

Assignment Overview In this assignment you will explore a couple of aspects of RNA-seq (with a small introduction to clustering). For this assignment you will have to generate some visualizations - we recommend R or Python, but use a language you are comfortable with! Make sure to show your work/code in your writeup!

As a reminder, any questions about the assignment should be posted to Piazza.

Question 1. Time Series [20 pts] This file contains normalized expression values for 100 genes over 10 time points. Most genes have a stable background expression level, but some special genes show increased expression over the timecourse and some show decreased expression.

Question 1a. Cluster the genes using an algorithm of your choice. Which genes show increasing expression and which genes show decreasing expression, and how did you determine this? What is the background expression level (numerical value) and how did you determine this? [Hint: K-means and hierarchical clustering are common clustering algorithms you could try.]

```
In [ ]:
Ans:
There are 12 genes that show decrease in expression over time and 20 genes that show increase in expression over time.

The genes that show decreasing in expression are:
gene_7, gene_14, gene_21, gene_28, gene_42, gene_49, gene_56, gene_63, gene_77, gene_84, gene_91, gene_98

The genes that show increasing in expression are:
gene_3, gene_10, gene_15, gene_20, gene_25, gene_30, gene_35, gene_40, gene_45,
gene_50, gene_55, gene_60, gene_65, gene_70, gene_75, gene_80, gene_85, gene_90, gene_95, gene_100

The genes were separated into clusters using kmeans clustering and also by hierarchical clustering.
The same genes with increasing and decreasing expression were then identified using time point regression for each gene in each cluster and then calculating the average slope.

Here cluster 3 corresponds to the background genes because it has an averaged slope of around -0.011
This cluster consists of 68 genes and the background expression level of these genes
mean(background$Q1_mean) == 55.0965

In [ ]:
getwd()
setwd("D:/Comparative_Genomics/Assignments/5/")

Expression_Data1 = read.csv("Expression1.txt", header = T, sep = "\t", row.names = 1)

#kmeans clustering
kmeans = kmeans(Expression_Data1, centers = 3)
kmeans$cluster

#hierarchical clustering
distance = dist(Expression_Data1[,c(1:10)], method = "manhattan")
clusters = hclust(distance, method = "average")
plot(clusters)

#time point regression
library(dplyr)
background = filter(Expression_Data1, Expression_Data1$clusters == 3)
increasing = filter(Expression_Data1, Expression_Data1$clusters == 1)
decreasing = filter(Expression_Data1, Expression_Data1$clusters == 2)

background_expres = t(background[,c(1:10)])
increasing_expres = t(increasing[,c(1:10)])
decreasing_expres = t(decreasing[,c(1:10)])

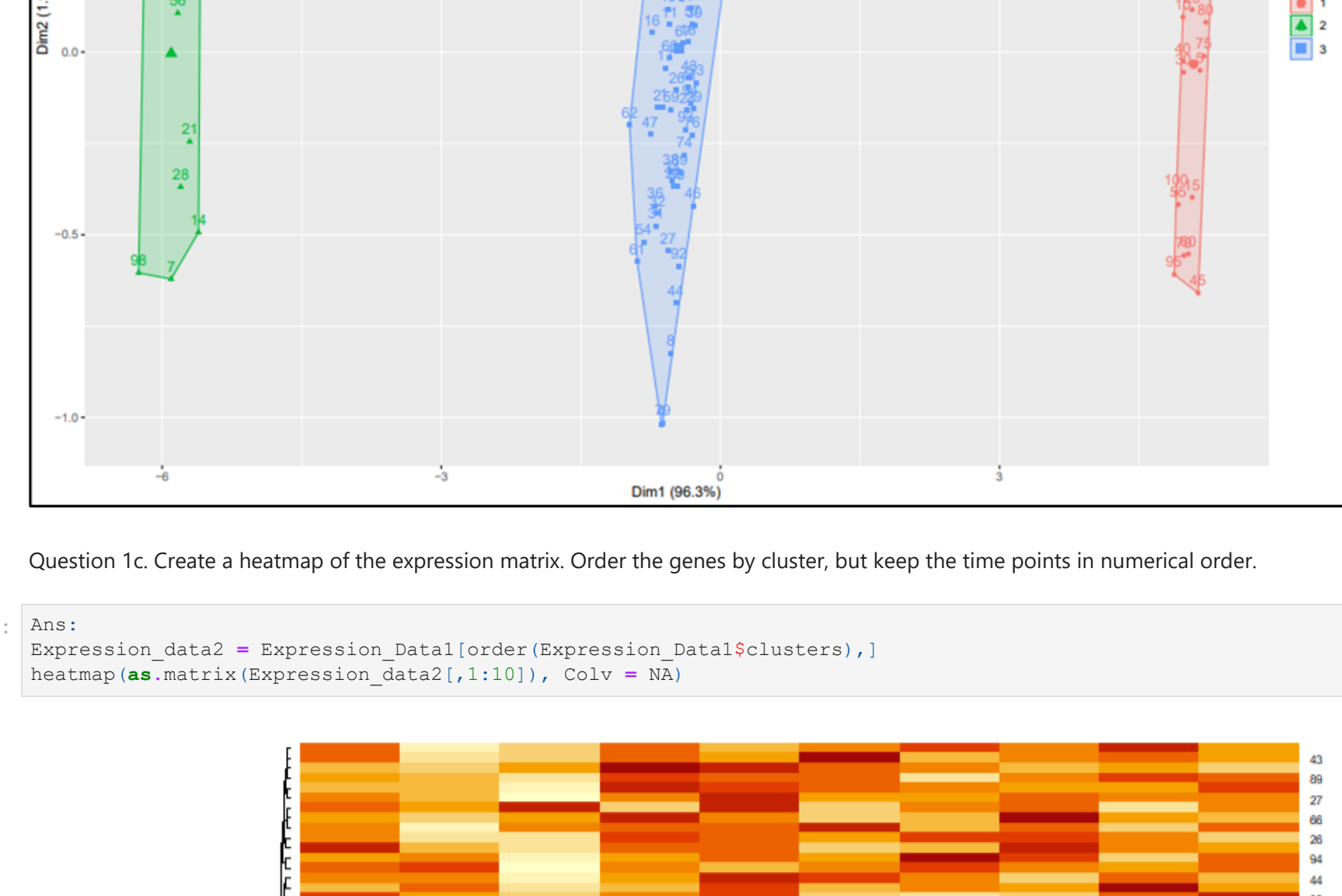
background_slopes = c()
increasing_slopes = c()
decreasing_slopes = c()

time_point = c(1:10)
my_rownames = c(1:112)
for (i in my_rownames){
  x = decreasing_expres[i,]
  a = linetime_point
  required = a$coefficients[2]
  decreasing_slopes = append(decreasing_slopes, required)
}

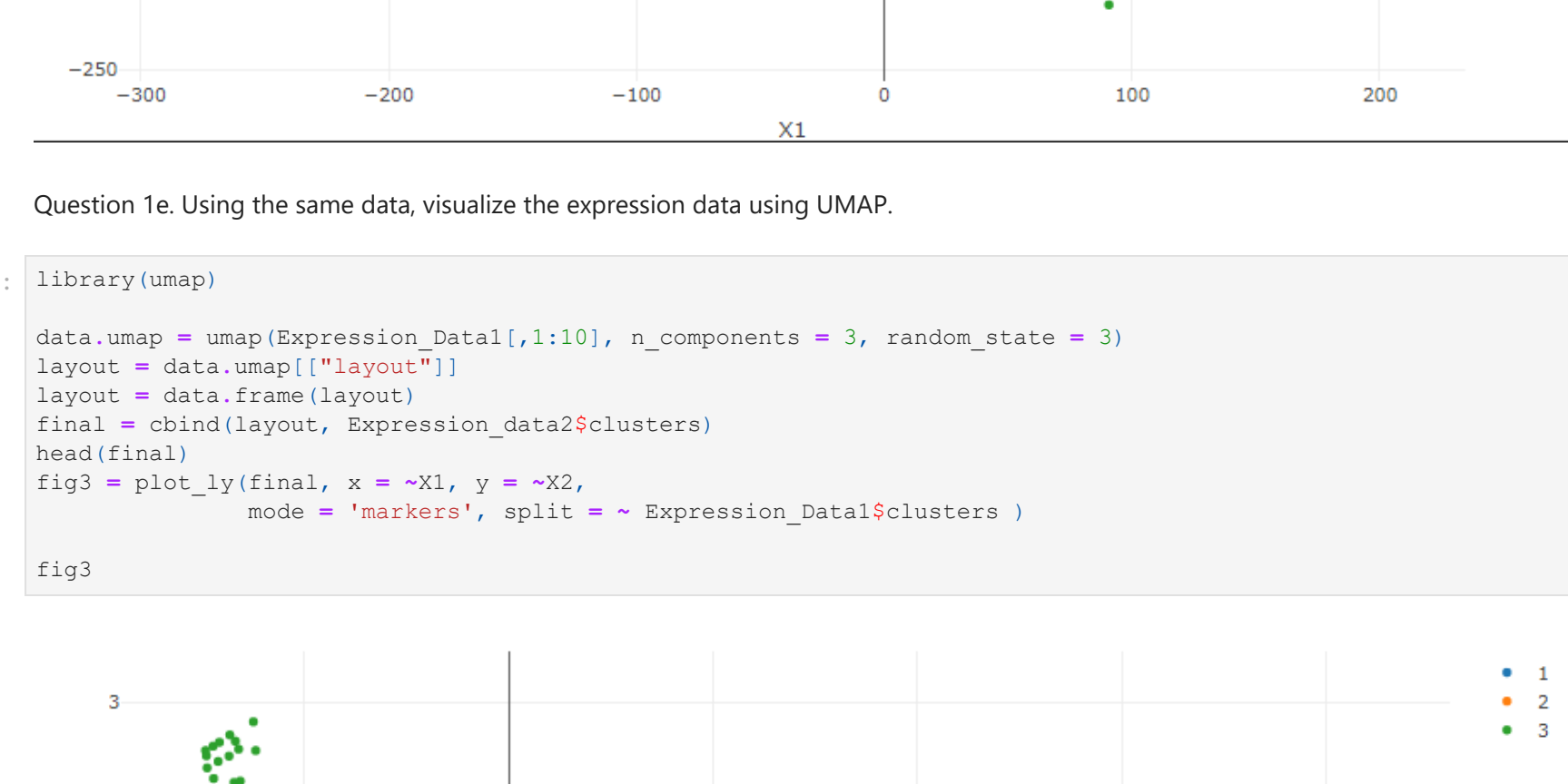
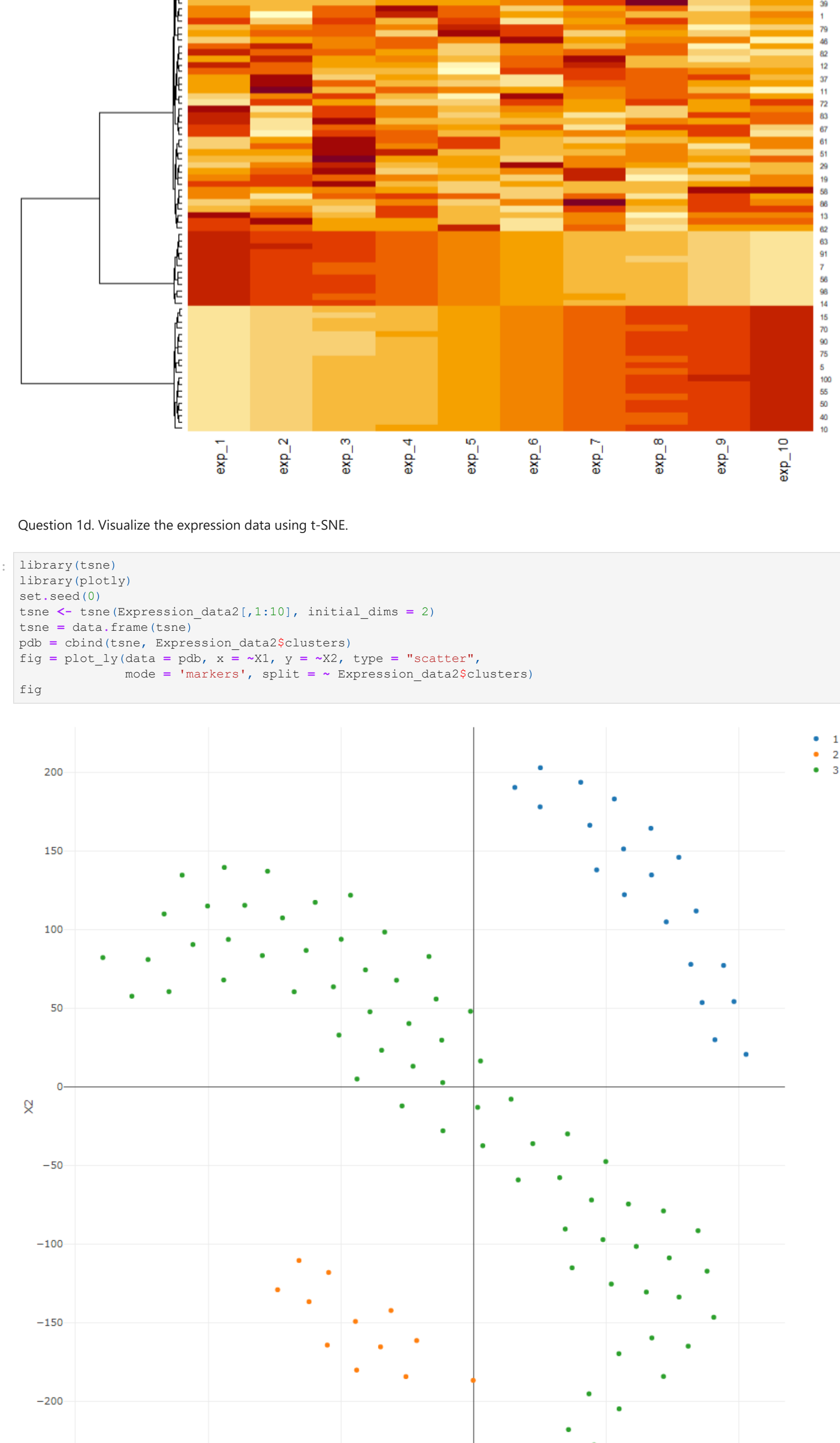
mean(background_slopes) # = -0.0122971
mean(increasing_slopes) # = 0.1001659
mean(decreasing_slopes) # = -0.09934661

mean(background$Q1_mean) # = 55.00665

#kmeans clustering result
#hierarchical clustering result
```



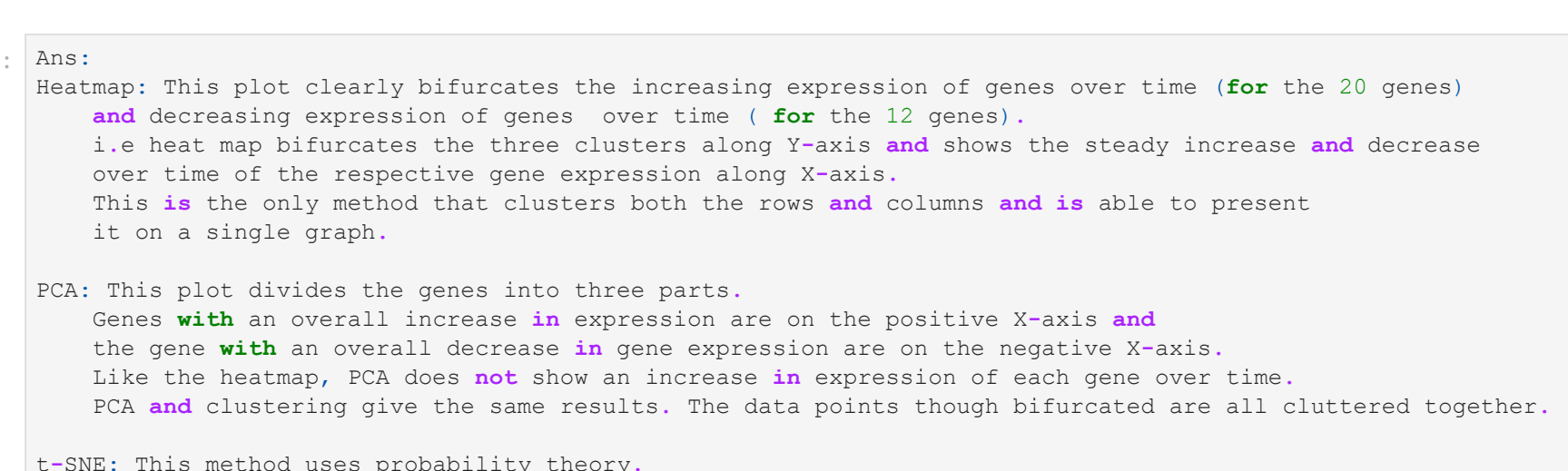
Question 1c. Create a heatmap of the expression matrix. Order the genes by cluster, but keep the time points in numerical order.



Question 1d. Visualize the expression data using t-SNE



Question 1e. Using the same data, visualize the expression data using UMAP.



Question 1f. In a few sentences, compare the (1) heatmap, (2) PCA, (3) t-SNE and (4) UMAP results. Be sure to comment on understandability, relative positioning of clusters, runtime, and any other significant factors that you see.

PCA: This plot clearly bifurcates the increasing expression of genes over time (for the 20 genes). The heat map bifurcates the three clusters along Y-axis and shows the steady increase and decrease over time of the respective gene expression along X-axis.

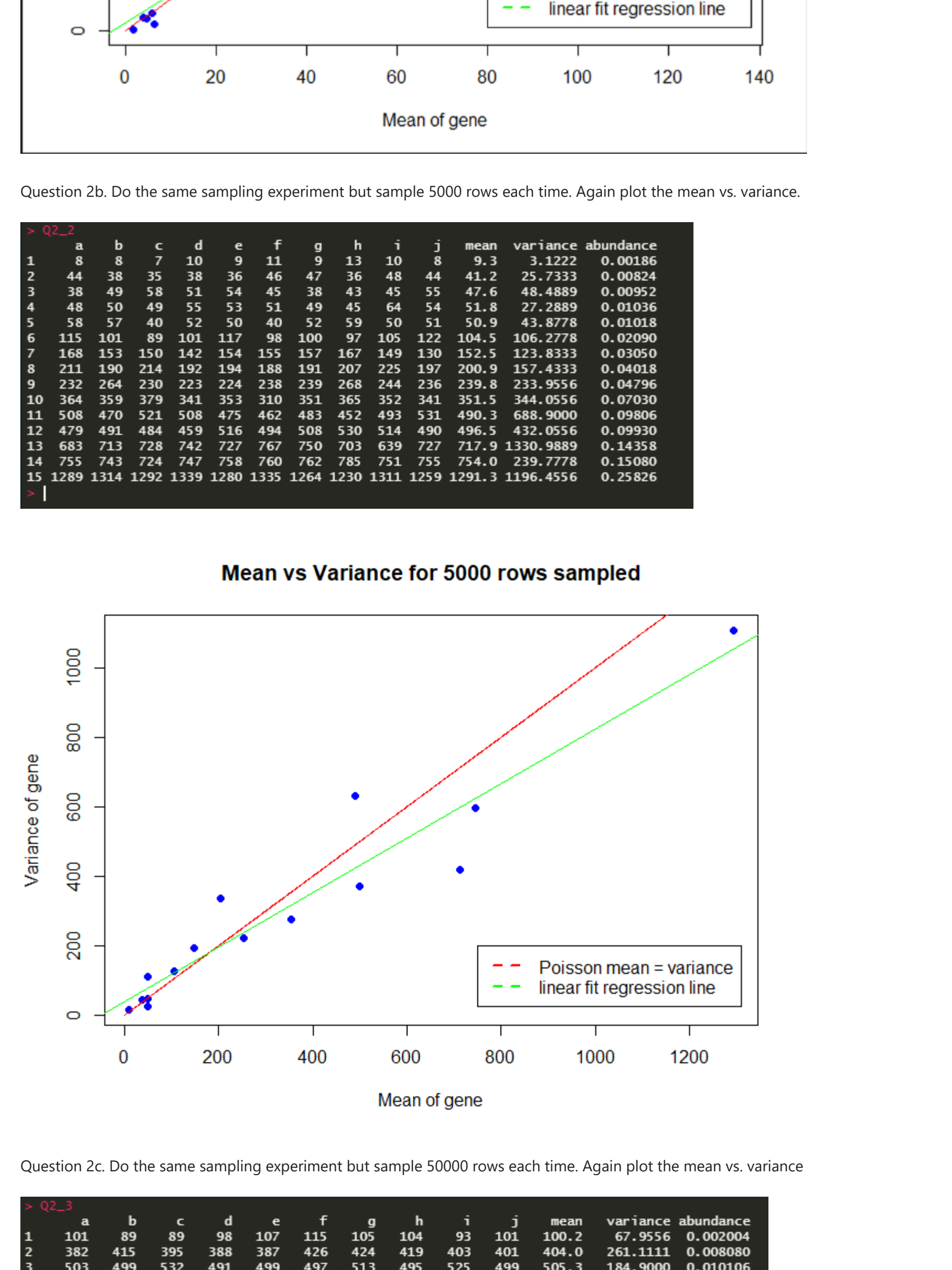
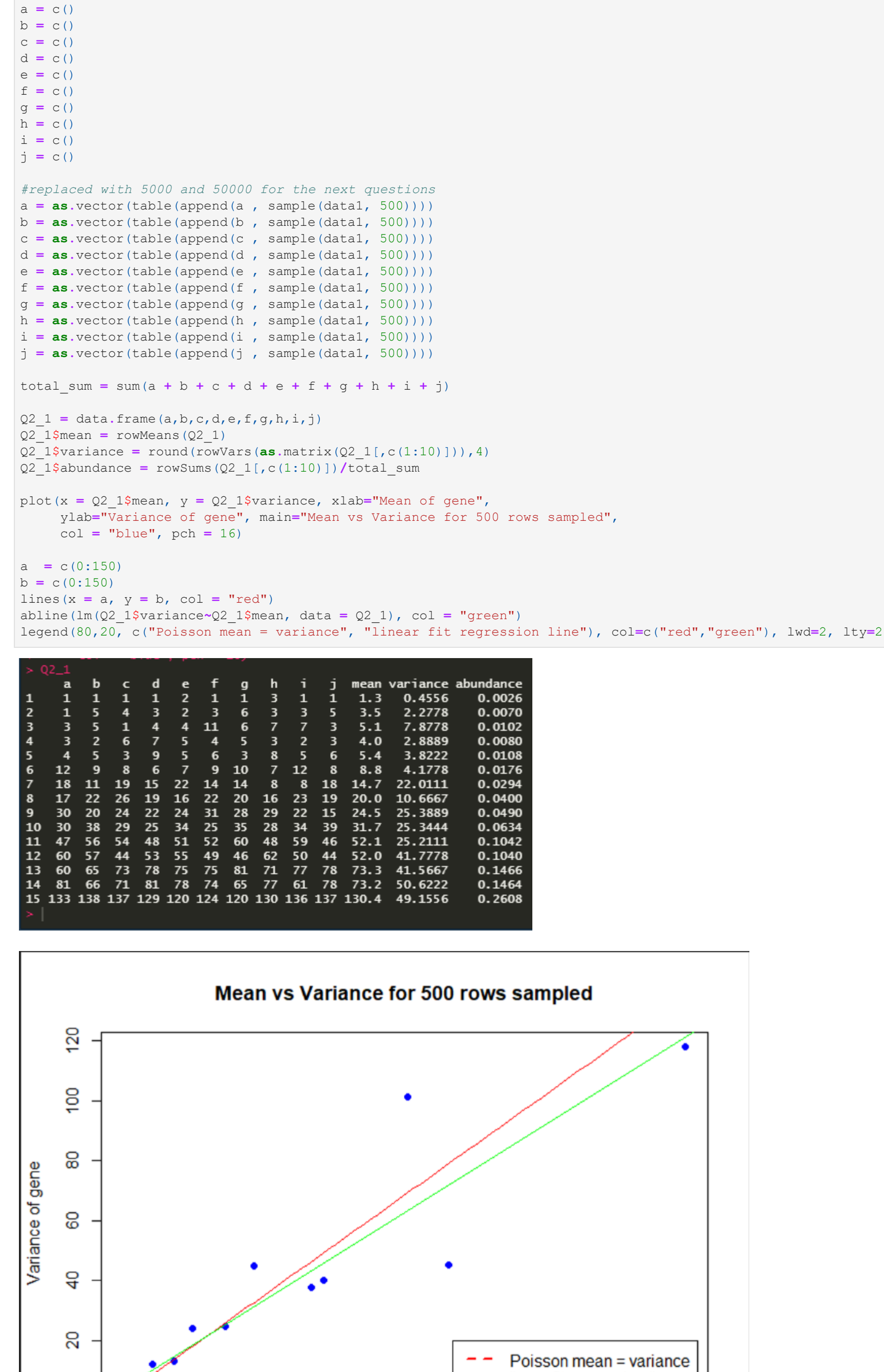
t-SNE: This clustering method uses probability theory. Compared to the PCA method, it takes significantly more time to converge, but presents a significantly better insights when visualised.

UMAP: This method is also used to visualize complex data into low dimensions. As the number of data points increase, UMAP becomes more time efficient as compared to the t-SNE.

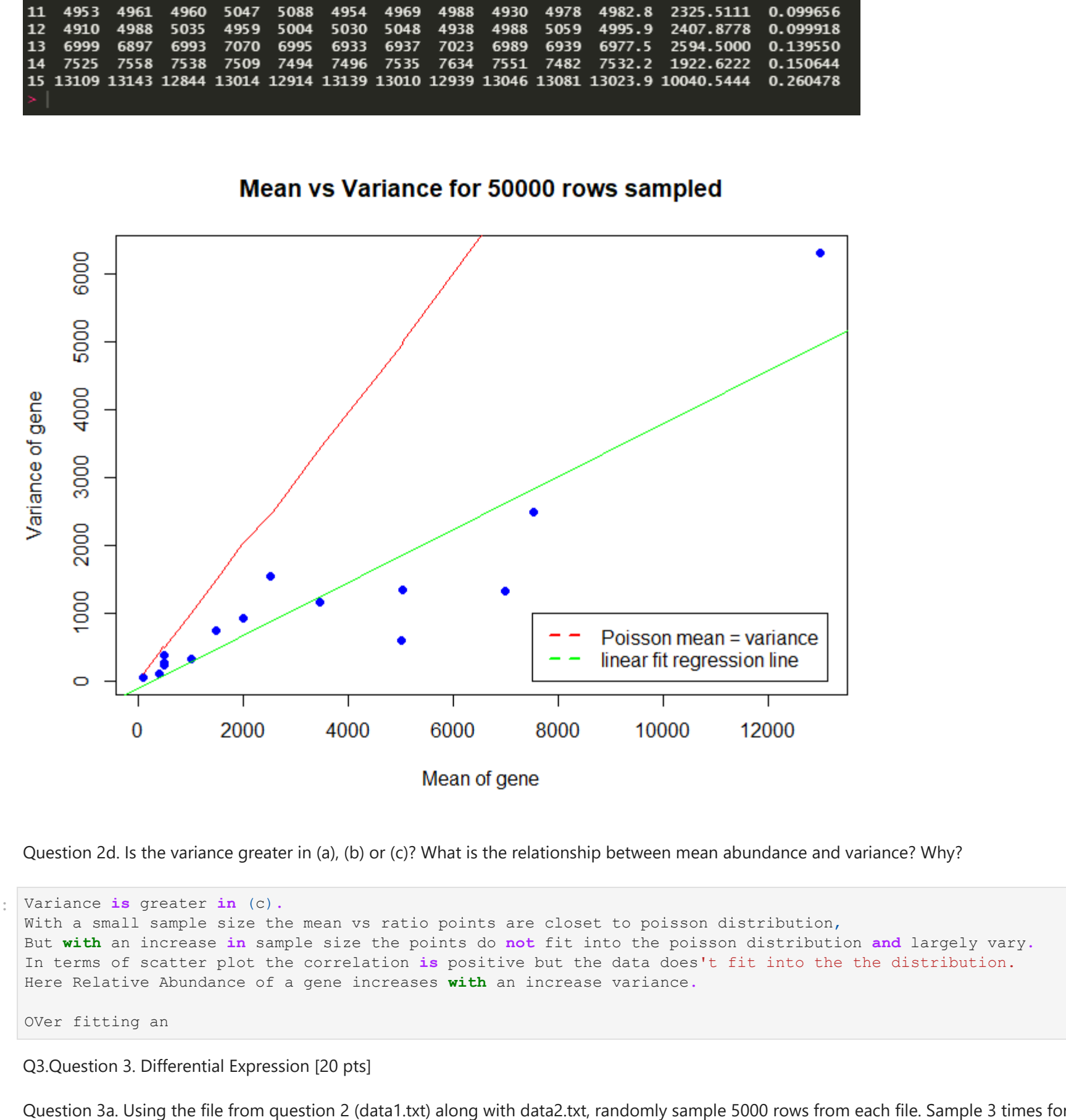
Question 2. Sampling Simulation [10 pts] A typical human cell has ~250,000 transcripts, and a typical bulk RNA-seq experiment may involve millions of cells. Consequently in an RNA-seq experiment you may start with trillions of RNA molecules, although your sequencer will only give a few tens of millions of reads. Therefore your RNA-seq experiment will be a small sampling of the full composition. We hope the sequences will be a representative sample of the total population, but if your sample is very uniquely or biased it may not represent the true distribution. We will explore this concept by sampling a small subset of transcripts (500 to 5000) out of a much larger set (1M) so that you can evaluate this bias.

In data1.txt with 100,000 lines we provide an abstraction of RNA-seq data where normalization has been performed and the number of times a gene name occurs corresponds to the number of transcripts in the sample.

Question 2a. Randomly sample 500 rows. Do this 10 times and record the relative abundance of each of the 15 genes. Make a scatterplot the mean vs. variance of each gene (x-axis=mean of gene, y-axis=variance of gene.)



Question 2b. Do the same sampling experiment but sample 5000 rows each time. Again plot the mean vs. variance.



Question 2c. Do the same sampling experiment but sample 50000 rows each time. Again plot the mean vs. variance.



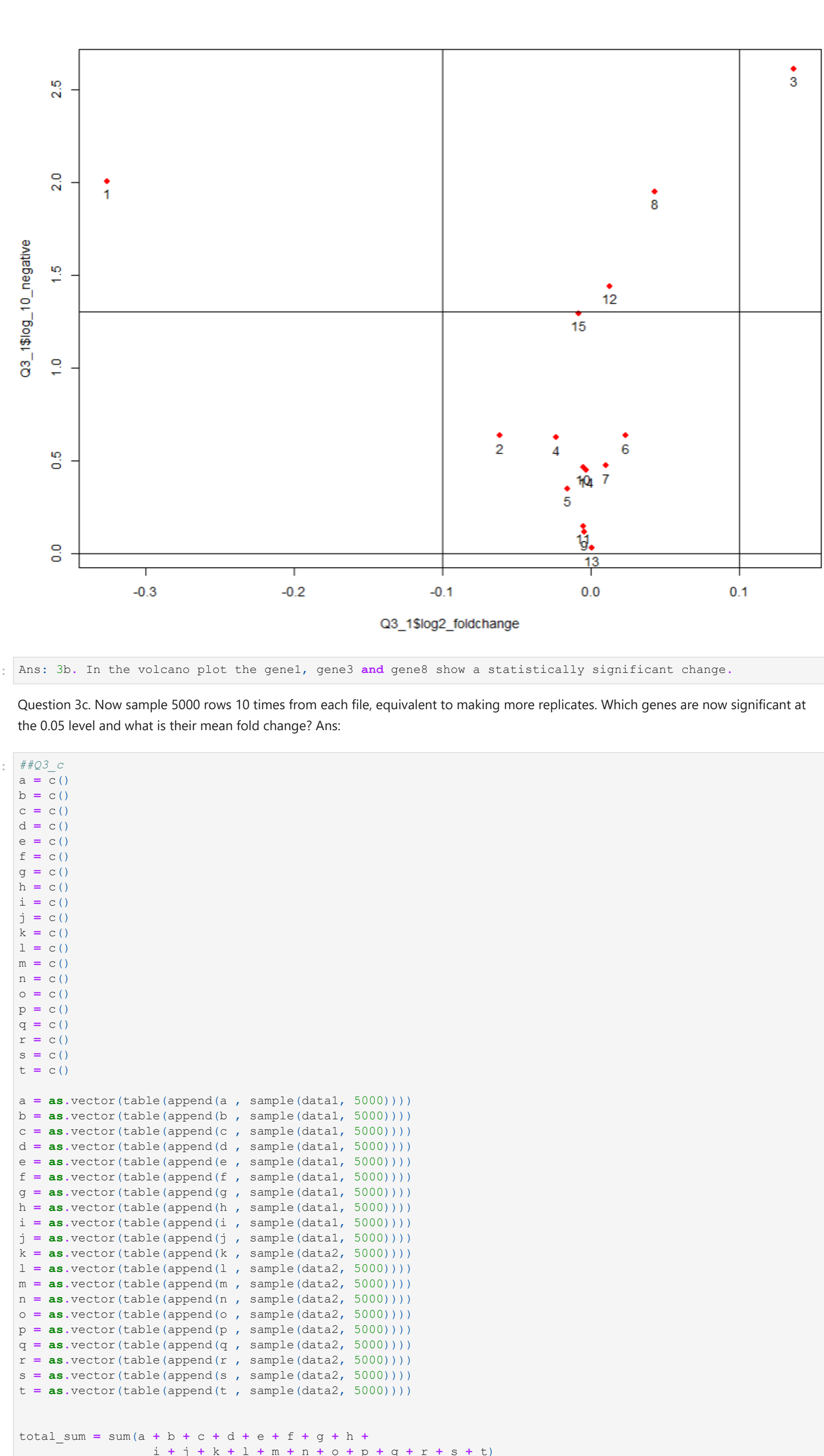
Question 2d. Is the variance greater in (a), (b) or (c)? What is the relationship between mean abundance and variance? Why?

Ans: Variance is greater in (c). With a small sample size the mean vs ratio points are closer to poisson distribution, but as n increases the points move further away from the poisson distribution and largely vary. In terms of scatter plot the correlation is positive but the data doesn't fit into the distribution. Here Relative Abundance of a gene increases with an increase variance.

Over fitting an

Question 3a. Differential Expression [20 pts]

Question 3a. Using the file from question 2 (data1.txt) along with data2.txt, randomly sample 5000 rows from each file. Sample 3 times for each file (this emulates making experimental replicates) and conduct a paired t-test for differential expression of each of the 15 genes. Which genes are significantly differentially expressed at the 0.05 level and what is their mean fold change? Ans: A difference in the mean values of the two files show the following results.



Question 3b. In the volcano plot the gene1, gene5 and gene8 show a statistically significant change.

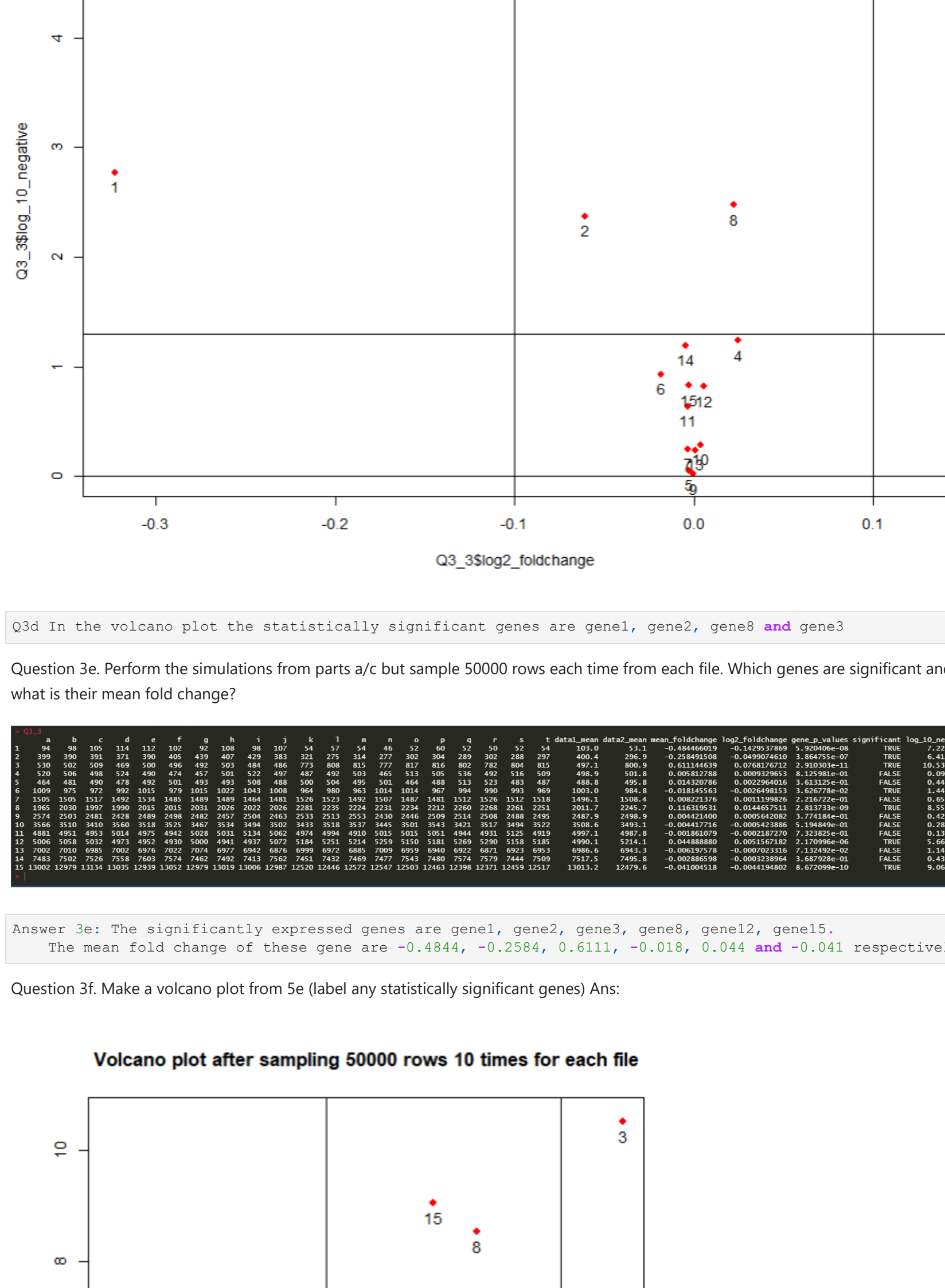
Question 3c. Now sample 5000 rows 10 times from each file, equivalent to making more replicates. Which genes are now significant at the 0.05 level and what is their mean fold change? Ans:



Question 3c. Ans: The most highly expressed genes are Gene1, Gene2, Gene3 and Gene8 and their mean fold changes are -0.543, -0.280, 0.614, 0.124 respectively

Question 3d. Make a volcano plot using the results from part c (label any statistically significant genes) Ans:





In [ ]: Q3d In the volcano plot the statistically significant genes are gene1, gene2, gene8 and gene3

Question 3e. Perform the simulations from parts a/c but sample 50000 rows each time from each file. Which genes are significant and what is their mean fold change?

Q3\_3log2\_foldchange

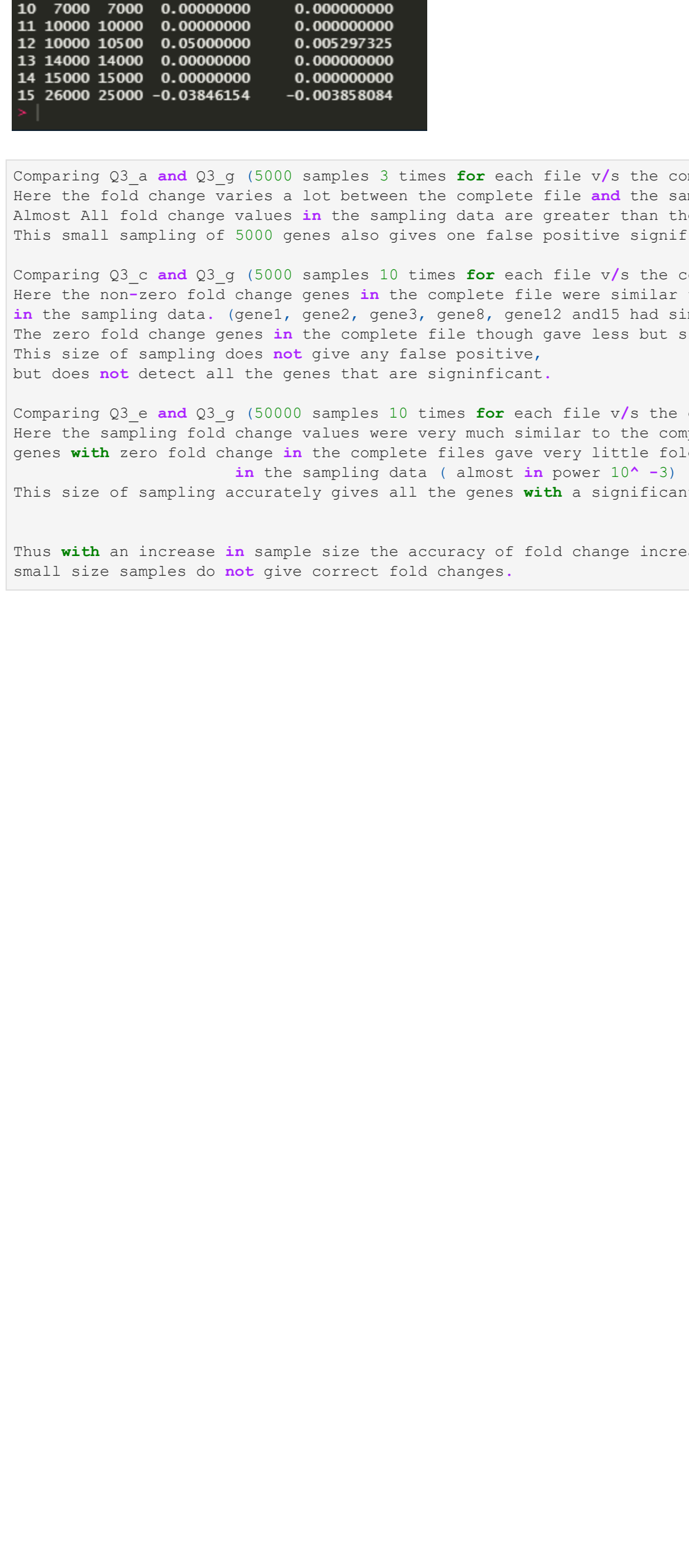
Q3\_3negLog10Pvalue

9  
13  
14

Answer 3f. In the volcano plot the statistically significant genes are gene1 and gene3

In [ ]: Answer 3e: The significantly expressed genes are gene1, gene2, gene3, gene8, gene12, gene15. The mean fold change of these gene are -0.4844, -0.2584, 0.6111, -0.018, 0.044 and +0.041 respectively

Question 3f. Make a volcano plot from 5e (label any statistically significant genes) Ans:



In [ ]: Answer 3f: In the volcano plot the statistically significant genes are gene1, and gene3

Question 3g. Now examine the complete files: compare the fold change in the complete files vs the different subsamples making sure to address replicates and the size of the random sample.

Ans:

|    | a     | b     | Foldchange  | log2_foldchange |
|----|-------|-------|-------------|-----------------|
| 1  | 200   | 100   | -0.50000000 | -1.10824021     |
| 2  | 800   | 600   | -0.25000000 | -0.043036467    |
| 3  | 1000  | 1600  | 0.60000000  | 0.068039994     |
| 4  | 1000  | 1000  | 0.00000000  | 0.000000000     |
| 5  | 1000  | 1000  | 0.00000000  | 0.000000000     |
| 6  | 2000  | 2000  | 0.00000000  | 0.000000000     |
| 7  | 3000  | 3000  | 0.00000000  | 0.000000000     |
| 8  | 4000  | 4500  | 0.12500000  | 0.014200908     |
| 9  | 5000  | 5000  | 0.00000000  | 0.000000000     |
| 10 | 7000  | 7000  | 0.00000000  | 0.000000000     |
| 11 | 10000 | 10000 | 0.00000000  | 0.000000000     |
| 12 | 10000 | 10500 | 0.05000000  | 0.005297315     |
| 13 | 14000 | 14000 | 0.00000000  | 0.000000000     |
| 14 | 15000 | 15000 | 0.00000000  | 0.000000000     |
| 15 | 26000 | 25000 | -0.03846154 | -0.003858084    |

In [ ]: Comparing Q3\_a and Q3\_g (5000 samples 3 times for each file v/s the complete file) Here the fold change varies a lot between the complete file and the sampling data.

Almost All fold change values in the sampling data are greater than the complete data. This small sampling of 5000 genes also gives one false positive significant genes

Comparing Q3\_c and Q3\_g (5000 samples 10 times for each file v/s the complete file) Here the non-zero fold change genes in the complete file were similar to the fold change values in the sampling data. (gene1, gene2, gene3, gene8, gene12 and15 had similar fold changes

The zero fold change genes in the complete file though gave less but significant fold change in comparison. This size of sampling does not give any false positive, but does not detect all the genes that are significant.

Comparing Q3\_e and Q3\_g (50000 samples 10 times for each file v/s the complete file) Here the sampling fold change values were very much similar to the complete files. genes with zero fold change in the complete files gave very little fold change in the sampling data ( almost in power 10^-3)

This size of sampling accurately gives all the genes with a significant fold change.

Thus with an increase in sample size the accuracy of fold change increases. small size samples do not give correct fold changes.