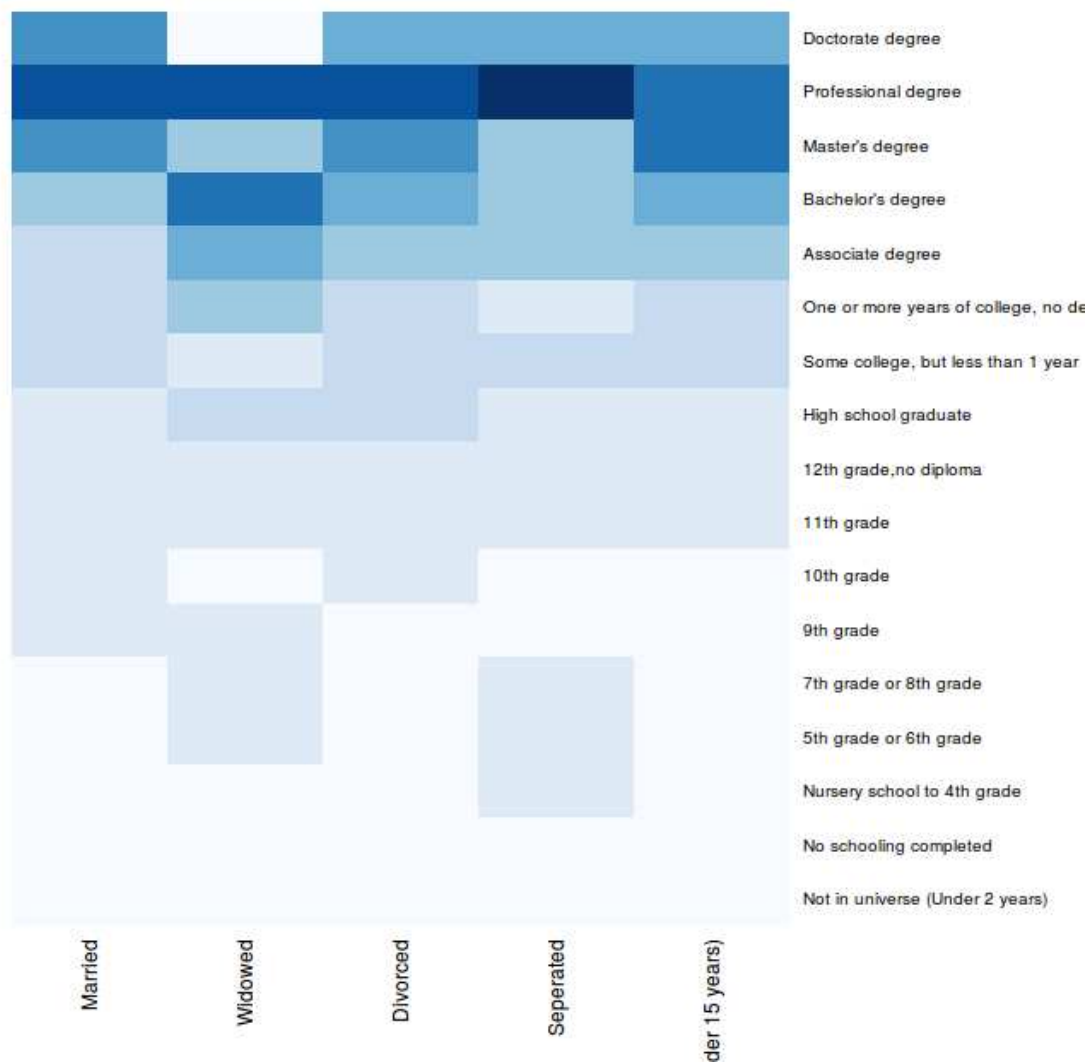


Figure 7: peoples' Earnings based on gender and education level

Figure. 8 looked into the people's earning based on gender and marital status. The male with married status are earning more than single or divorced.

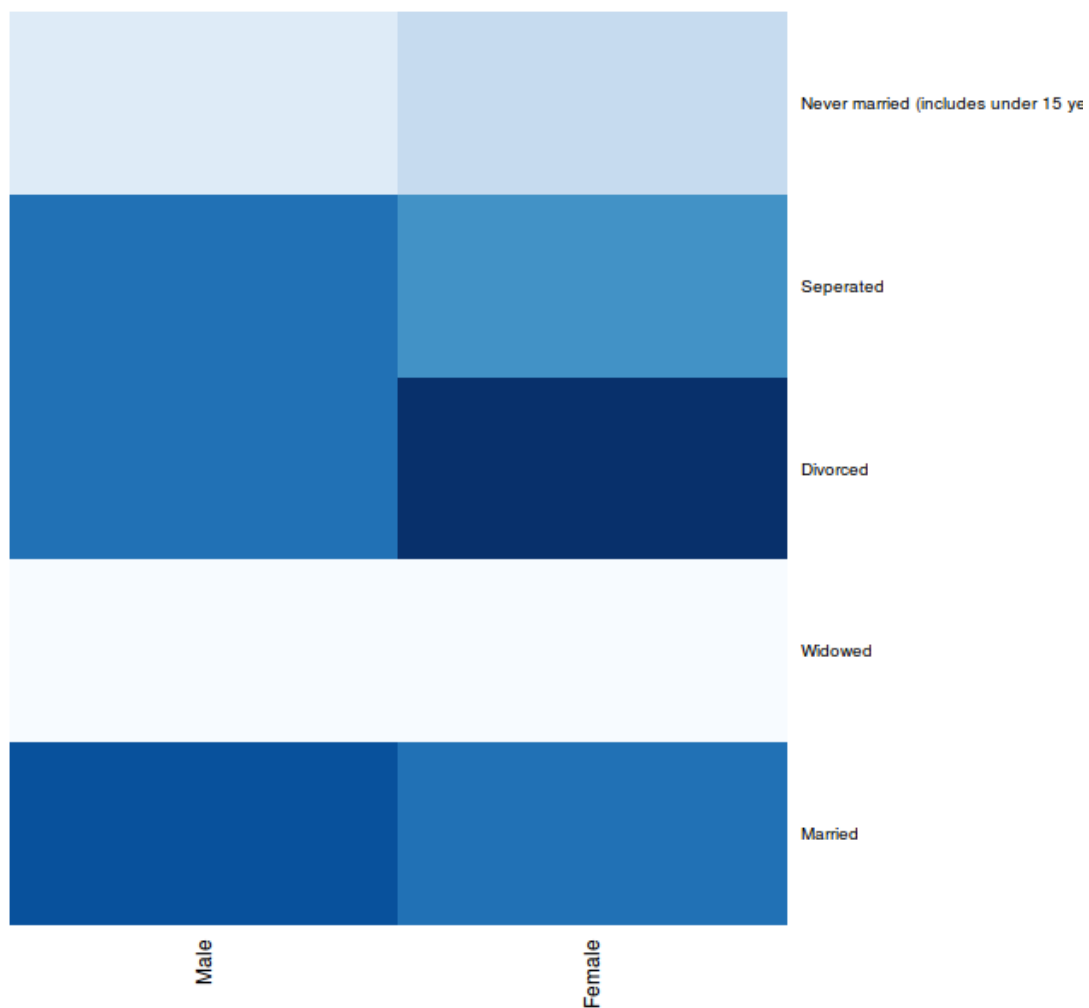


Figure 8: peoples' Earnings based on gender and education level

Question # 5

For this task , i want to present **Marital Status, Education, and Hours** as a data cube.

Slicing data can be sliced by differnt martical status values.

Dicing Using dice operation, i can pick a subcube of certain group of values (Hours, week, months etc.)

Drill up and Drill down We can drill up from hours to days and then to months. Similarly, we can drill down from months to days and then to hours.

Pivoting is about rotating the cube in space to see the data from different perspective.

Question #6

Below is the code that tried to recover the image using method “**repetitive_nn**” of TSP.

```
1 d_matrix = dist(img)
2 tsp = TSP(d_matrix)
3 x = solve_TSP(tsp, method="repetitive_nn")
4 img_order = x
```

Below is the code written in R to show the recovered image. The recovered image is shown in Figure. 9

```
1 order_img = matrix(NA, 1018, 4662)
2 for(i in c(1:1018))
3 {
4   order_img[i,] = img[img_order[i],]
5 }
6
7 img1 = array(NA, c(1018, 1554, 3))
8 for (i in c(1:1018))
9 {
10   row = order_img[i,]
11   y=1
12   for(pcol in seq(1, 4662, 3))
13   {
14     img1[i, y, 1] = as.numeric(row[pcol])/255
15     img1[i, y, 2] = as.numeric(row[pcol+1])/255
16     img1[i, y, 3] = as.numeric(row[pcol+2])/255
17     y = y + 1
18   }
19 }
20
21 require('jpeg')
22 res = dim(img1)[1:2]
23 plot(1, 1, xlim=c(1, res[1]), ylim=c(1, res[2]), asp=1, type='n', xaxs='i', yaxs='i', xaxt='n',
24       yaxt='n', xlab='', ylab='', bty='n')
25 rasterImage(img1, 1, 1, res[1], res[2])
```



Figure 9: Recovered image using method “repetitive_nn” of TSP

Figure. shows the recovered image using “**nearest_insertion**” method of TSP.



Figure 10: Recovered image using method“nearest_insertion” of TSP

Question #7

To import data in MySQL , we need similar kind of Table structure to save data. Below SQL command creates the schema to store census data.

```
1 CREATE TABLE tbl_Census (  
2 State varchar(50) DEFAULT NULL,  
3 House_id varchar(50) DEFAULT NULL,  
4 Weight varchar(50),  
5 House_relation varchar(50),  
6 Sex varchar(20),  
7 Aage varchar(50),  
8 Race varchar(50),  
9 Marriage varchar(50) DEFAULT NULL,  
10 Education varchar(50) DEFAULT NULL,  
11 Ancestry varchar(50) DEFAULT NULL,  
12 Language varchar(50) DEFAULT NULL,  
13 Employment_Status varchar(50) DEFAULT NULL,  
14 Traveltime varchar(50) DEFAULT NULL,  
15 Industry varchar(50) DEFAULT NULL,  
16 Occupation varchar(50) DEFAULT NULL,  
17 Hours varchar(50) DEFAULT NULL,  
18 Weeks varchar(50) DEFAULT NULL,
```

```

20 Salary varchar(50) DEFAULT NULL,
21 Income varchar(50) DEFAULT NULL,
22 Earnings int(20)DEFAULT NULL
23 );

```

Data is imported from large version of the csv file into table using following command.

```

1 LOAD DATA INFILE '/var/lib/mysql-files/extract_large.csv'
2 INTO TABLE tbl_Census
3 FIELDS TERMINATED BY ',';

```

In the end, I executed “group by” statement of SQL to generate the pivot table.

```

1 SELECT Sex, Education, avg(Earnings)
2 FROM tbl_Census
3 group by Sex, Education

```

Sex	Education	Avg(Earnings)
1	00	0.0000
1	01	3436.9163
1	02	1018.5143
1	03	5719.2300
1	04	4585.9244
1	05	7892.1493
1	06	8702.8604
1	07	10341.6579
1	08	16704.6764
1	09	22399.7355
1	10	26831.3657
1	11	29900.6087
1	12	35249.8174
1	13	52814.0881
1	14	63902.4224
1	15	97160.4380
1	16	69790.5412
2	00	0.0000
2	01	1358.9230
2	02	386.2795
2	03	1996.0910
2	04	1656.8459
2	05	3048.7202
2	06	3878.9171
2	07	4711.3237
2	08	7794.4182
2	09	10865.3802
2	10	15102.8411
2	11	16310.1780
2	12	20708.8457
2	13	27804.7783
2	14	35058.9334
2	15	44529.5096
2	16	44325.0832

Figure 11: pivot table using SQL group by statement